

В.М.Вержбицкий

---

# ЧИСЛЕННЫЕ МЕТОДЫ

Линейная алгебра  
и нелинейные  
уравнения

ВЫСШАЯ ШКОЛА

В.М.Вержбицкий

---

---

# ЧИСЛЕННЫЕ МЕТОДЫ

Линейная алгебра  
и нелинейные  
уравнения

*Рекомендовано  
Министерством образования  
Российской Федерации  
в качестве учебного пособия  
для студентов математических  
и инженерных специальностей  
высших учебных заведений*



Москва  
«Высшая школа» 2000

УДК 519.6  
ББК 22.19  
В31

**Рецензенты:**

кафедра математического моделирования систем и процессов Пермского государственного технического университета (зав. кафедрой доктор физ.-мат. наук, профессор П.В. Трусов);

кафедра вычислительной математики Удмуртского государственного университета (зав. кафедрой доктор физ.-мат. наук, профессор Г.Г. Исламов);

доктор физ.-мат. наук, профессор Н.В. Азбелев

ISBN 5-06-003654-5

© ГУПИ издательство «Высшая школа», 2000

Оригинал-макет данного издания является собственностью издательства «Высшая школа» и его репродуцирование (воспроизведение) любым способом без согласия издательства запрещено.

# Оглавление

Предисловие . . . . .	5
<b>Глава 1. Об учете погрешностей приближенных вычислений</b>	
1.1 Общая формула для оценки главной части погрешности . . . . .	7
1.2. Статистический и технический подходы к учету погрешностей действий . . . . .	11
1.3. Понятия о погрешностях машинной арифметики . . . . .	13
1.4. Примеры неустойчивых задач и методов . . . . .	19
1.5. Обусловленность линейных алгебраических систем . . . . .	21
1.6. Погрешности корней скалярных уравнений с приближенными коэффициентами . . . . .	26
Упражнения . . . . .	30
<b>Глава 2. Решение линейных алгебраических систем (прямые методы)</b>	
2.0. Введение . . . . .	32
2.1. Алгоритм решения СЛАУ методом Гаусса с постолбцовым выбором главного элемента . . . . .	34
2.2. Применение метода Гаусса к вычислению определителей и к обращению матриц . . . . .	38
2.3. LU-разложение матриц . . . . .	40
2.4. Решение линейных систем и обращение матриц с помощью LU-разложения . . . . .	43
2.5. Разложение симметричных матриц. Метод квадратных корней . . . . .	48
2.6. Метод прогонки решения систем с трехдиагональными матрицами коэффициентов . . . . .	51
2.7. Метод вращений решения линейных систем . . . . .	55
2.8. Два замечания к применению прямых методов . . . . .	58
Упражнения . . . . .	61
<b>Глава 3. Итерационные методы решения линейных алгебраических систем и обращения матриц</b>	
3.1. Решение СЛАУ методом простых итераций . . . . .	63
3.2. Метод Якоби . . . . .	69
3.3. Метод Зейделя . . . . .	72
3.4. Понятие о методе релаксации . . . . .	77
3.5. О других итерационных методах решения СЛАУ . . . . .	80
3.6. Быстросходящийся итерационный способ обращения матриц . . . . .	86
3.7. О роли ошибок округления в итерационных методах . . . . .	91
Упражнения . . . . .	93
<b>Глава 4. Методы решения алгебраических проблем собственных значений</b>	
4.1. Собственные пары матрицы и их простейшие свойства . . . . .	96
4.2. Степенной метод . . . . .	101
4.3. Обратные итерации . . . . .	111
4.4. Метод вращения Якоби решения симметричной полной проблемы собственных значений . . . . .	117
4.5. Понятие об LU-алгоритме для несимметричных задач . . . . .	124
4.6. QR-алгоритм . . . . .	129
Упражнения . . . . .	138

## **Глава 5. Методы решения нелинейных задач скалярных уравнений**

5.1. Локализация корней . . . . .	141
5.2. Метод дихотомии. Метод хорд . . . . .	147
5.3. Типы сходимостей итерационных последовательностей . . . . .	151
5.4. Метод Ньютона . . . . .	154
5.5. Применение метода Ньютона к вычислению значений функций . . . . .	163
5.6. Модификации метода Ньютона. Метод секущих . . . . .	166
5.7. Задача о неподвижной точке. Метод простых итераций . . . . .	176
5.8. Ускорение сходимости последовательных приближений . . . . .	186
а) D2-процесс Эйткена . . . . .	187
б) метод Вегстейна . . . . .	192
5.9. Нелинейные уравнения с параметром. Бифуркации . . . . .	195
5.10. О методах решения алгебраических уравнений. Метод Бернулли . . . . .	202
Упражнения . . . . .	209

## **Глава 6. Методы решения систем нелинейных уравнений**

6.1. Векторная запись нелинейных систем. Метод простых итераций . . . . .	212
6.2. Метод Ньютона, его реализации и модификации . . . . .	216
6.3. Метод Брауна . . . . .	223
6.4. О решении нелинейных систем методами спуска . . . . .	225
6.5. Численный пример . . . . .	229
6.6. Сходимость метода Ньютона и некоторых его модификаций . . . . .	231
Упражнения . . . . .	243

### **Приложение 1. Краткие сведения о нормах векторов и матриц.**

Сходимость в конечномерных пространствах . . . . .	245
--	-----

### **Приложение 2. Производные векторных функций . . . . .**

250

### **Приложение 3. Образцы постановок лабораторных заданий . . . . .**

254

Литература . . . . .

259

Предметный указатель . . . . .

263

## Предисловие

Предлагаемая вниманию читателя книга явилась результатом многолетнего опыта преподавания автором цикла дисциплин вычислительного характера студентам, обучающимся по специальности «Прикладная математика» и другим специальностям Ижевского государственного технического университета. Она отражает лишь небольшую часть материала основ численных методов, с которой, по мнению автора, должно начинаться их изучение. А именно, в книгу включены численные методы решения основных задач линейной алгебры (к которым, как правило, в итоге сводится решение других задач вычислительной математики, прикладной механики и т.п.), а также методы решения алгебраических и трансцендентных уравнений и их систем. Предполагается, что читатель знаком с линейной алгеброй и математическим (в частности, функциональным) анализом. Тем не менее, чтобы освежить некоторые базовые знания, а также чтобы расширить круг потенциальных пользователей книги, в ее конце приведены два приложения, посвященные нормам в конечномерных векторных пространствах и производным векторных функций векторных аргументов (последнее нужно лишь для облегчения чтения шестой главы); третье приложение содержит постановки заданий для лабораторных работ, выполнение которых автор считает необходимым элементом в освоении численных методов; фигурирующие в этих заданиях точностные константы могут быть заменены другими, в зависимости от конкретных параметров применяемого вычислительного инструментария. Написанию компьютерных программ, реализующих те или иные методы, должно способствовать наличие в книге во многих случаях достаточно детализированных алгоритмов.

Главное, что хотел бы автор от своей читательской аудитории, впервые изучающей численные методы, — это понимания основных идей (которых не так уж много) и умения довести их до числа, содержащего, по возможности, только верные знаки искомого результата. Автор сознательно позволяет себе «не замечать» двух очевидных барьеров на этом пути. Первый барьер связан с тем, что наряду с основными, наиболее важными методами, здесь отводится место и второстепенным, на первый взгляд, методам, т.е. как учебное пособие книга, возможно, несколько перегружена. Второй барьер — более «плотное», чем это принято в учебной литературе, цитирование. Ссылки на другие книги в разных случаях делаются из разных соображений: для указания на источники дополнительной информации и на возможные пути более глубокого изучения той или иной темы, для сопоставления различных терминов, означающих одно и то же, а также для того, чтобы «разделить ответственность» за некоторые приводимые без доказательства

результаты с авторами цитируемых работ. На журнальные статьи ссылки делаются лишь в исключительных случаях (в основном, когда речь идет о приоритетах).

Не вводя специально раздела, в котором были бы сосредоточены все встречающиеся в книге обозначения и сокращения (поскольку, в большинстве своем, они стандартны и должны быть читателю хорошо знакомы), автор обращает внимание на следующие особенности. Во-первых, автор считает удобным использование символа присваивания  $:=$  и в качестве знака "по определению", что позволяет сразу различать, особенно в цепи преобразований, где равенство фигурирует как утверждение, а где им определяется (вводится, обозначается) некоторая величина. Во-вторых, в некоторых местах трудно было обойтись без использования довольно редко применяемого, но очень удобного в вычислительной математике знака  $\leq$  означающего неравенство в смысле главных (линейных) частей сравниваемых выражений. В-третьих, учитывая распространенность в тексте стержневых для книги групп слов, таких как "система линейных алгебраических уравнений" и "метод простых итераций", автор часто использует вместо них аббревиатуры СЛАО и МПИ соответственно. Другие аббревиатуры носят локальный характер и не должны затруднять чтение книги. Наконец, автор старался выдержать общепринятую систему, в которой из написания одной и той же буквы, например,  $X$ ,  $x$  или  $x$ , сразу должно быть ясно, матрица это, вектор или скаляр; при этом для обозначения индексов применялись только буквы  $m, n, l, i, j, k$ , обычно по умолчанию принимающие только целые значения.

Автор выражает свою признательность коллегам по кафедре прикладной математики и информатики ИжГТУ, в частности, профессору В.Я. Дерру и доценту А.А. Айзиковичу за внимание и поддержку при подготовке рукописи и особенно профессору А.Л. Тептину, сделавшему своевременно ряд ценных замечаний. Автор благодарен А.П. Ананину, А.Н. Борисенкову, Р.М. Валееву, А.В. Дворникову, П.В. Куприянову, А.В. Лосеву, К.В. Перцеву, принимавшим участие в компьютерном наборе рукописи, и, более всего, С.В. Высоцкому, собравшему воедино и приведшему к определенному стандарту набранные в разное время разными исполнителями части рукописи и осуществившему окончательную подготовку оригинал-макета. Особо автор отмечает роль не одного поколения студентов, явно или неявно способствовавших появлению этой книги.

В настоящее время автор работает над продолжением этого издания – книгой «Численные методы (математический анализ и дифференциальные уравнения)».

Автор

# ГЛАВА 1 || ОБ УЧЕТЕ ПОГРЕШНОСТЕЙ ПРИБЛИЖЕННЫХ ВЫЧИСЛЕНИЙ

*Рассматривается круг вопросов, связанных с учетом погрешностей, появление которых неизбежно при численном анализе математических моделей, в частности, при решении линейных алгебраических систем и нелинейных скалярных уравнений. Обращается внимание на заведомо приближенный характер компьютерных операций над действительными числами. Приводятся примеры задач и методов, чрезмерно чувствительных к ошибкам исходных данных и к погрешностям арифметических действий.*

## 1.1. ОБЩАЯ ФОРМУЛА ДЛЯ ОЦЕНКИ ГЛАВНОЙ ЧАСТИ ПОГРЕШНОСТИ

При численном решении математических и прикладных задач почти неизбежно появление на том или ином этапе их решения погрешностей следующих трех типов.

а) **Погрешность задачи.** Она связана с приближенным характером исходной содержательной модели (в частности, с невозможностью учесть все факторы в процессе изучения моделируемого явления), а также ее математического описания, параметрами которого служат обычно приближенные числа (например, из-за принципиальной невозможности выполнения абсолютно точных измерений): Для вычислителя погрешность задачи следует считать *неустранимой* (безусловной), хотя постановщик задачи иногда может ее изменить.

б) **Погрешность метода.** Это погрешность, связанная со способом решения поставленной математической задачи и появляющаяся в результате подмены исходной математической модели другой или конечной последовательностью других, например, линейных моделей. При создании численных методов закладывается возможность отслеживания таких погрешностей и доведения их до сколь угодно малого уровня. Отсюда естественно отношение к погрешности метода как к *устранимой* (или условной).

в) **Погрешность округлений** (погрешность действий). Этот тип погрешностей обусловлен необходимостью выполнять арифметические операции над числами, усеченными до количества разрядов, зависящего от применяемой вычислительной техники (если, разумеется, не используются специальные программные средства, реализующие, например, арифметику рациональных чисел).



Все три описанных типа погрешностей в сумме дают *полную погрешность* результата решения задачи. Поскольку первый тип погрешностей не находится в пределах компетенции вычислителя, для него он служит лишь ориентиром точности, с которой следует рассчитывать математическую модель. Нет смысла решать задачу существенно точнее, чем это диктуется неопределенностью исходных данных. Таким образом, погрешность метода подчиняют погрешности задачи. Наконец, при выводе оценок погрешностей численных методов обычно исходят из предположения, что все операции над числами выполняются точно. Это означает, что погрешность округлений не должна существенно отражаться на результатах реализации методов, т.е. должна подчиняться погрешности метода. Влияние погрешностей округлений не следует упускать из виду ни на стадии отбора и алгоритмизации численных методов, ни при выборе вычислительных и программных средств, ни при выполнении отдельных действий и вычисления значений функций.

Рассмотрим некоторые возможные подходы к учету погрешностей действий ([1, 8, 10, 17, 21, 24, 34] и др.).

Пусть  $A$  и  $a$  – два "близких" числа; условимся считать  $A$  точным,  $a$  – приближенным.

Величина  $\Delta a := |A - a|$  называется *абсолютной погрешностью* приближенного числа  $a$ , а  $\delta a := \frac{\Delta a}{|a|}$  – его *относительной погрешностью* <sup>\*)</sup>.

Числа  $\Delta_a$  и  $\delta_a$  такие, что  $\Delta_a \geq \Delta a$  и  $\delta_a = \frac{\Delta_a}{|a|} \geq \delta a$ , называются *оценками* или *границами* абсолютной и относительной погрешностей соответственно (к  $\Delta_a$  и  $\delta_a$  часто применяют также термин "*предельные погрешности*"). Так как обычно истинные погрешности не известны, то там, где не может возникнуть недоразумений, будем иногда называть  $\Delta_a$  и  $\delta_a$  просто абсолютной и относительной погрешностями.

Поставим вопрос о грубом оценивании погрешности результата вычисления значения дифференцируемой функции  $u = f(x_1, x_2, \dots, x_n)$  приближенных аргументов  $x_1, x_2, \dots, x_n$ , если известны границы их абсолютных погрешностей  $\Delta_{x_1}, \Delta_{x_2}, \dots, \Delta_{x_n}$  соответственно. В этом случае точные значения аргументов  $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$  лежат соответственно на отрезках

---

<sup>\*)</sup> Возможно, чаще относительной погрешностью называют  $\frac{\Delta a}{|A|}$ .

Использование символа  $\Delta$  в дальнейшем может быть двояким: как для обозначения абсолютной погрешности, так и для обозначения приращения переменной, что будет ясно из контекста либо специально оговорено.

$[x_1 - \Delta_{x_1}, x_1 + \Delta_{x_1}]$ ,  $[x_2 - \Delta_{x_2}, x_2 + \Delta_{x_2}]$ , ...,  $[x_n - \Delta_{x_n}, x_n + \Delta_{x_n}]$ , а точная абсолютная погрешность результата  $u = f(x_1, x_2, \dots, x_n)$  есть

$$\Delta u = |f(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n) - f(x_1, x_2, \dots, x_n)|$$

– модуль полного приращения функции. Главной, т.е. линейной частью этого приращения является, как известно, полный дифференциал  $du$ . Таким образом, имеем:

$$\Delta u \approx |du| = \left| \sum_{i=1}^n \frac{\partial u}{\partial x_i} dx_i \right| \leq \sum_{i=1}^n \left| \frac{\partial u}{\partial x_i} \right| |\tilde{x}_i - x_i| \leq \sum_{i=1}^n \left| \frac{\partial u}{\partial x_i} \right| \cdot \Delta_{x_i},$$

т.е. за границу абсолютной погрешности результата приближенно может быть принята величина

$$\Delta_u = \sum_{i=1}^n \left| \frac{\partial u}{\partial x_i} \right| \cdot \Delta_{x_i}. \quad (1.1)$$

Отсюда легко получается формула приближенной оценки относительной погрешности значения  $u$ :

$$\delta_u = \frac{\Delta_u}{|u|} = \sum_{i=1}^n \left| \frac{\partial u}{\partial x_i} \right| \cdot \frac{\Delta_{x_i}}{|u|} = \sum_{i=1}^n \left| \frac{\partial u}{\partial x_i} \right| \cdot \Delta_{x_i} = \sum_{i=1}^n \left| \frac{\partial \ln u}{\partial x_i} \right| \cdot \Delta_{x_i}. \quad (1.2)$$

Как частные случаи формул (1.1), (1.2) (точных для функций, линейных относительно  $x_i$  или  $\ln x_i$  соответственно) можно получить известные правила оценивания погрешностей результатов арифметических действий.

Действительно, пусть  $u = \pm x_1 \pm x_2 \pm \dots \pm x_n$ . Тогда  $\left| \frac{\partial u}{\partial x_i} \right| = 1$  и

$\Delta_{\Sigma(\pm x_i)} = \sum_{i=1}^n 1 \cdot \Delta_{x_i} = \sum_{i=1}^n \Delta_{x_i}$ , т.е. при сложении и вычитании приближенных чисел их предельные абсолютные погрешности складываются.

Пусть теперь  $u = x_1 \cdot x_2 \cdot \dots \cdot x_n$ , где можно считать все сомножители положительными. Так как  $\ln u = \ln x_1 + \ln x_2 + \dots + \ln x_n$  и  $\frac{\partial \ln u}{\partial x_i} = \frac{1}{x_i}$ , то, согласно (1.2),

$$\delta_{\Gamma x_i} = \sum_{i=1}^n \frac{1}{x_i} \Delta_{x_i} = \sum_{i=1}^n \delta_{x_i}. \quad (1.3)$$

Если же  $u = \frac{x_1}{x_2}$ , где  $x_1, x_2 > 0$ , то  $\ln u = \ln x_1 - \ln x_2$ ,  $\frac{\partial \ln u}{\partial x_i} = \frac{1}{x_i}$  и, значит,

$$\delta_{x_1 x_2} = \frac{\Delta_{x_1}}{x_1} + \frac{\Delta_{x_2}}{x_2} = \delta_{x_1} + \delta_{x_2}.$$

Последнее вместе с (1.3) означает известный результат о сложении предельных относительных погрешностей при умножении и делении приближенных чисел.

Возвращаясь к сложению, рассмотрим относительную погрешность суммы  $n$  положительных приближенных чисел  $x_1, x_2, \dots, x_n$ , имеющих границы относительных погрешностей  $\delta_{x_1}, \delta_{x_2}, \dots, \delta_{x_n}$  соответственно:

$$\begin{aligned} \delta(x_1 + x_2 + \dots + x_n) &= \frac{\Delta(x_1 + x_2 + \dots + x_n)}{x_1 + x_2 + \dots + x_n} \leq \\ &\leq \frac{\Delta_{x_1 + x_2 + \dots + x_n}}{x_1 + x_2 + \dots + x_n} = \frac{\Delta_{x_1} + \Delta_{x_2} + \dots + \Delta_{x_n}}{x_1 + x_2 + \dots + x_n} = \\ &= \frac{x_1 \delta_{x_1} + x_2 \delta_{x_2} + \dots + x_n \delta_{x_n}}{x_1 + x_2 + \dots + x_n} \leq \frac{x_1 \delta^* + x_2 \delta^* + \dots + x_n \delta^*}{x_1 + x_2 + \dots + x_n} = \delta^*, \end{aligned}$$

где  $\delta^* = \max_i \delta_{x_i}$ . Полученное неравенство говорит о том, что относительная погрешность суммы  $n$  положительных приближенных чисел не превосходит максимальной относительной погрешности слагаемых.

С вычитанием приближенных чисел дело обстоит хуже: оценка

$$\delta_{x_1 - x_2} = \frac{\Delta_{x_1 - x_2}}{|x_1 - x_2|} = \frac{\Delta_{x_1} + \Delta_{x_2}}{|x_1 - x_2|}$$

относительной погрешности разности  $x_1 - x_2$  двух приближенных положительных чисел указывает на возможность сильного возрастания погрешности при  $x_1 - x_2 \rightarrow 0$ . В этом случае говорят о потере точности при вычитании близких чисел.

Часто возникает *обратная задача теории погрешностей*: какой точности данные нужно подать на вход, чтобы на выходе получить результат заданной точности? Применительно к поставленной выше прямой задаче оценивания погрешности результата вычисления значения функции при заданных оценках погрешностей аргументов обратная задача заключается в оценивании величин  $\Delta x_i$  (или  $\delta x_i$ ) по известной величине  $\Delta y$ . Для случая дифференцируемой функции одной переменной грубое решение обратной задачи тривиально: если  $y = f(x)$ , то  $\Delta y \approx |dy| = |f'(x)| \Delta x$ , откуда

$\Delta x \approx \frac{\Delta y}{|f'(x)|}$ . Для функции большего числа переменных обратная задача,

вообще говоря, некорректна. Нужны дополнительные условия. Например, **принцип равных влияний** состоит в предположении, что частные дифференциалы  $\left| \frac{\partial u}{\partial x_i} \right| \Delta_{x_i}$  в (1.1) одинаково влияют на погрешность значения функции; тогда

$$\Delta_u = n \left| \frac{\partial u}{\partial x_i} \right| \Delta_{x_i}, \quad \text{откуда} \quad \Delta_{x_i} = \frac{\Delta_u}{n \left| \frac{\partial u}{\partial x_i} \right|}.$$

В качестве другого довольно естественного допущения можно принять равенство относительных погрешностей всех аргументов, т.е. считать

$$\delta_{x_i} = \frac{\Delta_{x_i}}{|x_i|} = p \quad \text{при всех } i=1, 2, \dots, n. \quad \text{Тогда } \Delta_{x_i} = p|x_i| \quad \text{и, значит,}$$

$$\Delta_u = p \sum_{i=1}^n \left| \frac{\partial u}{\partial x_i} x_i \right|. \quad \text{Из последнего равенства получаем величину}$$

$$p = \frac{\Delta_u}{\sum_{i=1}^n \left| x_i \frac{\partial u}{\partial x_i} \right|} \quad (\text{характеризующую относительный уровень точности зада-$$

ния аргументов), на основе которой за границы абсолютных погрешностей аргументов принимаем  $\Delta_{x_i} = \frac{|x_i| \Delta_u}{\sum_{i=1}^n \left| x_i \frac{\partial u}{\partial x_i} \right|}$ . Имеются и другие, более сложные

подходы к решению обратной задачи (см., например, [6]).

## 1.2. СТАТИСТИЧЕСКИЙ И ТЕХНИЧЕСКИЙ ПОДХОДЫ К УЧЕТУ ПОГРЕШНОСТЕЙ ДЕЙСТВИЙ

Рассмотренный выше *аналитический* (или *классический*) способ учета погрешностей действий, предполагающий точное оценивание погрешностей, основанное либо на приведенных в предыдущем параграфе правилах подсчета погрешностей арифметических действий, либо на параллельной работе с верхними и нижними границами исходных данных, имеет два существенных недостатка. Во-первых, этот способ чрезвычайно громоздок и не может быть рекомендован при массовых вычислениях. Во-вторых, он учитывает крайние, наихудшие случаи взаимодействия погрешностей, которые допустимы, но маловероятны. Ясно, что, например, при суммировании нескольких приближенных чисел (полученных в результате измерений, округлений или каким-либо другим путем) среди них почти наверное будут слагаемые как с избытком, так и с недостатком, т.е. произой-

дет частичная компенсация погрешностей. При больших количествах однотипных вычислений вступают в силу уже *вероятностные* или *статистические законы* формирования погрешностей результатов действий. Например, методами теории вероятностей показывается, что математическое ожидание абсолютной погрешности суммы  $n$  слагаемых с одинаковым уровнем абсолютных погрешностей, при достаточно большом  $n$ , пропорционально  $\sqrt{n}$  ([8, 10, 21, ...]). В частности, если  $n > 10$  и все слагаемые округлены до  $m$ -го десятичного разряда, то для подсчета абсолютной погрешности суммы  $S$  применяют *правило Чеботарева*

$$\Delta S \approx \sqrt{3n} \cdot 0.5 \cdot 10^{-m}. \quad (1.4)$$

Различие в результатах классического и статистического подходов к оцениванию погрешности суммы рассмотрим на примере оценки погрешности среднего арифметического нескольких приближенных чисел.

Пусть  $x = \frac{1}{n}(x_1 + \dots + x_n)$  – среднее арифметическое  $n$  ( $> 10$ ) приближенных чисел (например, результатов измерений), имеющих одинаковый уровень абсолютных погрешностей  $\Delta_{x_i} = 0.5 \cdot 10^{-m}$ . Тогда классическая оценка абсолютной погрешности величины  $x$  есть

$$\Delta_x = \frac{1}{n}(\Delta_{x_1} + \dots + \Delta_{x_n}) = \frac{1}{n} \cdot n \cdot 0.5 \cdot 10^{-m} = 0.5 \cdot 10^{-m} = \Delta_{x_i},$$

т.е. такая же, как и у исходных данных. В то же время по формуле (1.4) имеем

$$\Delta_x \approx \frac{1}{n} \sqrt{3n} \cdot 0.5 \cdot 10^{-m} = \sqrt{\frac{3}{n}} \cdot 0.5 \cdot 10^{-m} = \sqrt{\frac{3}{n}} \cdot \Delta_{x_i}, \quad \xrightarrow{n \rightarrow \infty} 0.$$

Как видим, применение правила Чеботарева приводит к естественному выводу о том, что арифметическое усреднение результатов измерений или наблюдений увеличивает точность, чего нельзя сказать на основе классической теории погрешностей.

Прямое применение вероятностно-статистических оценок погрешностей также является достаточно сложным делом и вряд ли может быть рекомендовано при рядовых массовых вычислениях. Однако именно такие оценки подкрепляют практические правила работы с приближенными числами, составляющие основу так называемого *технического подхода*. Этот подход связывают с именем известного русского кораблестроителя, математика и механика академика А. Н. Крылова. Согласно *принципу А. Н. Крылова*, приближенное число должно записываться так, чтобы в нем все значащие цифры, кроме последней, были верными и лишь последняя была бы сомнительна, и притом в среднем\*) не более чем на одну единицу. Напомним, что *значащими цифрами* числа в его позиционной записи называются все его цифры, начиная с первой ненулевой слева. Знача-

\*) "В среднем" здесь понимается в вероятностном смысле.

щую цифру приближенного числа называют *верной*, если абсолютная погрешность числа не превосходит единицы разряда, в котором стоит эта цифра (или половины единицы; в этом случае иногда применяется термин *верная в узком смысле*).

Чтобы результаты арифметических действий, совершаемых над приближенными числами, записанными в соответствии с принципом А. Н. Крылова, также соответствовали этому принципу, нужно придерживаться следующих нехитрых правил [10, 17, 21].

1. При сложении и вычитании приближенных чисел в результате следует сохранять столько десятичных знаков, сколько их в приближенном данном с наименьшим количеством десятичных знаков.

2. При умножении и делении в результате следует сохранять столько значащих цифр, сколько их имеет приближенное данное с наименьшим числом значащих цифр.

3. Результаты промежуточных вычислений должны иметь один-два запасных знака (которые затем должны быть отброшены).

Таким образом, при техническом подходе к учету погрешностей приближенных вычислений предполагается, что в самой записи приближенного числа содержится информация о его точности. И хотя прямая выгода от применения приведенных правил работы с приближенными числами может быть получена лишь при ручном счете (не нужно оперировать с цифрами, не влияющими на информативную часть приближенного результата), их знание и понимание помогает правильной интерпретации компьютерных расчетов, а иногда и самой организации таковых.

### 1.3. ПОНЯТИЕ О ПОГРЕШНОСТЯХ МАШИННОЙ АРИФМЕТИКИ

Для представления вещественных чисел в ЭВМ применяют, в основном, два способа: с фиксированной и с плавающей запятой (точкой).

Пусть в основу запоминающего устройства машины положены однотипные физические устройства (базисные элементы), имеющие  $r$  устойчивых состояний (как правило,  $r = 2, 8, 16$  и т.п.), причем каждому числу ставится в соответствие одинаковое количество  $k$  этих элементов и, кроме того, с помощью таких или более простых элементов может фиксироваться знак. Упорядоченные элементы образуют разрядную сетку машинного слова: в каждом разряде может быть записано одно из базисных чисел  $0, 1, \dots, r-1$  (одна из  $r$  "цифр"  $r$ -ичной системы счисления) и в специальном разряде отображен знак  $+$  или  $-$ .

При записи числа с *фиксированной запятой* кроме упомянутых  $r$  параметров (основания системы счисления) и  $k$  (количества разрядов, отводимых под запись цифр числа) указывается еще количество  $\ell$  разря-

дов, выделяемых под дробную часть числа. Таким образом, положительное вещественное число  $a$ , представляющее собой в  $r$ -ичной системе бесконечную, вообще говоря, непериодическую дробь, здесь будет отображено конечной последовательностью

$$\alpha_1 \alpha_2 \dots \alpha_{k-\ell} \alpha_{k-\ell+1} \dots \alpha_{k-1} \alpha_k,$$

где  $\alpha_i \in \{0; 1; \dots; r-1\}$ , т.е. реализуется приближенное равенство<sup>\*)</sup>

$$a \approx \text{fix}(a) := \alpha_1 r^{k-\ell-1} + \alpha_2 r^{k-\ell-2} + \dots + \alpha_{k-\ell} r^0 + \\ + \alpha_{k-\ell+1} r^{-1} + \dots + \alpha_{k-1} r^{-(\ell-1)} + \alpha_k r^{-\ell}.$$

**Диапазон** представляемых таким способом чисел определяется числами с наибольшими цифрами во всех разрядах, т.е. наименьшим числом  $-(r-1)(r-1)\dots(r-1)$  и наибольшим  $(r-1)(r-1)\dots(r-1)$ , а **абсолютная точность представления** есть оценка величины  $|a - \text{fix}(a)|$ , зависящая от способа округления: это  $r^{-\ell}$  при простом отбрасывании "хвоста"  $\alpha_{k+1} r^{-(\ell+1)} + \alpha_{k+2} r^{-(\ell+2)} + \dots$  числа  $a$  и половина этой величины при **правильном округлении** (т.е. при увеличении  $\alpha_k$  на единицу, если  $\alpha_{k+1} > \frac{r}{2}$ ). Заметим, что абсолютная точность представления вещественных чисел с фиксированной запятой одинакова в любой части диапазона. В

то же время **относительная точность**, т.е. оценка величины  $\left| \frac{a - \text{fix}(a)}{a} \right|$

(или  $\left| \frac{a - \text{fix}(a)}{\text{fix}(a)} \right|$ ), очевидно, может значительно различаться в зависимости

от того, берется  $a$  близким к нулю или к границе диапазона. Иными словами, вещественные числа с фиксированной запятой имеют равномерную абсолютную плотность распределения на всем отрезке вещественной оси, определяемом границами диапазона, и неравномерную, возрастающую к границам отрезка, относительную плотность распределения.

В основе значительно чаще употребляемого представления с **плавающей запятой** лежит следующая экспоненциальная форма записи вещественного числа:

$$a = M \cdot r^p,$$

где  $r$  — **основание**,  $p$  — **порядок**, а  $M$  такое, что  $r^{-1} \leq |M| < 1$  ( $= r^0$ ) — **мантисса**. Если под мантиссу выделяется  $\ell$   $r$ -ичных элементов, а под порядок  $m$ , то в системе записи с плавающей запятой вещественное число  $a$

<sup>\*)</sup> Символ := нами используется в двух близких смыслах:

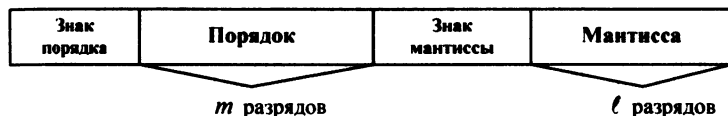
1) положить по определению, задать, выбрать;

2) присвоить, т.е. записать в ячейку памяти компьютера, определяемую идентификатором в левой части, значение, задаваемое правой частью.

представляется конечным числом  $\text{fl}(a)$  (от англ. floating — "плавающий") вида

$$a \approx \text{fl}(a) := \pm (\beta_1 r^{-1} + \beta_2 r^{-2} + \dots + \beta_\ell r^{-\ell}) \cdot r^\gamma,$$

где  $\gamma$  — целое число из промежутка  $[-r^m, r^m - 1]$ ;  $\beta_1 \in \{1; \dots; r-1\}$ ;  $\beta_i \in \{0; 1; \dots; r-1\}$  ( $i = 2, \dots, \ell$ ), т.е. *машинное слово* условно имеет структуру:



Числа  $\pm r^{r^m}$  определяют границы допустимого числового диапазона. Более информативно здесь нужно говорить о диапазоне представимости положительных вещественных чисел, составляющем промежутки  $[r^{-r^m}, r^{r^m - 1}]$ . Левую и правую границы этого отрезка называют соответственно *машинным нулем* и *машинной бесконечностью*, так как числа из промежутка  $[-r^{-r^m}, r^{-r^m}]$  машина заменяет нулем, а числа, лежащие за пределами промежутка  $[-r^{r^m - 1}, r^{r^m - 1}]$ , она не воспринимает (без специальных ухищрений).

Важной характеристикой является число  $\varepsilon$ , называемое *машинный эпсилон* и обозначаемое обычно идентификатором *macheps*. Эта характеристика определяется как расстояние между единицей и ближайшим следующим за ней числом системы машинных чисел с плавающей запятой. Так как

$$1 = (1 \cdot r^{-1} + 0 \cdot r^{-2} + \dots + 0 \cdot r^{-\ell} + \dots) \cdot r^1,$$

а следующее за 1 машинное число есть

$$(1 \cdot r^{-1} + 0 \cdot r^{-2} + \dots + 0 \cdot r^{-(\ell-1)} + 1 \cdot r^{-\ell}) \cdot r^1 = \text{fl}(1 + \varepsilon),$$

то за *macheps* можно принять величину

$$\varepsilon = 1 \cdot r^{-\ell} \cdot r^1 = r^{1-\ell}.$$

Это число непосредственно связано с относительной погрешностью представления чисел в системе с плавающей запятой. Имеем:



$$\left| \frac{a - \text{fl}(a)}{a} \right| = \frac{\beta_{\ell+1} r^{-(\ell+1)} + \beta_{\ell+2} r^{-(\ell+2)} + \dots}{\beta_1 r^{-1} + \beta_2 r^{-2} + \dots} \leq \frac{1 \cdot r^{-\ell}}{\beta_1 \cdot r^{-1}} \leq r^{1-\ell} = \varepsilon. \quad (1.5)$$

Таким образом, машинный эpsilon<sup>\*)</sup> служит мерой относительной точности представления вещественных чисел, причем эта точность одинакова в любой части числового диапазона и зависит лишь от числа  $r$ -ичных разрядов, отводимых под мантиссу числа. В то же время оценка абсолютной погрешности

$$|a - \text{fl}(a)| \leq |a| \cdot r^{1-\ell}$$

показывает, что расстояние между вещественными числами и конечными приближениями к ним в системе с плавающей запятой неодинаковы в разных частях числового диапазона: абсолютная плотность машинных чисел больше вблизи нуля при одинаковой относительной плотности их распределения.

Заметим, что величина *macheps* служит оценкой относительной точности представления вещественного числа  $a$  при условии, что  $|a| > r^{-r^m}$ . Если же  $a \in [-r^{-r^m}, r^{-r^m}]$ , то  $\text{fl}(a) \equiv 0$  и, значит, относительная погрешность

$$\left| \frac{a - \text{fl}(a)}{a} \right| \equiv 1,$$

т.е. является постоянной достаточно большой величиной, в то время как абсолютная погрешность не превосходит величины  $r^{-r^m}$ .

Приведем значения введенных выше теоретических параметров для нескольких типов отечественных ЭВМ 60-80 гг.<sup>\*\*) (реальные параметры, естественно, могут незначительно отличаться от приводимых) [1, 2, 47].</sup>

<sup>\*)</sup> Если трактовать *macheps* как минимальное положительное действительное число, прибавление которого к 1 дает следующее за 1 число с плавающей запятой, то очевидно, при правильном округлении значение *macheps* будет в два раза меньшим.

Действительно, полагая  $\varepsilon = \frac{1}{2} r^{1-\ell} = \frac{r}{2} \cdot r^{-(\ell+1)} \cdot r^1$ , получаем

$$1 + \varepsilon = \left( 1 \cdot r^{-1} + 0 \cdot r^{-2} + \dots + 0 \cdot r^{-\ell} + \frac{r}{2} \cdot r^{-(\ell+1)} \right) \cdot r^1$$

и, значит,  $\text{fl}(1 + \varepsilon) = 1 + r^{1-\ell} > 1$ . "Мера дискретности" множества машинных чисел, как видно, остается той же:  $r^{1-\ell}$ .

<sup>\*\*) В качестве упражнения читателю предлагается составить несложную программу, которая позволит получить соответствующие параметры используемого им современного компьютера.</sup>

Так, для записи числа в 48-разрядном машинном слове БЭСМ-6 40 двоичных разрядов выделяются под мантиссу, 6 – под порядок и 2 – под знаки мантиссы (т.е. числа) и порядка. Отсюда, принимая  $r=2$ ,  $\ell=40$ ,  $m=6$ , получаем, что точность представления чисел с плавающей запятой на БЭСМ-6 не хуже  $2^{-39} (\approx 10^{-12})$ , граница машинного нуля  $2^{-64} (\approx 10^{-19})$ , машинной бесконечности  $2^{63} (\approx 10^{19})$ .

Машинное слово СМ ЭВМ имеет 32 двоичных разряда, из которых под мантиссу выделяется 24, а под порядок 7. Зная параметры  $r=2$ ,  $\ell=24$ ,  $m=7$ , получаем  $macheps = 2^{-23} \approx 10^{-7}$ , машинные нуль и бесконечность  $\approx 10^{\mp 38}$ .

На ЕС ЭВМ используется представление вещественных чисел по основанию  $r=16$ . Эти машины имеют относительную точность представления  $\approx 10^{-7}$  и диапазон для положительных чисел  $\approx 10^{-77} \div 10^{76}$ .

Практически любая машина ЕС ЭВМ рассчитана на то, чтобы выделить под запись числа двойное машинное слово, что позволяет более чем вдвое увеличивать точность представления. Двойная точность предусматривается также многими языками программирования.

Обращаясь к арифметическим операциям над машинными числами, прежде всего заметим, что они утрачивают привычные свойства. Особенно это касается свойств ассоциативности и дистрибутивности, нарушаемых при выполнении арифметических операций на любых ЭВМ. (Сохранение или несохранение свойства коммутативности связывают со способом округления чисел.) Так, весьма утрированный пример сравнения

$$(r^{D/2} \cdot r^{3D/4}) \cdot r^{-D/2} \quad \text{с} \quad r^{D/2} \cdot (r^{3D/4} \cdot r^{-D/2}),$$

где  $D$  таково, что  $r^D$  – правая граница числового диапазона, показывает существенность способа расстановки скобок при умножении: в первом случае машина выдаст сообщение о переполнении, в силу

$$r^{D/2} \cdot r^{3D/4} > r^D,$$

а во втором случае будет получен правильный результат. Легко также представить ситуацию с тремя положительными числами  $a$ ,  $b$ ,  $c$ , когда расставляя по-разному скобки в выражении  $a + b - c$ , будем получать (или не получать вовсе) разные результаты.

Изучение погрешностей результатов арифметических операций над числами с плавающей запятой производится с помощью представления

$$fl(a) = a(1 + \delta), \quad \text{где} \quad |\delta| \leq macheps \quad (1.6)$$

(чтобы убедиться в справедливости (1.6), достаточно ввести  $\delta$  равенством

$\delta = \frac{fl(a) - a}{a}$ , равносильным фигурирующему в (1.6), и воспользоваться доказанным в (1.5) неравенством  $|\delta| \leq macheps$ ). Принимая во внимание, что

операции над двумя машинными числами  $a$  и  $b$  производятся точно (здесь используется двойная длина машинного слова), после чего производится округление, результат любой арифметической операции  $\otimes$  также может быть записан в виде

$$fl(a \otimes b) = (a \otimes b)(1 + \delta_1), \quad (1.7)$$

где  $|\delta_1| \leq \varepsilon$  (за исключением особых случаев).

Пусть складываются последовательно три положительных числа  $a_1, a_2, a_3$ . Тогда, согласно (1.7),

$$fl(a_1 + a_2) = (a_1 + a_2)(1 + \delta_1),$$

где  $|\delta_1| \leq \varepsilon$ ;

$$\begin{aligned} fl((a_1 + a_2) + a_3) &= ((a_1 + a_2)(1 + \delta_1) + a_3)(1 + \delta_2) = \\ &= (a_1 + a_2)(1 + \delta_1)(1 + \delta_2) + a_3(1 + \delta_2), \end{aligned}$$

где  $|\delta_i| \leq \varepsilon$  ( $i = 1, 2$ ).

Заменяя здесь  $\delta_i$  бóльшим значением  $\varepsilon$ , получим оценку абсолютной погрешности суммы трех слагаемых:

$$|fl(a_1 + a_2 + a_3) - (a_1 + a_2 + a_3)| \leq 2(a_1 + a_2)\varepsilon + a_3\varepsilon + (a_1 + a_2)\varepsilon^2.$$

Обращает на себя внимание неравноправность слагаемых в образовании погрешности суммы: меньшую роль в ней играет слагаемое, прибавляемое последним. Природа этого факта очевидна: первые слагаемые неявно (в просуммированном виде) участвуют в процессе каждого последующего сложения.

Если пренебречь степенями  $\varepsilon$  выше первой, то для суммы  $n$  положительных чисел  $a_i$  нетрудно получить [2] приближенную оценку абсолютной погрешности вида<sup>\*)</sup>

$$\left| fl\left(\sum_{i=1}^n a_i\right) - \sum_{i=1}^n a_i \right| \leq (n-1)(a_1 + a_2) + (n-2)a_3 + \dots + 2a_{n-1} + a_n \cdot \varepsilon$$

при последовательном суммировании, начинающемся с  $a_1$ . Очевидно, чтобы эта погрешность была минимальной, *последовательность чисел нужно суммировать в порядке возрастания членов*. Только за счет этого можно добиться уменьшения погрешности, как показано в [15], в  $n/\log_2 n$  раз.

На основе изучения погрешности произведения нескольких чисел строятся алгоритмы оптимального умножения. Пусть требуется перемножить  $n$  чисел  $a_i$  таких, что  $|a_1| \leq |a_2| \leq \dots \leq |a_n|$ . Погрешность произведения

<sup>\*)</sup> Здесь и далее знак  $\leq$  используется для обозначения неравенства в смысле главных (линейных) частей.

будет минимальной, если находить его по схеме: умножать  $a_1$  последовательно на  $a_n, a_{n-1}, \dots$  до тех пор, пока модуль частичного произведения не станет большим единицы, затем это частичное произведение умножить на  $a_2, a_3, \dots$  до тех пор, пока новое частичное произведение не станет по модулю меньше единицы, и так далее до исчерпывания всех сомножителей [15].

Естественно, что расплатой за выигрыш в точности при реализации таких алгоритмов будет проигрыш в скорости счета.

#### 1.4. ПРИМЕРЫ НЕУСТОЙЧИВЫХ ЗАДАЧ И МЕТОДОВ

В силу неизбежного появления погрешностей в исходных данных задачи (в процессе создания математической модели изучаемого объекта или явления), а также погрешностей округления при ее решении, следует иметь представление о том, насколько чувствительными могут оказаться сами задачи и методы их решения к таким погрешностям. Рассмотрим несколько примеров проявления чрезмерной чувствительности.

а) Пусть требуется найти вещественное решение уравнения

$$(x - a)^n = \varepsilon,$$

где  $\varepsilon$  — очень малое положительное число, а натуральное  $n$  достаточно велико. Тогда естественно заменить  $\varepsilon$  нулем и положить  $x \approx a$ . Так как точное решение данного уравнения есть  $x = a + \sqrt[n]{\varepsilon}$ , то абсолютная погрешность при таком подходе составит величину  $\sqrt[n]{\varepsilon}$ . Много это или мало — судить об этом можно, придавая  $\varepsilon$  и  $n$  численные значения. Например, взяв  $\varepsilon = 10^{-10}$ ,  $n=10$ , получим абсолютную погрешность значения корня  $x \approx a$ , равную 0.1. Относительная погрешность при этом может оказаться сколь угодно большой, если  $a$  взять сколь угодно малым.

б) Пример Уилкинсона [1, 2, 26, 53 и др.]. Многочлен

$$P_{20}(x) \equiv (x-1)(x-2)\dots(x-20) \equiv x^{20} - 210x^{19} + \dots + 20!$$

имеет 20 хорошо отделимых действительных корней

$$x_1 = 1, \quad x_2 = 2, \quad \dots, \quad x_{20} = 20.$$

Предположим, что только в одном его коэффициенте, а именно, при  $x^{19}$  сделана ошибка порядка *macheps*: вместо -210 в развернутый вид многочлена  $P_{20}(x)$  подставлено число  $-(210 + 2^{-23}) \approx -(210 + 10^{-7})$ . Полученный при этом так называемый *возмущенный многочлен* будет иметь следующие корни (ограничимся записью трех цифр после запятой):

$$x_1 \approx 1.000, \quad x_6 \approx 6.000, \quad x_{12,13} \approx 11.794 \pm 1.652i,$$

$$\begin{array}{lll}
 x_2 \approx 2.000, & x_7 \approx 7.000, & x_{14,15} \approx 13.992 \pm 2.519i, \\
 x_3 \approx 3.000, & x_8 \approx 8.007, & x_{16,17} \approx 16.731 \pm 2.813i, \\
 x_4 \approx 4.000, & x_9 \approx 8.917, & x_{18,19} \approx 19.502 \pm 1.940i, \\
 x_5 \approx 5.000, & x_{10,11} \approx 10.095 \pm 0.644i, & x_{20} \approx 20.847.
 \end{array}$$

Как видим, весьма малое возмущение (сопоставимое с точностью представления чисел в некоторых ЭВМ) всего лишь в одном коэффициенте даже качественно изменило набор корней данного многочлена: половина из них перестали быть действительными.

в) Линейная система

$$\begin{cases} x + 10y = 11, \\ 100x + 1001y = 1101 \end{cases}$$

имеет единственное решение  $x=1, y=1$ . Допустив абсолютную погрешность в 0.01 в правой части одного уравнения, получим *возмущенную систему*

$$\begin{cases} x + 10y = 11.01, \\ 100x + 1001y = 1101 \end{cases}$$

с единственным решением  $x=11.01, y=0$ . Последнее никак не назовешь близким к решению исходной системы.

г) Для вычисления определенных интегралов вида

$$I_n := \int_0^1 x^n e^x dx,$$

где  $n \in \mathbb{N}$ , с помощью метода интегрирования "по частям" легко вывести рекуррентную формулу

$$I_n = 1 - n I_{n-1}; \quad n = 1, 2, \dots; \quad I_0 = 1 - \frac{1}{e}. \quad (1.8)$$

В [4] приведена сводная таблица результатов подсчета  $I_n$  при  $n=0, 1, 2, \dots, 14$  по формуле (1.8), полученных в 60-х годах на разных вычислительных машинах в Чехословакии, Восточной Германии, России. Эта таблица наглядно демонстрирует рост разброса значений интеграла при увеличении значения  $n$ . Не приводя здесь полностью этих результатов, отметим лишь, что при  $n=14$  спектр значений  $I_{14}$  лежит в границах от  $-148$  до  $5356$ . Несмотря на то, что на некоторых ЭВМ вычисления велись с мантиссами, имеющими более десяти десятичных разрядов, ни один результат не имел ни одной верной цифры! Причина подобного явления – в численной неустойчивости схемы (1.8). Безупречная в теоретическом пла-

не, т.е. с точки зрения аналитической математики, она совершенно непригодна с позиций вычислительной математики, поскольку неизбежная погрешность стартового значения  $I_0$  при подсчете  $I_n$  увеличивается в  $n!$  раз, т.е. катастрофически нарастает.

Если примеры а)–в) указывали прямо на существование *неустойчивых задач*<sup>\*)</sup>, то в г) мы видим яркий пример *неустойчивого метода* вычислений. Последним, как правило, есть альтернативы. Например, для вычисления рассматриваемого в г) интеграла  $I_n$  при конкретном значении  $n$  можно применять квадратурные формулы, т.е. формулы, специально приспособленные для приближенного вычисления определенных интегралов. А можно воспользоваться и равенством (1.8), только переписать его в виде

$$I_{n-1} = \frac{1}{n}(1 - I_n). \quad (1.9)$$

Учитывая, что  $0 < I_{n+1} < I_n$  и  $I_n \xrightarrow{n \rightarrow \infty} 0$ , для подсчета  $I_n$  при некотором фиксированном  $n=k$  (и для меньших этого  $k$  индексов) можно задать значением  $I_N := 0$  и вести счет по формуле (1.9) при  $n = N, N-1, \dots, k+1$ . Так как начальная погрешность на каждом шаге теперь уменьшается в  $n$  раз, то такой алгоритм будет *численно устойчивым*. Значение  $N$  при этом может быть определено теоретически или экспериментально [2, 4].

Более подробно тема устойчивости и неустойчивости задач и методов развивается при изучении численных процессов решения дифференциальных уравнений. Для решения задач, требующих особого учета возмущений и обобщения самого понятия решения, в последние 10-20 лет рассматриваются особые подходы, приводящие к построению специальных устойчивых численных методов, называемых *методами регуляризации* [45, 49].

## 1.5. ОБУСЛОВЛЕННОСТЬ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ СИСТЕМ

Учитывая распространенность систем линейных алгебраических уравнений (ибо часто именно к ним сводится на определенном этапе процесс математического моделирования), попытаемся количественно охарактеризовать степень неопределенности этих задач. Знание таких характеристик позволяет обоснованно судить о корректности моделей, грамотно подбирать методы и строить алгоритмы, правильно трактовать полученные результаты.

<sup>\*)</sup> См. также пример неустойчивой задачи на собственные значения для  $n \times n$ -матрицы в п.4.5

Рассмотрим линейную алгебраическую систему, записанную в виде векторно-матричного уравнения

$$Ax = b, \quad (1.10)$$

где  $A$  – невырожденная  $n \times n$ -матрица коэффициентов данной системы;  $b$  – ненулевой  $n$ -мерный вектор свободных членов;  $x$  –  $n$ -мерный вектор неизвестных (решение, если трактовать (1.10) как верное равенство).

Пусть правая часть (1.10) получила приращение ("возмущение")  $\Delta b$ , т.е. вместо истинного вектора  $b$  используется приближенный вектор  $b + \Delta b$ . Реакцией решения  $x$  на возмущение  $\Delta b$  правой части будет вектор поправок  $\Delta x$ , т.е. если  $x$  – решение (1.10), то  $x + \Delta x$  – решение уравнения

$$A(x + \Delta x) = b + \Delta b. \quad (1.11)$$

Понимая под абсолютной погрешностью приближенного вектора норму<sup>\*)</sup> разности между точным и приближенным векторами, а под относительной погрешностью – отношение абсолютной погрешности к норме вектора (точного или приближенного), выясним связь между относительными погрешностями вектора свободных членов и вектора-решения. Иначе, получим оценку вида

$$\frac{\|\Delta x\|}{\|x\|} \leq (?) \frac{\|\Delta b\|}{\|b\|},$$

где  $\|\bullet\|$  – какая-либо векторная норма, а (?) – неизвестный пока коэффициент связи.

Подставляя (1.10) в (1.11), видим, что поправка  $\Delta x$  связана с возмущением  $\Delta b$  таким же, как и (1.10), равенством

$$A \Delta x = \Delta b,$$

из которого находим ее явное выражение

$$\Delta x = A^{-1} \Delta b. \quad (1.12)$$

Нормируя равенства (1.10) и (1.12), будем иметь

$$\|b\| \leq \|A\| \cdot \|x\| \quad \text{и} \quad \|\Delta x\| \leq \|A^{-1}\| \cdot \|\Delta b\|,$$

где матричная норма должна быть согласованной с выбранной векторной нормой. Эти два числовых неравенства одинакового смысла можно перемножить:

$$\|b\| \cdot \|\Delta x\| \leq \|A\| \cdot \|x\| \cdot \|A^{-1}\| \cdot \|\Delta b\|.$$

Из последнего делением на  $\|b\| \cdot \|x\|$  получаем искомую связь

$$\frac{\|\Delta x\|}{\|x\|} \leq \|A\| \cdot \|A^{-1}\| \cdot \frac{\|\Delta b\|}{\|b\|}. \quad (1.13)$$

<sup>\*)</sup> См. приложение 1

Положительное число  $\|A\| \cdot \|A^{-1}\|$  – коэффициент этой связи – называют **числом (мерой) обусловленности** матрицы  $A$  и обозначают  $\text{cond } A$  (от английского слова *conditioned* – "обусловленный"). Распространены также обозначения  $\nu(A)$  и  $\chi(A)$ .

Легко показать, что то же самое число  $\text{cond } A = \|A\| \cdot \|A^{-1}\|$  служит коэффициентом роста относительных погрешностей при неточном задании элементов матрицы  $A$  в (1.10). А именно, если матрица  $A$  получила возмущение  $\Delta A$  и  $x + \Delta x$  – решение возмущенной системы

$$(A + \Delta A)(x + \Delta x) = b,$$

то справедливы неравенства\*)

$$\frac{\|\Delta x\|}{\|x\|} \leq \text{cond } A \cdot \frac{\|\Delta A\|}{\|A + \Delta A\|} \quad \text{и} \quad \frac{\|\Delta x\|}{\|x + \Delta x\|} \leq \text{cond } A \cdot \frac{\|\Delta A\|}{\|A\|}. \quad (1.14)$$

Итак, неравенства (1.13) и (1.14) показывают, что *чем больше число обусловленности, тем сильнее сказывается на решении линейной системы ошибка в исходных данных*. Грубо говоря, если  $\text{cond } A = O(10^p)$  и исходные данные имеют погрешность в  $\ell$ -м знаке после запятой, то независимо от способа решения системы (1.10) в результате можно гарантировать не более  $\ell - p$  знаков после запятой.

Если число  $\text{cond } A$  велико, то система считается плохо обусловленной. Говорить о том, "что такое хорошо и что такое плохо", в отрыве от контекста решаемой задачи почти бессмысленно, так как здесь может играть роль размерность задачи, точность, с которой должно быть найдено ее решение, точность представления чисел в ЭВМ и т.п. Однако можно дать оценку снизу числа обусловленности. А именно, если используются подчиненные матричные нормы (для которых норма единичной матрицы есть единица), то, очевидно,

\*) Можно получить оценки, обобщающие в определенном смысле (1.13) и (1.14) в случае одновременного возмущения матрицы  $A$  системы (1.10) и ее правой части  $b$ . Более того, величина  $\nu(A) = \|A\| \cdot \|A^{-1}\|$  является также мерой обусловленности линейного обратимого оператора  $A$  в произвольном нормированном пространстве. Справедливо следующее утверждение [51].

Пусть  $Ax = b$  – данное, а  $\tilde{A}x = \tilde{b}$  – возмущенные линейные операторные уравнения с относительными уровнями возмущений  $\delta_A \geq \frac{\|A - \tilde{A}\|}{\|A\|}$  и  $\delta_b \geq \frac{\|b - \tilde{b}\|}{\|b\|}$ . Тогда если  $\delta_A \nu(A) < 1$ , то эти уравнения одновременно однозначно разрешимы и справедлива оценка относительной погрешности решения

$$\frac{\|x - \tilde{x}\|}{\|x\|} \leq \frac{\nu(A)}{1 - \delta_A \cdot \nu(A)} \cdot \delta_b + \frac{\nu^2(A)}{1 - \delta_A \cdot \nu(A)} \cdot \delta_A.$$



$$\text{cond } A = \|A\| \cdot \|A^{-1}\| \geq \|A \cdot A^{-1}\| = \|E\| = 1,$$

т.е. число обусловленности не может быть меньше 1. Можно также указать верхнюю границу для чисел обусловленности, превышение которой при решении линейных систем на конкретной ЭВМ может привести к заведомо ложным результатам. Так, решение считается ненадежным, если  $\text{cond } A \geq (\text{macheps})^{-1}$  или даже  $\text{cond } A \geq (\text{macheps})^{-0.5}$  [22]. При этом заметим, что масштабированием матрицы  $A$  путем умножения на скаляр  $\alpha$  ее обусловленность не улучшить, ибо

$$\text{cond}(\alpha A) = \|\alpha A\| \cdot \|(\alpha A)^{-1}\| = \|A\| \cdot \|A^{-1}\| = \text{cond } A.$$

Классическим примером плохо обусловленной матрицы является так называемая *матрица Гильберта*

$$H_n = \left( \frac{1}{i+j-1} \right)_{i,j=1}^n,$$

демонстрирующая катастрофическое возрастание числа обусловленности с ростом размерности [41, 55]. Так, уже при  $n=8$   $\text{cond } H_8 > 10^{10}$  и обратная матрица  $H_8^{-1}$ , полученная на машине с точностью представления чисел порядка  $10^{-8}$ , может не содержать ни одного верного знака.

Очевидно, число обусловленности зависит от выбора матричной нормы (индуцированной, как правило, той или иной векторной нормой, в терминах которой характеризуется относительная погрешность решения алгебраической системы). Однако нетрудно получить оценку числа обусловленности через собственные числа матрицы. Действительно, пусть собственные числа  $\lambda_i$  матрицы  $A$  упорядочены по модулю:

$$|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|,$$

т.е. *спектральный радиус* матрицы  $A$  есть  $\rho_A = |\lambda_1|$ . Тогда в силу известного неравенства  $\rho_A \leq \|A\|$  и соотношения между собственными числами прямой и обратной матриц, имеем

$$\|A\| \cdot \|A^{-1}\| \geq \rho_A \cdot \rho_{A^{-1}} = |\lambda_1| \cdot \frac{1}{|\lambda_n|}.$$

Таким образом, оценкой снизу меры обусловленности матрицы  $A$  может служить величина  $\left| \frac{\lambda_1}{\lambda_n} \right|$  (называемая иногда *числом обусловленности Тодда* [РЖМат, 1983, 10Б937]). Для симметричных матриц эта оценка и на самом деле является числом обусловленности, соответствующим спектральной норме матрицы (индуцированной евклидовой нормой вектора). Учитывая смысл собственных чисел матрицы, можно сказать, что число обусловленности показывает величину отношения наибольшего коэффи-

циента растяжения вектора посредством линейного преобразования  $A$  к наименьшему.

Следует отметить, что непосредственный подсчет чисел обусловленности, особенно при большой размерности матриц, является весьма дорогостоящим делом из-за необходимости обращать матрицы или находить их собственные значения. Поэтому зачастую о приемлемости порядка возможного роста относительной погрешности результата решения какой-либо алгебраической задачи относительно данной матрицы судят либо по каким-то достаточным признакам (например, по доминированию элементов главной диагонали матрицы), либо на основе теоретического изучения матрицы [41, 46], либо путем применения специальных алгоритмов приближенного оценивания  $\text{cond } A$  [22]. Исследование матриц на обусловленность может быть естественным образом связано со способом решения той или иной алгебраической задачи относительно данной матрицы.

Прокомментируем теперь пример неустойчивой системы, приведенной в примере в) п.1.4.

Матрица коэффициентов системы  $A = \begin{pmatrix} 1 & 10 \\ 100 & 1001 \end{pmatrix}$  имеет обратную

$A^{-1} = \begin{pmatrix} 1001 & -10 \\ -100 & 1 \end{pmatrix}$ . Следовательно, число обусловленности в матричной норме, индуцированной векторной нормой-максимум (иначе, нормой  $l_\infty$ ), есть

$$v_\infty(A) = 1101 \cdot 1011 = 1113111 > 10^6.$$

Учитывая, что в данном примере  $b = \begin{pmatrix} 11 \\ 1101 \end{pmatrix}$ ,  $\Delta b = \begin{pmatrix} 0.01 \\ 0 \end{pmatrix}$ , на основе (1.13) получаем оценку относительной погрешности решения в  $l_\infty$ -нормах

$$\delta_x \leq v_\infty(A) \cdot \delta_b = 1113111 \cdot \frac{0.01}{1101} = 10.11.$$

Так как норма-максимум решения  $x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$  равна 1, то оценка абсолютной погрешности  $\|\Delta x\|$  решения равна

$$\|\Delta x\| = \|x\| \cdot \delta_x \leq 10.11.$$

Как видим, решение  $x + \Delta x = \begin{pmatrix} 11.01 \\ 0 \end{pmatrix}$  возмущенной системы вписывается в оценку

$$\|x + \Delta x\| \leq \|x\| + \|\Delta x\| \leq 1 + 10.11 = 11.11.$$

Аналогичный результат может быть получен через число обусловленности Тодда. Решая характеристическое уравнение

$$\lambda^2 - 1002\lambda + 1 = 0,$$

находим собственные числа матрицы  $A$ :  $\lambda_1 \approx 1002$  и  $\lambda_2 \approx 0.000998$ , дающие оценку

$$\text{cond } A \geq \frac{\lambda_1}{\lambda_2} \approx 1004000 > 10^6.$$

На данном примере также можно наглядно убедиться в том, что малость невязки  $r = b - A\tilde{x}$  плохо обусловленной системы еще не говорит о близости приближенного решения  $\tilde{x}$  к точному  $x$ . Действительно, невязки  $r_1$  и  $r_2$  векторов  $\tilde{x}_1 = \begin{pmatrix} 11.01 \\ 0 \end{pmatrix}$  и  $\tilde{x}_2 = \begin{pmatrix} 1 \\ 1.1 \end{pmatrix}$  для основной (невозмущенной) системы из примера в) п.1.4 имеют нормы соответственно  $\|r_1\|_\infty = 0.01$  и  $\|r_2\|_\infty = 100.1$ . Вектор  $\tilde{x}_2$ , явно более близкий к точному решению  $x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ , чем  $\tilde{x}_1$ , имеет существенно большую невязку!

Двумерный случай допускает простую геометрическую трактовку понятия обусловленности. Плохая обусловленность системы двух уравнений с двумя неизвестными означает, что прямые, являющиеся геометрическими образами уравнений, пересекаются на координатной плоскости под очень острым углом. В этом случае небольшое искажение в данных, интерпретируемое как параллельный перенос (при возмущении свободного члена) или поворот прямых (при возмущении матрицы коэффициентов), приводит к значительному перемещению их точки пересечения, т.е. геометрического образа решения.

## 1.6. ПОГРЕШНОСТИ КОРНЕЙ СКАЛЯРНЫХ УРАВНЕНИЙ С ПРИБЛИЖЕННЫМИ КОЭФФИЦИЕНТАМИ

Пусть решается алгебраическое или трансцендентное уравнение вида

$$f(x) = 0. \quad (1.15)$$

Это уравнение, являясь нелинейным, может иметь один или несколько корней или не иметь их вовсе. Каждый корень так или иначе зависит от числовых данных уравнения (1.15). Эти числовые данные (будем иногда называть их коэффициентами уравнения) можно считать фиксированными значениями параметров, т.е. уравнение (1.15) целесообразно рассматривать в качестве представителя семейства уравнений

$$f(x, a_1, a_2, \dots, a_m) = 0 \quad (1.16)$$

и наличие корней у (1.15) связывать с существованием неявной функции  $x = \varphi(a_1, a_2, \dots, a_m)$  при данном наборе значений параметров (коэффициентов)  $a_1, a_2, \dots, a_m$ . Поскольку эти числа, как правило, точно не известны (грубость модели, неточность измерений, усечение чисел при вводе в ЭВМ

и т.п.), встает вопрос о том, как влияет погрешность коэффициентов уравнения (1.15) на погрешности его корней. Иначе, с какой точностью имеет смысл решать данное уравнение, если известно, что его коэффициенты не точны, но имеется информация об уровне их погрешностей?

Реально можно оценить лишь главную часть погрешности корня, принимая под ней, как и в п.1.1, модуль дифференциала. В данном случае речь идет о дифференциале неявной функции нескольких переменных. Подходящую основу для этого находим в математическом анализе: если, например, выполняются условия, при которых уравнение  $F(x, y, z) = 0$  неявно определяет дифференцируемую функцию  $z = z(x, y)$ , то выражение дифференциала  $dz$  этой функции можно получить из равенства  $dF = 0$ , а именно,

$$dz = -\frac{1}{F'_z} (F'_x dx + F'_y dy).$$

Аналогично, если  $x = \varphi(a_1, a_2, \dots, a_m)$  – корень уравнения (1.16), то линейная часть его изменения, соответствующая изменениям аргументов – коэффициентов  $a_1, a_2, \dots, a_m$ , равна

$$dx = -\frac{1}{\frac{\partial f}{\partial x}} \left( \frac{\partial f}{\partial a_1} da_1 + \frac{\partial f}{\partial a_2} da_2 + \dots + \frac{\partial f}{\partial a_m} da_m \right).$$

Переходя здесь к модулям и заменяя истинные абсолютные погрешности коэффициентов  $\Delta a_i (= da_i)$  их оценками  $\Delta a_i$ , получим *формулу для оценки абсолютной погрешности корня*:

$$|\Delta x| \leq \frac{1}{\left| \frac{\partial f}{\partial x} \right|} \left( \left| \frac{\partial f}{\partial a_1} \right| \Delta a_1 + \left| \frac{\partial f}{\partial a_2} \right| \Delta a_2 + \dots + \left| \frac{\partial f}{\partial a_m} \right| \Delta a_m \right). \quad (1.17)$$

Заметим, что при вычислении значений частных производных в (1.17) следует пользоваться фиксированными значениями коэффициентов  $a_1, a_2, \dots, a_m$  (такими, какими они используются при нахождении корня) и приближенным значением того корня, степень неопределенности которого устанавливается. Ясно, что для разных корней одного и того же уравнения значения этой величины могут сильно различаться.

Следуя [23], получаемую с помощью (1.17) оценку будем называть *безусловной абсолютной погрешностью* приближенного корня  $x$  уравнения (1.16) с приближенными коэффициентами и обозначать б.п.х. Это же  $x$  есть точный корень уравнения (1.15), в то время как при численном решении уравнения вместо  $x$  будет получено некоторое приближение к нему  $\bar{x}$ . Возникающую при этом остаточную погрешность или погрешность метода назовем *условной погрешностью* и обозначим у.п.х.

Итак, в ходе численного решения уравнения (1.15) в предположении, что его коэффициенты точны, находится приближенный корень  $\bar{x}$  с условной погрешностью у.п.х. Это означает, что точный корень  $x$  уравне-

ния (1.15) при данном предположении лежит на интервале  $(\bar{x} - \text{у.п.}x, \bar{x} + \text{у.п.}x)$ . Зная погрешности коэффициентов и приближенный корень  $\bar{x}$ , можно подсчитать безусловную погрешность б.п.х. Истинная величина корня  $x^*$  уравнения с приближенными коэффициентами – это потенциально любое число из интервала  $(x - \text{б.п.}x, x + \text{б.п.}x)$ . Ясно, что гипотетическое значение  $x^*$  – решение рассчитываемой модели – может отличаться от реально полученного значения  $\bar{x}$  на величину  $\text{у.п.}x + \text{б.п.}x$ , т.е. *полная погрешность корня уравнения с приближенными коэффициентами складывается из погрешностей условной и безусловной*. Так как величиной условной погрешности распоряжается вычислитель, причем ее уменьшение в разумных пределах не вызывает, как правило, больших затруднений, то обычно, чтобы не увеличивать полную погрешность корня приближенного уравнения, условную погрешность задают неравенством

$$\text{у.п.}x \leq \text{б.п.}x \quad (1.18)$$

или даже

$$\text{у.п.}x \leq 0.1 \text{ б.п.}x.$$

Если воспользоваться *правилом Ньютона*<sup>\*)</sup> оценки близости приближенного корня (нуля)  $\bar{x}$  дифференцируемой функции  $f(x)$  к точному корню  $x$

$$|\bar{x} - x| \leq \left| \frac{f(\bar{x})}{f'(\bar{x})} \right| \quad (1.19)$$

и оценкой (1.17), то можно на основе (1.18) получить простой критерий окончания процесса численного решения скалярного уравнения (1.15) с приближенными коэффициентами:

*уточнение корня  $\bar{x} \approx x$  ведется до тех пор, пока не выполнится условие*

$$|f(\bar{x})| \leq \left| \frac{\partial f}{\partial a_1} \right| \Delta a_1 + \left| \frac{\partial f}{\partial a_2} \right| \Delta a_2 + \dots + \left| \frac{\partial f}{\partial a_m} \right| \Delta a_m. \quad (1.20)$$

Как видим, большую роль в неравенстве (1.20) играют модули частных производных данной функции по параметрам, называемые *коэффициентами чувствительности*.

В случае, когда уравнение (1.15) – алгебраическое, т.е.

---

<sup>\*)</sup> По формуле конечных приращений Лагранжа при некотором  $\Theta \in (\bar{x}; x)$  (или  $(x; \bar{x})$ )

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

Предполагая, что  $x$  – корень,  $\bar{x} \approx x$ , а значит и  $\Theta \approx \bar{x}$ , имеем

$$f(\bar{x}) \approx f'(\bar{x})(\bar{x} - x),$$

откуда следует правило Ньютона (1.19).

$$f(x) \equiv P(x) := a_0 x^n + a_1 x^{n-1} + \dots + a_{n-1} x + a_n, \quad (1.21)$$

его корни – функции коэффициентов многочлена  $a_0, a_1, \dots, a_n$ , а коэффициенты чувствительности суть числа  $|x^n|, |x^{n-1}|, \dots, |x|, 1$ , получающиеся при подстановке сюда вместо  $x$  приближенных корней многочлена. Неравенство типа (1.20), закладываемое в процесс уточнения корня  $\bar{x}$  многочлена (1.21) с приближенными коэффициентами  $a_0(\pm \Delta_{a_0}), a_1(\pm \Delta_{a_1}), \dots, a_n(\pm \Delta_{a_n})$ , имеет вид

$$|P(\bar{x})| \leq |\bar{x}^n| \Delta_{a_0} + |\bar{x}^{n-1}| \Delta_{a_1} + \dots + |\bar{x}| \Delta_{a_{n-1}} + \Delta_{a_n}.$$

Легко заметить, какую роль играют погрешности тех или иных коэффициентов в зависимости от того, малы или велики модули корней.

Получим, наконец, связь между относительными погрешностями коэффициентов и корней многочлена.

На основе оценки (1.17) для многочлена (1.21) имеем

$$\delta_x = \left| \frac{\Delta x}{x} \right| \leq \frac{\sum_{i=0}^n \left| \frac{\partial P}{\partial a_i} \right| \Delta_{a_i}}{|x| \cdot \left| \frac{\partial P}{\partial x} \right|} = \frac{\sum_{i=0}^n |x^{n-i}| \cdot |a_i| \cdot \delta_{a_i}}{|x| \cdot \left| \sum_{i=0}^{n-1} (n-i) a_i x^{n-1-i} \right|} = \frac{\sum_{i=0}^n |a_i x^{n-i}| \delta_{a_i}}{\left| \sum_{i=0}^{n-1} (n-i) a_i x^{n-1-i} \right|}.$$

Если сделать допущение, что все коэффициенты многочлена имеют одинаковый уровень  $\delta$  относительных погрешностей, или положить  $\max_{i=0, \dots, n} \delta_{a_i} = \delta$ , то для граничной относительной погрешности  $\delta_x$  простого

ненулевого корня  $x$  будем иметь приближенную оценку

$$\delta_x \lesssim \frac{\sum_{i=0}^n |a_i x^{n-i}|}{\left| \sum_{i=0}^{n-1} (n-i) a_i x^{n-1-i} \right|} \cdot \delta. \quad (1.22)$$

Коэффициент

$$v_x = \frac{\sum_{i=0}^n |a_i x^{n-i}|}{\left| \sum_{i=0}^{n-1} (n-i) a_i x^{n-1-i} \right|}$$

связи (1.22) относительных погрешностей корней и коэффициентов (для каждого корня свой) по аналогии с терминологией предыдущего пункта

естественно назвать *числом* или *мерой обусловленности ненулевого простого корня*<sup>\*)</sup> многочлена с приближенными коэффициентами.

Обращаясь к примеру Уилкинсона в п.1.46), с помощью подсчета чисел обусловленности корней можно грубо объяснить картину разного влияния на разные корни внесенного малого возмущения. Если для первого корня число обусловленности  $\nu_x \approx 400$ , то с увеличением номера корня это число значительно возрастает, достигая, например, на пятнадцатом корне величины большей, чем  $10^{10}$ . Последний (двадцатый) корень имеет меньшее, чем предыдущий, число обусловленности, сравнимое с величиной  $1/macheps$ . (Значения несколько иначе определенных чисел обусловленности можно найти в [1].)

Отметим, что в числителе выражения для подсчета числа обусловленности  $\nu_x$  могут отсутствовать слагаемые, соответствующие теоретически точным и при этом точно реализуемым при вводе в ЭВМ значениям коэффициентов многочлена. Это видно из оценки  $\delta_x$ , где некоторые из  $\Delta_{a_i}$  могут быть тождественно равными нулю (т.е. соответствующие коэффициенты не считаются параметрами приближенного уравнения (1.16)).

## УПРАЖНЕНИЯ

1.1. Найти приближенное значение выражения

$$\frac{3.7894 \cdot 0.29^2}{5.63}$$

если известно, что входящие в него числа (за исключением показателя степени) – приближенные, записанные в соответствии с правилом Крылова.

Сделать теоретическую оценку главной части абсолютной погрешности результата.

1.2. Пусть машинное слово некой ЭВМ состоит из 16 четверичных разрядов.

\*) В случае нулевого или кратного корня  $x$  знаменатель дроби обращается в нуль.

В [26] мерой обусловленности простого корня  $x = x(a_1, a_2, \dots, a_n)$  многочлена (1.21) с  $a_0 = 1$  считается норма градиента  $\left( \frac{\partial x}{\partial a_1}, \frac{\partial x}{\partial a_2}, \dots, \frac{\partial x}{\partial a_n} \right)$ . Если все корни  $x_1, x_2, \dots, x_n$  – простые, то обусловленность многочлена характеризуется нормой матрицы Якоби  $\left( \frac{\partial x_i}{\partial a_j} \right)_{i,j=1}^n$  (см приложение 2)

а) Найти приближенно значения машинного нуля  $M_0$ , машинной бесконечности  $M_\infty$  и машинного эpsilon в системе представления с плавающей запятой, если под мантиссе выделяется 11 разрядов.

б) Оценить абсолютную погрешность и диапазон представимости чисел с фиксированной запятой, если под целую и дробную части числа выделяется поровну разрядов.

(Считается, что для отображения знаков  $+$ ,  $-$  и  $,$  используются те же четверичные элементы и что округление производится на основе простого отбрасывания).

1.3. Обращением матрицы коэффициентов решить систему:

$$\begin{cases} x_1 - x_3 = 2, \\ x_1 + x_2 + x_3 = 6, \\ x_2 + 3x_3 = 5. \end{cases}$$

Найдя число обусловленности, оценить возможные относительное и абсолютное отклонение от полученного решения, если допустить относительную ошибку в правой части порядка 0.01 (по некоторой норме).

1.4. Пусть коэффициенты уравнения

$$x^2 - 6.9x - 0.0212 = 0$$

являются приближенными числами, записанными в соответствии с принципом Крылова (коэффициент при  $x^2$  и правую часть считаем числами точными).

а) Найти корни уравнения с максимальной разумной точностью, диктуемой степенью неопределенности коэффициентов.

б) Подсчитать числа обусловленности для найденных корней и на этой основе дать приближенную оценку их относительных и абсолютных погрешностей.



## ГЛАВА 2 || РЕШЕНИЕ ЛИНЕЙНЫХ || АЛГЕБРАИЧЕСКИХ СИСТЕМ || (ПРЯМЫЕ МЕТОДЫ)

*Рассматриваются простые, часто применяемые численные методы решения систем линейных алгебраических уравнений с квадратными матрицами коэффициентов. Попутно решаются задачи обращения матриц и вычисления определителей. Наряду с широко известным методом Гаусса, изучаемым здесь в плане реальных вычислений, вниманию читателя предлагаются также метод вращений (более устойчивый к погрешностям арифметических действий, чем метод Гаусса), метод квадратных корней для решения систем с симметричными матрицами и метод прогонки для систем с трехдиагональными матрицами (т.е. для трехточечных разностных уравнений второго порядка).*

### 2.0. ВВЕДЕНИЕ

Как утверждается в книге известного американского математика Вальяха [11], 75% всех расчетных математических задач приходится на решение систем линейных алгебраических уравнений (в дальнейшем, СЛАУ); Это не удивительно, так как математические модели тех или иных явлений или процессов либо сразу строятся как линейные алгебраические, либо сводятся к таковому посредством дискретизации и/или линеаризации. Поэтому трудно переоценить роль, какую играет выбор эффективного (в том или ином смысле) способа решения СЛАУ. Современная вычислительная математика располагает большим арсеналом методов, а математическое обеспечение ЭВМ – многими пакетами прикладных программ, позволяющих решать различные возникающие на практике линейные системы. Чтобы ориентироваться в этом море методов и программ и в нужный момент сделать оптимальный выбор, нужно разбираться в основах построений методов и алгоритмов, учитывающих специфику постановок задач, знать их сильные и слабые стороны и границы применимости.

Все методы решения линейных алгебраических задач (наряду с задачей решения СЛАУ, это и вычисление определителей, и обращение матриц, и задачи на собственные значения) можно разбить на два класса: прямые и итерационные. Как явствует из заглавия, здесь будут рассмотрены только *прямые методы*, т.е. такие методы, которые приводят к решению



теоретическом плане и приемлемых на практике при небольших  $n$  (2, 3) **формул Крамера**

$$x_i = \frac{\det A_i}{\det A} \quad (i = 1, 2, \dots, n),$$

позволяющих находить неизвестные в виде дробей, знаменателем которых является определитель матрицы системы, а числителем – определители матриц  $A_i$ , полученные из  $A$  заменой столбца коэффициентов при вычисляемом неизвестном столбцом свободных членов. Если при реализации этих формул определители вычисляются понижением порядка на основе разложения по элементам какой-нибудь строки или столбца матрицы, то на вычисление определителя  $n$ -го порядка будет затрачиваться  $n!$  операций умножения. Факториальный рост количества арифметических операций (и вообще, очень быстрый рост) с увеличением размерности задачи называют "проклятием размерности". Что это такое, можно представить, зафиксировав, например,  $n = 100$ . Оценив величину  $100! \approx 10^{158}$  и прикинув потенциальные возможности развития вычислительной техники, приходим к выводу о том, что в обозримом будущем системы сотого порядка в принципе не могут быть решены по формулам Крамера [5, 33]. Заметим при этом, что, во-первых, метод Крамера будет неустойчив, т.е. погрешности округлений будут катастрофически нарастать, во-вторых, размерность  $n = 100$  для современных задач не так и велика: довольно часто решаются системы с сотнями и с тысячами неизвестных.

Если осуществлять вычисление обратной матрицы с помощью союзной матрицы, т. е. через алгебраические дополнения, то нахождение решения векторно-матричного уравнения (2.1') по формуле

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

фактически равнозначно применению формул Крамера и также практически непригодно по упомянутым выше причинам для вычислительных целей.

## **2.1. АЛГОРИТМ РЕШЕНИЯ СЛАУ МЕТОДОМ ГАУССА С ПОСТОЛБЦОВЫМ ВЫБОРОМ ГЛАВНОГО ЭЛЕМЕНТА**

Наиболее известным и популярным способом решения линейных систем вида (2.1) является метод Гаусса. Суть его проста – это последовательное исключение неизвестных. В отличие от курсов линейной алгебры, нас будут интересовать вычислительные аспекты метода Гаусса, а именно, технология получения вектора-решения  $x$  из исходных матрицы  $A$  и вектора  $b$ , причем, по возможности, минимизирующая влияние неизбежных ошибок округлений. С этой целью, работая с уравнениями системы (2.1), выведем сначала совокупность формул, позволяющих в итоге получить ис-



Продолжая этот процесс, на  $(n-1)$ -м этапе так называемого *прямого хода* метода Гаусса данную систему (2.1) приведем к треугольному виду:

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1, \\ a_{22}^{(1)}x_2 + \dots + a_{2n}^{(1)}x_n = b_2^{(1)}, \\ \dots \dots \dots \\ a_{nn}^{(n-1)}x_n = b_n^{(n-1)}. \end{cases} \quad (2.3)$$

На основе предыдущих рассуждений и формул легко убедиться, что коэффициенты этой системы могут быть получены из коэффициентов данной системы последовательным пересчетом по формулам

$$a_{ij}^{(k)} = a_{ij}^{(k-1)} - \frac{a_{ik}^{(k-1)}}{a_{kk}^{(k-1)}} a_{kj}^{(k-1)}, \quad b_i^{(k)} = b_i^{(k-1)} - \frac{a_{ik}^{(k-1)}}{a_{kk}^{(k-1)}} b_k^{(k-1)}, \quad (2.4)$$

где верхний индекс  $k$  (номер этапа) должен изменяться от 1 до  $n-1$ , нижние индексы  $i$  и  $j$  (в любой очередности) – от  $k+1$  до  $n$ ; по определению полагаем  $a_{ij}^{(0)} := a_{ij}$ ,  $b_i^{(0)} := b_i$ .

Треугольная, точнее, трапециевидная структура системы (2.3) позволяет последовательно одно за другим вычислять значения неизвестных, начиная с последнего:

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

Этот процесс последовательного вычисления значений неизвестных называют *обратным ходом* метода Гаусса. Очевидно, он определяется одной формулой

$$x_k = \frac{1}{a_{kk}^{(k-1)}} \left( b_k^{(k-1)} - \sum_{j=k+1}^n a_{kj}^{(k-1)} x_j \right), \quad (2.5)$$

где  $k$  полагают равным  $n, n-1, \dots, 2, 1$  и сумма по определению считается равной нулю, если нижний предел суммирования у знака  $\Sigma$  имеет значение больше верхнего.

Итак, решение СЛАУ вида (2.1) методом Гаусса сводится к последовательной реализации вычислений по формулам (2.4) и (2.5).

Учитывая цикличность выполняемых при этом операций, а также нецелесообразность хранения промежуточных результатов (пересчитывае-

мых коэффициентов промежуточного этапа), запишем простой алгоритм решения линейных систем (2.1) методом Гаусса [41]:

1) для  $k=1, 2, \dots, n-1$ ,

2) для  $i = k+1, \dots, n$ :

$$3) t_{ik} := a_{ik} / a_{kk},$$

$$4) b_i := b_i - t_{ik} b_k;$$

5) для  $j=k+1, \dots, n$ :

$$6) a_{ij} := a_{ij} - t_{ik} a_{kj}.$$

$$7) x_n := b_n / a_{nn};$$

8) для  $k = n-1, \dots, 2, 1$ :

$$9) x_k := \left( b_k - \sum_{j=k+1}^n a_{kj} x_j \right) / a_{kk}.$$

Подав на его вход квадратную матрицу  $(a_{ij})_{i,j=1}^n$  коэффициентов при неизвестных системы (2.1) и вектор  $(b_i)_{i=1}^n$  свободных членов и выполнив три вложенных цикла вычислений прямого хода (строки 1–6) и один цикл вычислений обратного хода (строки 7–9), на выходе алгоритма получим вектор-решение  $(x_k)_{k=1}^n$  (в обратном порядке), если, разумеется, ни один из знаменателей не обращается в нуль и все вычисления проводятся точно.

Так как реальные машинные вычисления производятся не с точными, а с усеченными числами, т.е. неизбежны ошибки округления, то анализируя, например, формулы (2.4), можно сделать вывод о том, что выполнение алгоритма может прекратиться или привести к неверным результатам, если знаменатели дробей на каком-то этапе окажутся равными нулю или очень маленькими числами. Чтобы уменьшить влияние ошибок округлений и исключить деление на нуль, на каждом этапе прямого хода уравнения системы (точнее, обрабатываемой подсистемы) обычно переставляют так, чтобы деление производилось на наибольший по модулю в данном столбце (обрабатываемом подстолбце) элемент. Числа, на которые производится деление в методе Гаусса, называются *ведущими* или *главными элементами*. Отсюда название рассматриваемой модификации метода, исключающей деление на нуль и уменьшающей вычислительные погрешности, *метод Гаусса с постолбцовым выбором главного элемента* (или, иначе, с *частичным упорядочиванием по столбцам*).

Частичное упорядочивание по столбцам требует внесения в алгоритм следующих изменений: между строками 1 и 2 нужно сделать вставку\*)

"Найти  $m \geq k$  такое, что  $|a_{mk}| = \max_{i \geq k} \{ |a_{ik}| \}$ ;

⊗ если  $a_{mk} = 0$ , остановить работу алгоритма ("однозначного решения нет"),

иначе поменять местами  $b_k$  и  $b_m$ ,  $a_{kj}$  и  $a_{mj}$  при всех  $j = k, \dots, n$ ."

Более разумным, наверное, является сравнение  $|a_{mk}|$  не с нулем, а с некоторым малым  $\varepsilon > 0$ , задаваемым вычислителем в зависимости от различных априорных соображений. Счет останавливается или берется под особый контроль, если окажется  $|a_{mk}| < \varepsilon$ . Заметим, что соответствующая частичному упорядочиванию вставка в алгоритм позволяет фактически в процессе его выполнения проводить алгоритмически исследование системы (2.1) на однозначную разрешимость.

Устойчивость алгоритма к погрешностям исходных данных и результатов промежуточных вычислений можно еще усилить, если выполнять деление на каждом этапе на элемент, наибольший по модулю во всей матрице преобразуемой на данном этапе подсистемы. Такая модификация метода Гаусса, называемая *методом главных элементов*, применяется довольно редко, поскольку сильно усложняет алгоритм. Усложнение связано как с необходимостью осуществления двумерного поиска главных элементов, так и с необходимостью запоминать номера столбцов, откуда берутся эти элементы (перестановка столбцов означает как бы переобозначение неизвестных, в связи с чем требуется обратная замена).

## 2.2. ПРИМЕНЕНИЕ МЕТОДА ГАУССА К ВЫЧИСЛЕНИЮ ОПРЕДЕЛИТЕЛЕЙ И К ОБРАЩЕНИЮ МАТРИЦ

Как было упомянуто во введении, решения линейных алгебраических систем можно получать с помощью определителей или обратных матриц. Однако нетрудно увидеть, что более эффективно поступать наоборот: вычислять определители и обращать матрицы в рамках метода Гаусса решения линейных систем.

Действительно, выполняемые в методе Гаусса преобразования прямого хода, приведшие матрицу  $A$  системы к треугольному виду (см. (2.3)), таковы, что они не изменяют определителя матрицы  $A$ . Учитывая, что оп-

---

\*) Чтобы частичное упорядочивание было более эффективным, перед этим целесообразно произвести *масштабирование* (уравновешивание) системы. например. разделить все числа каждой строки на наибольшее число строки [41]

ределитель треугольной матрицы равен произведению диагональных элементов, имеем

$$\det \mathbf{A} = \begin{vmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ 0 & a_{22}^{(1)} & \dots & a_{2n}^{(1)} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & a_{nn}^{(n-1)} \end{vmatrix} = a_{11} \cdot a_{22}^{(1)} \dots a_{nn}^{(n-1)} .$$

Таким образом,  $\det \mathbf{A}$  равен произведению всех ведущих элементов метода Гаусса.

При желании получить  $\det \mathbf{A}$  дополнительно к решению СЛАУ  $\mathbf{Ax} = \mathbf{b}$  алгоритм предыдущего пункта должен быть пополнен всего лишь одной строкой:

$$10) \quad \det \mathbf{A} = \prod_{k=1}^n a_{kk} .$$

Если метод Гаусса используется только для вычисления определителя, из алгоритма его реализации следует изъять строки 4 и 7–9.

Так как перестановка строк матрицы меняет знак определителя, то при постолбцовом выборе главного элемента, т. е. при включении в алгоритм вставки  $\otimes$ , нужно в результате учесть число  $p$  произведенных перестановок, точнее, четность этого числа. Это означает, что при вычислении  $\det \mathbf{A}$  алгоритмом Гаусса с частичным упорядочиванием вместо строки 10 должна быть включена строка

$$10') \quad \det \mathbf{A} = (-1)^p \prod_{k=1}^n a_{kk} .$$

Для получения матрицы  $\mathbf{A}^{-1}$ , обратной к матрице  $\mathbf{A} = (a_{ij})_{i,j=1}^n$ , будем исходить из того, что она является решением матричного уравнения

$$\mathbf{AX} = \mathbf{E}, \quad (2.6)$$

где  $\mathbf{E}$  – единичная матрица.

Представляя искомую матрицу  $\mathbf{X} = (x_{ij})_{i,j=1}^n$  как набор (вектор-строку) векторов-столбцов

$$\mathbf{x}_1 = \begin{pmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{pmatrix}, \quad \mathbf{x}_2 = \begin{pmatrix} x_{12} \\ x_{22} \\ \vdots \\ x_{n2} \end{pmatrix}, \quad \dots, \quad \mathbf{x}_n = \begin{pmatrix} x_{1n} \\ x_{2n} \\ \vdots \\ x_{nn} \end{pmatrix},$$



а единичную матрицу  $E$  как набор единичных векторов

$$e_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad e_2 = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}, \quad \dots, \quad e_n = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix},$$

матричное уравнение (2.6) в соответствии с правилами умножения матриц подменим эквивалентной системой не связанных между собой векторно-матричных уравнений

$$Ax_1 = e_1; \quad Ax_2 = e_2; \quad \dots; \quad Ax_n = e_n. \quad (2.7)$$

Каждое из последних уравнений имеет вид (2.1') и может быть решено методом Гаусса. При этом специфичным является то обстоятельство, что все СЛАУ (2.7) имеют одну и ту же матрицу коэффициентов, а это означает, что наиболее трудоемкая часть метода Гаусса – приведение матрицы системы к треугольному виду – общая для всех систем (2.7). Так что, если требуется приспособить рассмотренный выше алгоритм решения СЛАУ методом Гаусса к обращению матриц, целесообразно не просто применить его последовательно  $n$  раз к системам (2.7), а слегка подкорректировать: "размножить" строки 4 и 9 так, чтобы в роли вектора  $b$  оказались все единичные векторы  $e_1, e_2, \dots, e_n$ . Тогда в результате завершения работы алгоритма будут получаться столбец за столбцом (столбцы "перевернуты!") элементы обратной матрицы  $X = A^{-1}$ . При этом введение в алгоритм частичного упорядочивания, т.е. постолбцовый выбор главного элемента, не требует запоминаний и обратных замен.

### 2.3. LU - РАЗЛОЖЕНИЕ МАТРИЦ

Пусть  $A = (a_{ij})_{i,j=1}^n$  – данная  $n \times n$ -матрица, а  $L = (l_{ij})_{i,j=1}^n$  и  $U = (u_{ij})_{i,j=1}^n$  – соответственно нижняя (левая) и верхняя (правая) треугольные матрицы<sup>\*)</sup>. Справедлива следующая теорема.

**Теорема 2.1** [15, 47, 55] *Если все главные миноры квадратной матрицы  $A$  отличны от нуля, то существуют такие нижняя  $L$  и верхняя  $U$  треугольные матрицы, что  $A = LU$ . Если элементы диагонали одной из матриц  $L$  или  $U$  фиксированы (ненулевые), то такое разложение единственно.*

<sup>\*)</sup> Общепринятые обозначения  $L$  и  $U$  связаны с английскими словами lower (нижний) и upper (верхний). Существует другой стандарт обозначения:  $L$  и  $R$ , определяемый словами left (левый) и right (правый).

Вместо полного доказательства этой теоремы (см., например, [15]) получим формулы для фактического разложения матриц в случае фиксирования диагонали нижней треугольной матрицы  $L$ .

Будем находить  $l_{ij}$  (при  $i > j$ ) и  $u_{ij}$  (при  $i \leq j$ ) такие, чтобы

$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ l_{21} & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ l_{n1} & l_{n2} & \dots & 1 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ 0 & u_{22} & \dots & u_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & u_{nn} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix}.$$

Выполнив перемножение матриц, на основе поэлементного приравнивания левых и правых частей приходим к  $n \times n$ -матрице уравнений

$$\begin{aligned} u_{11} &= a_{11}, & u_{12} &= a_{12}, \dots, & u_{1n} &= a_{1n}, \\ l_{21}u_{11} &= a_{21}, & l_{21}u_{12} + u_{22} &= a_{22}, \dots, & l_{21}u_{1n} + u_{2n} &= a_{2n}, \\ \dots & \dots & \dots & \dots & \dots & \dots \\ l_{n1}u_{11} &= a_{n1}, & l_{n1}u_{12} + l_{n2}u_{22} &= a_{n2}, \dots, & l_{n1}u_{1n} + \dots + l_{n,n-1}u_{n-1,n} + u_{nn} &= a_{nn}; \end{aligned}$$

относительно  $n \times n$ -матрицы неизвестных

$$\begin{aligned} &u_{11}, u_{12}, \dots, u_{1n}, \\ &l_{21}, u_{22}, \dots, u_{2n}, \\ &\dots \dots \dots \\ &l_{21}, l_{2n}, \dots, u_{nn}. \end{aligned} \tag{2.8}$$

Специфика этой системы позволяет находить неизвестные (2.8) одно за другим в следующем порядке.

Из первой строки уравнений имеем

$$u_{1j} = a_{1j} \quad (j = 1, \dots, n);$$

из оставшейся части первого столбца уравнений

$$l_{i1} = \frac{a_{i1}}{u_{11}} \quad (i = 2, \dots, n);$$

из оставшейся части второй строки

$$u_{2j} = a_{2j} - l_{21}u_{1j} \quad (j = 2, \dots, n);$$

из оставшейся части второго столбца

$$l_{i2} = \frac{a_{i2} - l_{i1}u_{12}}{u_{22}} \quad (i = 3, \dots, n)$$

и т.д. Последним находится элемент

$$u_{nn} = a_{nn} - \sum_{k=1}^{n-1} l_{nk}u_{kn}.$$

Легко видеть, что все отличные от 0 и 1 элементы матриц  $L$  и  $U$  могут быть однозначно вычислены с помощью всего двух формул:

$$u_{ij} = a_{ij} - \sum_{k=1}^{i-1} l_{ik}u_{kj} \quad (i \leq j), \quad (2.9)$$

$$l_{ij} = \frac{1}{u_{jj}} \left( a_{ij} - \sum_{k=1}^{j-1} l_{ik}u_{kj} \right) \quad (i > j). \quad (2.10)$$

При практическом выполнении разложения<sup>\*)</sup> матрицы  $A$  нужно иметь в виду следующие два обстоятельства.

Во-первых, организация вычислений по формулам (2.9)–(2.10) должна предусматривать переключение счета с одной формулы на другую в соответствии с показанным выше процессом получения неизвестных, приведшим к этим формулам. Это удобно делать, ориентируясь на матрицу неизвестных (2.8) (ее, кстати, можно интерпретировать как  $n^2$ -мерный массив для компактного хранения LU-разложения в памяти ЭВМ), а именно, первая строка (2.8) вычисляется по формуле (2.9) при  $i = 1, j = 1, 2, \dots, n$ ; первый столбец (2.8) (без первого элемента) – по формуле (2.10) при  $j = 1, i = 2, \dots, n$  и т.д.

Во-вторых, препятствием для осуществимости описанного процесса LU-разложения матрицы  $A$  может оказаться равенство нулю диагональных элементов матрицы  $U$ , поскольку на них выполняется деление в формуле (2.10). Отсюда требование теоремы, накладываемое на главные миноры (напомним, что главными минорами матрицы  $A = (a_{ij})_{i,j=1}^n$  называются определители подматриц  $A_k = (a_{ij})_{i,j=1}^k$ , где  $k = 1, 2, \dots, n-1$ ). Заметим, что  $u_{11} = a_{11}$ , т.е. первый диагональный элемент матрицы  $U$  совпадает с

<sup>\*)</sup> Применяют также термины *факторизация*, *декомпозиция*.

первым главным минором  $A$  и по условию должен быть отличным от нуля. Второй диагональный элемент матрицы  $U$

$$u_{22} = a_{22} - l_{21}u_{12} = a_{22} - \frac{a_{21}}{a_{11}}a_{12} = \frac{\begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix}}{a_{11}}$$

не равен нулю, если отличен от нуля второй главный минор, и т.д. Ясно, что вместо проверки на равенство нулю главных миноров данной матрицы удобнее делать такую проверку для элементов  $u_{j,j}$  в процессе их вычисления, причем, чтобы уменьшить влияние погрешностей округлений, лучше сравнивать модули  $u_{j,j}$  с малой положительной константой (допуском).

Для определенных классов матриц требования теоремы о разложении заведомо выполняются. Это относится, например, к *матрицам с диагональным преобладанием*, т.е. к таким, для которых

$$|a_{ii}| > \sum_{\substack{j=1 \\ (j \neq i)}}^n |a_{ij}| \quad \forall i \in \{1, 2, \dots, n\}.$$

## 2.4. РЕШЕНИЕ ЛИНЕЙНЫХ СИСТЕМ И ОБРАЩЕНИЕ МАТРИЦ С ПОМОЩЬЮ LU-РАЗЛОЖЕНИЯ

Если матрица  $A$  исходной системы (2.1) разложена в произведение треугольных  $L$  и  $U$ , то, значит, вместо (2.1') мы можем записать эквивалентное (2.1) уравнение

$$LUx = b.$$

Введя вектор вспомогательных переменных  $y = (y_1, y_2, \dots, y_n)^T$ , последнее можно переписать в виде системы

$$\begin{cases} Ly = b, \\ Ux = y. \end{cases}$$

Таким образом, решение данной системы с квадратной матрицей коэффициентов свелось к последовательному решению двух систем с треугольными матрицами коэффициентов.

Получим сначала формулы для вычисления элементов  $y_i$ . Для этого запишем уравнение  $Ly = b$  в развернутом виде:

$$\begin{cases} y_1 & = b_1, \\ l_{21}y_1 + y_2 & = b_2, \\ \dots & \dots \\ l_{n1}y_1 + l_{n2}y_2 + \dots + l_{n,n-1}y_{n-1} + y_n & = b_n. \end{cases}$$



Вычисление определителя LU-факторизованной матрицы  $A$  опирается на свойство определителя произведения матриц и сводится к перемножению  $n$  чисел:

$$\det A = \det L \cdot \det U = u_{11} \cdot u_{22} \cdot \dots \cdot u_{nn}.$$

Для обращения LU-факторизацией матрицы  $A$  можно применить тот же прием, который рассмотрен в п. 2.2, т. е.  $n$ -кратно использовать формулы (2.11) и (2.13) для получения столбцов матрицы  $A^{-1}$ ; при этом в качестве  $b_i$  в (2.11) должны фигурировать только 0 или 1: для нахождения первого столбца  $A^{-1}$  полагаем  $b_1 = 1, b_2 = b_3 = \dots = b_n = 0$ , для второго —  $b_2 = 1, b_1 = b_3 = \dots = b_n = 0$ , и т.д. Можно однако вывести и специальные формулы для выражения элементов обратной матрицы через элементы матриц  $L$  и  $U$ .

Пусть матрицы  $A$  и  $U$  обратимы (матрица  $L$  обратима всегда). Тогда

$$A = L \cdot U \quad \Leftrightarrow \quad A^{-1} = U^{-1} \cdot L^{-1}.$$

Умножая последнее равенство поочередно на  $U$  слева и на  $L$  справа, будем иметь

$$UA^{-1} = L^{-1} \quad \text{и} \quad A^{-1}L = U^{-1}. \quad (2.14)$$

Обозначим, как и ранее, искомые элементы матрицы  $A^{-1}$  через  $x_{ij}$ . Учитывая, что треугольные матрицы при обращении сохраняют свою структуру, перепишем равенства (2.14) в виде

$$\begin{pmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ 0 & u_{22} & \dots & u_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & u_{nn} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nn} \end{pmatrix} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ * & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ * & * & \dots & 1 \end{pmatrix},$$

$$\begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nn} \end{pmatrix} \begin{pmatrix} 1 & 0 & \dots & 0 \\ l_{21} & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ l_{n1} & l_{n2} & \dots & 1 \end{pmatrix} = \begin{pmatrix} * & * & \dots & * \\ 0 & * & \dots & * \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & * \end{pmatrix};$$

где  $*$  — некоторые числа. Полученные матричные равенства можно рассматривать как систему  $2n^2$  уравнений с  $n^2$  неизвестными  $x_{ij}$  ( $i, j = 1, 2, \dots, n$ ). Из этих  $2n^2$  уравнений ровно  $n^2$  имеют известные пра-

вые части (0 или 1). Выпишем соответствующую им  $n \times n$ -матрицу уравнений:

$$\begin{aligned} u_{11}x_{11} + \dots + u_{1n}x_{n1} &= 1, & u_{11}x_{12} + \dots + u_{1n}x_{n2} &= 0, & \dots, & & u_{11}x_{1n} + \dots + u_{1n}x_{nn} &= 0, \\ x_{21} + x_{22}l_{21} + \dots + x_{2n}l_{n1} &= 0, & u_{22}x_{22} + \dots + u_{2n}x_{nn} &= 1, & \dots, & & u_{22}x_{2n} + \dots + u_{2n}x_{nn} &= 0, \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ x_{n1} + x_{n2}l_{21} + \dots + x_{nn}l_{n1} &= 0, & x_{n2} + \dots + x_{nn}l_{n2} &= 0, & \dots, & & u_{nn}x_{nn} &= 1. \end{aligned}$$

Короче все эти уравнения могут быть представлены следующими тремя типами связей\*):

$$\begin{aligned} \sum_{k=i}^n u_{ik}x_{kj} &= 0, & \text{если } i < j, \\ \sum_{k=i}^n u_{ik}x_{kj} &= 1 & \text{при } i = j \end{aligned}$$

и

$$x_{ij} + \sum_{k=j+1}^n x_{ik}l_{kj} = 0, \quad \text{если } i > j.$$

Отсюда можно выразить все элементы  $x_{ij}$  искомой обратной матрицы  $A^{-1}$ :

$$x_{jj} = \frac{1}{u_{jj}} \left( 1 - \sum_{k=j+1}^n u_{jk}x_{kj} \right); \quad (2.15)$$

$$x_{ij} = -\frac{1}{u_{ii}} \sum_{k=i+1}^n u_{ik}x_{kj} \quad (i < j); \quad (2.16)$$

$$x_{ij} = -\sum_{k=j+1}^n x_{ik}l_{kj} \quad (i > j). \quad (2.17)$$

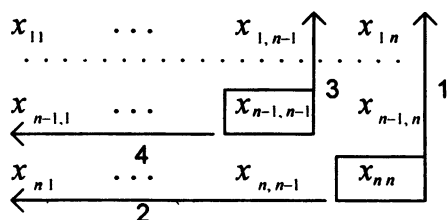
Формулы (2.15)–(2.17) позволяют эффективно обращаться LU-факторизованную матрицу, если соблюдать определенную технологию их использования. А именно, как видно из записанной выше матрицы уравнений, следует сначала из последнего столбца уравнений найти  $x_{nn}$ ,  $x_{n-1,n}, \dots, x_{2n}, x_{1n}$ , затем из оставшейся части последней строки уравнений найти  $x_{n,n-1}, \dots, x_{n2}, x_{n1}$ , потом переключиться на предпоследний

\* Используя символ Кронекера  $\delta_{ij}$ , первые две формулы можно совместить:

$$\sum_{k=i}^n u_{ik}x_{kj} = \delta_{ij}$$

Это относится и к формулам (2.15), (2.16)

столбец и т.д. Схематично последовательность вычислений элементов обратной матрицы можно изобразить пронумерованными стрелками следующим образом:



При этом стрелка 1 означает, что фиксируем  $j = n$  и ведем счет по формулам (2.15), (2.16) при  $i = n, n-1, \dots, 1$ ; стрелка 2 – счет по формуле (2.17) при  $i = n$  и  $j = n-1, n-2, \dots, 1$  и т.д.

Возвращаясь к началу п.2.3, заметим, что так же употребительно фиксирование единичной диагонали у правой треугольной матрицы, т.е. представление

$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} = \begin{pmatrix} l_{11} & 0 & \dots & 0 \\ l_{21} & l_{22} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{pmatrix} \begin{pmatrix} 1 & u_{12} & \dots & u_{1n} \\ 0 & 1 & \dots & u_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{pmatrix}.$$

В этом случае  $l_{ij}$  и  $u_{ij}$  находятся по формулам

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

$$u_{ij} = \frac{1}{l_{ii}} \left( a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj} \right) \quad (i < j),$$

где индексы фиксируются так, чтобы вычислялись поочередно столбец  $(l_{i1})_{i=1}^n$ , затем строка  $(u_{1j})_{j=2}^n$  и т.д.

Решение системы (2.1) с таким образом факторизованной матрицей коэффициентов получают по формулам

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

$$x_i = y_i - \sum_{k=i+1}^n u_{ik} x_k, \quad i = n, n-1, \dots, 1.$$



Детерминант матрицы  $A$  равен произведению  $l_{11}l_{22}\dots l_{nn}$ , а для подсчета элементов обратной матрицы используют совокупность формул

$$x_{ii} = \frac{1}{l_{ii}} \left( 1 - \sum_{k=i+1}^n x_{ik}l_{ki} \right),$$

$$x_{ij} = -\frac{1}{l_{jj}} \sum_{k=i+1}^n x_{ik}l_{kj} \quad (i > j),$$

$$x_{ij} = -\sum_{k=i+1}^n u_{ik}x_{kj} \quad (i < j)$$

с такой организацией вычислений, при которой сначала вычисляется последняя строка  $(x_{nj})$  при  $j = n, n-1, \dots, 1$ , затем последний столбец  $(x_{in})$  при  $i = n-1, \dots, 1$ , потом предпоследняя строка  $(x_{n-1,j})$  при  $j = n-1, \dots, 1$ , и т.д.

В отличие от рассмотренной в п.2.1 схемы единственного деления схема Холецкого менее удобна для усовершенствования с целью уменьшения влияния вычислительных погрешностей путем выбора подходящих ведущих элементов. Достоинством же ее можно считать то, что LU-разложение матрицы  $A$  играет роль обратной матрицы, может помещаться в память ЭВМ на место матрицы  $A$  и использоваться, например, при решении нескольких систем, имеющих одну и ту же матрицу коэффициентов и разные правые части.

## 2.5. РАЗЛОЖЕНИЕ СИММЕТРИЧНЫХ МАТРИЦ. МЕТОД КВАДРАТНЫХ КОРНЕЙ

Объем вычислений, требующихся для решения линейных алгебраических задач с симметрическими матрицами, можно сократить почти вдвое, если учитывать симметрию при треугольной факторизации матриц.

Пусть  $A = (a_{ij})_{i,j=1}^n$  — данная симметрическая матрица, т.е.  $a_{ij} = a_{ji}$ .

Будем строить ее представление в виде  $A = U^T U$ , где

$$U = \begin{pmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ 0 & u_{22} & \dots & u_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & u_{nn} \end{pmatrix}, \quad U^T = \begin{pmatrix} u_{11} & 0 & \dots & 0 \\ u_{12} & u_{22} & \dots & 0 \\ \dots & \dots & \dots & \dots \\ u_{1n} & u_{2n} & \dots & u_{nn} \end{pmatrix}.$$

Аналогично тому, как это делалось в п.2.3, составим систему  $\frac{n(n+1)}{2}$  уравнений относительно такого же количества неизвестных (элементов матрицы  $U$ ):

$$\begin{aligned}
 u_{11}^2 &= a_{11}, & u_{12}u_{11} &= a_{12}, \dots, & u_{1n}u_{11} &= a_{1n}, \\
 u_{12}^2 + u_{22}^2 &= a_{22}, \dots, & u_{12}u_{1n} + u_{22}u_{2n} &= a_{2n}, \\
 & \dots \dots \dots & & & & \\
 & & u_{1n}^2 + u_{2n}^2 + \dots + u_{nn}^2 &= a_{nn}.
 \end{aligned}$$

Из первой строки уравнений находим сначала  $u_{11} = \sqrt{a_{11}}$ , затем  $u_{1j} = \frac{a_{1j}}{u_{11}}$  при  $j = 2, \dots, n$ . Из второй —  $u_{22} = \sqrt{a_{22} - u_{12}^2}$ , затем  $u_{2j} = \frac{a_{2j} - u_{12}u_{1j}}{u_{22}}$  при  $j = 3, \dots, n$ , и т.д. Завершается процесс вычислением

$$u_{nn} = \sqrt{a_{nn} - \sum_{k=1}^{n-1} u_{kn}^2}.$$

Таким образом, матрица  $U$  может быть определена совокупностью формул

$$\begin{aligned}
 u_{ii} &= \sqrt{a_{ii} - \sum_{k=1}^{i-1} u_{ki}^2} && \text{при } i = 1, 2, \dots, n; \\
 u_{ij} &= \frac{a_{ij} - \sum_{k=1}^{i-1} u_{ki}u_{kj}}{u_{ii}} && \text{при } j = 2, \dots, n; \quad j > i \\
 &&& (u_{ij} = 0 \text{ при } j < i).
 \end{aligned} \tag{2.18}$$

Осуществимости вещественного  $U^T U$ -разложения вещественной симметрической матрицы  $A$  по этим формулам могут помешать два обстоятельства: обращение в нуль элемента  $u_{ii}$  при каком-либо  $i \in \{1, 2, \dots, n\}$  и отрицательность подкоренного выражения. Известно, что для важного в приложениях класса симметрических положительно определенных матриц разложение по формулам (2.18) выполнимо [15, 54 и др.]<sup>\*)</sup>.

<sup>\*)</sup> Более универсальным, чем  $U^T U$ -разложение Холецкого, является  $U^* D U$ -разложение, пригодное для эрмитовых матриц, частным случаем которых являются симметрические (см., например, [5, 16, 47]).

При наличии  $U^T U$ -разложения решение симметричной системы  $Ax = b$  сводится к последовательному решению двух треугольных систем

$$U^T y = b \quad \text{и} \quad Ux = y.$$

Первая из них имеет вид

$$\begin{cases} u_{11}y_1 & = b_1, \\ u_{12}y_1 + u_{22}y_2 & = b_2, \\ \dots & \dots \\ u_{1n}y_1 + u_{2n}y_2 + \dots + u_{nn}y_n & = b_n, \end{cases}$$

откуда получаем вспомогательные неизвестные  $y_1, y_2, \dots, y_n$  по формуле

$$y_i = \frac{b_i - \sum_{k=1}^{i-1} u_{ki}y_k}{u_{ii}}, \quad (2.19)$$

полагая в ней  $i = 1, 2, \dots, n$ . Из второй системы

$$\begin{cases} u_{11}x_1 + u_{12}x_2 + \dots + u_{1n}x_n = y_1, \\ u_{22}x_2 + \dots + u_{2n}x_n = y_2, \\ \dots \\ u_{nn}x_n = y_n \end{cases}$$

находим искомые значения  $x_i$  в обратном порядке, т.е. при  $i = n, n-1, \dots, 1$  по формуле

$$x_i = \frac{y_i - \sum_{k=i+1}^n u_{ik}x_k}{u_{ii}}. \quad (2.20)$$

Решение симметричных СЛАУ по формулам (2.18)–(2.20) называют *методом квадратных корней* или *схемой Холецкого*. В случае систем с положительно определенными матрицами можно ожидать хороших результатов применения такого метода (особенно, если в процессе решения делать проверку на немалость  $|u_{ii}|$  дабы избежать большого роста погрешностей). В противном случае нет, например, гарантий, что в процессе разложения не появятся чисто мнимые числа, что кстати может не отразиться на результатах, если в алгоритме реализации метода квадратных корней предусмотреть возможность появления мнимых чисел [54].

## 2.6. МЕТОД ПРОГОНКИ РЕШЕНИЯ СИСТЕМ С ТРЕХДИАГОНАЛЬНЫМИ МАТРИЦАМИ КОЭФФИЦИЕНТОВ

Часто возникает необходимость в решении линейных алгебраических систем, матрицы которых, являясь слабо заполненными, т.е. содержащими немного ненулевых элементов, имеют определенную структуру. Среди таких систем выделим системы с матрицами ленточной структуры, в которых ненулевые элементы располагаются на главной диагонали и на нескольких побочных диагоналях. Для решения систем с ленточными матрицами коэффициентов метод Гаусса можно трансформировать в более эффективные методы.

Рассмотрим наиболее простой случай *ленточных систем*, к которым, как увидим впоследствии, сводится решение задач сплайн-интерполяции функций, дискретизации краевых задач для дифференциальных уравнений методами конечных разностей, конечных элементов и др. А именно, будем искать решение такой системы, каждое уравнение которой связывает три "соседних" неизвестных:

$$b_i x_{i-1} + c_i x_i + d_i x_{i+1} = r_i, \quad (2.21)$$

где  $i = 1, 2, \dots, n$ ;  $b_1 = 0, d_n = 0$ . Такие уравнения называют *трехточечными разностными уравнениями второго порядка*<sup>\*)</sup>. Система (2.21) имеет трехдиагональную структуру, что хорошо видно из следующего, эквивалентного (2.21), векторно-матричного представления:

$$\begin{bmatrix} c_1 & d_1 & 0 & 0 & \dots & 0 & 0 & 0 \\ b_2 & c_2 & d_2 & 0 & \dots & 0 & 0 & 0 \\ 0 & b_3 & c_3 & d_3 & \dots & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & b_{n-1} & c_{n-1} & d_{n-1} \\ 0 & 0 & 0 & 0 & \dots & 0 & b_n & c_n \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} = \begin{bmatrix} r_1 \\ r_2 \\ r_3 \\ \vdots \\ r_{n-1} \\ r_n \end{bmatrix}.$$

Преследуя, как и в п.2.1, цель избавиться от ненулевых элементов в поддиагональной части матрицы системы, предположим, что существуют такие наборы чисел  $\delta_i$  и  $\lambda_i$  ( $i = 1, 2, \dots, n$ ), при которых

$$x_i = \delta_i x_{i+1} + \lambda_i, \quad (2.22)$$

<sup>\*)</sup> Чаще вместо принятой здесь записи системы (2.21), предполагающей как бы наличие фиктивных неизвестных  $x_0$  и  $x_{n+1}$  с нулевыми коэффициентами, считают в (2.21)  $i$  изменяющимся от 2 до  $n-1$ , выделяя первое и последнее уравнения системы соответственно

$$c_1 x_1 + d_1 x_2 = r_1 \quad \text{и} \quad b_n x_{n-1} + c_n x_n = r_n$$

в отдельные строки (так называемые *краевые условия* разностного уравнения).

т.е. трехточечное уравнение второго порядка (2.21) преобразуется в двухточечное уравнение первого порядка (2.22). Уменьшим в связи (2.22) индекс на единицу и полученное выражение  $x_{i-1} = \delta_{i-1}x_i + \lambda_{i-1}$  подставим в данное уравнение (2.21):

$$b_i \delta_{i-1} x_i + b_i \lambda_{i-1} + c_i x_i + d_i x_{i+1} = r_i,$$

откуда

$$x_i = -\frac{d_i}{c_i + b_i \delta_{i-1}} x_{i+1} + \frac{r_i - b_i \lambda_{i-1}}{c_i + b_i \delta_{i-1}}.$$

Последнее равенство имеет вид (2.22) и будет точно с ним совпадать, иначе говоря, представление (2.22) будет иметь место, если при всех  $i = 1, 2, \dots, n$  выполняются рекуррентные соотношения

$$\delta_i = -\frac{d_i}{c_i + b_i \delta_{i-1}}, \quad \lambda_i = \frac{r_i - b_i \lambda_{i-1}}{c_i + b_i \delta_{i-1}}. \quad (2.23)$$

Легко видеть, что, в силу условия  $b_1 = 0$ , процесс вычисления  $\delta_i, \lambda_i$  может быть начат со значений

$$\delta_1 = -\frac{d_1}{c_1}, \quad \lambda_1 = \frac{r_1}{c_1}$$

и продолжен далее по формулам (2.23) последовательно при  $i = 2, 3, \dots, n$ , причем при  $i = n$ , в силу  $d_n = 0$ , получим  $\delta_n = 0$ . Следовательно, полагая в (2.22)  $i = n$ , будем иметь

$$x_n = \lambda_n = \frac{r_n - b_n \lambda_{n-1}}{c_n + b_n \delta_{n-1}}$$

(где  $\lambda_{n-1}, \delta_{n-1}$  — уже известные с предыдущего шага числа). Далее по формулам (2.22) последовательно находятся  $x_{n-1}, x_{n-2}, \dots, x_1$  при  $i = n-1, n-2, \dots, 1$  соответственно.

Таким образом, решение уравнений вида (2.21) описываемым способом, называемым *методом прогонки*<sup>\*)</sup>, сводится к вычислениям по трем простым формулам: нахождение так называемых *прогночных коэффициентов*  $\delta_i, \lambda_i$  по формулам (2.23) при  $i = 1, 2, \dots, n$  (*прямая прогонка*) и затем получение неизвестных  $x_i$  по формуле (2.22) при  $i = n, n-1, \dots, 1$  (*обратная прогонка*).

Для успешного применения метода прогонки нужно, чтобы в процессе вычислений не возникало ситуаций с делением на нуль, а при больших размерностях систем не должно быть быстрого роста погрешностей округлений.

Будем называть прогонку *корректной*, если знаменатели прогночных коэффициентов (2.23) не обращаются в нуль, и *устойчивой*, если  $|\delta_i| < 1$  при всех  $i \in \{1, 2, \dots, n\}$ .

<sup>\*)</sup> Термин, характерный, в основном, для отечественной литературы по вычислительной математике, введён в 50-х годах (см., например, [1]).

Приведем простые достаточные условия корректности и устойчивости прогонки, которые во многих приложениях метода автоматически выполняются.

**Теорема 2.2.** Пусть коэффициенты  $b_i$  и  $d_i$  уравнения (2.21) при  $i = 2, 3, \dots, n-1$  отличны от нуля и пусть

$$|c_i| > |b_i| + |d_i| \quad \forall i = 1, 2, \dots, n. \quad (2.24)$$

Тогда прогонка (2.23), (2.22) корректна и устойчива (т.е.  $c_i + b_i \delta_{i-1} \neq 0$ ,  $|\delta_i| < 1$ ).

**Доказательство.** Воспользуемся методом математической индукции для установления обоих нужных неравенств одновременно.

При  $i = 1$ , в силу (2.24), имеем:

$$|c_1| > |d_1| \geq 0$$

– неравенство нулю знаменателя первой пары прогоночных коэффициентов, а также

$$|\delta_1| = \left| -\frac{d_1}{c_1} \right| < 1.$$

Предположим, что знаменатель  $(i-1)$ -х прогоночных коэффициентов не равен нулю и что  $|\delta_{i-1}| < 1$ . Тогда, используя свойства модулей, условия теоремы и индукционные предположения, получаем:

$$\begin{aligned} |c_i + b_i \delta_{i-1}| &\geq |c_i| - |b_i \delta_{i-1}| > |b_i| + |d_i| - |b_i| \cdot |\delta_{i-1}| = \\ &= |d_i| + |b_i|(1 - |\delta_{i-1}|) > |d_i| > 0, \end{aligned}$$

а с учетом этого

$$|\delta_i| = \left| -\frac{d_i}{c_i + b_i \delta_{i-1}} \right| = \frac{|d_i|}{|c_i + b_i \delta_{i-1}|} < \frac{|d_i|}{|d_i|} = 1.$$

Следовательно,  $c_i + b_i \delta_{i-1} \neq 0$  и  $|\delta_i| < 1$  при всех  $i \in \{1, 2, \dots, n\}$ , т.е. имеет место утверждаемая в данных условиях корректность и устойчивость прогонки. Теорема доказана.

Пусть  $A$  – матрица коэффициентов данной системы (2.21), удовлетворяющих условиям теоремы 2.2, и пусть

$$\delta_1 = -\frac{d_1}{c_1}, \quad \delta_i = -\frac{d_i}{c_i + b_i \delta_{i-1}} \quad (i = 2, 3, \dots, n-1), \quad \delta_n = 0$$

– прогоночные коэффициенты, определяемые первой из формул (2.23), а

$$\Delta_i = c_i + b_i \delta_{i-1} \quad (i = 2, 3, \dots, n)$$

– знаменатели этих коэффициентов (отличные от нуля согласно утверждению теоремы 2.2). Непосредственной проверкой легко убедиться, что имеет место представление  $\mathbf{A} = \mathbf{LU}$ , где

$$\mathbf{L} = \begin{bmatrix} c_1 & 0 & 0 & 0 & \dots & 0 & 0 & 0 \\ b_2 & \Delta_2 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & b_3 & \Delta_3 & 0 & \dots & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & b_{n-1} & \Delta_{n-1} & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & b_n & \Delta_n \end{bmatrix},$$

$$\mathbf{U} = \begin{bmatrix} 1 - \delta_1 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & 1 - \delta_2 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & 1 - \delta_3 & \dots & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & 0 & 1 - \delta_{n-1} \\ 0 & 0 & 0 & 0 & \dots & 0 & 0 & 1 \end{bmatrix},$$

единственное в силу утверждения теоремы 2.1. Как видим, LU-разложение трехдиагональной матрицы  $\mathbf{A}$  может быть выполнено очень простым алгоритмом, вычисляющим  $\Delta_i$  и  $\delta_i$  при возрастающих значениях  $i$ . При необходимости попутно может быть вычислен

$$\det \mathbf{A} = c_1 \prod_{i=2}^n \Delta_i.$$

В заключение этого пункта заметим, что, во-первых, имеются более слабые условия корректности и устойчивости прогонки, чем требуемое в теореме 2.2 условие строгого диагонального преобладания в матрице  $\mathbf{A}$  (см., например, [17, 47; 48]). Во-вторых, применяется ряд других, отличных от рассмотренной нами правой прогонки, методов подобного типа, решающих как поставленную здесь задачу (2.21) для систем с трехдиагональными матрицами (левая прогонка, встречная прогонка, немонотонная, циклическая, ортогональная прогонки и т.д.), так и для более сложных систем с матрицами ленточной структуры или блочно-матричной структуры (например, матричная прогонка). Выводы и исследование различных вариантов метода прогонки можно найти, например, в [27], [48].

## 2.7. МЕТОД ВРАЩЕНИЙ РЕШЕНИЯ ЛИНЕЙНЫХ СИСТЕМ

Вернемся к рассмотрению линейных систем общего вида (2.1).

Предположим, что методом Гаусса решается система  $Ax = b$  с матрицей коэффициентов

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

Приведение такой системы<sup>\*)</sup> к треугольному виду прямым ходом метода Гаусса равносильно следующей последовательности эквивалентных преобразований матрицы  $A$ :

$$A \sim \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & -1 & 1 & 2 \\ 0 & -1 & -1 & 2 \end{pmatrix} \sim \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 4 \\ 0 & 0 & -1 & 4 \end{pmatrix} \sim \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 4 \\ 0 & 0 & 0 & 8 \end{pmatrix}.$$

Очевидно, в случае  $n \times n$ -матрицы такого типа прямой ход метода Гаусса допускает рост элементов матрицы до  $2^{n-1}$ . При больших  $n$  это может привести если не к переполнению разрядной сетки ЭВМ, то к сильному влиянию погрешностей округлений, причем в данной ситуации не даст эффекта и постолбцовый выбор главного элемента.

Рассмотренный пример оправдывает поиск других подходов к построению прямых методов решения линейных систем (2.1), возможно, более сложных, чем метод Гаусса, но не допускающих большого роста элементов в процессе преобразований и как следствие численно более устойчивых.

Как и в методе Гаусса, цель прямого хода преобразований в новом методе – приведение системы к треугольному виду последовательным обнулением поддиагональных элементов сначала первого столбца, затем второго и т.д. Делается это следующим образом.

Пусть  $c_1$  и  $s_1$  – некоторые отличные от нуля числа. Умножим первое уравнение системы (2.1) на  $c_1$ , второе – на  $s_1$  и сложим их; полученным уравнением заменим первое уравнение системы. Затем первое уравнение исходной системы умножаем на  $-s_1$ , второе – на  $c_1$  и результатом их сло-

<sup>\*)</sup> Пример заимствован из [25].



жения заменяем второе уравнение. Таким образом, первые два уравнения системы (2.1) заменяются уравнениями

$$\begin{aligned} (c_1 a_{11} + s_1 a_{21})x_1 + (c_1 a_{12} + s_1 a_{22})x_2 + \dots + (c_1 a_{1n} + s_1 a_{2n})x_n &= c_1 b_1 + s_1 b_2, \\ (-s_1 a_{11} + c_1 a_{21})x_1 + (-s_1 a_{12} + c_1 a_{22})x_2 + \dots + (-s_1 a_{1n} + c_1 a_{2n})x_n &= -s_1 b_1 + c_1 b_2. \end{aligned}$$

На введенные два параметра  $c_1$  и  $s_1$  наложим два условия:

$$-s_1 a_{11} + c_1 a_{21} = 0$$

– условие обнуления (т.е. исключение  $x_1$  из второго уравнения) и

$$c_1^2 + s_1^2 = 1$$

– условие нормировки. Легко проверить, что за  $c_1$  и  $s_1$ , удовлетворяющие этим условиям, можно принять

$$c_1 = \frac{a_{11}}{\sqrt{a_{11}^2 + a_{21}^2}}, \quad s_1 = \frac{a_{21}}{\sqrt{a_{11}^2 + a_{21}^2}}. \quad (2.25)$$

Эти числа можно интерпретировать как косинус и синус некоторого угла  $\alpha_1$  (отсюда название *метод вращений*, так как один промежуточный шаг прямого хода такого метода может рассматриваться как преобразование вращения на угол  $\alpha_1$  расширенной матрицы системы в плоскости, определяемой индексами обнуляемого элемента<sup>\*)</sup>).

После фиксирования  $c_1$  и  $s_1$  способом (2.25) система (2.1) принимает вид

$$\begin{cases} a_{11}^{(1)}x_1 + a_{12}^{(1)}x_2 + \dots + a_{1n}^{(1)}x_n = b_1^{(1)}, \\ \quad \quad \quad a_{22}^{(1)}x_2 + \dots + a_{2n}^{(1)}x_n = b_2^{(1)}, \\ a_{31}x_1 + a_{32}x_2 + \dots + a_{3n}x_n = b_3, \\ \quad \quad \quad \dots \quad \quad \quad \dots \quad \quad \quad \dots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n, \end{cases} \quad (2.26)$$

где

$$\begin{aligned} a_{1j}^{(1)} &= c_1 a_{1j} + s_1 a_{2j} \quad (j = 1, 2, \dots, n), & b_1^{(1)} &= c_1 b_1 + s_1 b_2; \\ a_{2j}^{(1)} &= -s_1 a_{1j} + c_1 a_{2j} \quad (j = 2, 3, \dots, n), & b_2^{(1)} &= -s_1 b_1 + c_1 b_2. \end{aligned}$$

<sup>\*)</sup> Более подробно эта интерпретация, как и вообще идея ортогональных преобразований, будет рассмотрена позже применительно к решению полной алгебраической проблемы собственных значений (см. пп.4.4, 4.6). Там же, в замечании 4.10 описан другой способ численно устойчивого решения СЛАУ – метод отражений.



менными, то, значит, на любом этапе преобразований длина столбца будет одной и той же, т.е. не будет наблюдаться роста элементов (но какой ценой это достигается! коэффициенты первого уравнения пересчитываются  $n-1$  раз!).

Дальше точно так же за  $n-2$  промежуточных шага преобразуем подсистему

$$\begin{cases} a_{22}^{(1)}x_2 + \dots + a_{2n}^{(1)}x_n = b_2^{(1)}, \\ \dots \dots \dots \dots \dots \dots \dots \dots \\ a_{n2}^{(1)}x_2 + \dots + a_{nn}^{(1)}x_n = b_n^{(1)} \end{cases}$$

системы (2.27), создавая нули под элементом  $a_{22}^{(1)}$ , и т.д.

В результате  $n-1$  таких этапов прямого хода исходная система (2.1) будет приведена к треугольному виду (ср. с (2.3) в п.2.1):

$$\begin{cases} a_{11}^{(n-1)}x_1 + a_{12}^{(n-1)}x_2 + \dots + a_{1n}^{(n-1)}x_n = b_1^{(n-1)}, \\ a_{22}^{(n-1)}x_2 + \dots + a_{2n}^{(n-1)}x_n = b_2^{(n-1)}, \\ \dots \dots \dots \dots \dots \dots \dots \dots \\ a_{nn}^{(n-1)}x_n = b_n^{(n-1)}. \end{cases}$$

Нахождение отсюда неизвестных  $x_n, x_{n-1}, \dots, x_1$  не отличается от рассмотренного ранее обратного хода метода Гаусса.

## 2.8. ДВА ЗАМЕЧАНИЯ К ПРИМЕНЕНИЮ ПРЯМЫХ МЕТОДОВ

**а) О контроле точности и уточнении приближенного решения в рамках прямого метода.**

Как отмечалось во введении, прямые методы приводят к точному решению СЛАУ при точном выполнении предусматриваемых соответствующими алгоритмами арифметических операций (без округлений). Реальные же вычисления базируются на арифметике машинных (т.е. усеченных до определенного количества разрядов) чисел. Как отражается на результате решения системы подмена арифметики действительных чисел машинной арифметикой, зависит от самой решаемой системы, параметров применяемой ЭВМ, способов реализации алгоритмов (см. по этому поводу, в частности, гл. 1). В любом случае, практически вместо точного решения СЛАУ (2.1) прямой метод дает приближенное решение (обозначим его  $x^{(0)}$ ). Подставив  $x^{(0)}$  в выражение

$$\xi = b - Ax,$$

называемое *невязкой*, по малости полученного вектора-значения  $\xi^{(0)} = b - Ax^{(0)}$  можно с осторожностью судить о близости найденного решения  $x^{(0)}$  к точному решению  $x$ . Если, например,  $\|\xi^{(0)}\|$  — недостаточно малая величина, то следует искать вектор-поправку  $p$  такой, что  $x^{(0)} + p = x$  есть точное решение системы (2.1), т.е.

$$A(x^{(0)} + p) = b.$$

Последнее равносильно векторно-матричному уравнению

$$Ap = \xi^{(0)}.$$

Как видим, нахождение поправки сводится к решению такой же системы, как и (2.1), где в качестве вектора свободных членов должен быть взят вектор невязок. Поскольку матрица коэффициентов осталась той же, что и у исходной системы, нет надобности в выполнении прямого хода преобразований коэффициентов при неизвестных (иначе, LU-разложения); достаточно выполнить только действия, касающиеся новых свободных членов (решить две треугольные системы:  $Lz = \xi^{(0)}$  и  $Up = z$ ). Прибавив найденную поправку  $p = p^{(0)}$  к  $x^{(0)}$ , получаем уточненное приближенное решение  $x^{(1)} = x^{(0)} + p^{(0)}$ . В случае, если величина  $\|p^{(0)}\|$  (или  $\|p^{(0)}\|/\|x^{(1)}\|$ , если контролируется относительная, а не абсолютная погрешность) окажется недостаточно малой, процесс уточнения может быть повторен: ищется поправка  $p^{(1)}$  как приближенное решение уравнения  $Ap = \xi^{(1)}$ , где  $\xi^{(1)} = b - Ax^{(1)}$ ; тогда более точным должно быть решение  $x^{(2)} = x^{(1)} + p^{(1)}$ . Как аргументировано утверждается в [15], сходимость к нулю невязок в таком процессе уточнения решения может не наблюдаться<sup>\*)</sup>, т.е. следить нужно за установлением знаков самого решения. Обычно делают не более двух-трех шагов уточнения, причем рекомендуется производить вычисление невязок в режиме накопления. Если в этом процессе не происходит сближения  $x^{(k)}$  при  $k = 2, 3$ , то это говорит скорее всего о том, что данная система плохо обусловлена и ее решение не может быть найдено с требуемой точностью без привлечения дополнительной информации об исходной задаче. В таких случаях закономерно ставить вопрос о том, что понимать под точным решением системы и, возможно, обращаться к методам нахождения ее псевдорешений. (Достаточно глубокие исследования затронутых здесь вопросов, основанные на изучении поведения методов при введении ошибок в исходные данные, а также ряд других све-

<sup>\*)</sup> См. также пример поведения невязок у плохо обусловленной системы на приближенных решениях в гл. 1 (пп.1.4, 1.5).

дений, в частности, точностные характеристики прямых методов, можно найти в [15]).

Хотя описанный здесь контроль точности по невязкам и уточнение решений не требует больших вычислительных затрат, требуемая память ЭВМ должна быть увеличена вдвое, так как при этом нужно удерживать в памяти исходные данные.

### б) О вычислительных затратах.

Одним из факторов, предопределяющих выбор того или иного метода при решении конкретной задачи, является вычислительная эффективность метода. Особенностью прямых методов является то, что здесь можно точно подсчитать требуемое количество арифметических операций. Приведем пример такого подсчета для метода прогонки решения  $n$ -мерной системы с трехдиагональной матрицей коэффициентов (см. п.2.6). Необходимые операции и их число наглядно видны из следующей таблицы (вычитание отождествляется со сложением):

Таблица 2.1

Расчетные формулы метода прогонки	Умножений	Делений	Сложений
$\Delta_i = c_i + b_i \delta_{i-1} \quad (i = 2 \dots n)$	$n-1$		$n-1$
$\delta_i = -\frac{d_i}{c_i}; \delta_i = -\frac{d_i}{\Delta_i} \quad (i = 2 \dots n-1)$		$n-1$	
$\lambda_1 = \frac{r_1}{c_1}; \lambda_i = \frac{r_i - b_i \lambda_{i-1}}{\Delta_i} \quad (i = 2 \dots n)$	$n-1$	$n$	$n-1$
$x_n = \lambda_n; x_i = \delta_i x_{i+1} + \lambda_i \quad (i = n-1 \dots 1)$	$n-1$		$n-1$
Итого арифметических действий	$3(n-1) + 2n-1 + 3(n-1)$		

Учитывая, что часто операции сложения выполняются намного быстрее, чем умножения и деления, обычно ограничиваются подсчетом последних. Так, аналогично показанному нетрудно проверить, что для решения  $n$ -мерной СЛАУ методом Гаусса (без выбора главного элемента) требуется

$\frac{n^3}{3} + n^2 - \frac{n}{3}$  умножений и делений<sup>\*)</sup>, а методом квадратных корней —  $\frac{n^3}{6} + \frac{3}{2}n^2 + \frac{n}{3}$  плюс  $n$  операций извлечения корня (см. [8, 47, 54] и др.).

Метод вращений предполагает вчетверо больше операций умножения, чем метод Гаусса [25]. При больших значениях размерности  $n$  существенным является старший член выражения для подсчета числа арифметических операций. Можно сказать, что вычислительные затраты на операции ум-

<sup>\*)</sup> Несложно подсчитать и общее число операций, включая сложение [16, 17, 21], а также время решения  $n$ -мерной системы, если известна скорость выполнения различных операций [11].

ножения и деления в методе Гаусса составляют величину  $O\left(\frac{n^3}{3}\right)$ , в методе квадратных корней  $O\left(\frac{n^3}{6}\right)$ , в методе вращений  $O\left(\frac{4}{3}n^3\right)$ , в то время как прогонка требует всего  $O(5n)$  таких операций.

## УПРАЖНЕНИЯ

### 2.1. а) Решить систему

$$\begin{cases} 0.1x_1 + 2x_2 - 10x_3 = 0.6, \\ 0.3x_1 + 6.01x_2 - 25x_3 = 1.852, \\ 0.4x_1 + 8.06x_2 + 10.001x_3 = 2.91201 \end{cases}$$

методом Гаусса, пошагово выполняя предписания алгоритма из п.2.1.

б) Перемножением ведущих элементов метода Гаусса найти детерминант матрицы коэффициентов данной системы.

в) Выполнить задания а), б), имитируя работу модельного компьютера, в котором под запись мантиссы числа в режиме с плавающей запятой выделяется три десятичных разряда.

Проделать то же, проводя частичное упорядочивание по столбцам.

Сравнить результаты в) с результатами а) и б).

г) Подсчитав невязки приближенных решений данной системы, полученных в в), произвести итерационное уточнение этих решений в той же вычислительной среде (см. п.2.8а).

### 2.2. а) Выполняя LU-разложение матрицы

$$A = \begin{pmatrix} 4 & -3 & 2 \\ 8 & -8 & 7 \\ 12 & -5 & 5 \end{pmatrix}$$

по формулам (2.9)–(2.10), найти решение системы  $Ax = b$ , где  $b = (0; -12; 4)^T$ .

б) Используя полученное в а) LU-разложение, найти матрицу  $A^{-1}$  двумя способами: решая подсистемы  $AX_i = e_i$ , где  $X_i$  и  $e_i$  – столбцы со-

ответственно искомой и единичной матриц, и реализуя вычисления по формулам (2.15)–(2.17).

2.3. а) Методом квадратных корней решить систему

$$\begin{cases} 16x_1 - 8x_2 - 4x_3 & = -8, \\ -8x_1 + 13x_2 - 4x_3 - 3x_4 & = 7, \\ -4x_1 - 4x_2 + 9x_3 & = 6, \\ & - 3x_2 + 3x_4 = -3. \end{cases}$$

б) С помощью полученного в а)  $U^T U$ -разложения Холецкого найти детерминант матрицы коэффициентов данной системы и обратную ей матрицу.

в) Подсчитав число обусловленности, выяснить, какую относительную погрешность может иметь результат а), если в одной компоненте правой части данной системы допустить абсолютную ошибку 0.01.

2.4. При каких  $n \in \mathbb{N}$  можно гарантировать корректность и устойчивость метода прогонки для решения системы (2.21), где:

$$\begin{aligned} b_i &= 1+i && \text{при } i = 2, 3, \dots, n; \\ c_i &= 15+i && \text{при } i = 1, 2, \dots, n; \\ d_i &= -i && \text{при } i = 1, 2, \dots, n-1? \end{aligned}$$

2.5. Вывести формулы левой прогонки для решения системы (2.21) (т.е. такие, при которых неизвестные  $x_i$  вычислялись бы в порядке возрастания индексов).

Решить систему предыдущего упражнения при  $n = 7$ ,  $r_i = i^2 + 14i - 1$  (где  $i = 1, \dots, 6$ ),  $r_7 = 202$  по формулам левой и правой прогонки.

2.6. Пусть в (2.21)  $b_{i+1} = d_i$  при всех  $i = 1, 2, \dots, n-1$ . Записать для этого случая расчетные формулы метода квадратных корней. Подсчитать число требуемых арифметических операций и сравнить его с аналогичным результатом для метода прогонки (см. п.2.8б).

2.7. Решить систему из упражнения 2.1 методом вращений:

а) используя всю разрядную сетку калькулятора или компьютера;

б) работая, как и в упражнении 2.1в), с тремя значащими цифрами.

Проанализировать результаты, сравнивая их с результатами упражнений 2.1а) и 2.1в).

# ГЛАВА 3 ИТЕРАЦИОННЫЕ МЕТОДЫ РЕШЕНИЯ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ СИСТЕМ И ОБРАЩЕНИЯ МАТРИЦ

Рассматриваются итерационные способы решения систем линейных алгебраических уравнений и обращения матриц, служащие серьезной альтернативой прямым методам решения таких задач, по крайней мере в случаях, когда их размерность велика. Показывается логика построения нескольких наиболее важных итерационных процессов, таких как методы простых итераций, Якоби, Зейделя, релаксации, Шульца, и изучаются условия сходимости последовательностей приближений, получаемых этими методами, к искомым решениям; дается первое представление о методах установления.

## 3.1. РЕШЕНИЕ СЛАУ МЕТОДОМ ПРОСТЫХ ИТЕРАЦИЙ

Система

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

где  $A = (a_{ij})_{i,j=1}^n$  —  $n \times n$ -матрица, а  $x = (x_1, \dots, x_n)^T$  и  $b = (b_1, \dots, b_n)^T$  —  $n$ -мерные векторы-столбцы, тем или иным способом (таких способов существует бесконечное множество; некоторые из них будут рассмотрены ниже) может быть преобразована к эквивалентной ей системе вида

$$x = Bx + c, \quad (3.2)$$

где  $x$  — тот же вектор неизвестных, а  $B$  и  $c$  — некоторые новые матрица и вектор соответственно. Систему (3.2) можно трактовать как задачу о неподвижной точке линейного отображения  $B$  в пространстве  $R_n$  и по аналогии со скалярным случаем (более подробно изучаемым в гл.5) определить последовательность приближений  $x^{(k)}$  к неподвижной точке  $x^*$  рекуррентным равенством

$$x^{(k+1)} = Bx^{(k)} + c, \quad k = 0, 1, 2, \dots \quad (3.3)$$



Итерационный процесс (3.3), начинающийся с некоторого вектора  $x^{(0)} = (x_1^{(0)}, \dots, x_n^{(0)})^T$ , будем называть *методом простых итераций* (кратко МПИ).

Изучим комплекс вопросов о сходимости этого процесса. А именно:

- 1) какие нужно предъявить требования к  $B$ ,  $c$  и  $x^{(0)}$ , чтобы последовательность  $(x^{(k)})$  при  $k \rightarrow \infty$  имела пределом  $x^*$  – неподвижную точку задачи (3.2) (и значит, решение эквивалентной (3.2) исходной системы (3.1))?
- 2) с какой скоростью сходится этот процесс, т.е. каков закон убывания абсолютных погрешностей получаемых по формуле (3.3) приближений?
- 3) сколько нужно сделать итераций по формуле (3.3), чтобы при заданном начальном приближении найти решение задачи (3.2) с заданной точностью?

Ответы на подобные вопросы теории итерационных методов в  $R_n$  часто опираются на следующие два утверждения о сходимости степенных матричных рядов, точнее “матричной геометрической прогрессии”. Во втором из них, а также всюду далее под нормой матрицы понимается мультипликативная норма такая, что  $\|E\| = 1$  ( $E$  – единичная матрица).

**Лемма 3.1<sup>\*)</sup>.** *Условие, что все собственные числа матрицы  $B$  по модулю меньше 1, является необходимым и достаточным для того, чтобы:*

- 1)  $B^k \rightarrow 0$  при  $k \rightarrow \infty$  ( $k \in N$ );
- 2) матрица  $E - B$  имела обратную и

$$(E - B)^{-1} = E + B + B^2 + \dots + B^k + \dots$$

**Лемма 3.2<sup>\*\*)</sup>.** *Если  $\|B\| \leq q < 1$ , то матрица  $E - B$  имеет обратную матрицу  $(E - B)^{-1} = \sum_{k=0}^{\infty} B^k$  и при этом  $\|(E - B)^{-1}\| \leq \frac{1}{1 - q}$ .*

**Доказательство.** Рассмотрим матричный ряд

$$E + B + B^2 + \dots + B^k + \dots \quad (3.4)$$

<sup>\*)</sup> В некоторых литературных источниках лемма 3.1 называется *леммой Неймана* (см., например, [42]).

<sup>\*\*)</sup> В функциональном анализе для более общего случая, когда  $B$  – линейный оператор, действующий в полных нормированных пространствах, эту лемму называют *теоремой Банаха* ([29] и др.).

В силу условия леммы и вытекающего из мультипликативного свойства нормы неравенства  $\|B^k\| \leq \|B\|^k$ , этот ряд можно промажорировать сходящимся числовым рядом:

$$\|E\| + \|B\| + \|B^2\| + \dots + \|B^k\| + \dots \leq 1 + q + q^2 + \dots + q^k + \dots = \frac{1}{1-q}.$$

Следовательно, ряд (3.4) сходится, т.е. существует матрица

$$V = E + B + B^2 + \dots + B^k + \dots$$

такая, что  $\|V\| \leq \frac{1}{1-q}$ . Так как

$$\begin{aligned} (E - B)V &= (E - B)(E + B + B^2 + \dots + B^k + \dots) = \\ &= E + B + B^2 + \dots + B^k + \dots - B - B^2 - B^3 - \dots - B^{k+1} - \dots = E, \end{aligned}$$

то  $V = (E - B)^{-1}$ . Лемма доказана.

Доказательство леммы 1 более сложно. Его можно найти во многих учебных пособиях по вычислительной математике и функциональному анализу (см., например, [8, 21, 29, 34, 42]).

**Теорема 3.1.** *Необходимым и достаточным условием сходимости метода простых итераций (3.3) при любом начальном векторе  $x^{(0)}$  к решению  $x^*$  системы (3.2) является требование, чтобы все собственные числа матрицы  $B$  были по модулю меньше 1.*

**Доказательство. Достаточность.** Пусть  $\max \lambda_B < 1$ , тогда по лемме 3.1 общий член  $B^k$  ряда (3.4) стремится к нуль-матрице и существует матрица  $(E - B)^{-1}$ , являющаяся пределом частичных сумм  $(E + B + B^2 + \dots + B^k)$  при  $k \rightarrow \infty$ . Применяя рекурсию в равенстве (3.3), определяющем МПИ, получим:

$$\begin{aligned} x^{(k+1)} &= Bx^{(k)} + c = B^2x^{(k-1)} + (B + E)c = \dots = \\ &= B^{k+1}x^{(0)} + (E + B + B^2 + \dots + B^k)c. \end{aligned} \quad (3.5)$$

В силу сказанного выше, предел последнего выражения существует при любом фиксированном  $x^{(0)}$  и равен  $(E - B)^{-1}c$ . Следовательно, итерационный процесс (3.3) сходится и

$$x^* := \lim_{k \rightarrow \infty} x^{(k)} = (E - B)^{-1}c.$$

Подставляя  $\mathbf{x}^*$  в уравнение (3.2), преобразованное к виду  $(\mathbf{E} - \mathbf{B})\mathbf{x} = \mathbf{c}$ , имеем:

$$(\mathbf{E} - \mathbf{B})(\mathbf{E} - \mathbf{B})^{-1}\mathbf{c} = \mathbf{c},$$

т.е. вектор  $\mathbf{x}^* = \lim_{k \rightarrow \infty} \mathbf{x}^{(k)}$  удовлетворяет системе (3.2). (Заметим, что это  $\mathbf{x}^*$  – единственное решение (3.2). Действительно, допустив, что наряду с  $\mathbf{x}^*$  таким, что  $\mathbf{x}^* = \mathbf{B}\mathbf{x}^* + \mathbf{c}$ , имеется  $\mathbf{x}^{**}$ , удовлетворяющее такому же равенству  $\mathbf{x}^{**} = \mathbf{B}\mathbf{x}^{**} + \mathbf{c}$ , получаем  $\mathbf{x}^* - \mathbf{x}^{**} = \mathbf{B}(\mathbf{x}^* - \mathbf{x}^{**})$ . Последнее означает, что  $\lambda = 1$  по определению является собственным числом матрицы  $\mathbf{B}$ , что противоречит условию).

*Необходимость.* Как видно из представления общего члена итерационной последовательности  $(\mathbf{x}^{(k)})$  в форме (3.5), существование  $\lim_{k \rightarrow \infty} \mathbf{x}^{(k+1)}$  при любых векторах  $\mathbf{x}^{(0)}$  и  $\mathbf{c}$  (в том числе и нулевых, что гарантирует существование предела каждого слагаемого в правой части (3.5)) влечет сходимость матриц  $\mathbf{B}^{k+1}$  к нуль-матрице и сходимость ряда  $\sum_{k=0}^{\infty} \mathbf{B}^k$  к  $(\mathbf{E} - \mathbf{B})^{-1}$ . Согласно лемме 3.1, это равносильно условию  $\lambda_B < 1$  для каждого собственного числа матрицы  $\mathbf{B}$ .

Теорема доказана.

**Теорема 3.2.** Пусть  $\|\mathbf{B}\| \leq q < 1$ . Тогда при любом начальном векторе  $\mathbf{x}^{(0)}$  МПИ (3.3) сходится к единственному решению  $\mathbf{x}^*$  задачи (3.2) и при всех  $k \in \mathbb{N}$  справедливы оценки погрешности:

$$1) \quad \|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{q}{1-q} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \quad (\text{апостериорная});$$

$$2) \quad \|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{q^k}{1-q} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \quad (\text{априорная}).^*)$$

(Одно и то же обозначение  $\|\cdot\|$  здесь используется для матричных и векторных норм, согласованных между собой, т.е. таких, что  $\|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{x}\|$ ).

<sup>\*)</sup> Лат. a priori и a posteriori означают соответственно "до опыта" и "из опыта", т.е. априорной оценкой можно воспользоваться до начала счета, а апостериорной – лишь после проведения  $k$ -й итерации.

Доказательство. Вычитая из (3.3) равенство  $\mathbf{x}^{(k)} = \mathbf{B}\mathbf{x}^{(k-1)} + \mathbf{c}$ , имеем  $\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} = \mathbf{B}(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$ . Переходя в последнем к нормам, получаем неравенство

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq q \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|, \quad (3.6)$$

из которого видно, в силу условия  $q < 1$ , что элементы итерационной последовательности  $(\mathbf{x}^{(k)})$  сближаются с ростом номера  $k$ . С помощью (3.6) оценим разность между  $(k+m)$ -м и  $k$ -м членами этой последовательности при некотором  $m \in N$ :

$$\begin{aligned} \|\mathbf{x}^{(k+m)} - \mathbf{x}^{(k)}\| &= \|\mathbf{x}^{(k+m)} - \mathbf{x}^{(k+m-1)} + \mathbf{x}^{(k+m-1)} - \mathbf{x}^{(k+m-2)} + \\ &\quad + \mathbf{x}^{(k+m-2)} - \dots - \mathbf{x}^{(k+1)} + \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq \\ &\leq \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| + \|\mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)}\| + \dots + \|\mathbf{x}^{(k+m)} - \mathbf{x}^{(k+m-1)}\| \leq \\ &\leq q \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| + q^2 \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| + \dots + q^m \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| = \\ &= \frac{q(1-q^m)}{1-q} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \frac{q^k}{1-q} (1-q^m) \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|. \end{aligned}$$

Рассматривая итоговое неравенство при  $k \rightarrow \infty$  и фиксированном  $m$ , видим, что  $(\mathbf{x}^{(k)})$  является фундаментальной последовательностью и, в силу полноты пространства  $R_n$ , имеет предел. Обозначим его  $\mathbf{x}^*$ . Переходя к пределу в равенстве (3.3), получаем  $\mathbf{x}^* = \mathbf{B}\mathbf{x}^* + \mathbf{c}$ , т.е.  $\mathbf{x}^* = \lim_{k \rightarrow \infty} \mathbf{x}^{(k)}$  — решение уравнения (3.2). При этом  $\mathbf{x}^*$  — единственное решение (3.2), так как предположив существование другого решения  $\mathbf{x}^{**} \neq \mathbf{x}^*$  и нормировав равенство  $\mathbf{x}^* - \mathbf{x}^{**} = \mathbf{B}(\mathbf{x}^* - \mathbf{x}^{**})$ , приходим к противоречащему условию теоремы неравенству  $\|\mathbf{B}\| \geq 1$ .

Справедливость утверждаемых в теореме оценок погрешности видна из неравенств

$$\|\mathbf{x}^{(k+m)} - \mathbf{x}^{(k)}\| \leq \frac{q(1-q^m)}{1-q} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \frac{q^k}{1-q} (1-q^m) \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|,$$

если в них теперь зафиксировать  $k$  и перейти к пределу при  $m \rightarrow \infty$ . Теорема доказана.

**Замечание 3.1.** Последние неравенства говорят еще о том, что априорная оценка, как правило, грубее апостериорной.

**Замечание 3.2.** Теорема 3.2 могла быть доказана на основе леммы 3.2 и теоремы 3.1. В частности, сходимость  $(x^{(k)})$  к решению  $x^*$  системы (3.2) сразу следует из теоремы 3.1, в силу соотношений  $|\lambda_B| \leq \|B\| < 1$ . Из леммы 3.2 также легко вывести другую априорную оценку погрешности  $k$ -го приближения: вычитая из равенства

$$x^* = (E - B)^{-1}c = (E + B + \dots + B^k + \dots)c$$

равенство

$$x^{(k)} = B^k x^{(0)} + (E + B + \dots + B^{k-1})c$$

(см. (3.5)), имеем:

$$\begin{aligned} \|x^* - x^{(k)}\| &= \|(B^k + B^{k+1} + \dots)c - B^k x^{(0)}\| \leq \\ &\leq \|B^k\| \cdot \|(E - B)^{-1}c - x^{(0)}\| \leq q^k \left( \|x^{(0)}\| + \frac{\|c\|}{1-q} \right). \end{aligned}$$

**Замечание 3.3.** Априорная оценка позволяет подсчитывать заранее число итераций  $k$ , достаточное для получения решения  $x^*$  с заданной точностью  $\varepsilon$  (в смысле допустимого уровня абсолютных погрешностей) при выбранном начальном векторе  $x^{(0)}$ . Для этого нужно найти наименьшее целое решение неравенства

$$\frac{q^k}{1-q} \|x^{(1)} - x^{(0)}\| \leq \varepsilon$$

относительно переменной  $k$  (или неравенства  $q^k \left( \|x^{(0)}\| + \frac{\|c\|}{1-q} \right) \leq \varepsilon$  в соответствии с результатом предыдущего замечания). Апостериорной же оценкой удобно пользоваться непосредственно в процессе вычислений и останавливать этот процесс, как только выполнится неравенство

$$\|x^{(k)} - x^{(k-1)}\| \leq \frac{1-q}{q} \varepsilon.$$

Отметим, что неравенство  $\|x^{(k)} - x^{(k-1)}\| \leq \varepsilon$  будет гарантией того, что и

$$\|x^* - x^{(k)}\| \leq \varepsilon, \text{ только в том случае, когда } q \leq \frac{1}{2}.$$

**Замечание 3.4.** По поводу выбора начального приближения.

Как установлено выше, сходимость МПИ (3.3) гарантируется при любом начальном векторе  $x^{(0)}$ . Очевидно, итераций потребуется тем меньше, чем ближе  $x^{(0)}$  к  $x^*$ . Если нет никакой дополнительной информации о решении задачи (3.2) (например, может быть известным решение близкой задачи или грубое решение данной задачи), за  $x^{(0)}$  обычно принимают вектор с свободных членов системы (3.2). Мотивация этого может

быть такой: матрица  $B$  “мала”, значит вектор  $Bx$  “мал”, следовательно, и вектор  $x^*$  не должен сильно отличаться от вектора  $c$ . При выборе  $x^{(0)} = c$  фигурирующая в теореме 3.2 априорная оценка принимает вид

$$\|x^* - x^{(k)}\| \leq \frac{\|c\|}{1-q} q^{k+1} \quad \forall k \in N.$$

### 3.2. МЕТОД ЯКОБИ

Вернемся к рассмотрению задачи (3.1). После выяснения условия, которому должна удовлетворять матрица коэффициентов приведенной системы (3.2) для сходимости МПИ (3.3), следует осуществить приведение системы (3.1) к виду (3.2) так, чтобы это условие выполнялось. Рассмотрим один из способов такого приведения, достаточно эффективный в определенных случаях.

Представим матрицу  $A$  системы (3.1) в виде

$$A = L + D + R,$$

где  $D$  – диагональная, а  $L$  и  $R$  – соответственно левая и правая строго треугольные (т.е. с нулевой диагональю) матрицы. Тогда система (3.1) может быть записана в виде

$$Lx + Dx + Rx = b, \quad (3.7)$$

и если на диагонали исходной матрицы нет нулей, то эквивалентной (3.1) задачей вида (3.2) будет

$$x = -D^{-1}(L + R)x + D^{-1}b, \quad (3.8)$$

т.е. в (3.2) и (3.3) следует положить

$$B = -D^{-1}(L + R), \quad c = D^{-1}b.$$

Основанный на таком приведении системы (3.1) к виду (3.2) метод простых итераций (3.3) называют *методом Якоби*<sup>\*)</sup>. В векторно-матричных обозначениях он определяется формулой

$$x^{(k+1)} = -D^{-1}(L + R)x^{(k)} + D^{-1}b, \quad k = 0, 1, 2, \dots \quad (3.9)$$

Чтобы записать метод Якоби (3.9) решения системы (3.1) в развернутом виде, достаточно заметить, что обратной матрицей к матрице  $D = (a_{ii})_{i=1}^n$  служит диагональная матрица  $D^{-1}$  с элементами диагонали

<sup>\*)</sup> Карл Густав Якоб Якоби (1804–1851) – немецкий математик.



торной нормой-максимум) меньше единицы. Таким образом, существуют нормы, в которых к методу Якоби, рассматриваемому как метод простых итераций, применима теорема 3.2, т.е. метод Якоби сходится. Так как сходимость по одной норме в пространстве  $R_n$  означает сходимость по любой другой, тем самым теорема 3.3 доказана.

**Замечание 3.5.** Обратим внимание на то, что к методу Якоби при условии диагонального преобладания в матрице  $A$  относится полностью заключение теоремы 3.2, а также предыдущие замечания; нужно лишь учесть в них, что матрица  $B$  определяется с помощью (3.11), а вектор  $c$  – равенством

$$c = \left( \frac{b_1}{a_{11}}; \frac{b_2}{a_{22}}; \dots; \frac{b_n}{a_{nn}} \right)^T.$$

При этом матрица  $D$  заведомо обратима.

Следствием теоремы 1, устанавливающим необходимые и достаточные условия сходимости метода Якоби, является следующая теорема .

**Теорема 3.4.** *Метод Якоби (3.9) сходится к решению системы (3.1) в том и только в том случае, когда все корни уравнения*

$$\begin{vmatrix} a_{11}\lambda & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22}\lambda & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn}\lambda \end{vmatrix} = 0$$

*по модулю меньше единицы.*

Действительно, чтобы все собственные числа матрицы  $B = -D^{-1}(L + R)$  были по модулю меньше единицы, как этого требует теорема 3.1 для данного случая, нужно, чтобы меньше единицы были модули всех корней характеристического уравнения

$$\det(-D^{-1}(L + R) - \lambda E) = 0.$$

Последнее же эквивалентно уравнению

$$\det(L + R + \lambda D) = 0,$$

которое в записи через элементы исходной матрицы  $A$  и фигурирует в формулировке теоремы.







Доказательство. Применяя теорему 3.1 к МПИ (3.14), составляем характеристическое уравнение, определяющее собственные числа  $\lambda$  матрицы  $\mathbf{B} = -(\mathbf{L} + \mathbf{D})^{-1}\mathbf{R}$ :

$$\det\left(-(\mathbf{L} + \mathbf{D})^{-1}\mathbf{R} - \lambda\mathbf{E}\right) = 0.$$

Это уравнение равносильно уравнению

$$\det((\mathbf{L} + \mathbf{D})\lambda + \mathbf{R}) = 0,$$

которое с учетом  $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{R}$  совпадает с (3.16).

**Замечание 3.6.** Уравнение (3.16), а также метод (3.13), являющийся частным случаем более общей формы метода Зейделя (3.12), называют иногда соответственно *уравнением* и *методом Некрасова* [54]. Метод (3.12) называют еще и *методом Гаусса-Зейделя* [41, 42].

Прямым следствием теоремы 2 для метода Зейделя (3.13) является следующая теорема.

**Теорема 3.6.** Пусть  $\|(\mathbf{L} + \mathbf{D})^{-1}\mathbf{R}\| \leq t < 1$ . Тогда при любом начальном векторе  $\mathbf{x}^{(0)}$  метод Зейделя (3.13) сходится к решению  $\mathbf{x}^*$  системы (3.1) и справедливы оценки погрешности

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{t}{1-t} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \frac{t^k}{1-t} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|. \quad (3.17)$$

Ясно, что для непосредственного использования оценок (3.17) нужно предварительно выполнить обращение треугольной матрицы  $\mathbf{L} + \mathbf{D}$  и перемножить матрицы  $(\mathbf{L} + \mathbf{D})^{-1}$  и  $\mathbf{R}$ . В таком случае частично теряется смысл в поэлементной реализации метода Зейделя (3.13); вместо этого можно проводить итерации по формуле (3.14) до тех пор, пока не выполнится условие  $\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \frac{1-t}{t} \varepsilon$ , где  $\varepsilon > 0$  – требуемая точность. В частности, такой подход может быть рекомендован при решении СЛАУ методом Зейделя на ЭВМ с векторной обработкой информации.

Как видим, условия сходимости методов Зейделя и простых итераций, вообще говоря, различаются. Но некоторые достаточные условия можно применять к обоим методам одновременно.

**Теорема 3.7.** Если в матрице  $\mathbf{A}$  системы (3.1) имеет место диагональное преобладание, то метод Зейделя (3.13) сходится, причем быстрее, чем метод Якоби (3.9').

Доказательство. Вычитая тождественное (3.13) равенство (3.14) из равенства (3.15), рассматриваемого как верное равенство при подстановке в него решения  $\mathbf{x}^*$ , получаем

$$\mathbf{x}^* - \mathbf{x}^{(k+1)} = -(\mathbf{L} + \mathbf{D})^{-1} \mathbf{R}(\mathbf{x}^* - \mathbf{x}^{(k)}). \quad (3.18)$$

Введем в рассмотрение вектор ошибок

$$\Delta^{(k)} = \mathbf{x}^* - \mathbf{x}^{(k)}$$

с компонентами  $\delta_i^{(k)} = x_i^* - x_i^{(k)}$ . Через элементы матрицы  $\mathbf{A}$  исходной системы (3.1) равенство (3.18) теперь можно записать так (см. соответствие между (3.14) и (3.13)):

$$\delta_i^{(k+1)} = -\frac{1}{a_{ii}} \sum_{j=1}^{i-1} a_{ij} \delta_j^{(k+1)} - \frac{1}{a_{ii}} \sum_{j=i+1}^n a_{ij} \delta_j^{(k)},$$

где  $i = 1, 2, \dots, n$ ;  $k = 0, 1, 2, \dots$ . Переходя к модулям, отсюда имеем:

$$\begin{aligned} |\delta_i^{(k+1)}| &\leq \frac{1}{|a_{ii}|} \sum_{j=1}^{i-1} |a_{ij}| \cdot |\delta_j^{(k+1)}| + \frac{1}{|a_{ii}|} \sum_{j=i+1}^n |a_{ij}| \cdot |\delta_j^{(k)}| \leq \\ &\leq \left( \max_j |\delta_j^{(k+1)}| \right) \cdot \frac{1}{|a_{ii}|} \sum_{j=1}^{i-1} |a_{ij}| + \left( \max_j |\delta_j^{(k)}| \right) \cdot \frac{1}{|a_{ii}|} \sum_{j=i+1}^n |a_{ij}|. \end{aligned}$$

Обозначим  $\alpha_i = \frac{1}{|a_{ii}|} \sum_{j=1}^{i-1} |a_{ij}|$ ,  $\beta_i = \frac{1}{|a_{ii}|} \sum_{j=i+1}^n |a_{ij}|$  и  $\|\Delta^{(k)}\|_\infty = \max_i |\delta_i^{(k)}|$  (где

$\|\Delta^{(k)}\|_\infty$  может трактоваться как абсолютная погрешность  $k$ -го приближения по методу Зейделя<sup>\*)</sup>). В этих обозначениях последнее неравенство имеет вид<sup>^</sup>

$$|\delta_i^{(k+1)}| \leq \alpha_i \|\Delta^{(k+1)}\|_\infty + \beta_i \|\Delta^{(k)}\|_\infty. \quad (3.19)$$

Пусть  $m \in \{1, 2, \dots, n\}$  – значение индекса  $i$ , при котором реализуется равенство  $|\delta_m^{(k+1)}| = \max_i |\delta_i^{(k+1)}| = \|\Delta^{(k+1)}\|_\infty$ . Тогда из (3.19) следует

$$\|\Delta^{(k+1)}\|_\infty \leq \alpha_m \|\Delta^{(k+1)}\|_\infty + \beta_m \|\Delta^{(k)}\|_\infty,$$

<sup>\*)</sup> Индекс  $\infty$  у знака нормы использован согласно обозначению соответствующего частного случая  $l_p$ -нормы (нормы Гельдера [16], см. приложение 1).

т.е. при этом фиксированном  $i = m$  выполняется неравенство

$$\|\Delta^{(k+1)}\|_{\infty} \leq \frac{\beta_m}{1 - \alpha_m} \|\Delta^{(k)}\|_{\infty}. \quad (3.20)$$

Так как в условиях диагонального преобладания справедливо неравенство  $\alpha_i + \beta_i < 1$ , а это неравенство, в свою очередь, влечет неравенство  $\frac{\beta_i}{1 - \alpha_i} \leq \alpha_i + \beta_i$  (проверьте!), причем равенство в последнем случае имеет место лишь при  $i = 1$ , то абсолютная погрешность приближений по методу Зейделя (3.13), согласно (3.20), убывает со скоростью геометрической прогрессии, знаменатель которой, вообще говоря, меньше, чем для соответствующего этому случаю метода Якоби (3.9'). Такое мажорирование последовательности величин  $\|x^* - x^{(k+1)}\|$  позволяет сделать заключение о справедливости доказываемой теоремы.

**Замечание 3.7.** В соответствии с последней теоремой в методе Зейделя (3.13) вместо оценок (3.17), требующих дополнительных затрат на обращение треугольной матрицы, можно использовать оценки погрешности метода Якоби. Естественно, они заведомо грубее.

Остановимся еще на одном важном для приложений классе систем вида (3.1), для которых имеет место сходимость метода Зейделя (3.13).

**Определение 3.1** [21]. Система  $Ax = b$  называется *нормальной*, если матрица  $A$  – симметричная положительно определенная.

**Теорема 3.8.** Если система (3.1) – нормальная, то метод Зейделя (3.13) сходится.

Доказательство этой теоремы заключается в проверке того, что положительная определенность матрицы  $A = L + D + L^T$  влечет выполнение условия теоремы 3.5 (т.е. собственные числа матрицы  $-(L + D)^{-1}L^T$  по модулю меньше единицы). Это доказательство можно найти; например, в [8, 21].

Любая линейная система  $Ax = b$  легко может быть симметризована умножением на матрицу  $A^T$ . Более того, справедлива следующая теорема.

**Теорема 3.9** [21]. Пусть  $\det A \neq 0$ . Тогда система  $A^T Ax = A^T b$  — нормальная.<sup>\*)</sup>

Таким образом, если, например, известно, что система (3.1) однозначно разрешима, но в ее матрице коэффициентов нет диагонального преобладания, метод Зейделя типа (3.13) можно применять к системе  $A^T Ax = A^T b$ . Правда, при этом возникают трудности со своевременным окончанием процесса итерирования, обеспечивающим заданную точность приближенного решения, так как приведенные ранее оценки погрешности (см. теорему 3.6 и замечание 3.7) в этом случае “не работают”. Да и сходимость в этом случае может оказаться весьма медленной.

Наряду с рассмотренными, применяют и другие способы приведения систем (3.1) к виду (3.2) для их решения методами простых итераций и Зейделя. Достаточно общий подход к этой процедуре заключается в том, что эквивалентное (3.1) уравнение  $0 = b - Ax$  умножается на некоторую неособенную матрицу  $H$  (матричный параметр) и к обеим частям прибавляется вектор  $x$ . Полученное уравнение

$$x = x + H(b - Ax),$$

переписанное в виде

$$x = (E - HA)x + Hb,$$

имеет структуру (3.2). Проблема теперь заключается в подборе матрицы  $H$  такой, чтобы матрица  $B = E - HA$  обладала нужными свойствами для сходимости применяемых методов; для некоторых классов матриц  $A$  имеются определенные рекомендации [21, 34]. Заметим, что матрица  $H$  может быть как постоянной (в этом случае говорят о *стационарном* итерационном процессе), так и изменяющейся от шага к шагу. В последнем случае данное уравнение  $Ax = b$  подменяется последовательностью эквивалентных ему задач  $x = B_k x + c_k$ , и соответствующий итерационный процесс называется *нестационарным*.

### 3.4. ПОНЯТИЕ О МЕТОДЕ РЕЛАКСАЦИИ

В случаях, когда применение оценок погрешностей в методах простых итераций и Зейделя невозможно из-за отсутствия констант  $q < 1$  или  $t < 1$ , ограничивающих сверху какие-либо нормы матрицы итерирования соответствующего метода (см. теоремы 3.2 и 3.6), эти методы неэффективны и, более того, как будет показано в п.3.7, малонадежны ввиду медлен-

---

<sup>\*)</sup> Переход от системы  $Ax = b$  к системе  $A^T Ax = A^T b$  (или в более общем случае к  $A^* Ax = A^* b$ ) называют *симметризацией Гаусса*.

ной сходимости. Рассмотрим одно обобщение метода Зейделя, позволяющее иногда в несколько раз ускорить сходимость итерационной последовательности.

Пусть  $z_i^{(k)}$  – обозначение  $i$ -й компоненты  $k$ -го приближения к решению системы (3.1) по методу Зейделя, а обозначение  $x_i^{(k)}$  будем использовать для  $i$ -й компоненты  $k$ -го приближения новым методом. Этот метод определим равенством

$$x_i^{(k+1)} = x_i^{(k)} + \omega \left( z_i^{(k+1)} - x_i^{(k)} \right), \quad (3.21)$$

где  $i = 1, 2, \dots, n$ ;  $k = 0, 1, 2, \dots$ ;  $x_i^{(0)}$  – задаваемые начальные значения;  $\omega$  – числовой параметр, называемый *параметром релаксации*. Очевидно, при  $\omega = 1$  метод (3.21), называемый *методом релаксации (ослабления)*, совпадает с методом Зейделя<sup>\*)</sup>.

Конкретизируем метод релаксации для случая, когда исходная система (3.1) представляется в виде (3.7) и, следовательно, метод Зейделя имеет вид (3.13).

Пользуясь введенными здесь обозначениями, запишем на основе (3.13) дополнительное к (3.21) равенство для выражения компонент векторов  $z^{(k)} = \left( z_i^{(k)} \right)_{i=1}^n$  через компоненты векторов  $x^{(k)} = \left( x_i^{(k)} \right)_{i=1}^n$ :

$$z_i^{(k+1)} = \frac{1}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right). \quad (3.22)$$

Таким образом, метод релаксации можно понимать как поочередное применение формул (3.22) и (3.21) при каждом  $k = 0, 1, 2, \dots$ . Действительно, задав начальные значения  $x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)}$  и параметр  $\omega$ , при  $k = 0$ , полагая  $i = 1, 2, \dots, n$ , вычислим

$$z_1^{(1)}, x_1^{(1)}; z_2^{(1)}, x_2^{(1)}; \dots; z_n^{(1)}, x_n^{(1)},$$

при  $k = 1$ , так же полагая  $i = 1, 2, \dots, n$ , находим

$$z_1^{(2)}, x_1^{(2)}; z_2^{(2)}, x_2^{(2)}; \dots; z_n^{(2)}, x_n^{(2)},$$

и т.д. Но можно избавиться от вспомогательной последовательности  $\left( z^{(k)} \right)$ , подставив (3.22) в (3.21). Для  $i = 1, 2, \dots, n$  будем иметь:

$$x_i^{(k+1)} = (1 - \omega) x_i^{(k)} + \frac{\omega}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right). \quad (3.23)$$

<sup>\*)</sup> Метод Зейделя в качестве представителя семейства релаксационных методов называют иногда *методом полной релаксации*.

От формулы (3.23), объединяющей формулы (3.22) и (3.21) и пригодной для проведения покоординатных вычислений, мало отличающихся от вычислений по методу Зейделя, легко перейти к векторно-матричной записи процесса релаксации. С этой целью перепишем (3.23) в виде

$$a_{ii}x_i^{(k+1)} + \omega \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} = (1-\omega)a_{ii}x_i^{(k)} - \omega \sum_{j=i+1}^n a_{ij}x_j^{(k)} + \omega b_i$$

и далее, учитывая аддитивное представление матрицы  $\mathbf{A} = \mathbf{L} + \mathbf{D} + \mathbf{R}$ , получаем векторно-матричный итерационный процесс в неявной форме

$$(\mathbf{D} + \omega\mathbf{L})\mathbf{x}^{(k+1)} = (1-\omega)\mathbf{D}\mathbf{x}^{(k)} - \omega\mathbf{R}\mathbf{x}^{(k)} + \omega\mathbf{b}.$$

Умножив последнее равенство слева на матрицу  $(\mathbf{D} + \omega\mathbf{L})^{-1}$ , приходим к эквивалентному (3.23) методу простых итераций

$$\mathbf{x}^{(k+1)} = (\mathbf{D} + \omega\mathbf{L})^{-1}((1-\omega)\mathbf{D} - \omega\mathbf{R})\mathbf{x}^{(k)} + \omega(\mathbf{D} + \omega\mathbf{L})^{-1}\mathbf{b} \quad (3.24)$$

(подстановка сюда значения  $\omega = 1$  превращает (3.24) в МПИ (3.14), эквивалентный методу Зейделя (3.13)).

Представление метода релаксации (3.23) в виде (3.24) позволяет сделать для него некоторые утверждения о сходимости на основании соответствующих теорем о сходимости МПИ. Например, можно применить теоремы 3.1 и 3.2, полагая в них  $\mathbf{B} = (\mathbf{D} + \omega\mathbf{L})^{-1}((1-\omega)\mathbf{D} - \omega\mathbf{R})$ , правда, получаемые при этом результаты вряд ли будут вызывать интерес. Более глубокие результаты на этом пути получают, изучая спектральные свойства таких матриц  $\mathbf{B}$ . Так, установлено, что для сходимости процесса (3.23) необходимо, чтобы  $\omega \in (0; 2)$ . Для некоторых классов СЛАУ (3.1) это требование к параметру релаксации является и достаточным. Справедлива следующая теорема, обобщающая теорему 3.8.

**Теорема 3.10. (Островского-Рейча [41, 48]).** Для нормальной системы  $\mathbf{A}\mathbf{x} = \mathbf{b}$  метод релаксации (3.23) сходится при любом  $\mathbf{x}^{(0)}$  и любом  $\omega \in (0; 2)$ .

Поскольку итерационный процесс (3.23) содержит параметр, естественно распорядиться им так, чтобы сходимость последовательности  $(\mathbf{x}^{(k)})$  была наиболее быстрой. Очевидно, это достигается в том случае, когда спектральный радиус матрицы  $\mathbf{B} = (\mathbf{D} + \omega\mathbf{L})^{-1}((1-\omega)\mathbf{D} - \omega\mathbf{R})$  будет минимальным. В общем случае задача нахождения оптимального значения  $\omega = \omega_0$  не решена, и в практических расчетах применяют метод проб и ошибок. Однако для отдельных важных классов задач такие значения удастся выразить через собственные числа матрицы  $\mathbf{D}^{-1}(\mathbf{L} + \mathbf{R})$  (т.е. корни уравнения, фигурирующего в теореме 3.4) и даже оценить ускорение, дос-



тигаемое введением в процесс Зейделя оптимального параметра релаксации. Существенно отметить, что это оптимальное значение  $\omega_0 \in (1; 2)$ . При значениях  $\omega \in (1; 2)$  метод (3.23) называют *методом последовательной верхней релаксации* (сокращенно ПВР- или SOR-методом<sup>\*)</sup>. Ввиду неэффективности метода (3.23) при  $\omega \in (0; 1)$ , называемого в этом случае *методом нижней релаксации*, название метод ПВР в последнее время относят ко всему семейству методов (3.23), т.е. для любых  $\omega \in (0; 2)$ . При этом случай  $\omega \in (1; 2)$  называют *сверхрелаксацией*.

### 3.5. О ДРУГИХ ИТЕРАЦИОННЫХ МЕТОДАХ РЕШЕНИЯ СЛАУ

В основе построения и изучения или, по крайней мере, понимания многих итерационных методов лежит связь между системами алгебраических уравнений и методами дискретизации дифференциальных уравнений, их порождающими.

В простейшем абстрактном, но далеко не самом общем случае, легко установить такую связь между СЛАУ (3.1) и абстрактным дифференциальным уравнением

$$\frac{dy}{dt} + Ay(t) = b \quad (3.25)$$

с начальным условием  $y(0) = x^{(0)}$ , где  $t$  – абстрактная скалярная переменная, изменяющаяся на промежутке  $[0; +\infty)$ , а матрица  $A$  и вектор  $b$  те же, что и в уравнении (3.1).

Пусть постоянный вектор  $x$  и переменный вектор  $y = y(t)$  – решения задач (3.1) и (3.25) соответственно. Введем вектор

$$z(t) = x - y(t).$$

Учитывая равенство  $\frac{dz}{dt} = -\frac{dy}{dt}$ , из совместного рассмотрения (3.1) и (3.25) выясняем, что  $z(t)$  удовлетворяет однородному дифференциальному уравнению

$$\frac{dz}{dt} = -Az(t)$$

с начальным условием  $z(0) = x - x^{(0)}$ . Решением этой начальной задачи служит вектор

$$z(t) = e^{-At} \cdot z(0),$$

\* От англ. Successive over relaxation

и если спектр  $A$  лежит в правой полуплоскости (в частности, если, например, матрица  $A$  положительно определена), то  $z(t) \xrightarrow{t \rightarrow \infty} 0$  при любых  $z(0)$ . Таким образом, решение  $x$  системы (3.1) (стационарной задачи) может быть получено как предел решения  $y(t)$  задачи Коши (3.25) (эволюционной задачи) при  $t \rightarrow \infty$  с произвольным начальным вектором  $x^{(0)}$ .

Методы приближенного решения стационарных задач, основанные на нахождении решений нестационарных задач, асимптотически эквивалентных данным задачам, для достаточно больших значений искусственной скалярной переменной, называются *методами установления*<sup>\*)</sup>.

Будем далее считать параметром скалярную величину  $\tau_k$ , которую применительно к задаче (3.25) можно интерпретировать как шаг (вообще говоря, переменный), с которым на полуоси  $[0; +\infty)$  фиксируются точки

$$t_0 (= 0), t_1, t_2, \dots,$$

т.е.  $t_{k+1} = t_k + \tau_k$ , где  $k = 0, 1, 2, \dots$ .

При "замораживании"  $t = t_k$  уравнение (3.25) принимает вид

$$\left. \frac{dy}{dt} \right|_{t=t_k} = -Ay(t_k) + b. \quad (3.26)$$

Для производной в его левой части при малых  $\tau_k$  на основе определения можно записать приближенное равенство

$$\left. \frac{dy}{dt} \right|_{t=t_k} = \lim_{\tau_k \rightarrow 0} \frac{y(t_k + \tau_k) - y(t_k)}{\tau_k} \approx \frac{y(t_{k+1}) - y(t_k)}{\tau_k}.$$

Теперь ясно, что полагая  $x^{(k)} := y(t_k)$  (заметим, что  $y(t_0) = y(0) = x^{(0)}$ ), равенство (3.26) можно приближенно заменить равенством

$$\frac{x^{(k+1)} - x^{(k)}}{\tau_k} = -Ax^{(k)} + b, \quad (3.27)$$

---

<sup>\*)</sup> Иногда к дифференциальным уравнениям переходят не от исходной стационарной задачи, а от какого-то конкретного итерационного метода ее решения. Получающаяся при этом асимптотически эквивалентную дифференциальную задачу называют *непрерывным аналогом* соответствующего итерационного метода (см., например, [14, 18]).

которое можно рассматривать как некий *явный* итерационный процесс. Его называют *двухслойным*<sup>\*)</sup> *итерационным методом* [48] или *методом Ричардсона* [47].

Более общий вид семейства двухслойных итерационных методов примет, если ввести в (3.27) невырожденный матричный параметр  $\mathbf{B}_k$ :

$$\mathbf{B}_k \frac{\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}}{\tau_k} = -\mathbf{A}\mathbf{x}^{(k)} + \mathbf{b}. \quad (3.28)$$

Различные конкретные итерационные процессы решения СЛАУ (3.1) (в том числе и все рассмотренные выше) получаются из (3.28) фиксированием матриц  $\mathbf{B}_k$  и скаляров  $\tau_k$ . При этом, если  $\mathbf{B}_k$  и  $\tau_k$  не зависят от  $k$ , т.е. одни и те же на каждой итерации, то (3.28) определяет *стационарный метод*, в противном случае – *нестационарный*. В общем случае, за исключением  $\mathbf{B}_k \equiv \mathbf{E}$ , (3.28) – *неявный метод*.

Выбор параметров  $\mathbf{B}_k$ ,  $\tau_k$  в (3.28) осуществляют, добиваясь удовлетворения каких-то отдельных или совокупности нескольких, возможно в чем-то противоречивых требований таких, как простота, хорошая структура и легкая обрабатываемость матриц  $\mathbf{B}_k$ , и в то же время, как можно более быстрая сходимость последовательности  $(\mathbf{x}^{(k)})$  к решению  $\mathbf{x}^*$  системы (3.1). Разумеется, оптимальность или, скорее, квазиоптимальность некоторых методов рассматриваемого семейства удастся установить лишь при очень жестких ограничениях на решаемую систему (3.1).

Так, например, доказано [47, 48], что если система (3.1) – нормальная с известными границами  $\lambda_{\min} > 0$ ,  $\lambda_{\max} > 0$  спектра ее матрицы коэффициентов, то при заранее зафиксированном (максимальном в реализуемом процессе) числе итераций  $K$  метод (3.27) будет обеспечивать наименьшую погрешность, иначе, минимизировать  $\|\mathbf{x}^* - \mathbf{x}^{(K)}\|$  в том случае, когда параметры  $\tau_k$  вычисляются по формуле

$$\tau_k = \frac{2}{(\lambda_{\max} + \lambda_{\min}) + (\lambda_{\max} - \lambda_{\min})t_{k+1}}, \quad (3.29)$$

где  $k = 0, 1, \dots, K-1$ , а  $t_k = \cos \frac{(2k-1)\pi}{2K}$  – корни полинома Чебышева  $K$ -й степени. Совокупность формул (3.27), (3.29) называют *явным итерационным методом с чебышевским набором параметров*. Имеется обобщение приведенного утверждения и на неявный случай.

<sup>\*)</sup> Смысл термина “двухслойный” становится понятным при изучении численных процессов решения уравнений математической физики. Изучая же численное интегрирование систем ОДУ, обнаруживаем, что (3.27) есть не что иное, как явный метод Эйлера (с переменным шагом) для задачи (3.25).

Дальнейшее формальное развитие методы установления получают как сугубо неявные методы вида (3.28) с матрицами  $B_k$ , представляемыми в виде произведения простых легко обрабатываемых (например, ленточных) матриц, в связи с чем такие методы называются *методами расщепления*. Из методов расщепления наиболее известными являются *методы переменных направлений*<sup>\*)</sup> и *попеременно-треугольный метод*. Неформальное изучение этих методов более целесообразно по месту их применения: при численном решении многомерных задач математической физики.

Рассматриваются также *трехслойные итерационные методы* (в частности, с чебышевскими параметрами), связывающие уже не два, а три соседних приближения:  $x^{(k+1)}$ ,  $x^{(k)}$  и  $x^{(k-1)}$ . В отличие от предыдущих, такие методы являются *двухшаговыми*.

Другой большой класс методов итерационного решения СЛАУ (3.1) – это так называемые *методы вариационного типа*. К ним относятся методы минимальных невязок, минимальных поправок, минимальных итераций, наискорейшего спуска, сопряженных градиентов и т.п. Хорошего понимания и обоснования таких методов можно достигнуть лишь с привлечением теории оптимизации, ибо решение линейной алгебраической системы здесь подменяется решением эквивалентной экстремальной задачи.

А именно, пусть  $Ax = b$  – нормальная  $n$ -мерная система, т.е.  $A$  – положительно определенная симметричная матрица, и пусть  $(\cdot, \cdot)$  – скалярное произведение в пространстве  $R_n$ . Образует квадратичный функционал

$$\Phi(x) = (Ax, x) - 2(b, x) + c, \quad (3.30)$$

где  $c \in R_1$  – произвольная постоянная. Задача решения нормальной системы (3.1) и задача минимизации функционала (3.30) эквивалентны ([46] и др.). Действительно, нормальная система имеет и притом единственное решение; обозначим его  $x^*$ . Тогда при любом векторе  $x = x^* + \Delta$

$$\begin{aligned} \Phi(x) &= \Phi(x^* + \Delta) = (A(x^* + \Delta), x^* + \Delta) - 2(b, x^* + \Delta) + c = \\ &= (Ax^*, x^*) + (A\Delta, x^*) + (Ax^*, \Delta) + (A\Delta, \Delta) - 2(b, x^*) - 2(b, \Delta) + c = \\ &= \Phi(x^*) + (A\Delta, x^*) + (Ax^*, \Delta) - 2(Ax^*, \Delta) + (A\Delta, \Delta) = \\ &= \Phi(x^*) + (A\Delta, \Delta) > \Phi(x^*), \end{aligned}$$

в силу самосопряженности и положительности  $A$ , значит,

$$\Phi(x^*) = \min_{x \in R_n} \Phi(x).$$

<sup>\*)</sup> В зарубежной литературе используется аббревиатура ADI – Alternating Direction Implicit [16].

Теперь можно применять различные методы численной минимизации функционала  $\Phi(\mathbf{x})$  (функции  $n$  переменных  $x_1, x_2, \dots, x_n$ ).

Одним из наиболее популярных и хорошо разработанных методов подобного типа является *метод сопряженных градиентов*. Приведем без вывода алгоритм, быть может, недостаточно подробный, но вполне определенный, чтобы с его помощью можно было решать нормальные СЛАУ (3.1) таким способом [41].

- 1.1. Задать  $\mathbf{x}^{(0)}$  (начальный вектор) и число  $\varepsilon > 0$  (уровень допустимых погрешностей).
- 1.2. Вычислить вектор  $\xi^{(0)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}$  (невязка начального приближения).
- 1.3. Положить  $\mathbf{p}^{(0)} = \xi^{(0)}$ ,  $k = 0$  (номер итерации).
- 2.1. Вычислить скаляр  $\alpha_k = (\xi^{(k)}, \mathbf{p}^{(k)}) / (\xi^{(k)}, \mathbf{A}\mathbf{p}^{(k)})$ .
- 2.2. Вычислить вектор  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$  (очередное приближение).
- 2.3. Вычислить  $\xi^{(k+1)} = \xi^{(k)} - \alpha_k \mathbf{A}\mathbf{p}^{(k)}$  (невязка  $(k+1)$ -го приближения).
- 2.4. Проверить выполнение неравенства  $\|\xi^{(k+1)}\| \leq \varepsilon$ ; если "да", остановить работу алгоритма и вывести результаты.
- 3.1. Вычислить скаляр  $\beta_k = (\xi^{(k+1)}, \mathbf{A}\mathbf{p}^{(k)}) / (\mathbf{p}^{(k)}, \mathbf{A}\mathbf{p}^{(k)})$ .
- 3.2. Вычислить вектор  $\mathbf{p}^{(k+1)} = \xi^{(k+1)} - \beta_k \mathbf{p}^{(k)}$  (новое направление минимизации).
- 3.3. Положить  $k := k + 1$  и вернуться к шагу 2.1.

Интересно определить место, которое занимает этот метод в общей классификации методов решения линейных алгебраических систем. Дело в том, что метод сопряженных градиентов, являясь по форме итерационным, фактически должен быть отнесен к прямым методам, ибо доказано, что с его помощью минимум квадратичной функции (3.30) от  $n$  переменных, иначе, решение  $n$ -мерной линейной системы (3.1), достигается ровно за  $n$  шагов при любом начальном векторе  $\mathbf{x}^{(0)}$ . Применяют же метод сопряженных градиентов именно как итерационный метод (что видно и из при-

<sup>\*)</sup> Полагая  $\xi^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}$ , видим, что в силу 2.2.

$$\xi^{(k+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)} - \mathbf{A}\alpha_k \mathbf{p}^{(k)} = \xi^{(k)} - \alpha_k \mathbf{A}\mathbf{p}^{(k)}.$$

Такое выражение невязки  $\xi^{(k+1)}$  позволяет обходиться без вычисления вектора  $\mathbf{A}\mathbf{x}^{(k+1)}$ .

веденного алгоритма), имея в виду два обстоятельства. Во-первых, реальный вычислительный процесс может быть довольно далек от идеального и, вследствие неизбежных ошибок округления, на  $n$ -м шаге может быть не достигнута нужная точность. Во-вторых, если размерность  $n$  решаемой задачи велика, то число шагов, достаточное для получения решения системы с нужной точностью (т.е. выход по критерию 2.4), может оказаться значительно меньшим этой ( $n$ ) теоретической величины.

Простейший вариант *метода минимальных невязок* определяется совокупностью формул

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \tau_k \xi^{(k)}, \quad \xi^{(k)} = \mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}, \quad \tau_k = \frac{(\mathbf{A}\xi^{(k)}, \xi^{(k)})}{(\mathbf{A}\xi^{(k)}, \mathbf{A}\xi^{(k)})}.$$

Его можно рассматривать как явный двухслойный итерационный процесс (3.27), в котором параметр  $\tau_k$  на каждом итерационном шаге  $k = 0, 1, 2, \dots$  выбирается таким, чтобы минимизировалась евклидова норма невязки  $\xi^{(k+1)}$  получаемого приближения  $\mathbf{x}^{(k+1)}$ .

Действительно, вычтем из вектора  $\xi^{(k+1)} = \mathbf{A}\mathbf{x}^{(k+1)} - \mathbf{b}$  вектор  $\xi^{(k)}$ .

Имеем

$$\xi^{(k+1)} - \xi^{(k)} = \mathbf{A}\mathbf{x}^{(k+1)} - \mathbf{A}\mathbf{x}^{(k)} = \mathbf{A}\mathbf{x}^{(k)} - \mathbf{A}\tau_k \xi^{(k)} - \mathbf{A}\mathbf{x}^{(k)},$$

т.е.

$$\xi^{(k+1)} = \xi^{(k)} - \tau_k \mathbf{A}\xi^{(k)}.$$

Возводя последнее равенство в квадрат (в смысле скалярного умножения векторов), получаем

$$(\xi^{(k+1)}, \xi^{(k+1)}) = (\xi^{(k)}, \xi^{(k)}) - 2\tau_k (\mathbf{A}\xi^{(k)}, \xi^{(k)}) + \tau_k^2 (\mathbf{A}\xi^{(k)}, \mathbf{A}\xi^{(k)})$$

или, что то же,

$$\|\xi^{(k+1)}\|^2 = \|\xi^{(k)}\|^2 - 2\tau_k (\mathbf{A}\xi^{(k)}, \xi^{(k)}) + \tau_k^2 \|\mathbf{A}\xi^{(k)}\|^2.$$

Легко видеть, что минимум этой положительной квадратичной функции (значит, и величины  $\|\xi^{(k+1)}\|^2$ ) достигается именно при указанном в записи метода значении  $\tau_k$ .

В случае нормальной системы для метода минимальных невязок можно получить ту же оценку скорости сходимости, что и для метода простой итерации

$$\mathbf{x}^{(k+1)} = (\mathbf{E} - \tau \mathbf{A})\mathbf{x}^{(k)} + \tau \mathbf{b}$$

при оптимальном значении параметра  $\tau = \frac{2}{\lambda_{\min} + \lambda_{\max}}$  (в предположении, что известны границы  $\lambda_{\min}$  и  $\lambda_{\max}$  спектра матрицы  $\mathbf{A}$ ) [47].

Рассмотренные здесь методы далеко не исчерпывают все многообразие итерационных способов решения СЛАУ. В частности, нами совсем не затрагивалась проблема решения больших разреженных систем, где на первый план выходят блочные методы, максимально сохраняющие исходную разреженность матриц.

### 3.6. БЫСТРОСХОДЯЩИЙСЯ ИТЕРАЦИОННЫЙ СПОСОБ ОБРАЩЕНИЯ МАТРИЦ

Согласно леммам 3.1, 3.2 (см. п.3.1), если матрица  $\mathbf{B} = \mathbf{E} - \mathbf{A}$  мала (в смысле ее нормы или собственных значений), то обратная к  $\mathbf{A}$  матрица

$$\mathbf{A}^{-1} = (\mathbf{E} - \mathbf{B})^{-1} = \mathbf{E} + \mathbf{B} + \mathbf{B}^2 + \dots,$$

в принципе, может быть найдена сколь угодно точно приближенным суммированием данного матричного ряда. Однако такой непосредственный подход к вычислению имеет два очевидных недостатка: во-первых, его можно применить лишь для обращения матриц, близких к единичной, во-вторых, сходимость последовательностей частичных сумм этого ряда будет медленной даже при достаточно малых нормах матриц  $\mathbf{B}$ . Поэтому, пользуясь отмеченным фактом лишь как теоретической основой, построим итерационный процесс, определяющий существенно более быстро сходящуюся последовательность приближений к обратной для  $\mathbf{A}$  матрице  $\mathbf{A}^{-1}$ . Будем далее обозначать эти приближения, получаемые на  $k$ -м шаге, через  $\mathbf{U}_k$ , а их невязки  $\mathbf{E} - \mathbf{A}\mathbf{U}_k$  — через  $\Psi_k$ .

**Лемма 3.3.** *Если для матрицы  $\mathbf{A}$  найдется такая обратимая матрица  $\mathbf{U}_0$ , что модули всех собственных чисел матрицы  $\Psi_0 = \mathbf{E} - \mathbf{A}\mathbf{U}_0$  меньше единицы, то матрица  $\mathbf{A}$  обратима и для обратной матрицы справедливо представление*

$$\mathbf{A}^{-1} = \mathbf{U}_0(\mathbf{E} - \Psi_0)^{-1} = \mathbf{U}_0(\mathbf{E} + \Psi_0 + \Psi_0^2 + \dots). \quad (3.31)$$

**Доказательство.** Из равенства

$$\mathbf{A}\mathbf{U}_0 = \mathbf{E} - \Psi_0, \quad (3.32)$$

в силу обратимости  $\mathbf{U}_0$  и  $\mathbf{E} - \Psi_0$  (последнее по лемме 3.1), имеем

$$\mathbf{A} = (\mathbf{E} - \Psi_0)\mathbf{U}_0^{-1} = \left( \left( (\mathbf{E} - \Psi_0)^{-1} \right)^{-1} \mathbf{U}_0^{-1} \right) = \left( \mathbf{U}_0(\mathbf{E} - \Psi_0)^{-1} \right)^{-1},$$

т.е. матрица  $\mathbf{A}$  обратима и

$$\mathbf{A}^{-1} = \mathbf{U}_0(\mathbf{E} - \Psi_0)^{-1}.$$

Доказательство завершается разложением  $(\mathbf{E} - \Psi_0)^{-1}$  в матричный ряд (лемма 3.1).

Очевидным следствием лемм 3.2 и 3.3 является следующая лемма.

**Лемма 3.4.** Пусть матрица  $U_0$  обратима и  $\|\Psi_0\| < 1$ .

Тогда:

- 1) существует матрица  $A^{-1}$ ;
- 2) справедливо представление  $A^{-1}$  по формуле (3.31);
- 3) имеет место оценка  $\|A^{-1}\| \leq \frac{\|U_0\|}{1 - \|\Psi_0\|}$ .

Для построения итерационного процесса зафиксируем в разложении (3.31)  $m+1$  первых слагаемых и будем считать первым приближением к  $A^{-1}$  матрицу

$$U_1 = U_0(E + \Psi_0 + \dots + \Psi_0^m).$$

Найдем выражение невязки  $\Psi_1$  этого приближения через невязку  $\Psi_0$  предыдущего (в данном случае начального) приближения  $U_0$ :

$$\begin{aligned} \Psi_1 &= E - AU_1 = E - AU_0(E + \Psi_0 + \dots + \Psi_0^m) = \\ &= E - (E - \Psi_0)(E + \Psi_0 + \dots + \Psi_0^m) = E - (E - \Psi_0^{m+1}) = \Psi_0^{m+1}. \end{aligned} \quad (3.33)$$

Благодаря полученной связи между невязками, можно утверждать, что если выполняются условия лемм 3.3 или 3.4 по отношению к матрицам  $U_0, \Psi_0$ , то для матриц  $U_1, \Psi_1$  они будут выполнены и подавно. Следовательно, к матрицам  $U_1, \Psi_1$  можно применить все рассуждения, проведенные выше для  $U_0, \Psi_0$ . Таким образом, приходим к итерационному процессу

$$\begin{cases} \Psi_k = E - AU_k, \\ U_{k+1} = U_k(E + \Psi_k + \dots + \Psi_k^m). \end{cases} \quad (3.34)$$

где  $k = 0, 1, 2, \dots$  – номер итерации;  $U_0$  – задаваемая начальная матрица, близкая к  $A^{-1}$  в указанном выше смысле, а  $m \in \mathbb{N}$  – параметр метода.

Изучим сходимости этого процесса.

**Теорема 3.11.** Пусть квадратные матрицы  $A$  и  $U_0$  таковы, что матрица  $U_0$  обратима и  $\|\Psi_0\| < 1$ . Тогда существует обратная к  $A$  матрица  $A^{-1}$  и к ней сходится последовательность матриц  $U_k$ , определяемая итерационным процессом (3.34). При этом имеет место точное равенство

$$A^{-1} - U_k = (A^{-1} - U_0)\Psi_0^{(m+1)k-1} \quad (3.35)$$

и справедливы оценки погрешности:



$$1) \|A^{-1} - U_k\| \leq \frac{\|U_k \Psi_k\|}{1 - \|\Psi_k\|};$$

$$2) \|A^{-1} - U_k\| \leq \frac{\|U_0\|}{1 - \|\Psi_0\|} \cdot \|\Psi_0\|^{(m+1)^k}.$$

Доказательство. Существование  $A^{-1}$  следует из леммы 3.4. Упомянутая повторяемость рассуждений и выкладок, проведенных на первом итерационном шаге, позволяет считать очевидными равенства типа (3.31), (3.33) для  $k$ -й итерации:

$$A^{-1} = U_k (E - \Psi_k)^{-1} = U_k (E + \Psi_k + \Psi_k^2 + \dots), \quad (3.36)$$

$$\Psi_k = E - AU_k = \Psi_{k-1}^{m+1} = \Psi_{k-2}^{(m+1)^2} = \dots = \Psi_0^{(m+1)^k}. \quad (3.37)$$

Из (3.31) имеем

$$\begin{aligned} A^{-1} - U_0 &= U_0 (E + \Psi_0 + \Psi_0^2 + \dots) - U_0 = \\ &= U_0 (E + \Psi_0 + \Psi_0^2 + \dots) \Psi_0 = A^{-1} \Psi_0, \end{aligned} \quad (3.38)$$

а из (3.36) аналогично (с учетом (3.37)) получаем

$$A^{-1} - U_k = A^{-1} \Psi_k = A^{-1} \Psi_0^{(m+1)^k}. \quad (3.39)$$

Заменяя здесь в правой части  $A^{-1} \Psi_0$  на  $A^{-1} - U_0$  (см. (3.38)), получаем утверждаемое в теореме равенство (3.35). Переходя в нем к нормам, в соответствии с условием заключаем, что

$$\|A^{-1} - U_k\| \leq \|A^{-1} - U_0\| \cdot \|\Psi_0\|^{(m+1)^k - 1} \xrightarrow{k \rightarrow \infty} 0,$$

т.е. имеет место сходимость  $(U_k)_{k=1}^{\infty}$  к  $A^{-1}$  по норме, а значит, и поэлементная сходимость.

Для доказательства первой оценки (апостериорной) вычтем  $U_k$  из (3.36):

$$A^{-1} - U_k = U_k (E + \Psi_k + \Psi_k^2 + \dots) - U_k = U_k \Psi_k (E - \Psi_k)^{-1}.$$

Отсюда по лемме 3.2 с учетом (3.37) получаем требуемую оценку 1).

Вторая оценка (априорная) может быть найдена в результате закругления первой. Но можно вывести ее непосредственно из равенства (3.39), подставив в его правую часть вместо  $A^{-1}$  выражение  $U_0 (E - \Psi_0)^{-1}$  (см. (3.31)):

$$A^{-1} - U_k = U_0 (E - \Psi_0)^{-1} \Psi_0^{(m+1)^k}.$$

Переход к нормам в последнем равенстве и привлечение леммы 3.2 завершает доказательство теоремы.

Равенства (3.34) определяют фактически не один, а целое семейство итерационных методов обращения. Фиксированием параметра  $m = 1, 2, \dots$  можно получать конкретные процессы  $(m+1)$ -го порядка скорости сходимости<sup>\*)</sup>. Этот порядок может быть сколь угодно большим, однако обычно ограничиваются процессами второго ( $m = 1$ ) и третьего ( $m = 2$ ) порядков. Приоритет процесса второго порядка связан с его простотой и более ранним появлением: первая публикация об этом методе относится к 1933 г. и принадлежит Г. Шульцу [61], в связи с чем и все семейство (3.34) естественно называть *методом Шульца*<sup>\*\*)</sup>. Метод третьего порядка целесообразно использовать из тех соображений, что он, как показал М. Альтман [59], обладает свойством минимальности вычислительных затрат, требующихся для обращения матриц с заданной точностью методами семейства (3.34).

Отметим, что как сам быстросходящийся итерационный процесс (3.34), так и представленные теоремой 3.11 результаты можно без каких-либо особых дополнительных условий отнести к более общей задаче обращения линейных ограниченных операторов в полных нормированных пространствах.

Процесс (3.34) построения приближений к обратной матрице легко видоизменить подобно тому, как это было сделано с методом простых итераций решения СЛАУ, когда для более оперативного учета получаемой на текущей итерации информации перешли от него к методу Зейделя (см. п.3.3). Например, *зейделева модификация метода Шульца второго порядка* может быть определена равенствами

$$\begin{cases} \Psi_k = E - AU_k, \\ U_{k+1} = U_k + U_k \underline{\Psi}_k + U_{k+1} \overline{\Psi}_k. \end{cases} \quad (3.40)$$

где  $k = 0, 1, 2, \dots$ ;  $\Psi_k = \underline{\Psi}_k + \overline{\Psi}_k$ , а  $\underline{\Psi}_k$  и  $\overline{\Psi}_k$  — соответственно нижняя треугольная и строго верхняя треугольная матрицы [14]. При реализации этой модификации нужно либо расписывать формулы (3.40) поэлементно (чтобы не работать с заведомо нулевыми элементами), либо формировать матрицу  $U_{k+1}$  постепенным замещением старых элементов новыми, осуществляя на  $k$ -й итерации цикл присвоений

$$U := U + U\Psi,$$

<sup>\*)</sup> Определение порядка итерационного процесса см. далее в п.5.3

<sup>\*\*)</sup> В разных литературных источниках можно встретить и другие названия этого метода: *Хотеллинга*, *Бодевига* (Бодвига), а также *Нобо*.

где до начала цикла в правой части в двумерном массиве  $U$  должна содержаться матрица  $U_k$ , а в двумерном массиве  $\Psi$  – матрица  $\Psi_k$  (заполнение массивов новыми элементами производится по строкам). Процесс (3.40) при том же шаговом объеме вычислений и такой же простоте, что и в методе Шульца второго порядка, может дать определенный выигрыш в скорости сходимости. Это можно показать сравнением невязок первых приближений при тех или иных предположениях относительно начальной невязки, иначе, при тех или иных требованиях к начальной матрице  $U_0$  [14].

Вообще, проблема выбора начального приближения  $U_0$  в рассматриваемых здесь процессах итерационного обращения матриц не позволяет относиться к ним как к самостоятельным универсальным методам, конкурирующим с прямыми методами обращения, основанными, например, на LU-разложении матриц. Имеются некоторые рекомендации по выбору  $U_0$  (см. [8, 59] и др.), обеспечивающие выполнение условия  $\rho(\Psi_0) < 1$ , являющегося необходимым и достаточным для сходимости процесса (3.34). Однако при этом, во-первых, требуется знать оценку сверху спектра обращаемой матрицы  $A$  либо матрицы  $AA^T$  (а именно, если  $A$  – симметричная положительно определенная и  $\rho(A) \leq \beta$ , то можно взять  $U_0 = \alpha E$ , где

$\alpha \in \left(0, \frac{2}{\beta}\right)$ ; если же  $A$  – произвольная невырожденная матрица и

$\rho(AA^T) \leq \beta$ , то полагают  $U_0 = \alpha A^T$ , где также  $\alpha \in \left(0, \frac{2}{\beta}\right)$ ; можно, конечно, упростить ситуацию и, воспользовавшись тем, что  $\rho(AA^T) \leq \|AA^T\|$ ,

положить  $U_0 = \frac{A^T}{\|AA^T\|}$ ). Во-вторых, при таком задании начальной матрицы

нет гарантии, что  $\|\Psi_0\|$  будет малой (возможно, даже окажется  $\|\Psi_0\| > 1$ ), и высокий порядок скорости сходимости обнаружится далеко не сразу.

Все сказанное выше не означает, что подобные методы обращения матриц (и операторов) не имеют своей сферы применения. В частности, ниже (в п.б.2) будет рассматриваться способ решения систем нелинейных уравнений, базирующийся на методе Ньютона с приближенным обращением матриц Якоби по методу Шульца.

---

<sup>\*)</sup> Через  $\rho(\cdot)$  здесь обозначается спектральный радиус указанной в скобках матрицы.

### 3.7. О РОЛИ ОШИБОК ОКРУГЛЕНИЯ В ИТЕРАЦИОННЫХ МЕТОДАХ

Обратимся, наконец, к вопросам практической реализации итерационных методов решения линейных алгебраических задач.

Многие утверждения о сходимости итерационных процессов говорят о том, что решение поставленной задачи при определенных условиях может быть найдено этим процессом сколь угодно точно, причем погрешность каждого приближения может быть эффективно проконтролирована (см. теоремы 3.2, 3.6, 3.11, а также теорему 3.3 с замечанием 3.5 и теорему 3.7 с замечанием 3.7). Нетрудно понять, что все это справедливо на самом деле лишь до тех пор, пока на погрешность метода (остаточную погрешность) не наложится вычислительная погрешность (погрешность округлений), неизбежная при любых реальных компьютерных расчетах. Особенно существенное и даже пагубное влияние на результат решения задачи итерационным методом могут оказать ошибки округления в тех случаях, когда утверждения о сходимости метода не содержат эффективных оценок погрешности (теоремы 3.1, 3.4, 3.5, 3.8, 3.10).

Рассмотрим различие между реальным и идеальным итерационными процессами на простейшем объекте – на методе простой итерации.

Пусть на  $k$ -м итерационном шаге вычислений по методу (3.3) ошибки округлений составляют вектор  $\gamma^{(k)}$ . Тогда в отличие от идеального МПИ (3.3), генерирующего последовательность приближений  $x^{(k)}$  к решению  $x^*$  системы (3.1) такому, что

$$x^* = Bx^* + c, \quad (3.41)$$

реальный МПИ будет иметь вид

$$\tilde{x}^{(k+1)} = B\tilde{x}^{(k)} + c + \gamma^{(k)}. \quad (3.42)$$

Изучим поведение векторов

$$\mu_k := \tilde{x}^{(k)} - x^*$$

– ошибок приближений  $\tilde{x}^{(k)}$ , получаемых реальным МПИ (3.42).

Вычитая (3.41) из (3.42), имеем

$$\tilde{x}^{(k+1)} - x^* = B(\tilde{x}^{(k)} - x^*) + \gamma^{(k)},$$

т.е.

$$\begin{aligned} \mu_{k+1} &= B\mu_k + \gamma^{(k)} = B(B\mu_{k-1} + \gamma^{(k-1)}) + \gamma^{(k)} = \\ &= B^2(B\mu_{k-2} + \gamma^{(k-2)}) + B\gamma^{(k-1)} + \gamma^{(k)} = \dots = \\ &= B^{k+1}\mu_0 + (B^k\gamma^{(0)} + B^{k-1}\gamma^{(1)} + \dots + B\gamma^{(k-1)} + \gamma^{(k)}). \end{aligned} \quad (3.43)$$

Первое слагаемое в последнем выражении отвечает за погрешность идеального МПИ и может быть сделано сколь угодно малым в процессе итерирования при условии  $\rho(\mathbf{B}) < 1$  (см. лемму 3.1). Чтобы оценить второе слагаемое, предположим, что порог абсолютных погрешностей округлений, допускаемых на каждой итерации, есть  $\gamma$ , т.е.

$$\|\gamma^{(k)}\| \leq \gamma \quad \forall k \in N_0.$$

Тогда

$$\begin{aligned} & \|\mathbf{B}^k \gamma^{(0)} + \mathbf{B}^{k-1} \gamma^{(1)} + \dots + \mathbf{B} \gamma^{(k-1)} + \gamma^{(k)}\| \leq \\ & \leq \gamma \|\mathbf{E} + \mathbf{B} + \dots + \mathbf{B}^k\|, \end{aligned}$$

и, если  $\|\mathbf{B}\| \leq q < 1$ , то второе слагаемое в (3.43), хотя и не стремится к нулю, но ограничено по норме величиной

$$\gamma \frac{1-q^k}{1-q} < \frac{\gamma}{1-q}.$$

При условии же  $\rho(\mathbf{B}) < 1$ , теоретически обеспечивающем сходимость идеального МПИ (3.3), малость этого второго слагаемого отнюдь не гарантируется, что означает допустимость ситуаций, когда в ходе реальных итераций погрешность округлений будет накапливаться вплоть до переполнения множества чисел, представляемых используемой ЭВМ.

Более детальный анализ влияния ошибок округления на итерационный процесс с попыткой пролить свет на природу этого влияния можно найти, например, в [6]. Здесь же ограничимся напоминанием о том, что необходимо с осторожностью применять процессы, когда для них нет эффективных оценок погрешности, и по возможности, учитывать влияние ошибок округления, если такие оценки есть. Например, применительно к МПИ решения СЛАУ выше фактически доказана

**Теорема 3.12.**<sup>\*)</sup> Пусть  $\|\mathbf{B}\| \leq q < 1$  и приближения  $\tilde{\mathbf{x}}^{(k)}$  к решению  $\mathbf{x}^*$  системы (3.2) получаются посредством равенства (3.42), где  $\gamma^{(k)}$  – вектор ошибок округлений таких, что  $\|\gamma^{(k)}\| \leq \gamma$ . Тогда погрешность  $k$ -го приближения при любом  $k \in N$  можно оценить неравенством

$$\|\mathbf{x}^* - \tilde{\mathbf{x}}^{(k)}\| \leq \frac{q}{1-q} \|\tilde{\mathbf{x}}^{(k)} - \tilde{\mathbf{x}}^{(k-1)}\| + \frac{\gamma}{1-q}. \quad (3.44)$$

<sup>\*)</sup> См. также [1].

Действительно, для последовательности  $(x^{(k)})$ , получаемой МПИ (3.3), справедливо равенство

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

Следовательно, считая, что процессы (3.3) и (3.42) начинаются с одного начального приближения  $x^{(0)} = \tilde{x}^{(0)}$ , в идентичном (3.43) равенстве

$$x^* - \tilde{x}^{(k+1)} = B^{k+1}(x^* - \tilde{x}^{(0)}) - (B^k \gamma^{(0)} + B^{k-1} \gamma^{(1)} + \dots + B \gamma^{(k-1)} + \gamma^{(k)})$$

можно заменить  $B^{k+1}(x^* - \tilde{x}^{(0)})$  на  $x^* - x^{(k+1)}$ . Таким образом, погрешности  $(k+1)$ -х приближений реального (3.42) и идеального (3.3) методов различаются лишь слагаемым, оцененным выше по норме величиной  $\frac{\gamma}{1-q}$ ,

т.е. и для процесса (3.42) можно воспользоваться оценкой, выведенной в теореме 3.2.

Отметим, что как непосредственно видно из оценки (3.44) (при  $q$ , приближающихся к единице), роль ошибок округлений в образовании общей погрешности тем сильнее, чем медленнее сходимость итерационного процесса.

## УПРАЖНЕНИЯ

3.1. Записать итерационный процесс Якоби нахождения решения системы

$$\begin{cases} 5x_1 + 2x_2 - x_3 + x_4 = 9, \\ x_1 - 4x_2 + 2x_4 = 10, \\ 2x_1 + 3x_2 - 9x_3 - x_4 = -10, \\ 3x_1 + x_3 - 6x_4 = -5. \end{cases}$$

Каким должен быть критерий окончания процесса итерирования, чтобы максимальная из абсолютных погрешностей компонент приближенного решения не превышала заданного малого  $\varepsilon > 0$ ?

3.2. Сделать по пять итераций методов Якоби и Зейделя для системы

$$\begin{cases} 10x_1 + x_2 - 2x_3 = 10, \\ x_1 - 5x_2 + x_3 = 10, \\ 3x_1 - x_2 + 10x_3 = -5. \end{cases}$$

Сколько верных знаков можно гарантировать в приближенных решениях, полученных тем и другим способами?

3.3. Доказать, что при любом начальном векторе  $(x^{(0)}, y^{(0)}, z^{(0)})^T$  последовательности векторов  $(x_1^{(k)}, y_1^{(k)}, z_1^{(k)})^T$  и  $(x_2^{(k)}, y_2^{(k)}, z_2^{(k)})^T$ , определяемые при  $k = 0, 1, 2, \dots$  равенствами

$$\begin{cases} x_1^{(k+1)} = 0.1x_1^{(k)} + 0.2y_1^{(k)} - 3, \\ y_1^{(k+1)} = 0.2x_1^{(k)} - 0.1y_1^{(k)} + 0.1z_1^{(k)} + 2, \\ z_1^{(k+1)} = -0.3x_1^{(k)} + 0.2z_1^{(k)} - 1 \end{cases}$$

и

$$\begin{cases} x_2^{(k+1)} = (2y_2^{(k)} - 30)/9, \\ y_2^{(k+1)} = (2x_2^{(k)} + z_2^{(k)} + 20)/11, \\ z_2^{(k+1)} = -(3x_2^{(k)} + 10)/8, \end{cases}$$

сходятся, причем к одному и тому же предельному вектору  $(x^*, y^*, z^*)^T$ .

Записать линейную систему, решением которой служит этот предельный вектор. За сколько шагов итераций по данным формулам можно получить предельный вектор с точностью  $\varepsilon = 10^{-6}$  (по норме-максимум), если начать счет с нулевого вектора?

3.4. Пусть методом Якоби решение системы

$$\begin{aligned} b_i x_{i-1} + c_i x_i + d_i x_{i+1} &= r_i \\ (i = 1, 2, \dots, n; \quad b_1 = d_n = 0) \end{aligned}$$

с нужной точностью достигается за  $k$  шагов. Существуют ли такие  $k$  и  $n$ , при которых метод Якоби эффективнее метода прогонки по числу арифметических операций?

3.5. Для линейной системы

$$\begin{cases} x_1 + 2x_2 + 3x_3 & = 5, \\ 2x_1 - x_2 + & + 2x_4 = 8, \\ 3x_1 + x_2 - 2x_3 - x_4 & = -1, \\ 4x_1 - & x_3 + 2x_4 = 8 \end{cases}$$

записать метод Зейделя и обосновать его сходимость. Каковы расчетные формулы метода ПВР в этом случае?

### 3.6. Дана система

$$\begin{cases} 7x_1 + 5x_2 + x_3 = 2.2, \\ 5x_1 + 8x_2 + 2x_3 = 2.4, \\ x_1 + 2x_2 + 4x_3 = 1.6. \end{cases}$$

Найти четвертое приближение к ее решению по методу минимальных невязок, начиная итерационный процесс с нулевого вектора. За сколько итераций по методу Якоби достигается примерно такая же величина евклидовой нормы невязки? Сравнить вычислительные затраты, требующиеся для реализации одного шага каждого из этих методов.

3.7. Предположим, что некоторая  $n \times n$ -система вида  $x = Bx + c$  с  $\|B\| \approx 0.5$  решается методом простых итераций с уровнем абсолютных погрешностей арифметических операций порядка  $10^{-6}$ . Допустим, что при этом  $\|x^{(1)} - x^{(0)}\| \approx 1$ . Каким числом следует ограничить количество итераций, чтобы вычислительная погрешность не стала существенно превышать погрешность метода?

### 3.8. Даны матрицы

$$A = \begin{pmatrix} 1 & -2 & 3 \\ -1 & 1 & 2 \\ 2 & -1 & -1 \end{pmatrix} \quad \text{и} \quad U_0 = \begin{pmatrix} -0.1 & 0.6 & 0.9 \\ -0.4 & 0.9 & 0.6 \\ 0.1 & 0.4 & 0.1 \end{pmatrix}.$$

а) Подсчитав невязку  $\Psi_0 = E - AU_0$ , убедиться в существовании матрицы  $A^{-1}$  и оценить какую-либо ее норму.

б) Сделать по два приближения к  $A^{-1}$  методом Шульца второго и третьего порядков и оценить близость полученных приближений к  $A^{-1}$ .

в) Сравнить оценки погрешностей приближений с истинными ошибками, найдя  $A^{-1}$  каким-нибудь прямым методом.



# ГЛАВА 4 || МЕТОДЫ РЕШЕНИЯ АЛГЕБРАИЧЕСКИХ ПРОБЛЕМ СОБСТВЕННЫХ ЗНАЧЕНИЙ

Затрагивается наиболее сложная задача вычислительной линейной алгебры – нахождение собственных чисел и собственных векторов матриц. Рассматриваются современные подходы к решению спектральных задач для вещественных матриц умеренной размерности, базирующиеся на прямой и обратной итерациях (в том числе, со сдвигами), а также на приведении матриц к диагональной или треугольной формам ортогональными преобразованиями подобия. Показываются идеи методов, выводятся расчетные формулы, даются конкретные алгоритмы, позволяющие в оговоренных ситуациях решать частичные и полные проблемы собственных значений до конца.

## 4.1. СОБСТВЕННЫЕ ПАРЫ МАТРИЦЫ И ИХ ПРОСТЕЙШИЕ СВОЙСТВА

Пусть  $A$  – вещественная  $n \times n$ -матрица,  $y = y(t)$  –  $n$ -мерная векторная функция скалярного аргумента  $t$ , и пусть ищутся нетривиальные решения системы дифференциальных уравнений

$$\frac{dy}{dt} = Ay \quad (4.1)$$

в виде  $y = e^{\lambda t} x$ , где  $x \in C_n$ ,  $\lambda \in C$ . Подставляя  $y$  и  $\frac{dy}{dt}$  в (4.1), получаем

$$\lambda e^{\lambda t} x = A e^{\lambda t} x,$$

т.е. система (4.1) действительно будет иметь решения заданного вида в том и только том случае, если найдутся такие пары чисел  $\lambda$  и ненулевых векторов  $x$ , что

$$Ax = \lambda x. \quad (4.2)$$

Имеется ряд других примеров из областей, лежащих за пределами линейной алгебры, в которых также приходят к необходимости решать подобные (4.2) алгебраические задачи, называемые *задачами на собственные значения* (см., например, [30]). При этом различают *полную (алгебраическую или, иначе, матричную) проблему собственных значений*, предполагающую нахождение всех *собственных пар*  $\{\lambda, x\}$  матрицы  $A$ , и *частичные проблемы собственных значений*, состоящие, как правило, в нахождении одного или нескольких *собственных чисел*  $\lambda$  и, возможно,

соответствующих им *собственных векторов*)  $x$ . Чаще всего, в последнем случае речь идет о нахождении наибольшего и наименьшего по модулю собственных чисел; знание таких характеристик матрицы позволяет, например, делать заключения о сходимости тех или иных итерационных методов, оптимизировать параметры итерационных методов, учитывать влияние на результаты решения алгебраических задач погрешностей исходных данных и вычислительных погрешностей (потребность в таких числах неоднократно возникала в гл.3). Имеются и несколько иные постановки частичных проблем [5, 6, 17, 30].

Трактуя  $A$  в равенстве (4.2) как матрицу линейного преобразования в  $R_n$ , задачу на собственные значения можно сформулировать так: для каких ненулевых векторов  $x$  и чисел  $\lambda$  линейное преобразование вектора с помощью матрицы  $A$  не изменяет направления этого вектора в  $R_n$ , т.е. сводится к "растяжению" этого вектора в  $\lambda$  раз? Эта задача, очевидно, эквивалентна задаче исследования однородной СЛАУ<sup>\*)</sup> с параметром: при каких  $\lambda$  система

$$(A - \lambda E)x = 0 \quad (4.3)$$

имеет нетривиальные решения? Найти эти решения.

Теоретически эта задача легко решается: нужно найти корни так называемого *характеристического* или иначе "*векового*" уравнения

$$\det(A - \lambda E) = 0 \quad (4.4)$$

и, подставляя их поочередно в (4.3), получать из соответствующих переопределенных систем собственные векторы. Практическая реализация этого в сущности простого подхода сопряжена с рядом трудностей, возрастающих с ростом размерности решаемой задачи. Трудности эти обусловлены разрыванием "*векового*" определителя  $\det(A - \lambda E)$  и вычислением корней получающегося при этом многочлена  $n$ -й степени, а также поиском линейно независимых решений вырожденных СЛАУ. В связи с этим такой непосредственный подход к решению алгебраической проблемы собственных значений обычно применяют лишь при очень малых размерностях матриц  $A$  ( $n=2,3$ ); уже при  $n \geq 4$  на первый план выходят специальные численные методы решения таких задач. Ниже будут рассмотрены некоторые из этих методов так, чтобы можно было понять идеи, лежащие в их основе, и в то же время получить возможность решать поставленные задачи до конца для некоторых классов матриц (более полное и глубокое изложение этой темы см. в монографиях [26, 44, 53, 54] и в учебных пособиях [1, 5, 6, 8, 15, 21, 28, 34]).

<sup>\*)</sup> Собственное число и собственное значение в данном контексте – синонимы. В более общем случае, когда  $A$  в (4.2) – некоторый оператор, *собственные элементы*  $x$  могут иметь другую природу.

<sup>\*\*)</sup> Здесь, как и ранее, аббревиатура СЛАУ означает "система линейных алгебраических уравнений", а обозначение  $E$  зарезервировано за единичной матрицей.

Следует заметить, что если в недалеком прошлом численные методы решения задач на собственные значения опирались, как правило, на классический подход, т.е. на развертывание вековых определителей, в частности, в простейшем случае с помощью приведения матрицы  $A$  подходящим преобразованием к так называемой *сопровождающей матрице* [26, 53]

$$C = \begin{pmatrix} c_1 & c_2 & c_3 & \dots & c_{n-1} & c_n \\ 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 1 & \dots & 1 & 0 \end{pmatrix},$$

где в первой строке стоят коэффициенты уравнения (4.4), записанного в виде

$$(-1)^n (\lambda^n - c_1 \lambda^{n-1} - c_2 \lambda^{n-2} - \dots - c_{n-1} \lambda - c_n) = 0,$$

то современные методы решения полной проблемы ориентированы на алгоритмическое построение из матрицы  $A$  такой матрицы, определенные элементы которой являлись бы приближенными значениями собственных чисел  $A$ , причем параллельно формировались бы и ее собственные векторы.

Прежде чем приступить к изучению методов нахождения собственных чисел и векторов, вспомним некоторые простые их свойства, требующиеся в дальнейшем.

**Свойство 1.** Если  $\{\lambda, x\}$  – собственная пара матрицы  $A$ , а  $\alpha (\neq 0)$  – некоторое число, то  $\{\lambda, \alpha x\}$  также является собственной парой для  $A$ .

Действительно, умножив верное для данных  $\lambda$  и  $x$  равенство (4.2) на число  $\alpha$ , получаем верное равенство

$$A(\alpha x) = \lambda(\alpha x).$$

Оно означает, что каждому собственному числу  $\lambda$  соответствует бесчисленное множество собственных векторов, различающихся лишь скалярным множителем. Такие векторы задают одно направление в  $n$ -мерном пространстве; в соответствие этому направлению можно поставить нормированный вектор или орт. (Вообще говоря, одному собственному числу может соответствовать и несколько линейно независимых собственных векторов.)

**Свойство 2.** Пусть  $\{\mu, x\}$  – собственная пара матрицы  $A - pE$  при некотором  $p \in R$ . Тогда  $\{\lambda := \mu + p, x\}$  собственная пара матрицы  $A$ .

Чтобы убедиться в этом, заметим, что по условию

$$(A - pE)x = \mu x \quad (4.5)$$

при данных  $\mu$  и  $x$  – верное равенство. Рассмотрим равенство  $Ax = \lambda x$  при  $\lambda = \mu + p$ :

$$Ax = (\mu + p)x.$$

Оно равносильно (4.5), и значит, справедливо, с другой стороны, говорит о том, что  $\{\lambda, x\}$  – собственная пара  $A$ .

Как видим, прибавление к данной матрице  $A$  скалярной матрицы  $pE$  не изменяет ее собственных векторов и смещает спектр<sup>\*)</sup> исходной матрицы на число  $p$  (влево при  $p > 0$ ).

**Свойство 3.** Если  $\{\lambda, x\}$  – собственная пара обратимой матрицы  $A$ , то  $\left\{\frac{1}{\lambda}, x\right\}$  – собственная пара матрицы  $A^{-1}$ .

Справедливость этого свойства очевидна: умножив верное для данных  $\lambda$  и  $x$  равенство  $Ax = \lambda x$  слева на матрицу  $\frac{1}{\lambda}A^{-1}$ , получаем

$$\frac{1}{\lambda}x = A^{-1}x,$$

что и означает утверждаемое.

**Свойство 4.** Собственными числами диагональных и треугольных матриц являются их диагональные элементы.

Этот факт легко усматривается из очевидного представления характеристических уравнений (4.4) для таких матриц в виде

$$\prod_{i=1}^n (\lambda - a_{ii}) = 0.$$

Последнее равенство свидетельствует, что диагональные и треугольные вещественные матрицы имеют только вещественные собственные значения (ровно  $n$  с учетом возможной их кратности). Вещественность собственных чисел присуща и очень важному в приложениях классу симметричных матриц [1, 16].

---

<sup>\*)</sup> Напомним, что спектром матрицы называется множество всех ее собственных значений.

**Определение 4.1.** *Отношением Рэлея\** для  $n \times n$ -матрицы  $A$  называется функционал  $\rho(x) = \frac{(Ax, x)}{(x, x)}$ , определенный на множестве ненулевых  $n$ -мерных векторов  $x$ .

**Свойство 5.** Пусть  $x^*$  – собственный вектор матрицы  $A$ , тогда  $\rho(x^*)$  – ее собственное число.

Для доказательства этого утверждения обозначим через  $\lambda^*$  собственное число матрицы  $A$ , соответствующее вектору  $x^*$ . Подставляя  $Ax^* = \lambda^* x^*$  в вытекающее из определения 4.1 равенство

$$(Ax^*, x^*) = \rho(x^*) (x^*, x^*),$$

имеем

$$\lambda^* (x^*, x^*) = \rho(x^*) (x^*, x^*),$$

откуда после деления на  $(x^*, x^*) \neq 0$  получаем утверждаемое:  $\lambda^* = \rho(x^*)$ .

Отношение Рэлея обладает рядом других ценных свойств. Например, если матрица  $A$  – симметричная положительно определенная со спектром  $\lambda_1 > \lambda_2 > \dots > \lambda_n$ , то  $\max \rho(x) = \lambda_1$ ,  $\min \rho(x) = \lambda_n$ ,  $\rho(x) \in [\lambda_n, \lambda_1]$  при любых  $n$ -мерных  $x \neq 0$  и, кроме того,  $\text{grad } \rho(x) = 0$  тогда и только тогда, когда  $x$  – собственный вектор матрицы  $A$  (см. [1, 44]). Эти свойства служат основой для некоторых способов локализации собственных значений и построения градиентных методов их вычисления.

В дальнейшем (п.4.2, п.4.3) будет полезно следующее экстремальное свойство отношения Рэлея.

**Свойство 6.** Минимум евклидовой нормы вектора  $\xi(\lambda) := Ax - \lambda x$  для любого фиксированного ненулевого вектора  $x$  достигается при  $\lambda = \rho(x)$ .

Смысл этого факта в следующем: если некоторый вектор  $x$  считать приближением к собственному вектору матрицы  $A$  (а значит,  $\xi$  – его невязкой), то отношение Рэлея  $\rho(x)$  будет наилучшим приближением

---

\*) Лорд Рэлей (до получения титула лорда – Стретт Джон Уильям, 1842–1919) – английский физик, один из основоположников теории колебаний.

к соответствующему этому вектору собственному числу в смысле евклидовой метрики.

Доказательство свойства 4.6 для симметричных  $A$  и вещественных  $x$  весьма просто. Действительно, рассмотрим квадрат евклидовой нормы невязки

$$\begin{aligned}\|\xi(\lambda)\|_2^2 &= (Ax - \lambda x, Ax - \lambda x) = \\ &= (Ax, Ax) - 2\lambda(Ax, x) + \lambda^2(x, x) = q(\lambda)(x, x),\end{aligned}$$

где  $q(\lambda) = \lambda^2 - 2\lambda\rho(x) + \frac{(Ax, Ax)}{(x, x)}$ . Очевидно, квадратный трехчлен  $q(\lambda)$  всегда имеет минимум при  $\lambda = \rho(x)$ , а поскольку  $(x, x) > 0$ , это значение  $\lambda$  доставляет минимум величине  $\|\xi(\lambda)\|_2^2$  и, следовательно, величине  $\|\xi(\lambda)\|_2$ .

## 4.2. СТЕПЕННОЙ МЕТОД

Рассмотрим простейший метод решения частичных проблем собственных значений, который вряд ли может быть отнесен к широко применяемым методам решения таких задач, но много значащий для понимания и построения других, более эффективных методов.

Пусть о вещественной  $n \times n$ -матрице  $A$  известно, что она является *матрицей простой структуры*, т.е. имеет ровно  $n$  линейно независимых собственных векторов (базис):

$$x_1 = \begin{pmatrix} x_{11} \\ x_{12} \\ \dots \\ x_{1n} \end{pmatrix}, x_2 = \begin{pmatrix} x_{21} \\ x_{22} \\ \dots \\ x_{2n} \end{pmatrix}, \dots, x_n = \begin{pmatrix} x_{n1} \\ x_{n2} \\ \dots \\ x_{nn} \end{pmatrix}. \quad (4.6)$$

Пусть нумерация этих векторов отвечает упорядочению соответствующих им собственных чисел по убыванию модулей (где первое из неравенств — строгое):

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|. \quad (4.7)$$

Ставим задачу приближенного вычисления наибольшего по модулю собственного числа  $\lambda_1$  (вещественного, в силу предположения о строгом доминировании его модуля) и соответствующего ему собственного вектора  $x_1$  данной матрицы  $A$ .

Возьмем произвольный ненулевой вектор  $y^{(0)}$  и запишем его разложение по базису из собственных векторов  $x_1, x_2, \dots, x_n$ :

$$y^{(0)} = c_1 x_1 + c_2 x_2 + \dots + c_n x_n. \quad (4.8)$$

При этом без ограничения общности можно считать, что  $c_1 \neq 0$ , так как в противном (маловероятном) случае можно взять другой начальный вектор  $y^{(0)}$ .

Выполним первую итерацию вектора  $y^{(0)}$  умножением (4.8) слева на матрицу  $A$ :

$$y^{(1)} = Ay^{(0)} = c_1 Ax_1 + c_2 Ax_2 + \dots + c_n Ax_n.$$

Так как  $\{\lambda_i, x_i\}$  при всех  $i \in \{1, 2, \dots, n\}$  по предположению являются собственными парами матрицы  $A$ , то, в силу (4.2), последнее можно переписать в виде

$$y^{(1)} = c_1 \lambda_1 x_1 + c_2 \lambda_2 x_2 + \dots + c_n \lambda_n x_n.$$

Для второй итерации по тому же принципу получаем

$$\begin{aligned} y^{(2)} &= Ay^{(1)} = A^2 y^{(0)} = \\ &= c_1 \lambda_1^2 x_1 + c_2 \lambda_2^2 x_2 + \dots + c_n \lambda_n^2 x_n = \\ &= c_1 \lambda_1^2 x_1 + c_2 \lambda_2^2 x_2 + \dots + c_n \lambda_n^2 x_n. \end{aligned}$$

Очевидно,  $k$ -я итерация вектора  $y^{(0)}$  с помощью матрицы  $A$  дает вектор

$$y^{(k)} = Ay^{(k-1)} = A^k y^{(0)} = c_1 \lambda_1^k x_1 + c_2 \lambda_2^k x_2 + \dots + c_n \lambda_n^k x_n \quad (4.9)$$

или, с учетом представления  $x_1, x_2, \dots, x_n$  в исходном базисе (см. (4.6)),

$$y^{(k)} = \begin{pmatrix} y_1^{(k)} \\ y_2^{(k)} \\ \dots \\ y_n^{(k)} \end{pmatrix} = c_1 \lambda_1^k \begin{pmatrix} x_{11} \\ x_{12} \\ \dots \\ x_{1n} \end{pmatrix} + c_2 \lambda_2^k \begin{pmatrix} x_{21} \\ x_{22} \\ \dots \\ x_{2n} \end{pmatrix} + \dots + c_n \lambda_n^k \begin{pmatrix} x_{n1} \\ x_{n2} \\ \dots \\ x_{nn} \end{pmatrix}.$$

Беря отношения компонент *итерированного вектора*<sup>\*)</sup>  $y^{(k)}$  к соответствующим компонентам предыдущего вектора  $y^{(k-1)}$ , будем иметь:

$$\begin{aligned} \frac{y_i^{(k)}}{y_i^{(k-1)}} &= \frac{c_1 \lambda_1^k x_{1i} + c_2 \lambda_2^k x_{2i} + \dots + c_n \lambda_n^k x_{ni}}{c_1 \lambda_1^{k-1} x_{1i} + c_2 \lambda_2^{k-1} x_{2i} + \dots + c_n \lambda_n^{k-1} x_{ni}} = \\ &= \lambda_1 \cdot \frac{1 + \frac{c_2}{c_1} \cdot \frac{x_{2i}}{x_{1i}} \left(\frac{\lambda_2}{\lambda_1}\right)^k + \dots + \frac{c_n}{c_1} \cdot \frac{x_{ni}}{x_{1i}} \left(\frac{\lambda_n}{\lambda_1}\right)^k}{1 + \frac{c_2}{c_1} \cdot \frac{x_{2i}}{x_{1i}} \left(\frac{\lambda_2}{\lambda_1}\right)^{k-1} + \dots + \frac{c_n}{c_1} \cdot \frac{x_{ni}}{x_{1i}} \left(\frac{\lambda_n}{\lambda_1}\right)^{k-1}}. \end{aligned} \quad (4.10)$$

\*) Термин взят из [31].

Предел дроби в последнем равенстве при сделанных допущениях равен 1 в процессе  $k \rightarrow \infty$ , и значит,  $y_i^{(k)}/y_i^{(k-1)} \xrightarrow[k \rightarrow \infty]{} \lambda_i$  для каждого  $i \in \{1, 2, \dots, n\}$ , при котором  $x_{1i} \neq 0$  (заметим, что числа  $x_{11}, x_{12}, \dots, x_{1n}$  не могут быть одновременно нулями, так как  $x_1$  – базисный вектор и поэтому не может быть нулевым).

Представляя вектор  $y^{(k)}$  на основе (4.9) в виде

$$y^{(k)} = c_1 \lambda_1^k \left[ x_1 + \frac{c_2}{c_1} \left( \frac{\lambda_2}{\lambda_1} \right)^k x_2 + \dots + \frac{c_n}{c_1} \left( \frac{\lambda_n}{\lambda_1} \right)^k x_n \right], \quad (4.11)$$

можно сделать вывод, что при тех же исходных допущениях, в силу  $\left| \frac{\lambda_i}{\lambda_1} \right|^k \xrightarrow[k \rightarrow \infty]{(i \neq 1)} 0$ , в фигурирующей в скобках выражения (4.11) линейной комбинации векторов  $x_1, x_2, \dots, x_n$  с ростом  $k$  начнет доминировать первое слагаемое. Это означает, что вектор  $y^{(k)}$  от итерации к итерации будет давать все более хорошие приближения к собственному вектору  $x_1$  по направлению, т.е. с точностью до скалярного множителя  $c_1 \lambda_1^k$  (см. п.4.1, свойство 1).

Таким образом, как показывают приведенные рассуждения, метод нахождения "старшей" собственной пары матрицы простой структуры, называемый *степенным методом*<sup>\*)</sup>, в своей основе весьма примитивен и состоит в следующем: берется произвольный вектор  $y^{(0)} (\neq 0)$ , простыми итерациями  $y^{(k)} = Ay^{(k-1)}$  строится последовательность векторов  $y^{(k)}$  и параллельно рассматриваются последовательности отношений соответствующих компонент векторов  $k$ -й и  $(k-1)$ -й итераций (отношения с чрезвычайно малыми по модулю знаменателями следует игнорировать). Как только установятся несколько первых цифр во всех этих отношениях (что выясняется проверкой выполнения приближенных равенств

$$\frac{y_i^{(k)}}{y_i^{(k-1)}} \approx \frac{y_i^{(k-1)}}{y_i^{(k-2)}},$$

так можно считать, что найдено наибольшее по модулю

собственное число с точностью, определяемой последним установившимся в отношениях знаком, и соответствующий ему собственный вектор, за который принимается последний итерированный вектор  $y^{(k)}$ .

Для практической реализации такая схема нахождения старшей собственной пары мало пригодна по многим причинам и требует определен-

<sup>\*)</sup> Этимология данного термина совершенно ясна. Можно встретить и другие названия: *счет на установление* [28], *итерационный метод фон Мизеса* [32]. Иногда применяют латинскую аббревиатуру РМ (от англ. Power method) [44].



ной доводки. Рассмотрим некоторые из этих причин и соответственно пути модификации вышеописанного простейшего алгоритма.

Анализируя выражение  $y^{(k)}$  в форме (4.11), видим, что при достаточно большом числе итераций  $k$  за счет множителя  $\lambda_1^k$  в процессе счета может произойти либо превышение допустимых для ЭВМ чисел, если  $|\lambda_1| > 1$ , либо пропадание значащих цифр итерированных векторов, если  $|\lambda_1| < 1$ . Устранить это явление можно достаточно легко, введя в итерационный процесс нормировку итерированных векторов (т.е. приведение к единичной длине по той или иной метрике) на каждой итерации или через некоторое фиксированное число итерационных шагов.

Так, пошаговая нормировка векторов порождает следующий

#### **PM-алгоритм.**

Шаг 1. Ввести  $n \times n$ -матрицу  $A$ , задать  $n$ -мерный вектор  $y^{(0)}$ , вычислить  $\|y^{(0)}\|$  и вектор  $x^{(0)} := y^{(0)} / \|y^{(0)}\|$ ; положить  $k=1$ .

Шаг 2. Вычислить вектор  $y^{(k)} = Ax^{(k-1)}$ .

Шаг 3. Вычислить  $\|y^{(k)}\|$  и  $x^{(k)} := y^{(k)} / \|y^{(k)}\|$ .

Шаг 4. Вычислить отношения  $\lambda_i^{(k)} = y_i^{(k)} / x_i^{(k-1)}$  (координат векторов  $y^{(k)}$  и  $x^{(k-1)}$ ) при  $i=1, 2, \dots, n$  таких, что  $|x_i^{(k-1)}| > \delta$ , где  $\delta > 0$  — некоторое задаваемое малое число (допуск).

Шаг 5. Подвергнуть числа  $\lambda_i^{(k)}$  тесту на сходимость.

Если обнаруживается совпадение требуемого числа знаков в  $\lambda_i^{(k)}$  и  $\lambda_i^{(k-1)}$  ( $\lambda^{(0)}$  можно задавать произвольно), то работу алгоритма прекратить и за старшее собственное число  $\lambda_i$  принять усредненное (по  $i$ ) значение  $\lambda_i^{(k)}$ , а за нормированный старший собственный вектор  $x_1$  — вектор  $x^{(k)}$ .

В противном случае — вернуться к шагу 2.

Слабым местом данного алгоритма, очевидно, является последний шаг, т.е. решение проблемы своевременного останова работы алгоритма. Этот шаг описан из рациональных соображений и не может гарантировать во всех случаях (даже при сделанных допущениях) получения собственной пары  $\{\lambda_1, x_1\}$  с наперед заданной точностью, поскольку при разработке метода не было получено никаких оценок погрешности.

Относительно характера сходимости степенного метода можно утверждать (см. формулы (4.10) и (4.11)), что в указанных условиях итераци-

онный процесс является линейным<sup>\*)</sup>, т.е. сходится со скоростью геометрической прогрессии, знаменатель которой определяется в основном величиной отношения  $\left| \frac{\lambda_2}{\lambda_1} \right|$ . Это означает, что сходимость будет тем лучше и, как

следствие, критерий останова в шаге 5 тем надежнее, чем сильнее доминирует в спектре матрицы  $A$  собственное число  $\lambda_1$ . Подмеченный факт вкупе со свойством 2 п.4.1 позволяет существенно ускорить нахождение наибольшего по модулю собственного числа матрицы  $A$  путем удачного смещения ее спектра, чему могут способствовать какие-либо априорные сведения об исходной задаче<sup>\*\*)</sup>.

То же свойство 2 собственных пар позволяет применять непосредственно степенной метод для нахождения наименьшего по модулю собственного числа  $\lambda_n$  знакоопределенной матрицы  $A$  в случае, когда наибольшее  $\lambda_1$  уже найдено. Для этого достаточно найти наибольшее по модулю собственное число  $\Lambda$  матрицы  $A - \lambda_1 E$ ; соответствующий ему собственный вектор этой матрицы и число  $\lambda_n = \Lambda + \lambda_1$  будут образовывать искомую собственную пару.

Действительно, вычитая из верного для собственной пары  $\{\lambda_n, x_n\}$  равенства  $Ax_n = \lambda_n x_n$  тождество  $\lambda_1 x_n = \lambda_1 x_n$ , получаем верное равенство

$$(A - \lambda_1 E)x_n = (\lambda_n - \lambda_1)x_n,$$

означающее, что  $\Lambda = \lambda_n - \lambda_1$  и  $x_n$  служат собственной парой матрицы  $A - \lambda_1 E$ . Так как для знакоопределенной матрицы справедливо неравенство  $|\lambda_n - \lambda_i| \geq |\lambda_i - \lambda_1|$  при любом  $i \in \{1, 2, \dots, n\}$ , то  $\Lambda$  — наибольшее по модулю собственное число матрицы  $A - \lambda_1 E$  и может быть найдено степенным методом.

Знание старшего собственного числа  $\lambda_1$  матрицы  $A$  простой структуры, получаемого в процессе прямых итераций по формулам (4.9), (4.10), в предположении, что

$$|\lambda_1| > |\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|,$$

<sup>\*)</sup> См. определение 5.1 в п.5.3.

<sup>\*\*)</sup> См. по этому поводу [53, 54]. В [53] приведен пример, наглядно показывающий эффективность подходящего сдвига: если некая матрица  $A$  шестого порядка имеет собственные числа  $\lambda_i = 21 - i$ , то непосредственное применение степенного метода к вычислению  $\lambda_1$  порождает итерационный процесс, сходящийся со скоростью порядка

$\left(\frac{19}{20}\right)^k$ , в то время как сдвиг на величину  $p = 17$ , оставляющий соответствующее  $\lambda_1$  число  $\mu_1 = \lambda_1 - p = 3$  старшим в спектре матрицы  $A - pE$ , позволяет найти его степенным методом, сходящимся уже со скоростью порядка  $\left(\frac{2}{3}\right)^k$ .

позволяет без больших дополнительных затрат найти приближенное значение второго по модулю собственного числа  $\lambda_2$ . Это можно сделать по формуле

$$\lambda_2 \approx \frac{y_i^{(k+1)} - \lambda_1 y_i^{(k)}}{y_i^{(k)} - \lambda_1 y_i^{(k-1)}} \quad (4.12)$$

вычисляя фигурирующие в правой части отношения для достаточно больших  $k$  и для всех  $i=1,2,\dots,n$ , при которых абсолютная величина знаменателя не меньше некоторого порогового значения, и затем усредняя результат. Понятно, что при этом неизбежна потеря точности.

Для обоснования<sup>\*)</sup> приближенного равенства (4.12) подставим в его правую часть выражения компонент  $(k+1)$ -го,  $k$ -го и  $(k-1)$ -го итерированных векторов в соответствии с представлением (4.9) в исходном базисе. После взаимного уничтожения по паре первых членов в числителе и в знаменателе будем иметь:

$$\begin{aligned} \frac{y_i^{(k+1)} - \lambda_1 y_i^{(k)}}{y_i^{(k)} - \lambda_1 y_i^{(k-1)}} &= \frac{c_2 \lambda_2^{k+1} x_{2i} - c_2 \lambda_1 \lambda_2^k x_{2i} + \dots + c_n \lambda_n^{k+1} x_{ni} - c_n \lambda_1 \lambda_n^k x_{ni}}{c_2 \lambda_2^k x_{2i} - c_2 \lambda_1 \lambda_2^{k-1} x_{2i} + \dots + c_n \lambda_n^k x_{ni} - c_n \lambda_1 \lambda_n^{k-1} x_{ni}} = \\ &= \frac{c_2 \lambda_2^{k+1} x_{2i} \left( 1 - \frac{\lambda_1}{\lambda_2} + \sum_{j=3}^n \frac{\lambda_j^{k+1} - \lambda_1 \lambda_j^k}{\lambda_2^{k+1}} \cdot \frac{c_j}{c_2} \cdot \frac{x_{ji}}{x_{2i}} \right)}{c_2 \lambda_2^k x_{2i} \left( 1 - \frac{\lambda_1}{\lambda_2} + \sum_{j=3}^n \frac{\lambda_j^k - \lambda_1 \lambda_j^{k-1}}{\lambda_2^k} \cdot \frac{c_j}{c_2} \cdot \frac{x_{ji}}{x_{2i}} \right)} \xrightarrow{k \rightarrow \infty} \lambda_2. \end{aligned}$$

так как  $\frac{\lambda_j^{k+1} - \lambda_1 \lambda_j^k}{\lambda_2^{k+1}} = \left( \frac{\lambda_j}{\lambda_2} \right)^{k+1} - \frac{\lambda_1}{\lambda_2} \cdot \left( \frac{\lambda_j}{\lambda_2} \right)^k \xrightarrow{k \rightarrow \infty} 0$  при всех  $j=3,\dots,n$ .

Вернемся к вопросу о недостатках степенного метода нахождения наибольшего по модулю собственного числа и путях их устранения. При этом далее ограничимся рассмотрением класса симметричных положительно определенных матриц. Известно, что такие матрицы имеют положительный вещественный спектр  $\lambda_1, \lambda_2, \dots, \lambda_n$ , ортонормированный базис из собственных векторов  $x_1, x_2, \dots, x_n$  и, естественно, являются матрицами простой структуры.

Обсудим шаг 4 предложенного выше РМ-алгоритма.

Вычисление на каждом итерационном шаге отношений в с е х пар соответствующих компонент векторов  $x$  и  $y=Ax$ , да еще с определенными проверками, при больших размерностях  $n$  требует значительных вычисли-

<sup>\*)</sup> Другое обоснование см. например, в [21]. Там же показано, что за соответствующий  $\lambda_2$  собственный вектор  $x_2$  можно принять нормированный вектор  $y^{(k+1)} - \lambda_1 y^{(k)}$ .

тельных затрат, хотя и дает о старшем собственном числе  $\lambda_1$  дополнительную информацию: как утверждается в [32], значение  $\lambda_1$  заключено между наименьшим и наибольшим из этих отношений, т.е. имеются двусторонние оценки  $\lambda_1$  на каждой итерации.

Чтобы упростить соответствующую шагу 4 РМ-алгоритма процедуру, проведем следующие рассуждения.

Пусть  $R_n$  – евклидово пространство,  $A$  – симметричная положительно определенная матрица и последовательность итерированных векторов  $y^{(k)}$  строится, как и ранее, по формулам (4.9).

Рассмотрим скалярные произведения  $(y^{(k)}, y^{(k)})$  и  $(y^{(k)}, y^{(k-1)})$ . Выполнив умножение правых частей (4.9) по правилам умножения многочленов и учитывая ортонормированность собственных векторов, т.е. условие  $(x_i, x_j) = \delta_{ij}$  при  $i, j \in \{1, 2, \dots, n\}$ , имеем:

$$\begin{aligned} (y^{(k)}, y^{(k)}) &= c_1^2 \lambda_1^{2k} + c_2^2 \lambda_2^{2k} + \dots + c_n^2 \lambda_n^{2k}, \\ (y^{(k)}, y^{(k-1)}) &= c_1^2 \lambda_1^{2k-1} + c_2^2 \lambda_2^{2k-1} + \dots + c_n^2 \lambda_n^{2k-1}. \end{aligned}$$

Отношение этих чисел

$$\frac{(y^{(k)}, y^{(k)})}{(y^{(k)}, y^{(k-1)})} = \lambda_1 \cdot \frac{1 + \left(\frac{c_2}{c_1}\right)^2 \left(\frac{\lambda_2}{\lambda_1}\right)^{2k} + \dots + \left(\frac{c_n}{c_1}\right)^2 \left(\frac{\lambda_n}{\lambda_1}\right)^{2k}}{1 + \left(\frac{c_2}{c_1}\right)^2 \left(\frac{\lambda_2}{\lambda_1}\right)^{2k-1} + \dots + \left(\frac{c_n}{c_1}\right)^2 \left(\frac{\lambda_n}{\lambda_1}\right)^{2k-1}} \quad (4.13)$$

в оговоренных выше условиях при  $k \rightarrow \infty$  имеет пределом наибольшее собственное число  $\lambda_1$ , причем скорость сходимости к пределу будет больше, чем в степенном методе, опирающемся на отношения (4.10)

$$\left( O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^{2k}\right) \text{ против } O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right) \right).$$

Базирующаяся на таком подходе модификация степенного метода называется *методом скалярных произведений*. Реализовать ее можно, например, в виде следующего *SP-алгоритма*<sup>\*)</sup>.

1. Ввести: данную симметричную  $n \times n$ -матрицу  $A$ , произвольный  $n$ -мерный начальный вектор  $y^{(0)} (\neq 0)$ , малое число  $\varepsilon > 0$  (определяющее допустимую абсолютную погрешность искомого собственного числа  $\lambda_1$ ), число  $\lambda^{(0)}$  для начального сравнения (например, 0). Положить  $k=1$  (счетчик итераций).

<sup>\*)</sup> SP от англ. Scalar product

2. Вычислить скаляры  $s^{(0)} = (y^{(0)}, y^{(0)})$ ,  $\|y^{(0)}\| = \sqrt{s^{(0)}}$  и вектор  $x^{(0)} = y^{(0)} / \|y^{(0)}\|$ .
3. Вычислить  $y^{(k)} = Ax^{(k-1)}$  (итерация нормированного вектора).
4. Вычислить:  $s^{(k)} = (y^{(k)}, y^{(k)})$  и  $t^{(k)} = (y^{(k)}, x^{(k-1)})$  (скалярные произведения),  $\|y^{(k)}\| = \sqrt{s^{(k)}}$ ,  $x^{(k)} = y^{(k)} / \|y^{(k)}\|$  (приближение к нормированному собственному вектору),  $\lambda^{(k)} = s^{(k)} / t^{(k)}$  (приближение к собственному числу  $\lambda_1$ ).
5. Если  $|\lambda^{(k)} - \lambda^{(k-1)}| > \varepsilon$ , положить  $k := k+1$  и вернуться к шагу 3, иначе завершить работу алгоритма, считая  $\lambda_1 \approx \lambda^{(k)}$ ,  $x_1 \approx x^{(k)}$ .

**Замечание 4.1.** Данный алгоритм позволяет более быстро (т.е. за меньшее число итераций), чем РМ-алгоритм, найти с нужной точностью наибольшее собственное число симметричной матрицы, но при этом точность приближенного равенства  $x_1 \approx x^{(k)}$  для соответствующего собственного вектора может оказаться недостаточной (почему?).

**Замечание 4.2.** Очевидно, в методе скалярных произведений вместо отношения (4.13), стремящегося к  $\lambda_1$  при  $k \rightarrow \infty$ , можно с тем же успехом взять отношение  $(y^{(k+1)}, y^{(k)})$  к  $(y^{(k)}, y^{(k)})$ , имеющее тот же предел. Последнее же есть не что иное, как отношение Рэля:

$$\frac{(y^{(k+1)}, y^{(k)})}{(y^{(k)}, y^{(k)})} = \frac{(Ay^{(k)}, y^{(k)})}{(y^{(k)}, y^{(k)})} = \rho(y^{(k)});$$

отсюда другое название метода скалярных произведений – *метод частных Рэля*. В соответствии со свойствами 5, 6 предыдущего пункта можно сказать, что этим методом на каждом итерационном шаге ищется наилучшее для вычисленного вектора  $y^{(k)}$  приближение к собственному числу  $\lambda_1$  в смысле евклидовой нормы невязки.

Наличие ортонормированного базиса из собственных векторов  $x_1, x_2, \dots, x_n$  матрицы  $A$  позволяет применять степенной метод (метод скалярных произведений) для последовательного вычисления собственных пар  $\{\lambda_i, x_i\}$  при  $i \geq 2$  более совершенными, чем определяемый формулой (4.12), способами. Рассмотрим один из них.

Пусть первая (старшая) собственная пара  $\{\lambda_1, x_1\}$  уже найдена, причем  $\|x_1\| = \sqrt{(x_1, x_1)} = 1$ . Возьмем произвольный ненулевой вектор  $z^{(0)}$  и образуем вектор

$$y^{(0)} = z^{(0)} - (z^{(0)}, x_1)x_1. \quad (4.14)$$

Так как

$$(y^{(0)}, x_1) = (z^{(0)}, x_1) - (z^{(0)}, x_1)(x_1, x_1) = 0,$$

то вектор  $y^{(0)}$  ортогонален  $x_1$ , т.е. его проекция на первый базисный вектор системы  $x_1, x_2, \dots, x_n$  равна нулю. Значит, разложение (4.8) вектора  $y^{(0)}$  по этому базису имеет вид

$$y^{(0)} = c_2 x_2 + c_3 x_3 + \dots + c_n x_n$$

и, соответственно, степенные итерации этого вектора типа (4.9) порождают векторы

$$y^{(k)} = c_2 \lambda_2^k x_2 + c_3 \lambda_3^k x_3 + \dots + c_n \lambda_n^k x_n. \quad (4.15)$$

Легко видеть (ср. с (4.13)), что если  $|\lambda_2| > |\lambda_i|$  при всех  $i \in \{3, \dots, n\}$ , то

$$\frac{(y^{(k)}, y^{(k)})}{(y^{(k)}, y^{(k-1)})} \xrightarrow{k \rightarrow \infty} \lambda_2 \quad \text{со скоростью} \quad O\left(\left|\frac{\lambda_3}{\lambda_2}\right|^{2k}\right) \quad \text{и}$$

$$x^{(k)} = \frac{y^{(k)}}{\|y^{(k)}\|} \xrightarrow{k \rightarrow \infty} x_2 \quad \text{со скоростью} \quad O\left(\left|\frac{\lambda_3}{\lambda_2}\right|^k\right).$$

Следующая собственная пара  $\{\lambda_3, x_3\}$  может быть найдена приближенно тем же методом, если за начальный вектор последовательности  $(y^{(k)})$  принять вектор

$$y^{(0)} = z^{(0)} - (z^{(0)}, x_1)x_1 - (z^{(0)}, x_2)x_2,$$

ортогональный одновременно  $x_1$  и  $x_2$  при любом  $z^{(0)}$ , и т.д.

Известны и другие способы последовательного нахождения собственных пар, опирающиеся на непосредственное применение степенного метода. При этом имеются возможности понижения размерности при нахождении каждой последующей собственной пары.

**Замечание 4.3.** В реальных расчетах, в силу неизбежных ошибок округлений, в представлении (4.15) итерированного вектора  $y^{(k)}$  при вычислении второй собственной пары появится малое, но растущее с увеличением номера  $k$  слагаемое, соответствующее проекции  $y^{(k)}$  на первый собст-

венный вектор  $x_1$ . Поэтому реальный алгоритм должен предусматривать возврат к началу процесса итерирования, т.е. проведение операции ортогонализации по формуле (4.14) с  $z^{(0)} := y^{(m)}$  через некоторое число итераций  $k=m$ .

**Замечание 4.4.** Не всегда бывает известным, выполняются ли оговоренные выше условия, при которых изучался степенной метод. В таких ситуациях при его применении нужно принимать особые меры предосторожности. Целесообразно, например, контролировать, сближаются ли члены последовательности  $(\lambda^{(k)})$  посредством проверки неравенств

$$|\lambda^{(k+1)} - \lambda^{(k)}| < |\lambda^{(k)} - \lambda^{(k-1)}|$$

(прием Гарвика [1, 28]), а также осуществлять итерационный процесс с разных начальных векторов (в случае кратности находимого собственного числа это просто необходимо для вычисления степенным методом всех соответствующих ему собственных векторов).

**Замечание 4.5.** Как отмечалось выше, степенной метод сходится линейно, точнее, имеет лишь асимптотическую скорость сходимости геометрической прогрессии. При слабом доминировании модуля вычисляемого собственного числа эта сходимость может оказаться чрезвычайно медленной. Ускорения итерационного процесса можно достигнуть за счет быстрого накопления степеней матриц по схеме

$$A \cdot A = A^2, \quad A^2 \cdot A^2 = A^4$$

и т.д., что позволяет производить не последовательное, пошаговое, а скачкообразное построение последовательности  $(y^{(k)})$  с помощью равенств вида

$$y^{(k)} = A^{k-m} y^{(m)}$$

при фиксированных  $m \in \{0, 1, \dots, k-1\}$  (таких, при которых  $k-m$  является некоторой целой степенью двойки) и нормированием сразу после очередного скачка. Здесь, правда, нужно особенно внимательно относиться к риску выхода за границы диапазона компьютерных чисел в процессе счета.

Если не требуется находить собственный вектор  $x_1$ , то более быстро вычислять максимальное по модулю собственное число  $\lambda_1$  можно на основе соотношения [21]

$$\lambda_1^k + \lambda_2^k + \dots + \lambda_n^k = \text{Sp } A^k \quad \forall k \in N.$$

Вычислив  $A^k$  по закону удвоения степеней, а затем  $A^{k+1} = A^k \cdot A$ , находим отношение следов (сумм диагональных элементов) этих матриц:

$$\frac{\text{Sp } \mathbf{A}^{k+1}}{\text{Sp } \mathbf{A}^k} = \frac{\lambda_1^{k+1} \left( 1 + \left( \frac{\lambda_2}{\lambda_1} \right)^{k+1} + \dots + \left( \frac{\lambda_n}{\lambda_1} \right)^{k+1} \right)}{\lambda_1^k \left( 1 + \left( \frac{\lambda_2}{\lambda_1} \right)^k + \dots + \left( \frac{\lambda_n}{\lambda_1} \right)^k \right)} \xrightarrow{k \rightarrow \infty} \lambda_1.$$

**Замечание 4.6.** Более популярный способ улучшения сходимости степенного метода – это применение  $\Delta^2$ -процесса Эйткена<sup>\*)</sup>. Считается, что если  $\lambda^{(k-1)}$ ,  $\lambda^{(k)}$ ,  $\lambda^{(k+1)}$  являются тремя последовательными приближениями к собственному числу, полученными степенным методом, то число

$$\tilde{\lambda} := \lambda^{(k-1)} - \frac{(\lambda^{(k)} - \lambda^{(k-1)})^2}{\lambda^{(k+1)} - 2\lambda^{(k)} + \lambda^{(k-1)}}$$

ближе к пределу этой последовательности, чем каждое из них. Этот факт может быть использован в реальных алгоритмах либо через несколько итерационных шагов (например, через два на третий), либо на завершающем этапе вычислений. Для искомого собственного вектора такое ускорение может производиться по координатно.

### 4.3. ОБРАТНЫЕ ИТЕРАЦИИ

В предыдущем пункте было показано, что при определенных условиях наименьшее по модулю собственное число  $\lambda_n$  может быть найдено степенным методом, когда уже известно наибольшее  $\lambda_1$ . Если же проблема состоит лишь в нахождении младшей собственной пары матрицы  $\mathbf{A}$ , то можно обойтись и без вычисления  $\lambda_1$ , применяя степенной метод к матрице  $\mathbf{A}^{-1}$ .

В самом деле, если данная матрица  $\mathbf{A}$  имеет собственные пары

$$\{\lambda_1, \mathbf{x}_1\}, \{\lambda_2, \mathbf{x}_2\}, \dots, \{\lambda_{n-1}, \mathbf{x}_{n-1}\}, \{\lambda_n, \mathbf{x}_n\},$$

то по свойству 3 п.4.1 собственными парами матрицы  $\mathbf{A}^{-1}$  будут

$$\left\{ \frac{1}{\lambda_1}, \mathbf{x}_1 \right\}, \left\{ \frac{1}{\lambda_2}, \mathbf{x}_2 \right\}, \dots, \left\{ \frac{1}{\lambda_{n-1}}, \mathbf{x}_{n-1} \right\}, \left\{ \frac{1}{\lambda_n}, \mathbf{x}_n \right\}.$$

При этом упорядочиванию спектра  $\mathbf{A}$

$$|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_{n-1}| > |\lambda_n|$$

<sup>\*)</sup> Более подробно об этом процессе см. в п.5.8а).



соответствует цепочка неравенств

$$\left| \frac{1}{\lambda_n} \right| > \left| \frac{1}{\lambda_{n-1}} \right| \geq \dots \geq \left| \frac{1}{\lambda_2} \right| \geq \left| \frac{1}{\lambda_1} \right|$$

для собственных чисел  $\gamma_1 := \frac{1}{\lambda_n}$ ,  $\gamma_2 := \frac{1}{\lambda_{n-1}}$ , ...,  $\gamma_{n-1} := \frac{1}{\lambda_2}$ ,  $\gamma_n := \frac{1}{\lambda_1}$

матрицы  $A^{-1}$ . Это значит, что наименьшим по модулю собственным числом данной матрицы  $A$  является величина, обратная наибольшему по модулю собственному числу матрицы  $A^{-1}$ . Последнее же может быть получено прямыми итерациями произвольного начального вектора  $y^{(0)}$  посредством матрицы  $A^{-1}$  по аналогичной (4.9) формуле

$$y^{(k)} = A^{-1}y^{(k-1)}, \quad k=1, 2, \dots \quad (4.16)$$

При достаточно больших  $k \in N$  последовательность отношений одноименных координат векторов  $y^{(k)}$  и  $y^{(k-1)}$  должна давать приближенное значение  $\frac{1}{\lambda_n}$ , а вектор  $y^{(k)}$  (желательно его нормирование) можно принять за собственный вектор  $x_n$ .

Вместо прямых итераций (4.16), требующих предварительного обращения исходной матрицы  $A$ , обычно предпочитают строить ту же последовательность векторов  $(y^{(k)})$ , решая при  $k=1, 2, 3, \dots$  линейные системы

$$Ay^{(k)} = y^{(k-1)}. \quad (4.17)$$

Так как все эти системы имеют одну и ту же матрицу коэффициентов, то самая трудоемкая часть метода Гаусса для их решения – LU-факторизация матрицы  $A$  – может быть выполнена лишь один раз.

Построение последовательности векторов, приближающих собственный вектор  $x_n$  по неявной формуле (4.17), называют *обратными итерациями*, а процесс решения частных проблем собственных значений на этой основе – *методом обратных итераций*<sup>\*)</sup>.

Применение обратных итераций к нахождению младшей собственной пары матрицы  $A$  не требует написания специального алгоритма, достаточно лишь заменить один шаг в алгоритмах предыдущего пункта. А именно, наполнение шага 2 в РМ-алгоритме для матриц простой структуры и шага 3 в SP-алгоритме для симметричных матриц должно быть следующим:

$$\text{решить уравнение } Ay^{(k)} = y^{(k-1)}.$$

Полученный алгоритм называют *INVIT-алгоритмом*<sup>\*\*) [44]</sup>.

<sup>\*)</sup> Другое название *обратный степенной метод* [11].

<sup>\*\*)</sup> От англ. Inverse iteration.

Метод обратных итераций, а точнее, *обратные итерации со сдви-гами* часто применяют в тех случаях, когда нужно с большой точностью найти собственный вектор, отвечающий какому-либо собственному числу из спектра заданной матрицы при условии, что известно приближенное значение этого числа. При этом, очевидно, прямое решение однородной системы (4.3) заведомо неприменимо, так как подстановка в нее значения  $\lambda$ , хоть сколько-нибудь отличного от собственного, сделает систему однозначно разрешимой, т.е. допускающей только тривиальное решение.

Пусть для собственного числа  $\lambda_j$  матрицы простой структуры  $A$  известно его приближение  $\sigma$  такое, что

$$|\lambda_j - \sigma| < |\lambda_i - \sigma| \quad \forall i \neq j. \quad (4.18)$$

Начиная с вектора  $x^{(0)}$  такого, что  $\|x^{(0)}\| = 1$ , образуем последовательность нормированных векторов  $(x^{(k)})$  по формулам

$$(A - \sigma E)y^{(k)} = x^{(k-1)}; \quad (4.19)$$

$$x^{(k)} = y^{(k)} / \|y^{(k)}\|, \quad k = 1, 2, \dots \quad (4.20)$$

Изучим поведение этой последовательности, для чего запишем разложение векторов  $y^{(k)}$  и  $x^{(k-1)}$  по базису из собственных векторов  $x_1, x_2, \dots, x_n$  с некоторыми коэффициентами  $c_i^{(k)}$  и  $b_i^{(k-1)}$  соответственно:

$$y^{(k)} = c_1^{(k)} x_1 + c_2^{(k)} x_2 + \dots + c_n^{(k)} x_n, \quad (4.21)$$

$$x^{(k-1)} = b_1^{(k-1)} x_1 + b_2^{(k-1)} x_2 + \dots + b_n^{(k-1)} x_n.$$

Подставляя это в (4.19) и учитывая, что по определению

$$Ax_i = \lambda_i x_i \quad \forall i \in \{1, 2, \dots, n\},$$

имеем

$$\begin{aligned} (\lambda_1 - \sigma)c_1^{(k)} x_1 + (\lambda_2 - \sigma)c_2^{(k)} x_2 + \dots + (\lambda_n - \sigma)c_n^{(k)} x_n = \\ = b_1^{(k-1)} x_1 + b_2^{(k-1)} x_2 + \dots + b_n^{(k-1)} x_n, \end{aligned}$$

откуда, в силу единственности разложения вектора по базису, следует

$$(\lambda_i - \sigma)c_i^{(k)} = b_i^{(k-1)} \quad \forall i \in \{1, 2, \dots, n\}.$$

Анализируя получающиеся отсюда выражения

$$c_i^{(k)} = \frac{b_i^{(k-1)}}{\lambda_i - \sigma} \quad (4.22)$$

коэффициентов разложения вектора  $y^{(k)}$  по базису из собственных векторов, видим, что вследствие малости модуля знаменателя  $\lambda_j - \sigma$  по сравнению с другими знаменателями  $\lambda_i - \sigma$  (см. (4.18)), можно рассчитывать на преимущественное возрастание коэффициентов  $c_j^{(k)}$  именно при собственном векторе  $x_j$  с ростом  $k$ . Значит, чем сильнее неравенство в (4.18), тем сильнее (быстрее) будет доминировать составляющая собственного вектора  $x_j$  в представлении (4.21) вектора  $y^{(k)}$ , а значит, и вектора  $x^{(k)}$ , получаемого из  $y^{(k)}$  нормированием (4.20). Последнее же говорит о том, что каков бы ни был начальный вектор  $x^{(0)}$  ( $\neq 0$ ), быстрое доминирование  $c_j^{(k)}$  среди остальных коэффициентов  $c_i^{(k)}$  происходит еще и за счет числителей дробей (4.22).

Следует заметить, что обратные итерации со сдвигами (4.19), (4.20) позволяют не только найти собственный вектор  $x_j$ , но и служат основой для уточнения приближенного равенства  $\lambda_j \approx \sigma$ .

Действительно, формулы (4.19), (4.20) определяют не что иное, как метод обратных итераций для нахождения наименьшего по модулю собственного числа матрицы  $A - \sigma E$ , и, если  $\sigma$  существенно ближе к  $\lambda_j$ , чем к любому другому собственному числу  $\lambda_i$  матрицы  $A$ , то уточняющие  $\lambda_j$  значения, согласно свойству 2 п.4.1, можно получать при  $k=1,2,\dots$  по формуле

$$\lambda_j^{(k)} = \sigma + \left\langle \frac{x_i^{(k-1)}}{y_i^{(k)}} \right\rangle, \quad (4.23)$$

где  $x_i^{(k-1)}$  и  $y_i^{(k)}$  — координаты векторов  $x^{(k-1)}$  и  $y^{(k)}$  соответственно, а  $\langle \cdot \rangle$  — знак усреднения по всем  $i$ , при которых  $y_i^{(k)} \neq 0$  [28].

Как показывает практика вычислений, сходимость процесса обратных итераций со сдвигами характеризуется высокой скоростью (по сравнению с обычным степенным методом). Но еще более высокая скорость сходимости может быть получена введением переменных сдвигов, определяемых какой-нибудь последовательностью чисел  $\sigma_0, \sigma_1, \sigma_2, \dots$ , сходящейся к находимому собственному числу. Не вызывает сомнений целесообразность

---

<sup>\*)</sup> Лишь бы при его выборе не попасть на ортогональный  $x_j$  вектор: зачастую берут вектор  $x^{(0)}$  с равными координатами.

использования в роли таких чисел приближений  $\lambda_j^{(k)}$  к собственному числу  $\lambda_j$ , получаемых по формуле (4.23).

Таким образом, *обратные итерации с переменными сдвигами* можно определить совокупностью равенств

$$\begin{aligned} (A - \lambda_j^{(k-1)} E) y^{(k)} &= x^{(k-1)}, \\ x^{(k)} &= y^{(k)} / \|y^{(k)}\|, \\ \lambda_j^{(k)} &= \lambda_j^{(k-1)} + \left\langle \frac{x_i^{(k-1)}}{y_i^{(k)}} \right\rangle, \end{aligned} \quad (4.24)$$

где  $k=1,2,\dots$ , а число  $\lambda_j^{(0)} \approx \lambda_j$  и нормированный вектор  $x^{(0)}$  задаются.

Скорость сходимости такого процесса – квадратичная [16, 28], в то время как в случае постоянных сдвигов – лишь линейная, хотя и с малыми, как правило, знаменателями геометрической прогрессии. Зачастую бывает достаточно сделать 2-3 итерации, чтобы получить заданную собственную пару с реально возможной точностью. Нужно только видеть разницу в цене реализации обратных итераций с постоянными и с переменными сдвигами: в первом случае при каждом  $k$  решаются линейные системы с одной и той же матрицей коэффициентов (как это было и при обратных итерациях (4.17) без сдвигов), во втором случае на разных шагах приходится решать совершенно различные системы.

В методе (4.24) так же, как и в предыдущем, неясно, как подбирать начальный сдвиг  $\sigma = \lambda_j^{(0)}$ , за исключением случаев, когда решается частичная проблема заведомо в такой постановке, при которой требуется найти собственное число, ближайшее к заданному значению, и соответствующий ему собственный вектор.

Более определенной в этом смысле, к тому же более быстро сходящейся является следующая модификация метода (4.24) – *обратные итерации с отношениями Рэдея*, применяемые для решения симметричных задач на собственные значения.

Ее основу составляет **RQI-алгоритм**<sup>\*)</sup>:

0. Задать вектор  $x^{(0)}$  такой, что  $\|x^{(0)}\| = 1$ .

Для  $k=1,2,\dots$ :

1. Вычислить  $\rho_{k-1} = (Ax^{(k-1)}, x^{(k-1)}) / (x^{(k-1)}, x^{(k-1)})$ .

<sup>\*)</sup> RQI – Rayleigh quotient iteration (англ.).

2. Найти  $y^{(k)}$  из уравнения  $(A - \rho_{k-1}E)y^{(k)} = x^{(k-1)}$ .

3. Нормировать  $y^{(k)}$ , т.е. положить  $x^{(k)} = y^{(k)} / \|y^{(k)}\|$ .

4. Проверить  $\rho_k, x^{(k)}$  на сходимость.

После "штатного" останова работы алгоритма<sup>\*)</sup> при некотором  $k = k_0$  в качестве собственной для данной матрицы  $A$  объявляется пара  $\{\rho_{k_0}, x^{(k_0)}\}$ .

Сдвиги на отношения Рэля при наличии ортонормированного базиса из собственных векторов  $x_1, x_2, \dots, x_n$  обеспечивают асимптотически кубическую скорость сходимости *последовательности Рэля*  $x^{(0)}, x^{(1)}, x^{(2)}, \dots$  к некоторому из векторов этого базиса [16, 44]. К какому именно, зависит от выбора начального вектора этой последовательности; беря различные линейно независимые векторы  $x^{(0)}$ , можно получать разные собственные пары данной симметричной матрицы  $A$ . При этом, правда, без дополнительных условий (типа (4.18) применительно к  $\rho_0$  в роли  $\sigma$ ) нельзя гарантировать, что найденное как предел последовательности  $\rho_0, \rho_1, \rho_2, \dots$  собственное число будет ближайшим к числу  $\rho_0$ .

Чтобы если не обосновать, то хотя бы осмыслить RQI-алгоритм, нужно вспомнить свойство б п.4.1, согласно которому, при выбранном векторе  $x^{(0)}$  вычисленное на первом шаге при  $k=1$  отношение Рэля  $\rho_0 = (Ax^{(0)}, x^{(0)}) / (x^{(0)}, x^{(0)})$  можно считать некоторым приближением к собственному числу, связанному с заданным в  $R_n$  направлением  $x^{(0)}$ . С этим начальным приближением к какому-то собственному числу  $\lambda_j$  далее выполняются обратные итерации с переменными сдвигами, как и в (4.24), только приближения к  $\lambda_j$  находятся не через отношения координат векторов  $y^{(k)}$  и  $x^{(k-1)}$ , а через отношения Рэля (как в методе скалярных произведений, см. замечание 4.2 в п.4.2), причем поскольку здесь  $\rho_k$  приближает собственное число данной матрицы  $A$ , а не "сдвинутой", то нет необходимости корректировать получаемое значение на величину смещения спектра, что имело место в формуле (4.23) и в последней из формул (4.24).

**Замечание 4.7.** RQI-алгоритм допускает использование любых векторных норм. Более естественно здесь применение евклидовой нормы; в

<sup>\*)</sup> Один из вариантов останова:  $\|y^{(k)}\| > C$ , где  $C > 0$  – большая константа [44].

таком случае, в силу  $\|x^{(k-1)}\|_2 = \sqrt{(x^{(k-1)}, x^{(k-1)})} \doteq 1$ , вычисление  $\rho_{k-1}$  можно производить по формуле

$$\rho_{k-1} = (Ax^{(k-1)}, x^{(k-1)}).$$

**Замечание 4.8.** Ясно, что применение переменных сдвигов в методе обратных итераций сильно ухудшает от шага к шагу обусловленность матриц решаемых там СЛАУ (они быстро приближаются к вырожденным). Однако, как показали проведенные Уилкинсоном [53] исследования, это не сказывается на достижимой точности получаемых таким методом результатов. Более того, Парлетт [44] аргументировано утверждает полезность плохой обусловленности (редкий случай!) матриц линейных систем в методах обратных итераций с хорошими сдвигами; объяснением этому парадоксальному явлению служит сосредоточение ошибок округлений именно в направлении искомого собственного вектора, что только ускоряет доминирование нужной составляющей.

#### 4.4. МЕТОД ВРАЩЕНИЙ ЯКОБИ РЕШЕНИЯ СИММЕТРИЧНОЙ ПОЛНОЙ ПРОБЛЕМЫ СОБСТВЕННЫХ ЗНАЧЕНИЙ

Дальнейшее изучение методов решения алгебраических проблем собственных значений существенно опирается на матричное *преобразование подобия*. Напомним, что *подобными называются матрицы*  $A$  и  $B = C^{-1}AC$ , где  $C$  – произвольная невырожденная матрица.

Пополним перечень простейших свойств собственных пар (см. п.4.1) еще двумя свойствами.

**Свойство 7.** Пусть  $\{\lambda, x\}$  – собственная пара матрицы  $B = C^{-1}AC$ . Тогда  $\{\lambda, Cx\}$  – собственная пара матрицы  $A$ .

Чтобы убедиться в справедливости этого свойства, достаточно подставить выражение  $B = C^{-1}AC$  в верное для пары  $\{\lambda, x\}$  равенство  $Bx = \lambda x$ : имеем  $C^{-1}ACx = \lambda x$ , откуда после умножения слева на матрицу  $C$  получаем равенство  $ACx = \lambda Cx$ , означающее истинность утверждения.

Как видим, *преобразование подобия сохраняет неизменным спектр любой матрицы.*

**Свойство 8.** Пусть  $A$  —  $n$ -мерная матрица простой структуры (см. п.4.2), а матрицы  $\Lambda = \text{diag}(\lambda_i)$  и  $X = (x_1, x_2, \dots, x_n)$  образованы из ее собственных чисел и собственных векторов соответственно. Тогда  $\Lambda = X^{-1}AX$ .

Действительно, то, что  $\{\lambda_i, x_i\}$  являются собственными парами матрицы  $A$ , означает, что

$$Ax_i = \lambda_i x_i \quad \forall i \in \{1, 2, \dots, n\}.$$

Эти  $n$  равенств могут быть записаны в виде одного матричного равенства

$$AX = X\Lambda. \quad (4.25)$$

В силу простой структуры  $A$ , все ее собственные векторы, т.е. столбцы матрицы  $X$ , линейно независимы, поэтому матрица  $X$  обратима. Умножив равенство (4.25) слева на матрицу  $X^{-1}$ , получим нужное представление  $\Lambda = X^{-1}AX$ .

Так как для диагональной матрицы  $\Lambda$ , образованной из собственных чисел, собственными векторами могут служить единичные векторы исходного базиса ( $\Lambda e_i = \lambda_i e_i \quad \forall i \in \{1, 2, \dots, n\}$ ), то применяя к последнему случаю свойство 7 с  $C=X$  и с  $x = e_i$  (т.е. с  $Cx = Xe_i = x_i$ ), приходим к другой формулировке свойства 8:

если  $\{\lambda_i, e_i\}$  — собственные пары матрицы  $\Lambda = \text{diag}(\lambda_i) = X^{-1}AX$ ,  
то  $\{\lambda_i, x_i\}$  — собственные пары матрицы  $A$

(обозначения те же, что и в свойстве 8).

Далее (в пределах данного пункта) будем рассматривать только симметричные вещественные матрицы. Пользуясь известным фактом о наличии у таких матриц полной ортонормированной системы собственных векторов, т.е. тем, что заявленная выше матрица  $X$  из собственных векторов в этом случае является ортогональной ( $X^{-1} = X^T$ ), запишем как следствие свойства 8 равенство

$$\Lambda = X^T A X. \quad (4.26)$$

Значит, для всякой симметричной матрицы  $A$  найдется диагональная матрица  $\Lambda$ , ей ортогонально подобная. Вопрос теперь состоит в том, как реализовать хотя бы приближенно равенство (4.26), которое позволило бы найти сразу все собственные числа матрицы  $A$  (элементы диагонали матрицы  $\Lambda$ ) и все соответствующие им собственные векторы (столбцы матрицы  $X$ )? Один из возможных ответов на этот вопрос состоит в применении к  $A$  последовательности однопипных преобразований, сохраняющих спектр и приводящих в пределе данную матрицу к диагональному виду.

Для этих целей будем использовать преобразования с помощью так

называемой *матрицы плоских вращений*

$$T_{ij} = \begin{pmatrix} & i & & j & \\ 1 & : & & : & 0 \\ \dots & c & \dots & -s & \dots \\ & : & & : & \\ \dots & s & \dots & c & \dots \\ 0 & : & & : & 1 \end{pmatrix} \begin{matrix} i \\ j \end{matrix} \quad (4.27)$$

Она получается из единичной матрицы заменой двух единиц и двух нулей на пересечениях  $i$ -х и  $j$ -х строк и столбцов числами  $c$  и  $(\pm s)$ , как показано в (4.27), такими, что

$$c^2 + s^2 = 1. \quad (4.28)$$

Условие нормировки (4.28) позволяет интерпретировать числа  $c$  и  $s$  как косинус и синус некоторого угла  $\alpha$ , и, так как умножение любой матрицы на матрицу  $T_{ij}$  изменяет у нее только две строки и два столбца по формулам поворота на угол  $\alpha$  в плоскости, определяемой выбранной парой индексов  $i$  и  $j$ , то это полностью оправдывает название матрицы  $T_{ij}$ .

Матрица  $T_{ij}$  ортогональна при любых  $i, j \in \{1, 2, \dots, n\}$  (проверьте!), и значит, матрица

$$B = T_{ij}^T A T_{ij} \quad (4.29)$$

подобна  $A$ , т.е. имеет тот же набор собственных чисел, что и матрица  $A$ .

*Классический итерационный метод вращений*, предложенный Якоби (1846 г.), предполагает построение последовательности матриц

$$B_0 (= A), B_1, B_2, \dots, B_k, \dots$$

с помощью преобразований типа (4.29)

$$B_k = T_{ij}^T B_{k-1} T_{ij} \quad (4.30)$$

такой, что на  $k$ -м шаге обнуляется максимальный по модулю элемент матрицы  $B_{k-1}$  предыдущего шага (а значит, и симметричный ему элемент). Эта стратегия определяет способ фиксирования пары индексов  $i, j$ , задающих позиции  $(i, i), (j, j), (i, j), (j, i)$  существенных элементов в матрице вращения  $T_{ij}$ , и угол поворота  $\alpha$ , конкретизирующий значения этих элементов  $c = \cos \alpha$  и  $\pm s = \pm \sin \alpha$ . На каждом шаге таких преобразований пересчитываются только две строки (или два столбца, что неважно в силу симметрии) матрицы предыдущего шага. Хотя, к сожалению, нельзя рассчитывать, что таким путем за конечное число шагов можно точно найти диагональную матрицу  $\Lambda$ , ибо полученные на некотором этапе преобразований нулевые элементы на следующем этапе станут, вообще говоря, ненулевыми, но нужное предельное поведение

$$B_k \xrightarrow{k \rightarrow \infty} \Lambda,$$

как будет показано ниже, есть.



Определив идею метода вращений, рассмотрим теперь его несколько подробней.

Пусть  $A = (a_{ml})_{m,l=1}^n$  — исходная симметричная матрица, а  $B = (b_{ml})_{m,l=1}^n$  — матрица, получающаяся после одного шага преобразования по формуле (4.29). Обозначим соответственно через  $\tilde{A}$  и  $\tilde{B}$  двумерные подматрицы этих матриц<sup>\*)</sup>, определяемые фиксированием позиции  $(i, j)$  некоторого элемента  $a_{ij}$  матрицы  $A$ :

$$\tilde{A} = \begin{pmatrix} a_{ii} & a_{ij} \\ a_{ij} & a_{jj} \end{pmatrix}, \quad \tilde{B} = \begin{pmatrix} b_{ii} & b_{ij} \\ b_{ji} & b_{jj} \end{pmatrix},$$

а через  $\tilde{T}$  — такую же подматрицу матрицы  $T_{ij}$ :

$$\tilde{T} = \begin{pmatrix} c & -s \\ s & c \end{pmatrix}.$$

Очевидно, что равенство (4.29), записанное для матриц  $A, B, T_{ij}$ , будет верным и для их подматриц  $\tilde{A}, \tilde{B}, \tilde{T}$ . Пользуясь этим, подсчитаем элементы матрицы  $\tilde{B}$ , выполняя умножение в правой части двумерного аналога (4.29):

$$\begin{aligned} \tilde{B} &= \tilde{T}^T \tilde{A} \tilde{T} = \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} ca_{ii} + sa_{ij} & ca_{ij} - sa_{ii} \\ ca_{ij} + sa_{jj} & ca_{jj} - sa_{ij} \end{pmatrix} = \\ &= \begin{pmatrix} c^2 a_{ii} + 2csa_{ij} + s^2 a_{jj} & c^2 a_{ij} - csa_{ii} + csa_{jj} - s^2 a_{ij} \\ c^2 a_{ij} - csa_{ii} + csa_{jj} - s^2 a_{ij} & c^2 a_{jj} - 2csa_{ij} + s^2 a_{ii} \end{pmatrix}. \end{aligned}$$

Отсюда видим, что  $b_{ij} = b_{ji} = 0$ , если

$$(c^2 - s^2)a_{ij} - cs(a_{ii} - a_{jj}) = 0,$$

т.е. если

$$\frac{cs}{c^2 - s^2} = \frac{a_{ij}}{a_{ii} - a_{jj}}.$$

Учитывая тригонометрическую интерпретацию чисел  $c = \cos \alpha$  и  $s = \sin \alpha$ , в соответствии с чем можно считать

$$cs = \frac{\sin 2\alpha}{2}, \quad c^2 - s^2 = \cos 2\alpha,$$

приходим к выводу, что матрица  $B$  будет иметь нулевые внедиагональные элементы  $b_{ij} = b_{ji}$ , если использовать преобразование плоского вращения

<sup>\*)</sup> В  $\tilde{A}$  элемент  $a_{jj}$  сразу заменяем равным ему элементом  $a_{ii}$ , с равенством же  $b_{ii} = b_{jj}$  в  $\tilde{B}$  пока не торопимся.

по формуле (4.29) на угол  $\alpha$  такой, что

$$\operatorname{tg} 2\alpha = \frac{2a_{ij}}{a_{ii} - a_{jj}}$$

(для определенности считают  $\alpha \in \left(-\frac{\pi}{4}, \frac{\pi}{4}\right)$ ).

Ясно, что нет необходимости находить непосредственно угол  $\alpha$ , поскольку нужные для выполнения преобразований числа  $c$  и  $s$  можно получить через значение  $\operatorname{tg} 2\alpha$  по формулам тригонометрии. При этом сразу отметим, что наибольшие требования к точности в описываемом методе предъявляются именно на стадии вычисления  $c$  и  $s$ , так как здесь возможны наибольшие потери точности, а искажение  $c$  и  $s$  нарушает ортогональность матриц  $T$ , что ведет к неустраняемым погрешностям (метод вращений, итерационный по форме, не является итерационным по существу: ему не присуща самоисправляемость методов последовательных приближений).

Проделав соответствующие элементарные выкладки и возвратившись к  $n$ -мерному случаю, запишем теперь совокупность формул, определяющую один шаг метода вращений Якоби. Для того чтобы не перегружать эти формулы лишними индексами, будем считать, что преобразуется матрица  $A$  в матрицу  $B$  согласно (4.29), хотя на самом деле на  $k$ -м шаге должно применяться преобразование (4.30) к матрице  $B_{k-1} = \left(b_{ml}^{(k-1)}\right)$  с результатом  $B_k = \left(b_{ml}^{(k)}\right)$ .

Итак, пусть  $a_{ij}$  — *ключевой элемент* преобразуемой матрицы  $A$ . Матрица  $B$ , подобная  $A$ , формируется следующим образом:

1. Вычисляют  $p := 2a_{ij}$ ,  $q := a_{ii} - a_{jj}$ ,  $d := \sqrt{p^2 + q^2}$ .
2. Если  $q \neq 0$ , то  $r := |q|/(2d)$ ,  $c := \sqrt{0.5 + r}$ ,  $s := \sqrt{0.5 - r} \cdot \operatorname{sign}(pq)$  (если  $|p| \ll |q|$ , то лучше  $s := |p| \cdot \operatorname{sign}(pq)/(2cd)$ ),  
если же  $q = 0$ , то  $c = s = \sqrt{2}/2$ .
3. Вычисляют:  
 $b_{ii} := c^2 a_{ii} + s^2 a_{jj} + 2csa_{ij}$ ,  
 $b_{jj} := s^2 a_{ii} + c^2 a_{jj} - 2csa_{ij}$   
 (новые диагональные элементы).
4. Полагают  $b_{ij} = b_{ji} := 0$   
 (или для контроля вычисляют  $b_{ij} = b_{ji} := (c^2 - s^2)a_{ij} + cs(a_{jj} - a_{ii})$ );

5. При  $m = 1, 2, \dots, n$  таких, что  $m \neq i$ ,  $m \neq j$ , вычисляют

$$\begin{aligned} b_{im} &= b_{mi} := ca_{mi} + sa_{mj}, \\ b_{jm} &= b_{mj} := -sa_{mi} + ca_{mj}. \end{aligned} \quad (4.31)$$

6. Для всех остальных пар индексов  $m, l$  принимают  $b_{ml} := a_{ml}$ .

Конечно, в реальных вычислениях, если это считать основой алгоритма, не все записанное здесь следует выполнять, а именно, не нужно делать последних переприсвоений, а также должна учитываться симметрия получающейся матрицы  $\mathbf{B}$ .

Убедимся теперь, что если в качестве ключевого или, в иной терминологии [44], *обремененного элемента* на каждом шаге преобразований подобия по указанным формулам брать максимальный по модулю элемент преобразуемой матрицы, то в пределе получится диагональная матрица.

Для доказательства этого, т.е. сходимости последовательности  $(\mathbf{B}_k)$  к  $\Lambda$ , используется (см. приложение 1) норма Фробениуса

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j=1}^n a_{ij}^2}.$$

Проследим за поведением норм матриц (точнее, квадратов этих норм), получающихся из матриц  $\mathbf{B}_k$  заменой диагональных элементов нулями. Такие матрицы, определяемые внедиагональными элементами матриц  $\mathbf{B}_k$ , будем обозначать  $\mathbf{V}_{\mathbf{B}_k}$ . При этом опять для упрощения записей будем пока рассматривать переход от  $\mathbf{A}$  к  $\mathbf{B}$ .

Найдем выражение суммы квадратов внедиагональных элементов матрицы

$$\mathbf{V} = \begin{pmatrix} a_{11} & \dots & b_{1i} & \dots & b_{1j} & \dots & a_{1n} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ b_{i1} & \dots & b_{ii} & \dots & 0 & \dots & b_{in} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ b_{j1} & \dots & 0 & \dots & b_{jj} & \dots & b_{jn} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ a_{n1} & \dots & b_{ni} & \dots & b_{nj} & \dots & a_{nn} \end{pmatrix},$$

принадлежащих изменяющимся по сравнению с  $\mathbf{A}$  строке  $i$  и столбцу  $j$ , через элементы матрицы  $\mathbf{A}$ . Имеем (с учетом (4.31), (4.28) и условия обнуления элементов  $b_{ij} = b_{ji}$ ):

$$\begin{aligned} \sum_{m \neq i} b_{im}^2 + \sum_{m \neq j} b_{mj}^2 &= \sum_{m \neq i, j} (c^2 a_{mi}^2 + 2csa_{mi}a_{mj} + s^2 a_{mj}^2) + \\ &+ b_{ij}^2 + \sum_{m \neq j, i} (s^2 a_{mi}^2 - 2csa_{mi}a_{mj} + c^2 a_{mj}^2) + b_{ji}^2 = \\ &= \sum_{m \neq i, j} [(c^2 + s^2) a_{mi}^2 + (c^2 + s^2) a_{mj}^2] = \sum_{m \neq i, j} (a_{mi}^2 + a_{mj}^2). \end{aligned}$$

Аналогичные суммы квадратов  $j$ -й строки и  $i$ -го столбца, в силу симметрии, дадут такое же выражение. Это означает, что если полученное равенство удвоить и дополнить левую и правую части суммой квадратов всех остальных внедиагональных элементов матрицы  $\mathbf{A}$  (служащих соответствующими элементами и матрицы  $\mathbf{B}$ ), то в левой части будет стоять сумма квадратов всех внедиагональных элементов матрицы  $\mathbf{B}$ , а в правой части – сумма квадратов всех внедиагональных элементов матрицы  $\mathbf{A}$ , кроме  $a_{ji}^2$  и  $a_{ij}^2$ . Следовательно, справедливо равенство

$$\|\mathbf{V}_\mathbf{B}\|_F^2 = \|\mathbf{V}_\mathbf{A}\|_F^2 - 2a_{ij}^2, \quad (4.32)$$

говорящее об убывании сумм квадратов внедиагональных элементов в рассматриваемом процессе преобразований подобия.

Пусть  $|a_{ij}| = \max_{m \neq l} |a_{ml}|$ . Тогда можно считать, что  $a_{ij}^2 = \max_{m \neq l} a_{ml}^2$  не меньше, чем среднее значение множества из  $n^2 - n$  квадратов всех внедиагональных элементов  $n$ -мерной матрицы  $\mathbf{A}$ , т.е. величины

$$\frac{1}{n(n-1)} \sum_{m \neq l} a_{ml}^2 = \frac{1}{n(n-1)} \|\mathbf{V}_\mathbf{A}\|_F^2.$$

Подставляя полученную оценку  $a_{ij}^2 \geq \frac{1}{n(n-1)} \|\mathbf{V}_\mathbf{A}\|_F^2$  в (4.32), приходим к неравенству

$$\|\mathbf{V}_\mathbf{B}\|_F^2 \leq \left(1 - \frac{2}{n(n-1)}\right) \|\mathbf{V}_\mathbf{A}\|_F^2.$$

На основании этого для последовательности матриц  $\mathbf{B}_k$ , представляемых в виде

$$\mathbf{B}_k = \text{diag}\left(b_{ii}^{(k)}\right) + \mathbf{V}_{\mathbf{B}_k},$$

можно записать:

$$\|\mathbf{V}_{\mathbf{B}_k}\|_F^2 \leq \left(1 - \frac{2}{n(n-1)}\right) \|\mathbf{V}_{\mathbf{B}_{k-1}}\|_F^2 \leq \dots \leq \left(1 - \frac{2}{n(n-1)}\right)^k \|\mathbf{V}_\mathbf{A}\|_F^2 \xrightarrow{k \rightarrow \infty} 0.$$

Значит, при указанном способе выбора ключевого элемента последовательность подобных матриц  $\mathbf{B}_k$  сходится к диагональной матрице  $\mathbf{\Lambda}$  из собственных значений, по крайней мере, со скоростью геометрической прогрессии.

Описанный выше классический вариант метода вращений Якоби, как показывают более тонкие оценки, на самом деле имеет асимптотически квадратичную скорость сходимости [15, 44, 53]. Однако при больших размерах  $n$  его реализация наталкивается на существенные потери машинных ресурсов, связанные с поиском наибольшего по модулю ключевого элемента. Поэтому чаще применяется более медленно, но все-таки тоже асимптотически квадратично сходящийся *циклический метод Якоби* с



Обозначим<sup>\*)</sup>  $A_1 := UL$ , тогда  $U = A_1 L^{-1}$ . Подставив это выражение матрицы  $U$  в равенство  $A = LU$ , получаем новое представление  $A$ :

$$A = LA_1 L^{-1}, \quad (4.33)$$

которое говорит о подобии матриц  $A$  и  $A_1$ , т.е. о равенстве их собственных чисел  $\lambda_A$  и  $\lambda_{A_1}$ .

Если матрица  $A_1$  может быть, как и  $A$ , представлена в виде произведения нижней  $L_1$  и верхней  $U_1$  треугольных, т.е.  $A_1 = L_1 U_1$ , то, положив  $A_2 := U_1 L_1$  и выразив отсюда  $U_1 = A_2 L_1^{-1}$ , аналогично предыдущему получаем

$$A_1 = L_1 A_2 L_1^{-1}. \quad (4.34)$$

Следовательно,  $A_1$  подобна  $A_2$  и, значит,  $\lambda_{A_1} = \lambda_{A_2}$ .

Суперпозиция этих двух преобразований, т.е. подстановка (4.34) в (4.33), дает выражение  $A$  через  $A_2$ :

$$A = LL_1 A_2 L_1^{-1} L^{-1} = LL_1 A_2 (LL_1)^{-1},$$

непосредственно утверждающее равенство собственных чисел  $\lambda_A$  и  $\lambda_{A_2}$ .

Такой процесс построения теоретически бесконечной последовательности подобных матриц и составляет основу *LU*- (иначе, *LR*-) *алгоритма*<sup>\*\*)</sup>. Он определяется фактически двумя формулами:

$$A_k = L_k U_k, \quad A_{k+1} = U_k L_k, \quad (4.35)$$

где  $A_0 := A$ ,  $k = 0, 1, 2, \dots$ , причем первая из этих формул означает процедуру треугольной факторизации матрицы  $A_k$  на  $k$ -м шаге, а вторая - простое умножение верхней треугольной матрицы на нижнюю.

Доказано [41, 53], что при ряде ограничений на данную матрицу  $A$  (простейшим из которых является, в частности, требование, чтобы все ее собственные числа были различны по модулю) итерационный процесс (4.35) осуществим, и формируемая им последовательность  $(A_k)$  сходится к треугольной матрице вида

$$\begin{pmatrix} \lambda_1 & * & * & \dots & * \\ 0 & \lambda_2 & * & \dots & * \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \lambda_n \end{pmatrix} \quad \text{или} \quad \begin{pmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ * & \lambda_2 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ * & * & * & \dots & \lambda_n \end{pmatrix}$$

<sup>\*)</sup> Напомним, что произведение матриц, вообще говоря, некоммукативно.

<sup>\*\*)</sup> Алгоритм предложен Рутисхаузером (1958 г.).

в зависимости от того, фиксируется единичная диагональ при LU-факторизации у матрицы  $L$  или у  $U$  соответственно. К сожалению, эти ограничения трудно назвать конструктивными, и реализующие LU-алгоритм программы больше опираются на эмпирику. Осуществимости, устойчивости и ускорения сходимости процесса (4.35) обычно добиваются (если это в принципе возможно) путем подходящих сдвигов матриц и перестановок их элементов; соответствующие исследования и рекомендации по этому поводу можно найти в [53].

**Пример 4.1.** Рассмотрим, как ведет себя LU-алгоритм (4.35), примененный к нахождению собственных чисел матрицы  $A = \begin{pmatrix} 2 & 1 \\ 6 & 1 \end{pmatrix}$ .

Выполнив LU-разложение<sup>\*)</sup>, получим

$$A_0 := A = L_0 U_0 = \begin{pmatrix} 2 & 0 \\ 6 & -2 \end{pmatrix} \begin{pmatrix} 1 & 0.5 \\ 0 & 1 \end{pmatrix}.$$

Перемножая матрицы  $L_0$  и  $U_0$  в обратном порядке, строим матрицу

$$A_1 := U_0 L_0 = \begin{pmatrix} 5 & -1 \\ 6 & -2 \end{pmatrix}.$$

Факторизуя эту матрицу аналогично предыдущему, имеем

$$A_1 = L_1 U_1 = \begin{pmatrix} 5 & 0 \\ 6 & -0.8 \end{pmatrix} \begin{pmatrix} 1 & -0.2 \\ 0 & 1 \end{pmatrix},$$

откуда

$$A_2 := U_1 L_1 = \begin{pmatrix} 3.8 & 0.16 \\ 6 & -0.8 \end{pmatrix}.$$

Следующий шаг дает

$$A_2 = L_2 U_2 = \begin{pmatrix} 3.8 & 0 \\ 6 & -1.0526\dots \end{pmatrix} \begin{pmatrix} 1 & 0.0421\dots \\ 0 & 1 \end{pmatrix},$$

$$A_3 := U_2 L_2 = \begin{pmatrix} 4.0526\dots & 0.0443\dots \\ 6 & -1.0526\dots \end{pmatrix}.$$

Как видим, диагональные элементы матрицы  $A_2$  отличаются от точных значений собственных чисел  $\lambda_1 = 4$ ,  $\lambda_2 = -1$  на 0.2, а матрица  $A_3$  позволяет указать значения  $\lambda_1$  и  $\lambda_2$  с погрешностью  $\approx 0.05$ .

---

<sup>\*)</sup> По формулам из п 2.3  $u_{ij} = a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj}$  при  $i \leq j$ ,  $l_{ij} = \frac{1}{u_{jj}} \left( a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \right)$  при  $i > j$ , используемым поочередно

Если в этом же примере фиксировать единичную диагональ у матриц  $L_k$ , то процесс (4.35) будет развиваться следующим образом<sup>\*)</sup>:

$$\begin{aligned}
 A &= \tilde{L}_0 \tilde{U}_0 = \begin{pmatrix} 1 & 0 \\ 3 & 1 \end{pmatrix} \begin{pmatrix} 2 & 1 \\ 0 & -2 \end{pmatrix}, & \tilde{A}_1 &:= \tilde{U}_0 \tilde{L}_0 = \begin{pmatrix} 5 & 1 \\ -6 & -2 \end{pmatrix}; \\
 \tilde{A}_1 &= \tilde{L}_1 \tilde{U}_1 = \begin{pmatrix} 1 & 0 \\ -1.2 & 1 \end{pmatrix} \begin{pmatrix} 5 & 1 \\ 0 & -0.8 \end{pmatrix}, & \tilde{A}_2 &:= \tilde{U}_1 \tilde{L}_1 = \begin{pmatrix} 3.8 & 1 \\ 0.96 & -0.8 \end{pmatrix}; \\
 \tilde{A}_2 &= \tilde{L}_2 \tilde{U}_2 = \begin{pmatrix} 1 & 0 \\ 0.2526\dots & 1 \end{pmatrix} \begin{pmatrix} 3.8 & 1 \\ 0 & -1.0526\dots \end{pmatrix}, \\
 \tilde{A}_3 &:= \tilde{U}_2 \tilde{L}_2 = \begin{pmatrix} 4.0526\dots & 1 \\ -0.2659\dots & -1.0526\dots \end{pmatrix}.
 \end{aligned}$$

Диагонали матриц  $A_k$  и  $\tilde{A}_k$ , несущие приближения к собственным числам  $A$ , при одних и тех же значениях  $k$  полностью совпадают, но во втором случае не так заметно стремление к нулю поддиагональных элементов, хотя относительная скорость убывания модулей наддиагональных элементов  $A_k$  и поддиагональных элементов  $\tilde{A}_k$  примерно одинакова.

Одним из серьезных факторов, ограничивающих сферу применения LU-алгоритмов, является их недостаточно хорошая численная устойчивость (улучшение этого параметра путем перестановок строк и столбцов сильно отражается на экономичности метода). Этот фактор может играть особенно существенную роль на фоне возможной неустойчивости самой несимметричной проблемы собственных значений.

Ярким примером матрицы, для которой задача нахождения собственных чисел является неустойчивой, служит  $n \times n$ -матрица [41]

$$A \doteq \begin{pmatrix} 0 & 1 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 0 & 1 \\ 0 & 0 & 0 & \dots & 0 & 0 \end{pmatrix},$$

имеющая число 0 собственным значением  $n$ -й кратности. Введем возмущение  $\epsilon$  в левый нижний элемент матрицы  $A$ . Характеристическим урав-

<sup>\*)</sup> Для LU-разложения здесь поочередно используются формулы (см. п.2.3)

$$l_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \quad (i \geq j), \quad u_{kj} = \frac{1}{l_{jj}} \left( a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj} \right) \quad (i < j).$$



нением для возмущенной матрицы

$$A_\varepsilon = \begin{pmatrix} 0 & 1 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 0 & 1 \\ \varepsilon & 0 & 0 & \dots & 0 & 0 \end{pmatrix}$$

будет уравнение

$$\begin{vmatrix} -\lambda & 1 & 0 & \dots & 0 & 0 \\ 0 & -\lambda & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & -\lambda & 1 \\ \varepsilon & 0 & 0 & \dots & 0 & -\lambda \end{vmatrix} = 0.$$

Раскрывая определитель по элементам первого столбца, получаем:

$$-\lambda \cdot \begin{vmatrix} -\lambda & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & -\lambda & 1 \\ 0 & 0 & \dots & 0 & -\lambda \end{vmatrix} + (-1)^{n+1} \varepsilon \cdot \begin{vmatrix} 1 & 0 & \dots & 0 & 0 \\ -\lambda & 1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & -\lambda & 1 \end{vmatrix} = 0 \Leftrightarrow$$

$$\Leftrightarrow -\lambda(-\lambda)^{n-1} + (-1)^{n+1} \varepsilon = 0 \Leftrightarrow \lambda^n - \varepsilon = 0.$$

Следовательно, матрица  $A_\varepsilon$  имеет  $n$  различных, в общем, комплексных собственных значений  $\lambda_i = \sqrt[n]{\varepsilon}$ ,  $i = 1, 2, \dots, n$ . Если взять, например,  $n = 100$ , а  $\varepsilon = 10^{-100}$ , то  $|\lambda_i| = 0.1$ , т.е. чрезвычайно малое, неощутимое для вычислительной машины искажение всего одного элемента данной специфической матрицы приводит к существенному изменению ее спектра.

Разумеется, большинство важных в приложениях задач на собственные значения не так плохи. Однако, обозначив и, возможно, намеренно утрировав проблему, призовем читателя к осторожности в применениях уже рассмотренных методов и интерпретации их результатов, а также к пониманию необходимости построения более устойчивых методов численного решения несимметричных спектральных алгебраических задач (см. следующий пункт). Численная устойчивость всех описываемых здесь методов подробно изучается в [53].

## 4.6. QR-АЛГОРИТМ

В идейном плане от схематично описанного выше LU-алгоритма мало чем отличается так называемый **QR-алгоритм**<sup>\*)</sup> [1, 6, 15, 16, 26, 41, 44, 53]. При  $k=0, 1, 2, \dots$ , начиная с  $A_0 := A$ , здесь строят последовательность матриц  $(A_k)$  по формулам

$$A_k = Q_k R_k, \quad A_{k+1} = R_k Q_k, \quad (4.36)$$

первая из которых означает разложение матрицы  $A_k$  в произведение ортогональной  $Q_k$  и правой треугольной  $R_k$  (такое разложение существует для любой квадратной матрицы [6]), а вторая – перемножение полученных в результате факторизации  $A_k$  матриц  $Q_k$  и  $R_k$  в обратном порядке.

Аналогично предыдущему (см. п.4.5) на основе свойства ортогональных матриц  $Q_k^T = Q_k^{-1}$  в соответствии с (4.36) можно записать представление данной матрицы  $A$  в виде

$$A = Q_0 Q_1 \dots Q_{k-1} Q_k A_{k+1} Q_k^T Q_{k-1}^T \dots Q_1^T Q_0^T$$

или, иначе,

$$A = (Q_0 Q_1 \dots Q_{k-1} Q_k) A_{k+1} (Q_0 Q_1 \dots Q_k)^{-1}. \quad (4.37)$$

Следовательно, любая из матриц последовательности  $(A_k)$  ортогонально подобна матрице  $A$ .

При определенных ограничениях, одним из которых опять выступает требование, чтобы матрица  $A$  не имела равных по модулю собственных значений, генерируемая процессом (4.36) последовательность матриц  $(A_k)$  сходится к матрице правой треугольной формы с диагональю из собственных чисел. Скорость обнуления поддиагональных частей матриц  $A_k$  линейна и зависит, как и во многих ранее рассмотренных методах, от отношений  $\frac{|\lambda_i|}{|\lambda_j|}$  при  $i > j$  (по-прежнему считаем  $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$ ). Наличие

комплексно сопряженных пар собственных чисел у данной вещественной матрицы  $A$  не является, вообще говоря, препятствием для применения QR-алгоритма; просто в этом случае предельной матрицей для последовательности  $(A_k)$  будет матрица квазитреугольного (иначе, блочно-треугольного) вида. Каждой комплексной паре собственных чисел в такой матрице будет соответствовать диагональный  $2 \times 2$ -блок, причем сходимость здесь наблюдается по форме матрицы, а не поэлементно (т.е. элементы внутри этих блоков могут изменяться без видимой зависимости от  $k$  при сохранении неизменными их собственных чисел).

<sup>\*)</sup> Этот метод предложен почти одновременно российским математиком Кублановской (1961 г.) и англичанином Фрэнсисом (1962 г.).

Обычно QR-алгоритм (4.36) применяют не к исходной матрице  $A$ , а к подобной ей *правой почти треугольной матрице*  $B$ , называемой также *матрицей Хессенберга*<sup>1)</sup>, вида

$$B = \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1,n-1} & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2,n-1} & b_{2n} \\ 0 & b_{32} & \dots & b_{3,n-1} & b_{3n} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & b_{n-1,n-1} & b_{n-1,n} \\ 0 & 0 & \dots & b_{n,n-1} & b_{nn} \end{pmatrix}.$$

В основе преобразования  $A$  к виду  $B$  (определяемому условием  $b_{ij} = 0$  при  $j < i - 1$ ) лежит *преобразование Хаусхолдера* или, иначе, *преобразование отражения*, осуществляемое с помощью *матрицы отражения* (Хаусхолдера)

$$H = E - 2WW^T,$$

где  $W$  – произвольный вектор-столбец, но такой, что его евклидова норма равна единице. В силу этого требования к вектору  $W$ ,

$$\|W\|_2^2 = (W, W) = W^T W = 1,$$

и с учетом симметричности матрицы  $H$ , вытекающей из симметричности матрицы

$$WW^T = \begin{pmatrix} W_1 \\ W_2 \\ \dots \\ W_n \end{pmatrix} \cdot (W_1, W_2, \dots, W_n) = \begin{pmatrix} W_1^2 & W_1 W_2 & \dots & W_1 W_n \\ W_2 W_1 & W_2^2 & \dots & W_2 W_n \\ \dots & \dots & \dots & \dots \\ W_n W_1 & W_n W_2 & \dots & W_n^2 \end{pmatrix},$$

имеем:

$$HH^T = H^2 = E - 4WW^T + 4WW^T WW^T = E.$$

Следовательно, матрица отражения ортогональна, и, значит, матрицы  $A$  и  $B$ , связанные соотношением

$$B = HAH \quad (= HAH^T = HAH^{-1}),$$

являются подобными<sup>2)</sup>.

Теперь нетрудно сообразить, как нужно распорядиться свободой задания элементов векторов  $W$  при построении матриц отражения, чтобы за

<sup>1)</sup> Герхард Хессенберг (1874–1925) – немецкий математик. Иногда почти треугольной формой матрицы называют упомянутые матрицы блочно-треугольного вида [41].

<sup>2)</sup> Матрица отражения обладает рядом других интересных свойств; в частности, ее название связано с тем, что линейное преобразование, осуществляемое такой матрицей, оставляет без изменений векторы, ортогональные вектору  $W$ , а коллинеарные ему векторы переводит в противоположные ("отражает").

конечное число шагов преобразований Хаусхолдера произвольно заданную матрицу  $A$  привести к форме Хессенберга  $B$ .

А именно, можно показать, что начатым с  $B_1 = A$  процессом

$$B_{m+1} = H_m B_m H_m, \quad m=1, 2, \dots, n-2, \quad (4.38)$$

где  $H_m = E - 2W_m W_m^T$ , данная  $n \times n$ -матрица  $A$  за  $n-2$  шага будет приведена к виду  $B$ , т.е. матрица  $B := B_{n-1}$  подобна  $A$ , если задающие матрицы Хаусхолдера  $H_m$  векторы  $W_m$  по данной матрице  $A$  строить следующим образом<sup>\*)</sup>.

При  $m=1$  вектор  $W_1$  определяется равенством

$$W_1^T = \mu_1(0, a_{21} - s_1, a_{31}, \dots, a_{n1}), \quad (4.39)$$

где  $s_1 = \text{sign}(-a_{21}) \cdot \sqrt{\sum_{i=2}^n a_{i1}^2}$ ,  $\mu_1 = \frac{1}{\sqrt{2s_1(s_1 - a_{21})}}$ . Такое задание  $W_1$  обеспечивает ортогональность симметричной матрицы

$$H_1 = E - 2W_1 W_1^T$$

и одновременное получение с ее помощью нужных  $n-2$  нулей в первом столбце матрицы

$$B_2 = H_1 B_1 H_1 \quad (= H_1 A H_1).$$

Вектор  $W_2$  по матрице  $B_2$  строится совершенно аналогично, только фиксируются нулевыми не одна, а две первые его координаты, и определяющую роль играют теперь не первый, а второй столбец матрицы  $B_2$  и его третий элемент. При этом у матрицы  $B_3 = H_2 B_2 H_2$  окажется  $n-3$  нулевых элемента во втором столбце и сохранятся полученные на предыдущем шаге нули в первом столбце.

Этот процесс очевидным образом может быть продолжен до исчерпания и без особого труда может быть описан общими формулами типа формул (4.39), для чего нужно лишь ввести обозначения для элементов последовательности матриц  $B_m$  (например, с помощью верхних индексов  $m$ ).

Рассмотрим на простом числовом примере, как приводится матрица к форме Хессенберга, когда для этого требуется только один шаг преобразований Хаусхолдера.

**Пример 4.2.** Дана матрица  $A = \begin{pmatrix} 5 & 1 & -3 \\ 3 & 0 & -2 \\ -4 & -1 & 1 \end{pmatrix}$ . Найти матрицу  $B$ , подобную матрице  $A$  и имеющую форму Хессенберга.

<sup>\*)</sup> Логику таких построений, основанную на сохранении евклидовой нормы столбца, см., например, в [53].

Решение проводим по формулам (4.38), (4.39) при  $n = 3$ , которые для данного случая можно записать так (в естественном для выполнения порядка):

$$s = \text{sign}(-a_{21}) \cdot \sqrt{a_{21}^2 + a_{31}^2}; \quad \mu = \frac{1}{\sqrt{2s(s - a_{21})}};$$

$$\mathbf{W}^T = \mu(0, a_{21} - s, a_{31}); \quad \mathbf{H} = \mathbf{E} - 2\mathbf{W}\mathbf{W}^T; \quad \mathbf{B} = \mathbf{H}\mathbf{A}\mathbf{H}.$$

Имеем:

$$s = -\sqrt{3^2 + (-4)^2} = -5; \quad \mu = \frac{1}{\sqrt{2(-5)(-5-3)}} = \frac{1}{4\sqrt{5}};$$

$$2\mathbf{W}\mathbf{W}^T = \frac{1}{40} \begin{pmatrix} 0 \\ 8 \\ -4 \end{pmatrix} (0, 8, -4) = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1.6 & -0.8 \\ 0 & -0.8 & 0.4 \end{pmatrix};$$

$$\mathbf{H} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & -0.6 & 0.8 \\ 0 & 0.8 & 0.6 \end{pmatrix}; \quad \mathbf{H}\mathbf{A} = \begin{pmatrix} 5 & 1 & -3 \\ -5 & -0.8 & 2 \\ 0 & -0.6 & -1 \end{pmatrix};$$

$$\mathbf{B} = \begin{pmatrix} 5 & -3 & -1 \\ -5 & 2.08 & 0.56 \\ 0 & -0.44 & -1.08 \end{pmatrix}.$$

**Замечание 4.9.** При решении симметричных проблем собственных значений методом вращений на первой стадии решения также часто применяют преобразования Хаусхолдера. Абсолютно те же формулы (4.38), (4.39) приведут симметричную матрицу  $\mathbf{A}$  к подобной ей матрице  $\mathbf{B}$  трехдиагонального вида (частному случаю формы Хессенберга), что значительно повышает эффективность последующих преобразований плоских вращений Якоби.

**Замечание 4.10.** Если начать процесс построения матриц отражения  $\mathbf{H}_m$  не по формулам (4.39), а по аналогичным им, но с ключевым элементом  $a_{11}$ , т.е. полагая

$$\mathbf{W}_1^T = \mu_1(a_{11} - s, a_{21}, \dots, a_{n1}),$$

$$s_1 = \text{sign}(-a_{11}) \cdot \sqrt{\sum_{i=1}^n a_{i1}^2}, \quad \mu_1 = \frac{1}{\sqrt{2s_1(s_1 - a_{11})}},$$

то за  $n-1$  шаг ортогональных преобразований

$$\mathbf{R}_{m+1} = \mathbf{H}_m \mathbf{R}_m; \quad m = 1, 2, \dots, n-1; \quad \mathbf{R}_1 := \mathbf{A}$$

можно получить разложение матрицы  $A$  в произведение ортогональной  $H := H_{n-1} \dots H_2 H_1$  и правой треугольной  $R := R_n$ , так как результирующее равенство  $R = HA$  равносильно равенству  $A = HR$  в силу ортогональности и симметричности  $H$ . Приведение матрицы  $A$  такими преобразованиями к треугольному виду составляет основу *метода отражений решения линейных алгебраических систем*. Из равносильности равенств  $Ax=b$ ,  $HAx=Hb$  и  $Rx=Hb$  легко понять, что для решения СЛАУ этим методом нужно над вектором свободных членов  $b$  выполнять те же преобразования, что и над матрицей коэффициентов  $A$ , после чего нужно будет сделать только обратный ход, как в методе Гаусса. Метод отражений решения СЛАУ в полтора раза экономичнее метода вращений (см. п.2.7) и практически не уступает последнему по устойчивости к накоплению ошибок округлений.

Вообще говоря, весь QR-алгоритм (4.36) от начала и до конца может быть построен на базе описанной выше процедуры преобразований Хаусхолдера, направленной сразу на триангуляризацию в соответствии с замечанием 4.10. Однако такой подход значительно сужает границы применимости алгоритма и ухудшает его скоростные качества.

Обычно на втором этапе применяют другое ортогональное преобразование – *преобразование плоских вращений Гивенса*. Определяющая эти преобразования матрица  $G_{ij} = (g_{ml})_{m,l=1}^n$  при фиксированных  $i, j$  – индексах ключевого элемента преобразуемой матрицы – имеет точно такую же структуру, как и матрица плоских вращений Якоби  $T_{ij}$  (ср.(4.27)), только здесь, следуя [41], двумерную подматрицу из элементов, стоящих на пересечении  $i$ -х и  $j$ -х строк и столбцов, возьмем в виде

$$\hat{G}_{ij} = \begin{pmatrix} s & c \\ -c & s \end{pmatrix}.$$

Как и прежде, числа  $s$  и  $c$  связываем соотношением  $s^2+c^2=1$  (это позволяет интерпретировать их как синус и косинус некоторого угла  $\theta$ ), обеспечивающим ортонормированность матриц  $G_{ij}$ .

Первый полный шаг преобразования Гивенса, применяемого к матрице Хессенберга  $B$   $n$ -го порядка в рамках QR-алгоритма (4.36), состоит из  $n-1$  элементарных подшагов, имеющих целью последовательное обнуление поддиагональных элементов в столбцах от первого до  $(n-1)$ -го. В результате этого получается разложение матрицы  $B$  в произведение ортогональной и треугольной, что требуется первой формулой (4.36) при  $k=1$ ,  $A_1 := B$ .

Чтобы определить  $s$  и  $c$  на первом промежуточном шаге, рассмотрим произведение матриц<sup>\*)</sup>

$$\mathbf{G}_1 \mathbf{B} = \begin{pmatrix} s & c & 0 & \dots & 0 \\ -c & s & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} & b_{13} & \dots & b_{1n} \\ b_{21} & b_{22} & b_{23} & \dots & b_{2n} \\ 0 & b_{32} & b_{33} & \dots & b_{3n} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & b_{nn} \end{pmatrix} =$$

$$= \begin{pmatrix} sb_{11}+cb_{21} & sb_{12}+cb_{22} & sb_{13}+cb_{23} & \dots & sb_{1n}+cb_{2n} \\ -cb_{11}+sb_{21} & -cb_{12}+sb_{22} & -cb_{13}+sb_{23} & \dots & -cb_{1n}+sb_{2n} \\ 0 & b_{32} & b_{33} & \dots & b_{3n} \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & b_{nn} \end{pmatrix}.$$

Беря  $c = \cos \theta$  и  $s = \sin \theta$  такими, что  $\operatorname{tg} \theta = \frac{b_{11}}{b_{21}}$ , будем иметь

$b_{11} \cos \theta = b_{21} \sin \theta$ , т.е.  $-c b_{11} + s b_{21} = 0$ . Значит, результат первого промежуточного шага – матрица  $\mathbf{B}_1 = \mathbf{G}_1 \mathbf{B}$ , получающаяся при таких  $c$  и  $s$ , не будет содержать ненулевых элементов под диагональю в первом столбце.

Второй промежуточный шаг совершается аналогично: матрица  $\mathbf{B}_2 = \mathbf{G}_2 \mathbf{B}_1$  получается из предыдущей  $\mathbf{B}_1$  с помощью матрицы Гивенса

$\mathbf{G}_2$ , отличающейся от  $\mathbf{G}_1$  тем, что подматрица  $\begin{pmatrix} s & c \\ -c & s \end{pmatrix}$  смещается на одну

позицию вдоль диагонали и угол поворота подбирается так, чтобы в матрице  $\mathbf{B}_2$  обнулить элемент  $b_{32}^{(2)}$ .

Продолжая этот процесс преобразований Гивенса далее, в итоге получим правую треугольную матрицу

$$\mathbf{B}_{n-1} = \mathbf{G}_{n-1} \mathbf{G}_{n-2} \dots \mathbf{G}_2 \mathbf{G}_1 \mathbf{B}.$$

Последнее равенство можно переписать в виде

$$(\mathbf{G}_{n-1} \mathbf{G}_{n-2} \dots \mathbf{G}_2 \mathbf{G}_1)^{-1} \mathbf{B}_{n-1} = \mathbf{B},$$

который позволяет считать выполненным требуемое в (4.36) при  $k=1$  разложение

$$\mathbf{B} = \mathbf{Q}_1 \mathbf{R}_1,$$

где  $\mathbf{Q}_1 := (\mathbf{G}_{n-1} \mathbf{G}_{n-2} \dots \mathbf{G}_2 \mathbf{G}_1)^{-1} = \mathbf{G}_1^T \dots \mathbf{G}_{n-1}^T$  – ортогональная, а

<sup>\*)</sup> Так как в каждом столбце матрицы Хессенберга нужно "убивать" только по одному элементу, то в данном случае матрицы элементарных вращений Гивенса можно помечать только одним индексом

$\mathbf{R}_1 := \mathbf{V}_{n-1}$  – правая треугольная матрица. При этом матрица

$$\mathbf{A}_2 := \mathbf{R}_1 \mathbf{Q}_1 = (\mathbf{G}_{n-1} \dots \mathbf{G}_1) \mathbf{B} (\mathbf{G}_{n-1} \dots \mathbf{G}_1)^{-1}, \quad (4.40)$$

являющаяся результатом первого полного шага QR-алгоритма (примененного к  $\mathbf{B}$ ), сохраняет не только спектр данной матрицы, но и форму Хессенберга [16, 41], благодаря чему приведение исходной матрицы  $\mathbf{A}$  к почти треугольному виду  $\mathbf{B}$  достаточно сделать только один раз.

Очевидно, скалярные параметры  $c_j = \cos \theta_j$  и  $s_j = \sin \theta_j$  матриц Гивенса  $\mathbf{G}_j$ , с помощью которых осуществляется переход от матрицы Хессенберга  $\mathbf{B} = (b_{ij})_{i,j=1}^n$  "транзитом" через матрицы Хессенберга

$\mathbf{B}_j = (b_{im}^{(j)})_{i,m=1}^n$  к матрице Хессенберга  $\mathbf{A}_2$ , можно вычислять на  $j$ -м промежуточном шаге ( $j=1, 2, \dots, n-1$ ) по формулам

$$c_j = \frac{1}{\sqrt{1+t_j^2}}, \quad s_j = t_j c_j,$$

где

$$t_j = \frac{b_{jj}^{(j-1)}}{b_{j+1,j}^{(j-1)}} (= \operatorname{tg} \theta_j), \quad b_{ij}^{(0)} := b_{ij}$$

(если знаменатель в выражении  $t_j$  равен нулю или по модулю меньше некоторого существенно малого порогового значения, то можно считать  $c_j=0$ ,  $s_j=1$ , т.е.  $\mathbf{G}_j := \mathbf{E}$ ).

**Пример 4.3.** Преобразованиями Гивенса выполним один шаг

QR-алгоритма для матрицы  $\mathbf{B} = \begin{pmatrix} 5 & -3 & -1 \\ -5 & 2.08 & 0.56 \\ 0 & -0.44 & -1.08 \end{pmatrix}$ , полученной в результате преобразований Хаусхолдера в предыдущем примере.

При  $j=1$  последовательно находим:

$$t_1 = \frac{5}{-5} = -1, \quad c_1 = \frac{1}{\sqrt{1+(-1)^2}} = \frac{1}{\sqrt{2}}, \quad s_1 = -\frac{1}{\sqrt{2}};$$

$$\mathbf{G}_1 = \begin{pmatrix} -0.5\sqrt{2} & 0.5\sqrt{2} & 0 \\ -0.5\sqrt{2} & -0.5\sqrt{2} & 0 \\ 0 & 0 & 1 \end{pmatrix},$$

$$\mathbf{B}_1 = \mathbf{G}_1 \mathbf{B} = \begin{pmatrix} -5\sqrt{2} & 2.54\sqrt{2} & 0.78\sqrt{2} \\ 0 & 0.46\sqrt{2} & 0.22\sqrt{2} \\ 0 & -0.44 & -1.08 \end{pmatrix}.$$



При  $j=2$  вычисляем (округляя до  $10^{-6}$ ):

$$t_2 = \frac{0.46\sqrt{2}}{-0.44} = -1.478496, \quad c_2 = 0.560248, \quad s_2 = -0.828325;$$

значит,

$$G_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & -0.828325 & 0.560248 \\ 0 & -0.560248 & -0.828325 \end{pmatrix},$$

$$R_1 := B_2 = G_2 B_1 = \begin{pmatrix} -7.071068 & 3.592102 & 1.103087 \\ 0 & -0.785366 & -0.862782 \\ 0 & 0 & 0.720283 \end{pmatrix},$$

Следовательно,

$$Q_1 = G_1^T G_2^T = \begin{pmatrix} -0.707107 & -0.585714 & -0.396155 \\ 0.585714 & -0.396155 & 0.396155 \\ 0 & 0.560248 & -0.828325 \end{pmatrix},$$

и, согласно (4.40), матрица

$$A_2 = R_1 Q_1 = \begin{pmatrix} 7.103946 & 3.336597 & 3.310554 \\ -0.460000 & -0.172245 & 0.403537 \\ 0 & 0.403537 & -0.596628 \end{pmatrix}$$

есть искомый результат первого полного шага QR-алгоритма. Она имеет те же собственные числа, что  $B$  и  $A$ , сохраняет форму Хессенберга и модули ее поддиагональных элементов меньше, чем у матрицы  $B$ , т.е. она более близка к подобной  $B$  матрице треугольного вида, на диагонали которой должны быть собственные числа данной матрицы ( $\lambda_1 \approx 7.693$ ,  $\lambda_2 \approx -1.205$ ,  $\lambda_3 \approx -0.435$ ).

**Замечание 4.11.** Весь QR-алгоритм можно было бы построить на базе одних только преобразований Гивенса, т.е. не приводя исходную матрицу  $A$  к форме Хессенберга (или к трехдиагональному виду, если  $A$  симметрична) другими преобразованиями. В таком случае стала бы заметной разница между преобразованиями Якоби и Гивенса. Суть этой разницы в следующем: если для преобразований Якоби понятия "ключевой элемент" и "обреченный элемент" совпадают, то для преобразований Гивенса это, вообще говоря, не так. В общем случае при вращениях Гивенса угол поворота  $\theta$  в фиксированной индексными  $i, j$  плоскости вращения подбирается так, чтобы аннулировать какой-нибудь элемент, стоящий либо в одном столбце, либо в одной строке с ключевым элементом  $a_{ij}$ . Такие преобразования теряют свойство минимальности суммы квадратов внедиагональных элементов, имевшее место в преобразованиях Якоби для симметричных матриц, но позволяют (Гивенс, 1954г. [53]) привести симметричную мат-

рицу к трехдиагональному виду существенно быстрее, чем это требуется для выполнения одного цикла преобразований в методе вращений Якоби<sup>1)</sup>. Приведение несимметричных матриц к форме Хессенберга методом Гивенса требует большего числа арифметических операций, чем это нужно для такого приведения методом Хаусхолдера, поэтому обычно для этих целей отдают предпочтение последнему.

Приведенных выше сведений, в принципе, вполне достаточно, чтобы находить QR-алгоритмом все хорошо отделимые вещественные собственные числа вещественных матриц, реализуя равенство (4.37) при некотором  $k$ , хотя здесь и нет должного для этого обоснования<sup>2)</sup>. Однако в такой непосредственной форме QR-алгоритм не применяется ввиду его медленной сходимости. Для ускорения сходимости в процесс (4.36) вводят *сдвиги*, о роли которых немало говорилось при изучении метода обратных итераций (см. п.4.3). При выборе параметров сдвигов в QR-алгоритме учитывается подмеченная ранее целесообразность использования для этого приближений к собственным числам. В данном случае принимается во внимание, что формируемая QR-алгоритмом последовательность матриц  $A_k$  в пределе дает правую треугольную матрицу с диагональю из собственных чисел (когда все собственные числа матрицы  $A$  вещественны). Следовательно, есть основания утверждать, что последовательность элементов  $a_{nn}^{(k)}$  может рассматриваться при  $k=1, 2, \dots$  как последовательность приближений к какому-то определенному собственному числу матрицы  $A$  и служить соответствующей последовательностью параметров переменных сдвигов, ускоряющих процесс обнуления поддиагональных элементов.

Таким образом, по крайней мере, в случае, когда данная матрица  $A$  имеет только вещественные собственные значения, QR-алгоритм будет сходиться более быстро (квадратично), если при каждом  $k$  преобразования Гивенса применять не к матрице  $A_k$ , а к матрице  $\tilde{A}_k := A_k - a_{nn}^{(k)} E$ . При этом каждый раз спектр матрицы смещается на величину произведенного сдвига (см. свойство 2 в п.4.1), что может учитываться двояко: либо параметры сдвигов суммируются и затем сумма прибавляется к найденным в итоге значениям, либо каждый раз делается обратный сдвиг, т.е. вместо формул (4.36) используются формулы:

$$A_k - a_{nn}^{(k)} E = Q_k R_k$$

(QR-факторизация матрицы  $A_k - a_{nn}^{(k)} E$ ),

<sup>1)</sup> Для вычисления собственных значений трехдиагональных матриц разработан не рассматриваемый здесь довольно эффективный *метод бисекций* [11, 15, 53], который также можно считать одним из способов локализации собственных чисел.

<sup>2)</sup> Доказательство сходимости QR-алгоритма и его модификаций см. в [26, 53]. В [44] изучается сходимость принятого там за основу QL-алгоритма.

$$A_{k+1} = R_k Q_k + a_{nn}^{(k)} E$$

(перемножение в обратном порядке и обратный сдвиг).

Важно учесть, что большая экономия вычислительных затрат при реализации QR-алгоритма получается в результате включения сюда *процедуры исчерпывания*. Нетрудно показать, что если на каком-то шаге  $k=k_0$  последняя строка матрицы  $n$ -го порядка  $A_{k_0}$  стала нулевой (с некоторой точностью), то это позволяет считать  $a_{nn}^{(k_0)}$  собственным числом, а другие ее собственные значения находить, работая далее только с подматрицей  $(n-1)$ -го порядка, получающейся из  $A_{k_0}$  отбрасыванием последних строки и столбца.

Если у данной матрицы возможны комплексные собственные значения, применяют более сложный процесс двойных сдвигов [26, 41, 53], позволяющий преодолеть ситуацию, когда в конце диагоналей матриц  $A_k$  "прорисовывается"  $2 \times 2$ -блок, соответствующий паре комплексно сопряженных собственных чисел.

Нахождение собственных векторов в рамках QR-алгоритма (и других методов, основанных на асимптотической триангуляризации) не является такой простой задачей, как это было в методе вращений Якоби для симметричных матриц, диагонализированных в процессе ортогональных преобразований. Однако при известном собственном числе соответствующий ему собственный вектор эффективно может быть найден рассмотренным в п.4.3 методом обратных итераций (см. формулы (4.19), (4.20)). При этом обратные итерации обычно применяются не к исходной матрице  $A$ , а к матрице  $B$ , подобной  $A$  и имеющей форму Хессенберга. Если приведение  $A$  к виду  $B$  выполнялось преобразованиями Хаусхолдера (4.38), то  $B = HAH$ , где  $H$  – результирующая матрица  $n-2$  элементарных вращений, и значит, согласно свойству 7 п.4.4, найдя собственный вектор  $u$  матрицы  $B$ , искомый собственный вектор  $x$  данной матрицы  $A$  получаем равенством  $x = Hu$ .

## УПРАЖНЕНИЯ

4.1. Дана матрица  $A = \begin{pmatrix} 30 & -12 & 53 \\ -42 & 19 & -78 \\ -28 & 12 & -51 \end{pmatrix}$ .

а) Степенным методом найти несколько последовательных приближений к доминирующему собственному числу матрицы  $A$  и к соответствующему собственному вектору. Зная, что искомое собственное число есть

$\lambda_1 = -5$ , проверить, насколько эффективно здесь применение  $\delta^2$ -процесса Эйткена (см. замечание 4.6).

б) Методом обратных итераций найти младшую собственную пару  $\{\lambda_3, x_3\}$  данной матрицы  $A$ . Можно ли утверждать, что  $\lambda_3 = \Lambda + \lambda_1$ , где  $\Lambda$  – наибольшее по модулю собственное число матрицы  $A - \lambda_1 E$ ?

4.2. Найдя грубые приближения к собственным числам матрицы

$$A = \begin{pmatrix} 5 & 2 & -3 \\ 4 & 5 & -4 \\ 6 & 4 & -4 \end{pmatrix}$$

степенным методом, уточнить эти значения обратными итерациями со сдвигами (см. (4.24)).

4.3.а) Проанализировать сходимость степенного метода в случае, когда  $\lambda_1$  – кратное вещественное наибольшее по модулю собственное число  $n$ -мерной матрицы простой структуры (см. формулы (4.10), (4.11)). Как можно найти все соответствующие ему собственные векторы в зависимости от показателя кратности?

б) Что можно сказать о поведении последовательности отношений (4.10), если  $\lambda_1 = -\lambda_2$ , и  $|\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|$  ( $\lambda_i \in R$ )?

в) Рассмотреть и объяснить поведение степенного метода в случае, когда данная матрица  $A$  – диагональная.

4.4. Найти все собственные пары матрицы

$$A = \begin{pmatrix} 4 & 2 & -1 \\ 2 & 4 & 1 \\ -1 & 1 & 3 \end{pmatrix} :$$

а) методом скалярных произведений (для нахождения второй собственной пары использовать формулы (4.14), (4.15));

б) RQI-алгоритмом, начиная его с различных векторов.

4.5. Сравнить два подхода к нахождению всех собственных чисел матрицы

$$A = \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}$$

методом вращений Якоби:

а) применяя его непосредственно к данной матрице;

б) предварительно приведя ее к трехдиагональному виду преобразованиями Хаусхолдера.

4.6. Для нахождения собственных пар симметричных положительно определенных матриц построить LU - алгоритм на базе  $U^T U$  (или  $LL^T$ )-разложения Холецкого (см. главу 2). Опробовать его на матрицах

$$A = \begin{pmatrix} 5 & -2 \\ -2 & 2 \end{pmatrix} \quad \text{и} \quad B = \begin{pmatrix} 6 & 1 & 0 \\ 1 & 5 & -2 \\ 0 & -2 & 1 \end{pmatrix}.$$

Сохраняют ли получаемые на каждом шаге такого алгоритма подобные  $B$  матрицы трехдиагональную структуру?

4.7. Дана матрица  $A = \begin{pmatrix} 3 & -1 \\ -2 & 2 \end{pmatrix}$ . Сделать по три шага:

- LU-алгоритма;
- QR-алгоритма на основе преобразований Гивенса;
- QR-алгоритма на основе преобразований Хаусхолдера (с учетом замечания 4.10).

Сравнить полученные приближения к собственным числам матрицы  $A$  по точности (найдя сначала ее точные значения с помощью характеристического уравнения) и по вычислительным затратам.

4.8. Матрицу  $A$  из упражнения 4.5 привести к трехдиагональному виду преобразованием Гивенса (отличным от преобразования Якоби, см. замечание 4.11).

4.9. Для  $n \times n$ -матрицы сделать приблизительный подсчет количества арифметических операций, приходящихся на:

- один шаг степенного метода;
- один шаг метода обратных итераций (без сдвигов);
- один полный цикл метода вращений Якоби;
- полный цикл приведения матрицы к форме Хессенберга преобразованиями Хаусхолдера.

4.10. Методом отражений решить систему

$$\begin{cases} x + 2y + z = 0, \\ 2x - y - 2z = 6, \\ 2x + y - z = 9 \end{cases}$$

(см. замечание 4.10).

4.11. Вывести формулы для приведения  $n \times n$ -матрицы Хессенберга к треугольной форме плоскими поворотами на углы  $\varphi_j$  ( $j=1,2,\dots,n-1$ ) с помощью матриц  $T_j := T_{j,j+1}$  вида (4.27). Как отличаются углы  $\varphi_j$  от углов  $\theta_j$ , неявно присутствующих в представлении (4.40)?

# ГЛАВА 5 || МЕТОДЫ РЕШЕНИЯ НЕЛИНЕЙНЫХ СКАЛЯРНЫХ УРАВНЕНИЙ

Рассматривается задача вычисления действительных корней алгебраических и трансцендентных уравнений. Изучаются как универсальные методы (дихотомии, Ньютона, секущих) решения нелинейных уравнений, так и специальные подходы к нахождению корней многочленов. Большое внимание уделено методу простых итераций и базирующимся на нем  $\Delta^2$ -процессу Эйткена и методу Вегстейна. На примере логистического уравнения (с параметром) исследуется возможное поведение итерационных последовательностей в случае нарушения одного из достаточных условий сходимости, дается представление о бифуркациях решений и циклов.

## 5.1. ЛОКАЛИЗАЦИЯ КОРНЕЙ

Будем рассматривать задачу *приближенного нахождения нулей функции одной переменной*, иначе, задачу нахождения корней уравнения вида

$$f(x)=0, \quad (5.1)$$

где  $f: R_1 \rightarrow R_1$  – алгебраическая или трансцендентная функция. Такие уравнения называют *скалярными, числовыми, конечными* и т.п. Методам их решения посвящена обширная литература; кроме перечисленных здесь учебных и научных изданий, эту тему можно встретить почти в любом учебнике математического анализа или высшей математики, а во многих учебных пособиях по программированию даются “рецепты” решения таких задач<sup>\*)</sup>. Однако не все авторы учебников по вычислительной математике считают нужным включать рассматриваемую тему самостоятельным разделом. И в этом есть резон. Действительно, многие наиболее фундаментальные методы решения скалярных уравнений можно рассматривать как частные случаи соответствующих методов решения систем нелинейных уравнений со многими неизвестными и даже, более того, нелинейных операторных уравнений (как правило, в банаховых пространствах). Такой путь изучения методов более короток, но и более труден. Учитывая, что одномерный случай более прост и легко интерпретируется геометрически

---

<sup>\*)</sup> Можно понять преподавателей, обучающих алгоритмическим языкам: методы решения скалярных уравнений предоставляют для этого благодатный материал. К сожалению, программирование формул без понимания их сути часто лишь компрометирует вычислительную математику.

(что немаловажно для нелинейных задач), а также то, что обычно теоретические результаты, полученные для скалярных уравнений, затем переносятся (не без потерь) на системы и операторные уравнения, изучаемые методы будут гораздо лучше поняты, если в них сначала как следует разобраться на объектах вида (5.1).

В общем случае, если речь идет не об отдельных достаточно узких классах уравнений, например, изучавшихся в школьном курсе математики, можно говорить лишь о приближенном вычислении корней уравнений (5.1), т.е. таких значений аргумента  $x=\xi$ , при которых равенство

$$f(\xi) = 0$$

истинно. При этом под близостью приближенного значения  $\bar{x}$  к корню  $\xi$  уравнения (5.1), как правило, понимают выполнение неравенства

$$|\xi - \bar{x}| < \varepsilon$$

при малых  $\varepsilon > 0$ , хотя часто бывает важным контролировать не абсолютную погрешность приближенного равенства  $\bar{x} \approx \xi$ , а относительную, т.е. величину  $|\xi - \bar{x}|/|\bar{x}|$ , например, когда величина  $|\bar{x}|$  близка к нулю.

Нелинейная функция  $f(x)$  в своей области определения  $D(f) \subseteq R_1$  может иметь конечное или бесконечное количество нулей или не иметь их вовсе. Большинство же методов нахождения нулей требует знания промежутков (возможно, малых), где имеется и притом единственный нуль функции<sup>\*)</sup>. Если такие конкретные промежутки не предоставляются постановщиком задачи, то на первый план выходят: выявление ситуации с наличием и количеством корней уравнения (5.1); нахождение области их расположения, получение отрезков, на которых имеется точно по одному корню. Иными словами, ставятся *подзадачи существования и единственности, нахождения границ и локализации корней*. Эти подзадачи обычно решаются в комплексе средствами математического анализа. Но и численные методы здесь часто выступают в помощь математическому анализу: как будет видно из дальнейшего, многие теоремы сходимости итерационных методов можно считать локальными теоремами существования и единственности.

Для функций общего вида нет универсальных способов решения поставленных подзадач. Так что здесь, как говорится, все средства хороши.

Если функция  $f(x)$  такова, что без особого труда можно построить ее график, этим следует воспользоваться, чтобы представить ситуацию с количеством и расположением нулей  $f(x)$ , выделяя те промежутки оси абсцисс, где график  $y=f(x)$  пересекает  $Ox$ . (Знание графика много дает и для понимания поведения тех или иных процессов вычисления приближе-

---

<sup>\*)</sup> Перефразируя Конфуция, можно сказать: трудно найти корень на бесконечном промежутке, особенно когда его там нет.

ний к корням.) Может оказаться, что построение графика  $y=f(x)$  вызывает затруднения, но исходное уравнение (5.1) очевидным образом представляется в виде

$$f_1(x) = f_2(x)$$

и функции  $f_1(x)$  и  $f_2(x)$  таковы, что легко строятся графики  $y=f_1(x)$  и  $y=f_2(x)$ . Тогда задача определения количества корней и областей их единственности решается отслеживанием точек пересечения этих графиков и выделением на оси абсцисс тех промежутков, которым принадлежат проекции таких точек. Описанный прием называют *графическим способом локализации* (иначе, *отделения, изоляции*) корней.

**Пример 5.1.** Найдем промежутки изоляции корней уравнения

$$x^2 - \sin x - 1 = 0.$$

Представив это уравнение в виде  $x^2 - 1 = \sin x$ , строим схематично графики функций  $y = x^2 - 1$  и  $y = \sin x$  (рис. 5.1).

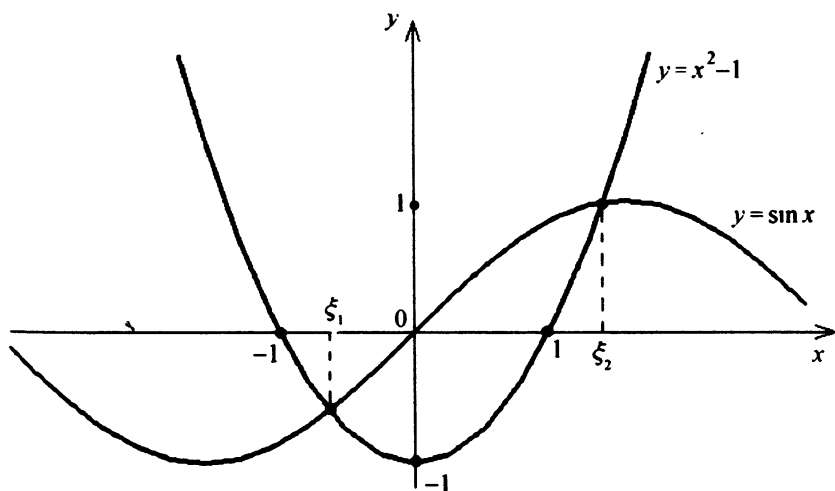


Рис. 5.1. Графическая локализация корней уравнения  $x^2 - \sin x - 1 = 0$

Совместное рассмотрение графиков позволяет сделать заключение, что данное уравнение имеет два корня:  $\xi_1 \in [-1; 0]$  и  $\xi_2 \in [1; \pi]$ .



Убедиться в том, что на данном отрезке  $[a; b]$  (например, грубо определенном графическим способом) действительно имеется нуль непрерывной функции  $f(x)$ , можно *аналитическим способом*, в основе которого лежит известное утверждение математического анализа.

**Теорема 5.1 (Больцано-Коши).** *Если непрерывная на отрезке  $[a; b]$  функция  $f(x)$  на концах его имеет противоположные знаки, т.е.*

$$f(a)f(b) < 0, \quad (5.2)$$

*то на интервале  $(a; b)$  она хотя бы один раз обращается в нуль*<sup>\*)</sup>.

Очевидна слабость теоремы 1 при ее применении к поставленной задаче: она не дает ответа на вопрос о количестве корней на отрезке  $[a; b]$  в случае выполнения условия (5.2) (их может быть нечетное число) и не позволяет утверждать, что на отрезке  $[a; b]$  нет корней, когда условие (5.2) не выполнено, так как в этом случае их может оказаться на  $[a; b]$  четное количество (нетрудно представить себе соответствующие “картинки” возможных поведений графиков функций на отрезке).

Результат, сформулированный в виде теоремы 1, можно значительно усилить, если требование непрерывности функции  $f(x)$  на  $[a; b]$  дополнить требованием монотонности ее на этом отрезке.

**Теорема 5.2.** *Непрерывная строго монотонная функция  $f(x)$  имеет и притом единственный нуль на отрезке  $[a; b]$  тогда и только тогда, когда на его концах она принимает значения разных знаков.*

Последняя теорема позволяет не только принимать, но и отвергать те или иные промежутки из области определения  $D(f)$  данной функции на предмет дальнейшего поиска ее нулей, если известно о ее монотонном поведении на этих промежутках и определены знаки значений функции на их концах.

Реально установить монотонность на данном отрезке можно для дифференцируемой функции, потребовав знакопостоянства ее производной на всем отрезке. Для таких функций основой решения задачи локализации корней уравнения (5.1) может служить следующая теорема.

---

<sup>\*)</sup> Допустимо считать, что  $a$  и/или  $b$  могут принимать значения  $-\infty$  и/или  $+\infty$  соответственно. В таком случае  $f(a)$  и/или  $f(b)$  следует понимать в предельном смысле, также допуская возможность бесконечных предельных значений  $f$  определенного знака.

**Теорема 5.3.** Пусть  $f \in C^1_{[a;b]}$ . Тогда если  $f'(x)$  не меняет знак на  $(a;b)$ , то условие (5.2) является необходимым и достаточным для того, чтобы уравнение (5.1) имело и притом единственный корень на  $[a;b]$ .

Так как производная может менять знак только в точках, где она равна нулю или не существует, а также в граничных точках области определения функции, то в случаях (к сожалению, редких), когда уравнение  $f'(x) = 0$  легко решается, вопрос о числе и расположении корней уравнения (5.1) не вызывает трудностей.

**Пример 5.2.** Выясним, сколько корней у уравнения  $x^2 e^x = \pi$  и где они расположены.

Обозначим  $f(x) = x^2 e^x - \pi$ . Тогда  $f'(x) = x(x+2)e^x$ . Очевидно,  $f'(x) = 0$  только при  $x = 0$  и  $x = -2$ . Поскольку  $f(x)$  и  $f'(x)$  определены и непрерывны на всей числовой оси, точки  $-2$  и  $0$  — единственные на  $Ox$  такие, в которых может происходить смена убывания функции  $y = f(x)$  на возрастание или наоборот. Поэтому, найдя знаки значений (в том числе, бесконечных)  $\lim_{x \rightarrow -\infty} f(x)$ ,  $f(-2)$ ,  $f(0)$  и  $\lim_{x \rightarrow +\infty} f(x)$ , т.е. заполнив таблицу знаков

$$f : \begin{array}{c|c|c|c} -\infty & -2 & 0 & +\infty \\ \hline - & - & - & + \end{array},$$

можно на основании теоремы 5.3 утверждать, что данное уравнение имеет единственный корень, и этот корень положителен. Выяснив еще знаки  $f(x)$  в точках  $x = 1$  (минус) и  $x = 2$  (плюс), область поиска корня данного уравнения с бесконечного промежутка  $[0; +\infty)$  сужаем до промежутка единичной длины  $[1; 2]$ .

В ситуациях, далеких от рассмотренной идеальной, часто поступают следующим образом. Всю область определения (если она конечна) или какую-нибудь ее часть, вызывающую по тем или иным соображениям интерес, разбивают на отрезки точками  $x_i$ , расположенными на условно небольшом расстоянии  $h$  одна от другой (сюда включаются также граничные точки области определения). Вычислив значения  $f(x)$  во всех этих точках (или только определив знаки  $f(x_i)$ ), сравнивают их в соседних точках, т.е. проверяют, не выполняется ли на отрезке  $[x_{i-1}; x_i]$  условие  $f(x_{i-1})f(x_i) \leq 0$ . Если заранее известно количество корней в исследуемой области, то, измельчая шаг поиска  $h$ , таким процессом можно либо все их

локализовать, либо довести процесс до состояния, позволяющего утверждать, что возможно наличие пар корней, не различимых с точностью  $h = \varepsilon$ . Этот хорошо приспособленный для вычислительных машин *способ перебора* является дорогим в смысле затрат на получение многочисленных пробных значений функции и не дает гарантий выявления количества и локализации всех корней в общем случае, что ограничивает сферу его применения.

Одна из проблем, в которую упирается решение задачи локализации корней, — это практическая невозможность точного вычисления значений функций. Она уже обсуждалась в главе 1. Вернемся к ней еще раз.

Из-за ограниченности разрядной сетки ЭВМ даже алгебраические функции вычисляются приближенно, вычисление же трансцендентных функций составляет тему отдельного разговора; ясно, что их значения за редкими исключениями по определению могут записываться лишь приближенно. Погрешности, с которыми вычисляются значения функции, порождают блуждания (иначе, флуктуации, случайные отклонения) ординат графика около средних значений. Если рассматривать небольшой участок графика какой-либо функции, построенной путем изображения всех ее точек, вычисленных на компьютере при дискретном изменении аргумента с очень мелким шагом  $h$ , с последующим соединением этих точек отрезками прямых, то окажется, что этот реальный график имеет пилообразную форму (рис. 5.2).

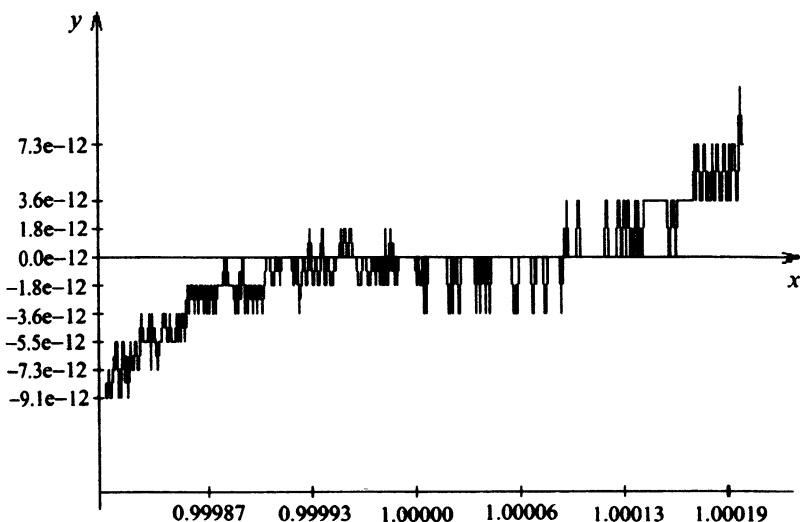


Рис. 5.2. Участок графика функции  $y = x^3 - 3x^2 + 3x - 1$ , полученный на компьютере DX5-133 вычислением ее значений с точностью  $O(10^{-12})$  с шагом  $O(10^{-6})$

Эти вычислительные погрешности ограничивают точность, с которой можно находить корень. В результате их влияния та или иная процедура сужения промежутка локализации корня рано или поздно приведет к такому промежутку, называемому *промежутком (отрезком, интервалом) неопределенности корня*, что любое число из этого промежутка с одинаковым успехом может быть принято за корень уравнения (ср. с промежутком  $(\xi - \text{б.п. } \xi, \xi + \text{б.п. } \xi)$ ), который фигурировал в п.1.6 при рассмотрении затронутой проблемы в несколько иной постановке).

Наличие флуктуаций значений функций, порождающих невозможность точного вычисления корней нелинейных уравнений, следует всегда иметь в виду, применяя те или иные численные методы. Если не с самого начала итерационного процесса, то с приближением момента, когда погрешность метода становится сравнимой с погрешностью вычисления значений функций или выражений, входящих в расчетные формулы метода, целесообразно подключение приема Гарвика, упоминавшегося ранее (см. замечание 4.4).

## 5.2. МЕТОД ДИХОТОМИИ. МЕТОД ХОРД

Пусть функция  $f(x)$  определена и непрерывна при всех  $x \in [a; b]$  и на  $[a; b]$  меняет знак, т.е.  $f(a)f(b) < 0$ . Тогда согласно теореме 5.1 уравнение (5.1) имеет на  $(a; b)$  хотя бы один корень. Возьмем произвольную точку  $c \in (a; b)$ . Будем называть в этом случае отрезок  $[a; b]$  *промежутком существования корня*, а точку  $c$  – *пробной точкой*. Поскольку речь здесь идет лишь о вещественнозначных функциях вещественной переменной, то вычисление значения  $f(c)$  приведет к какой-либо одной из следующих взаимоисключающих ситуаций:

- а)  $f(a)f(c) < 0$ ; б)  $f(c)f(b) < 0$ ; в)  $f(c) = 0$ .

Применительно к рассматриваемой задаче их можно интерпретировать так<sup>\*)</sup>:

- а) корень находится на интервале  $(a; c)$ ;  
 б) корень находится на интервале  $(c; b)$ ;  
 в) точка  $c$  является искомым корнем.

Таким образом, одно вычисление значения функции позволяет уменьшить промежуток  $[a; b]$  существования корня (ситуация а) или б)) или указать его значение (ситуация в), маловероятная в смысле “прямого попадания” пробной точкой  $c$  в корень, но вполне реальная в смысле выполнения при-

<sup>\*)</sup> Из-за допустимости неединственности корня в этой интерпретации уже нет взаимоисключаемости ситуаций.

ближенного равенства  $f(c) \approx 0$ , когда длина промежутка существования корня близка к длине промежутка его неопределенности). Ясно, что в зависимости от того, имеет место ситуация а) или б), описанная процедура одного шага сужения промежутка существования нуля непрерывной функции  $f(x)$  может быть применена к промежутку  $[a; c] \subset [a; b]$  или к  $[c; b] \subset [a; b]$  соответственно и далее повторяться циклически. Такой простой и легко программируемый процесс называется *методом дихотомии* (от греческого слова, означающего деление на две части), *методом бисекции*, *методом вилки*, *методом проб*. Если способ задания пробных точек  $c$  определен так, что последовательность длин получающихся в этом процессе промежутков существования корня стремится к нулю, то методом дихотомии можно найти какой-либо корень уравнения (5.1) с наперед заданной точностью.

Наиболее употребительным частным случаем метода дихотомии является *метод половинного деления*, реализующий самый простой способ выбора пробной точки – деление промежутка существования корня пополам. Выполнить приближенное вычисление с точностью  $\varepsilon$  корня уравнения (5.1) методом половинного деления при условии, что  $f(x)$  непрерывна на  $[a; b]$  и  $f(a)f(b) < 0$ , можно, например, по следующей схеме:

Шаг 0. Задать концы отрезка  $a$  и  $b$ , функцию  $f$ , малое число  $\varepsilon > 0$  (допустимую абсолютную погрешность корня или полудлину его промежутка неопределенности), малое число  $\delta > 0$  (допуск, связанный с реальной точностью вычисления значений данной функции);  
вычислить (или ввести)  $f(a)$ .

Шаг 1. Вычислить  $c := 0.5(a + b)$ .

Шаг 2. Если  $b - a < 2\varepsilon$ , положить  $\xi := c$  ( $\xi$  – корень) и остановиться.

Шаг 3. Вычислить  $f(c)$ .

Шаг 4. Если  $f(c) < \delta$ , положить  $\xi := c$  и остановиться.

Шаг 5. Если  $f(a)f(c) < 0$ , положить  $b := c$  и вернуться к шагу 1;  
иначе положить  $a := c$ ,  $f(a) := f(c)$  и вернуться к шагу 1.

**Замечание 5.1.** В упрощенных вариантах схем реализации метода половинного деления обходится без введения допуска  $\delta$ . В таком случае в шаге 4 вместо неравенства  $|f(c)| < \delta$  используют равенство  $f(c) = 0$ . Тогда разветвление алгоритма, диктуемое шагами 4 и 5, можно производить сравнением с нулем (с тремя исходами) произведения  $f(a)f(c)$ .

За один шаг метода половинного деления промежутков существования корня сокращается ровно вдвое. Поэтому, если за  $k$ -е приближение этим методом к корню  $\xi$  уравнения (5.1) примем точку  $x_k$ , являющуюся серединой полученного на  $k$ -м шаге отрезка  $[a_k; b_k]$  в результате последовательного сужения данного отрезка  $[a; b]$ , полагая  $a_1 = a$ ,  $b_1 = b$ , то придем к неравенству

$$|\xi - x_k| < \frac{b-a}{2^k} \quad \forall k \in N \quad (5.3)$$

(априори,  $\xi$  — любая точка интервала  $(a_k; b_k)$ ), и расстояние от нее до середины этого интервала не превосходит половины его длины. Это как раз и видим в (5.3) при  $k = 1$ ).

Неравенство (5.3), с одной стороны, позволяет утверждать, что последовательность  $(x_k)$  имеет предел — искомый корень  $\xi$  уравнения (5.1); с другой стороны, являясь априорной оценкой абсолютной погрешности приближенного равенства  $x_k \approx \xi$ , дает возможность подсчитать число шагов (итераций) метода половинного деления, достаточное для получения корня  $\xi$  с заданной точностью  $\varepsilon$ , для чего нужно лишь найти наименьшее натуральное  $k$ , удовлетворяющее неравенству

$$\frac{b-a}{2^k} < \varepsilon.$$

Используемый в методе половинного деления способ фиксирования пробной точки можно охарактеризовать как пассивный, ибо он осуществляется по заранее жестко заданному плану и никак не учитывает вычисляемые на каждом шаге значения функции. Логично предположить, что в семействе методов дихотомии можно достичь несколько лучших результатов, если отрезок  $[a; b]$  делить точкой  $c$  на части не пополам, а пропорционально величинам ординат  $f(a)$  и  $f(b)$  графика данной функции  $f(x)$ . Это означает, что точку  $c$  есть смысл находить как абсциссу точки пересечения оси  $Ox$  с прямой, проходящей через точки  $A(a; f(a))$  и  $B(b; f(b))$ , иначе, с хордой  $AB$  дуги  $A\xi B$  (рис. 5.3).

Запишем уравнение прямой, проходящей через две данные точки  $A$  и  $B$ :

$$\frac{y - f(a)}{f(b) - f(a)} = \frac{x - a}{b - a}.$$

Отсюда, полагая  $y = 0$  (уравнение оси  $Ox$ ),  $x = c$  (обозначение искомой точки пересечения прямой  $AB$  с осью  $Ox$ ), находим

$$c = a - \frac{f(a)(b-a)}{f(b) - f(a)}. \quad (5.4)$$

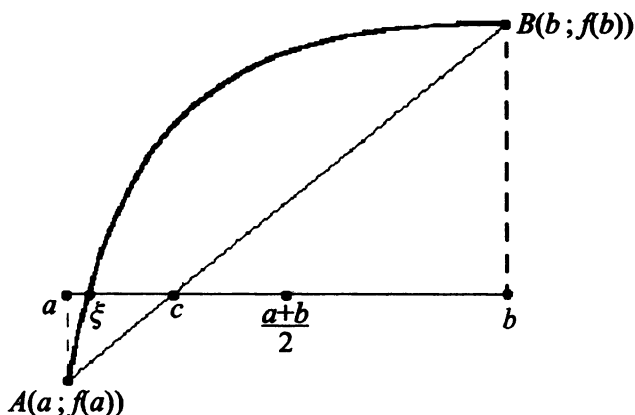


Рис. 5.3. Дуга графика функции  $f(x)$  и стягивающая ее хорда

Метод, получающийся из метода дихотомии таким фиксированием пробной точки, называют *методом хорд*, *методом пропорциональных частей*, *методом линейной интерполяции*. Все названия метода вполне естественны и отражают различные подходы к его выводу или интерпретации. Иногда (в последнее время реже) используют еще и название *правило ложного положения* или *regula falsi* [1, 23, 43].

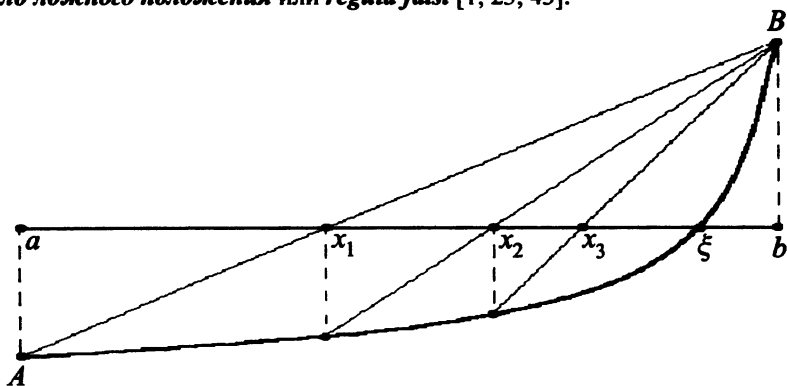


Рис. 5.4. Приближения к корню нелинейного уравнения по методу хорд

Существует несколько версий реализации метода хорд. Одна из них – подсчет значений  $c$  по формуле (5.4) в рамках алгоритма типа рассмотренного выше алгоритма половинного деления, где следует положить  $f(b) := f(c)$  при  $f(a)f(c) < 0$ . Длина промежутка локализации корня при этом может не стремиться к нулю, поэтому обычно счет ведется до совпа-

дения значений  $s$  на двух соседних итерациях с точностью  $\epsilon$  (лучше с точностью  $\frac{m\epsilon}{M-m}$ , если  $0 < m \leq |f'(x)| \leq M \quad \forall x \in [a; b]$ ).

Так как для линейной функции  $f(x)$  метод хорд дает корень  $\xi$  точно за один шаг при любой длине отрезка  $[a; b]$ , то можно рассчитывать на его довольно быструю сходимость, если  $f(x)$  близка к линейной. При определенных достаточно жестких условиях можно доказать соответствующие утверждения о монотонной сходимости, получить более точные оценки и упростить алгоритм (см., например, [21]). Однако в общем случае, если на функцию  $f(x)$  не накладывать дополнительных ограничений, может оказаться, что метод хорд будет проигрывать в скорости методу половинного деления; чтобы убедиться в этом, достаточно взглянуть на рис. 5.4, демонстрирующий возможное поведение нескольких приближений по методу хорд.

### 5.3. ТИПЫ СХОДИМОСТЕЙ ИТЕРАЦИОННЫХ ПОСЛЕДОВАТЕЛЬНОСТЕЙ

Чтобы более объективно судить о скорости сходимости тех или иных итерационных методов, введем следующие понятия [22 и др.].

Пусть некоторый итерационный процесс генерирует последовательность  $(x_k)_{k=0}^{\infty}$ , имеющую пределом  $x^*$ .

**Определение 5.1.** Сходимость последовательности  $(x_k)$  к  $x^*$  называется *линейной* (соответственно, *итерационный процесс – линейно сходящимся*), если существует такая постоянная  $C \in (0; 1)$  и такой номер  $k_0$ , что

$$|x^* - x_{k+1}| \leq C |x^* - x_k| \quad \forall k \geq k_0, \quad (5.5)$$

и *сверхлинейной*, если существует такая положительная последовательность  $(C_k)_{k=0}^{\infty}$ , что  $C_k \rightarrow 0$  и

$$|x^* - x_{k+1}| \leq C_k |x^* - x_k| \quad \forall k \in N_0. \quad (5.6)$$

**Определение 5.2.** Говорят, что последовательность  $(x_k)$  сходится к  $x^*$  по меньшей мере с  $p$ -м порядком (соответственно, итерационный процесс имеет по меньшей мере  $p$ -й порядок), если найдутся такие константы  $C > 0$  и  $p \geq 1$ , что



$$|x^* - x_{k+1}| \leq C |x^* - x_k|^p \quad (5.7)$$

при всех  $k \in N_0$ , начиная с некоторого  $k = k_0$ .

Фиксируя в определении 5.2 значение  $p = 1$ , видим, что линейно сходящийся процесс можно называть *процессом первого порядка*; значению  $p = 2$  в (5.7) соответствует *квадратично сходящийся процесс*,  $p = 3$  означает *кубическую сходимость*.\*)

К линейной сходимости применяют также термин *сходимость со скоростью геометрической прогрессии*. Объяснение ему можно найти в том, что определяющее линейную сходимость неравенство (5.5) между абсолютными погрешностями  $(k+1)$ -го и  $k$ -го приближений к предельной точке  $x^*$  означает существование последовательности положительных чисел  $\varepsilon_k$ , мажорирующих эти погрешности и связанных соотношением  $\varepsilon_{k+1} = C\varepsilon_k$ , т.е. являющихся членами геометрической прогрессии со знаменателем  $C = \varepsilon_{k+1}/\varepsilon_k$  ( $=\text{const}$ ). Отсюда следует также естественность в определении 5.1 условия  $C < 1$ , чтобы последовательность погрешностей была убывающей, иначе и речи не может быть о сходимости (в определении 5.2 для предельного случая  $p = 1$  также следует ограничить  $C$  единицей; при  $p > 1$  в этом, вообще говоря, нет необходимости; проанализируйте, почему?).

Если требуемой в неравенстве (5.5) константы  $C$  не удастся найти, но установлено неравенство (5.6) с  $C_k \rightarrow C \in (0; 1)$ , то в этом случае говорят об *асимптотически линейной сходимости* (пример такой сходимости дает степенной метод в п.4.4 или рассматриваемый далее в п.5.10 метод Бернулли). Аналогично можно определить *асимптотически  $p$ -й порядок*.

Имеются и другие понятия и термины, позволяющие более тонко классифицировать сходимость итерационных последовательностей и рассматривать ее как бы под разными углами зрения (см., например, [42, 43, 50, 52]). Как правило, они вводятся в более общем случае конечномерных или бесконечномерных нормированных пространств. Ничто не мешает и данные здесь определения распространить на многомерный случай, достаточно лишь элементы последовательности  $(x_k)$  и предельный элемент  $x^*$  считать  $n$ -мерными векторами ( $n \geq 1$ ) или матрицами, а вместо модуля использовать норму.

---

\*) В качестве известных примеров методов высоких порядков в пространстве  $R_n$  следует обратить внимание на изученный в п.3.6 класс итерационных процессов обращения матриц.

Среди нескольких способов охарактеризовать скорость сходимости итерационных последовательностей наиболее четко оформились два способа: опирающиеся на *q-сходимость* (от англ. quotient – частное) и на *r-сходимость* (от англ. root – корень). Происхождение этих терминов можно связать соответственно с признаками Даламбера (через отношение) и Коши (через арифметический корень), применяемыми для установления абсолютной сходимости ряда

$$x_0 + (x_1 - x_0) + (x_2 - x_1) + \dots + (x_k - x_{k-1}) + \dots \quad (5.8)$$

что равносильно установлению сходимости данной последовательности  $(x_k)$ , так как сходимость ряда (5.8) означает существование предела его частичных сумм

$$S_1 = x_0, \quad S_2 = x_1, \quad S_3 = x_2, \quad \dots, \quad S_{k+1} = x_k, \dots$$

Четкие определения и различия между этими двумя типами сходимостей можно найти в [42]: более существенно эти различия проявляются в многомерном случае. Здесь же главное понимать, что представление о порядке сходимости того или иного метода важно для возможности сравнить его с другими: более точно для этого подходит знание порядка *q-сходимости* (что и определено выше), порядок же *r-сходимости* говорит лишь о наличии такой последовательности положительных чисел, которая, сходясь к нулю с этим порядком, мажорирует последовательность величин  $|x^* - x_k|$ . Не вникая в тонкости, можно сказать, что обычно, изучая итерационный метод, устанавливают факт сходимости итерационной последовательности  $(x_k)$  к искомому элементу  $x^*$  и получают апостериорные и априорные оценки погрешности, а о порядке метода (в том или ином смысле, не всегда уточняя, в каком) судят или на основе неравенства типа (5.7), или по априорной оценке погрешности вида

$$|x^* - x_k| \leq C v^p, \quad (5.9)$$

где  $C > 0$ ,  $v \in (0;1)$  – некоторые константы, а  $p \geq 1$  – порядок метода, или по неравенству вида

$$|x_{k+1} - x_k| \leq C |x_k - x_{k-1}|^p, \quad (5.10)$$

показывающему скорость сближения членов итерационной последовательности и являющемуся ключевым для установления сходимости и получения оценок погрешностей. Чаще всего разные способы приводят к одному и тому же значению  $p$ , хотя это и не гарантировано.

Коснемся еще одного аспекта понятия сходимости итерационного метода. В приведенных выше определениях отождествлялись сходимость итерационного процесса и сходимость итерационной последовательности, порождаемой этим процессом; при этом негласно считалось, что последовательность  $(x_k)$  уже как бы фиксирована заданием начальной точки  $x_0$ .

Большой же интерес представляет сходимость множества всевозможных итерационных последовательностей, генерируемых итерационным методом при варьировании  $x_0$  в границах некоторой области. Итерационные методы, дающие в пределе решение данной задачи при любом начальном приближении  $x_0$ , называются *глобально сходящимися*. Если же сходимость итерационной последовательности  $(x_k)$  к искомому элементу  $x^*$  имеет место лишь при задании  $x_0$  из некоторой, вообще говоря, достаточно малой окрестности  $x^*$ , то соответствующий итерационный метод называют *локально сходящимся*.

Так, рассмотренные выше методы дихотомии можно отнести к глобально сходящимся методам, так как с их помощью всегда можно получить какой-нибудь из корней уравнения  $f(x) = 0$ , если начать итерационный процесс с отрезка  $[a; b]$  любой длины, входящего в область непрерывности  $f(x)$ , лишь бы было выполнено условие  $f(a)f(b) < 0$ . При этом, как видно из оценки (5.3), имеющей форму (5.9), метод половинного деления нужно считать линейно сходящимся методом, т.е. он сходится со скоростью геометрической прогрессии со знаменателем  $q = 0.5$  и имеет [56] *среднюю скорость сходимости*  $-\ln q = \ln 2$ . Метод хорд также является методом первого порядка и в зависимости от свойств  $f(x)$  может иметь как большую, так и меньшую, чем  $\ln 2$ , среднюю скорость (часто большую).

## 5.4. МЕТОД НЬЮТОНА

Одним из популярнейших итерационных методов решения нелинейных уравнений, что связано с его идейной простотой и быстрой сходимостью, является *метод Ньютона*<sup>1)</sup>. Правило построения итерационной последовательности  $(x_k)$  здесь получают или из геометрических соображений, откуда другое название этого метода – *метод касательных*, говорящее само за себя, или из аналитических путем подмены данной нелинейной функции ее линейной моделью на основе формулы конечных приращений Лагранжа или формулы Тейлора, в связи с чем метод Ньютона также называют *методом линеаризации*. В любом случае говорить о нахождении нуля функции  $f(x)$  методом Ньютона можно лишь в предположении, что данная функция обладает достаточной гладкостью.

Для простоты будем считать, что функция  $f(x)$  дважды дифференцируема на отрезке  $[a; b]$ , содержащем корень  $\xi$  уравнения (5.1).

---

<sup>1)</sup> В зарубежной литературе его часто называют *методом Ньютона-Рафсона* [22. 38. 43. 50]

Пусть  $x_k \in [a; b]$  – уже известный член последовательности приближений к  $\xi$ , полученный конструируемым методом (или заданное начальное приближение  $x_0$  при  $k=0$ ). Для любого  $x$  из  $[a; b]$  можно записать формальное представление  $f(x)$  по формуле Тейлора

$$f(x) = f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2}f''(\Theta_k)(x - x_k)^2, \quad (5.11)$$

где  $\Theta_k \in [a; b]$  – некоторая точка между  $x$  и  $x_k$ .

Так как корень  $\xi$  – потенциально произвольная точка отрезка  $[a; b]$ , то разложение (5.11) справедливо и для  $x = \xi$ , т.е. существует точка  $\Theta_k = \overline{\Theta}_k$  такая, что

$$f(\xi) = f(x_k) + f'(x_k)(\xi - x_k) + \frac{1}{2}f''(\overline{\Theta}_k)(\xi - x_k)^2.$$

Но  $f(\xi) = 0$ , и если точка  $\overline{\Theta}_k$  известна, то корень  $\xi$  можно точно найти из квадратного уравнения

$$f(x_k) + f'(x_k)(\xi - x_k) + \frac{1}{2}f''(\overline{\Theta}_k)(\xi - x_k)^2 = 0. \quad (5.12)$$

Считая, что значение  $x_k$  близко к  $\xi$ , т.е. разность  $(\xi - x_k)$  по модулю достаточно мала, можно рассчитывать, что величина  $(\xi - x_k)^2$  будет тем более малой. На этом основании отбросим в (5.12) последнее слагаемое и подменим квадратное уравнение (5.12) линейным уравнением. Естественно, что при этом будет найден не корень  $\xi$ , а некоторая другая точка, которую обозначим  $x_{k+1}$ .

Таким образом, итерационный процесс Ньютона определяется линейным уравнением

$$f(x_k) + f'(x_k)(x_{k+1} - x_k) = 0 \quad (5.13)$$

или в явном виде формулой

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad (5.14)$$

где  $k=0, 1, 2, \dots$ , и предполагается, что по крайней мере на элементах последовательности  $(x_k)$  первая производная данной функции в нуль не обращается\*).

Если в равенстве (5.13) фиксированную точку  $x_{k+1}$  заменить переменной  $x$ , а 0 в правой части – переменной  $y$ , то в полученном легко уз-

\* Как видим, процесс вычислений по методу Ньютона не требует знания второй производной. Можно обойтись без нее и при выводе (другим способом), но при этом значительно усложнилось бы изучение метода.

нать уравнение касательной к кривой  $y = f(x)$ , проведенной к ней в точке  $(x_k, f(x_k))$ . Отсюда геометрический смысл метода Ньютона: приближения к корню  $\xi$  совершаются по абсциссам точек пересечения касательных к графику данной функции, проводимых в точках, соответствующих предыдущим приближениям (рис. 5.5).

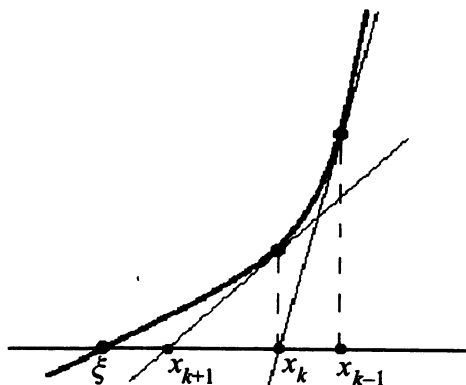


Рис. 5.5. Приближения к корню нелинейного уравнения методом касательных

Изучение сходимости метода Ньютона (осуществимость процесса вычисления элементов последовательности  $(x_k)$  по формуле (5.14) в пределах заданного отрезка, сходимость к корню  $\xi$  данного уравнения, порядок метода, оценки погрешности и критерии окончания процесса построения приближений, условия на выбор начального элемента последовательности) проводится при более ограничительных требованиях к данной функции  $f(x)$ .

Интуитивно ясно (из вида формулы (5.14), рассуждений при ее выводе и из ее геометрического смысла), что сходимость  $(x_k)$  к  $\xi$  будет тем быстрее и говорить о ней можно тем увереннее, чем ближе функция  $f(x)$  к линейной и чем круче ее график пересекает ось абсцисс; так что есть смысл потребовать от  $f(x)$ , чтобы по модулю вторая ее производная была ограничена сверху, а первая снизу. При этих условиях обратимся сначала к исследованию быстроты сходимости итерационного метода Ньютона (5.14) в предположении, что факт его осуществимости и сходимости к корню  $\xi$  сомнений не вызывает.

**Теорема 5.4.** Пусть функция  $f(x)$  удовлетворяет условиям

$$(A) \quad \begin{cases} |f'(x)| \geq \alpha > 0 \\ |f''(x)| \leq \beta < \infty \end{cases} \quad \forall x \in [a, b].$$

Тогда если члены последовательности  $(x_k)$ , определяемые методом Ньютона (5.14), при любом фиксированном  $k \in N_0$  принадлежат отрезку  $[a; b]$  и эта последовательность сходится на  $[a; b]$  к корню  $\xi$  уравнения (5.1), то справедливы неравенства ( $\forall k \in N_0$ ):

$$|\xi - x_{k+1}| \leq \frac{\beta}{2\alpha} |\xi - x_k|^2; \quad (5.15)$$

$$|\xi - x_{k+1}| \leq \frac{\beta}{2\alpha} |x_{k+1} - x_k|^2. \quad (5.16)$$

(Первое из этих неравенств в соответствии с определением 5.2 позволяет считать (5.14) *методом второго порядка*, а второе, являясь апостериорной оценкой погрешности, может служить в качестве критерия останова процесса вычислений).

**Доказательство.** Отметим, что требование принадлежности точек  $x_k$  отрезку  $[a; b]$  дает право пользоваться оговоренным условиями (A) поведением данной функции в этих точках и их малых окрестностях.

Подставляя в правую часть формулы (5.12) вместо нуля левую часть равенства (5.13) (оба равенства истинны для рассматриваемых точек  $x_k, x_{k+1}$ ), имеем:

$$f'(x_k)\xi - f'(x_k)x_k + \frac{1}{2}f''(\bar{\Theta}_k)(\xi - x_k)^2 = f'(x_k)x_{k+1} - f'(x_k)x_k.$$

Это равенство можно записать в виде точной связи между ошибками  $k$ -го и  $(k+1)$ -го приближений<sup>\*)</sup>:

$$\xi - x_{k+1} = -\frac{f''(\bar{\Theta}_k)}{2f'(x_k)}(\xi - x_k)^2,$$

из которой, переходя к модулям и привлекая условия (A), получаем первое из доказываемых неравенств.

Для доказательства второго неравенства сначала установим связь между невязкой  $(k+1)$ -го приближения и разностью соседних  $(k$ -го и  $(k+1)$ -го) приближений. С этой целью в формулу Тейлора (5.11) подставим  $x = x_{k+1}$ . Имеем

<sup>\*)</sup> Для элемента  $x_{k+1}$  последовательности  $(x_k)$  приближений к корню  $\xi$  уравнения  $f(x) = 0$  величину  $\xi - x_{k+1}$  называют *ошибкой*,  $x_{k+1} - x_k$  — *поправкой*, а  $f(x_{k+1})$  — *невязкой*.

$$f(x_{k+1}) = f(x_k) + f'(x_k)(x_{k+1} - x_k) + \frac{1}{2}f''(\tilde{\Theta}_k)(x_{k+1} - x_k)^2,$$

и воспользуемся тем, что согласно (5.13), первые два слагаемых правой части равенства в совокупности дают нуль. Таким образом,

$$f(x_{k+1}) = \frac{1}{2}f''(\tilde{\Theta}_k)(x_{k+1} - x_k)^2,$$

и значит,

$$|f(x_{k+1})| \leq \frac{\beta}{2}|x_{k+1} - x_k|^2. \quad (5.17)$$

Применим теперь формулу Лагранжа к разности  $f(\xi) - f(x_{k+1})$ . Согласно этой формуле, между точками  $\xi$  и  $x_{k+1}$  найдется точка  $\tau_{k+1}$  такая, что

$$f(\xi) - f(x_{k+1}) = f'(\tau_{k+1})(\xi - x_{k+1}).$$

Принимая во внимание, что  $f(\xi) = 0$ , а также условия (A), получаем неравенство между абсолютными величинами невязок и ошибок приближений:

$$|\xi - x_{k+1}| \leq \frac{1}{\alpha}|f(x_{k+1})|. \quad (5.18)$$

Усилив последнее с помощью неравенства (5.17), приходим к доказываемой оценке (5.16).

Теорема полностью доказана.

**Замечание 5.2.** Неравенство (5.18), справедливое для приближений  $x_{k+1}$  к нулю  $\xi$  дифференцируемой функции  $f(x)$  независимо от способа их получения, обосновывает контроль точности по невязкам, а именно, оправдывает критерий окончания итерационного процесса

$$|f(x_k)| < \varepsilon \Rightarrow \xi := x_k \text{ с точностью } \varepsilon,$$

если в окрестности  $\xi$  выполняется требование  $|f'(x)| > 1$ , и говорит о необходимости учитывать множитель  $\frac{1}{\alpha}$ , т.е. добиваться выполнения неравенства

$$|f(x_k)| < \alpha \varepsilon, \quad \text{если } \alpha < 1.$$

Чтобы выяснить, при каких условиях на выбор начального приближения  $x_0$  начинающаяся с него последовательность  $(x_k)$ , генерируемая методом Ньютона (5.14), будет сходиться к корню  $\xi$ , проитерируем формально неравенство (5.15). Имеем:

$$\text{при } k = 0 \quad |\xi - x_1| \leq \frac{\beta}{2\alpha}|\xi - x_0|^2;$$

$$\text{при } k = 1 \quad |\xi - x_2| \leq \frac{\beta}{2\alpha} |\xi - x_1|^2 \leq \frac{\beta}{2\alpha} \left(\frac{\beta}{2\alpha}\right)^2 |\xi - x_0|^2{}^2;$$

$$\text{при } k = 2 \quad |\xi - x_3| \leq \frac{\beta}{2\alpha} |\xi - x_2|^2 \leq \frac{\beta}{2\alpha} \left(\frac{\beta}{2\alpha}\right)^2 \left(\frac{\beta}{2\alpha}\right)^2 |\xi - x_0|^2{}^3;$$

далее по индукции получаем

$$\begin{aligned} |\xi - x_k| &\leq \left(\frac{\beta}{2\alpha}\right)^{1+2+2^2+\dots+2^{k-1}} |\xi - x_0|^{2^k} = \\ &= \left(\frac{\beta}{2\alpha}\right)^{\frac{2^k-1}{2-1}} |\xi - x_0|^{2^k} = \frac{2\alpha}{\beta} \left(\frac{\beta}{2\alpha}\right)^{2^k} |\xi - x_0|^{2^k}, \end{aligned}$$

т.е. при любых  $k \in N$

$$|\xi - x_k| \leq \frac{2\alpha}{\beta} \left(\frac{\beta}{2\alpha}\right)^{2^k} |\xi - x_0|^{2^k}. \quad (5.19)$$

Отсюда следует, что

$$x_k \rightarrow \xi, \quad \text{если } \frac{\beta}{2\alpha} |\xi - x_0| < 1.$$

Таким образом, появилась возможность судить о том, насколько далеко от корня  $\xi$  можно брать начальное приближение  $x_0$  в зависимости от свойств данной функции  $f(x)$ , и по априорной оценке (5.19) заранее подсчитывать число итераций, достаточное для вычисления корня с заданной точностью, если есть оценка близости  $x_0$  к  $\xi$ .

**Теорема 5.5.** Пусть для функции  $f(x)$  на отрезке  $[a; b]$  выполнены условия (A).

Тогда если интервал  $J := \left(\xi - \frac{2\alpha}{\beta}; \xi + \frac{2\alpha}{\beta}\right)$  содержится в  $[a; b]$ ,

то при произвольном выборе  $x_0$  из  $J$  для определяемой методом Ньютона (5.14) последовательности  $(x_k)$ :

- 1)  $x_k \in J \quad \forall k \in N$ ;
- 2)  $\exists \lim_{k \rightarrow \infty} x_k = \xi$  и  $f(\xi) = 0$ ;
- 3) справедливо утверждение теоремы 5.4 и оценка (5.19).

**Доказательство.** Прежде всего заметим, что условие  $x_0 \in J$  равносильно неравенству

$$v := \frac{\beta}{2\alpha} |\xi - x_0| < 1, \quad (5.20)$$

к которому мы пришли выше, анализируя неравенство (5.19) (рис. 5.6).



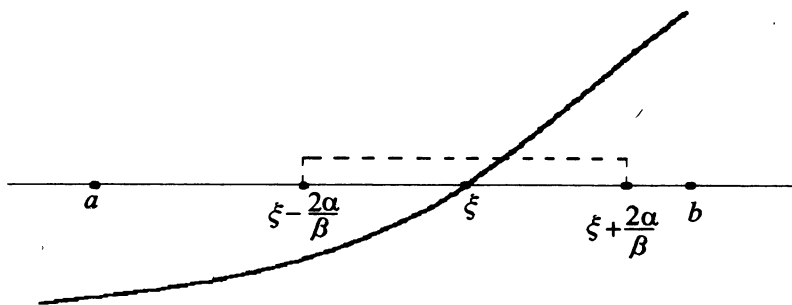


Рис. 5.6. Промежуток выбора начальной точки в методе Ньютона

Покажем осуществимость процесса (5.14) в  $J$ .

Так как  $x_0 \in J \subset [a; b]$ , то, согласно условиям (А),  $f'(x_0) \neq 0$ , и значение  $x_1$  по формуле (5.14) может быть получено. Более того, справедливо неравенство (5.15) при  $k = 0$ . Следовательно,

$$|\xi - x_1| \leq \frac{\beta}{2\alpha} |\xi - x_0|^2 = \frac{2\alpha}{\beta} v^2 < \frac{2\alpha}{\beta},$$

т.е.  $x_1 \in J$ . Аналогично из индукционного предположения, что  $x_k \in J$  при некотором  $k \in N$ , т.е. что  $|\xi - x_k| < \frac{2\alpha}{\beta}$ , используя (5.15), получаем:

$$|\xi - x_{k+1}| \leq \frac{\beta}{2\alpha} |\xi - x_k|^2 < \frac{\beta}{2\alpha} \cdot \left(\frac{2\alpha}{\beta}\right)^2 = \frac{2\alpha}{\beta} \Rightarrow x_{k+1} \in J.$$

Итак, согласно принципу математической индукции, все элементы последовательности  $(x_k)$  лежат в  $J \subset [a; b]$ , и значит, можно воспользоваться неравенством (5.19), из которого тут же следует, что  $\xi = \lim x_k$ , и заключение предыдущей теоремы.

То, что точка  $\xi$  – середина заявленного в теореме интервала  $J$  – есть корень уравнения (5.1), если заранее это неизвестно, устанавливается переходом к пределу в формуле (5.14) (с учетом условия  $f'(x) \neq 0 \forall x \in [a; b]$ ).

Слабым местом только что доказанной теоремы 5.5 является то, что точкой отсчета при построении промежутка применимости метода Ньютона служит корень  $\xi$ , который как раз и неизвестен. Более естественно за центр такого промежутка принимать некоторую конкретную точку  $x_0$  – начальное приближение, полагая по тем или иным соображениям (например, геометрическим), что в ее окрестности должен быть корень. Имеется

ряд теорем подобного типа (см. [34, 35, 40, 43, 47] и др.). Как правило, им присуща некоторая громоздкость и трудная проверяемость условий.

Заменяя условия (A) на требование знакопостоянства первой и второй производных данной функции, означающих монотонность и определенную выпуклость ее графика, докажем простую теорему несколько иного плана.

**Теорема 5.6** [21]. Пусть на отрезке  $[a; b]$  функция  $f(x)$  имеет первую и вторую производные постоянного знака и пусть

$$f(a)f(b) < 0.$$

Тогда если точка  $x_0$  выбрана на  $[a; b]$  так, что

$$f(x_0)f''(x_0) > 0, \quad (5.21)$$

то начатая с нее последовательность  $(x_k)$ , определяемая методом Ньютона (5.14), монотонно сходится к корню  $\xi \in (a; b)$  уравнения (5.1).

Доказательство опирается на теорему Вейерштрасса о сходимости монотонной ограниченной последовательности.

Положим, для определенности, что

$$f'(x) > 0 \quad \text{и} \quad f''(x) > 0 \quad \forall x \in [a; b].$$

Тогда  $f(a) < 0$ ,  $f(b) > 0$ , и в качестве начальной точки  $x_0$ , удовлетворяющей условию (5.21),<sup>\*)</sup> можно взять любую точку из промежутка  $(\xi; b]$  (наличие корня  $\xi$ , единственного в  $(a; b)$ , условиями теоремы обеспечено), и при этом  $x_0$  из  $(\xi; b]$  будет  $f(x_0) > 0$ .

Покажем ограниченность последовательности  $(x_k)$  снизу и ее монотонное убывание. Из равенства (5.12) при  $k = 0$ , т.е. из

$$f(x_0) + f'(x_0)(\xi - x_0) + \frac{1}{2}f''(\Theta_0)(\xi - x_0)^2 = 0,$$

в силу положительности последнего слагаемого, следует, что

$$f(x_0) + f'(x_0)(\xi - x_0) < 0.$$

Отсюда, учитывая положительность  $f'(x)$ , имеем

$$\xi < x_0 - \frac{f(x_0)}{f'(x_0)}.$$

Это неравенство показывает, что

$$\xi < x_1 = x_0 - \frac{f(x_0)}{f'(x_0)} < x_0$$

(значит,  $f(x_1) > 0$ ). По индукции устанавливаем, что

$$\xi < x_{k+1} < x_k \quad \forall k \in N_0,$$

<sup>\*)</sup> Называемому иногда *условием Фурье*.

т.е. последовательность  $(x_k)$  монотонно убывает и ограничена снизу самим корнем  $\xi$ , следовательно, имеет предел. Сходимость  $(x_k)$  именно к корню, как и в предыдущей теореме, получаем переходом к пределу в равенстве (5.14).

Остальные комбинации знаков производных рассматриваются аналогично. Теорема доказана.

**Замечание 5.3.** Нарушение условия Фурье (5.21) на выбор начального приближения  $x_0$  при выполненных требованиях к знакопостоянству производных может отразиться лишь на первом приближении: его перебросит через корень  $\xi$  на другую часть отрезка  $[a; b]$  (рис. 5.7). Если при этом  $x_1$  не окажется за пределами отрезка  $[a; b]$ , то далее итерационный процесс Ньютона пойдет монотонно (сформулируйте самостоятельно аналог теоремы 5.6 без условия Фурье).

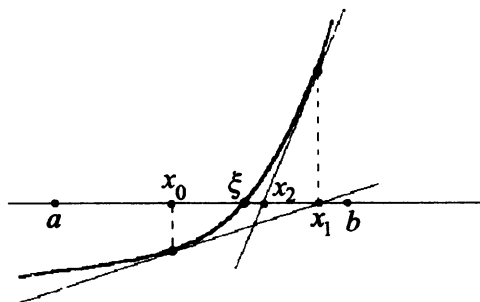


Рис.5.7. Поведение первых приближений по методу Ньютона при нарушении условия Фурье на выбор  $x_0$

Обратимся вновь к зафиксированным теоремами 5.4 и 5.5 квадратичной сходимости метода Ньютона, которая имеет место и в условиях теоремы 5.6. О поведении такого процесса можно получить наглядное представление из оценки (5.19), схематично записанной в виде

$$|\xi - x_k| \leq C \cdot v^{2^k},$$

где  $v \in (0; 1)$  (см. (5.20)), а  $C \left( = \frac{2\alpha}{\beta} \right)$  — некоторая положительная постоянная.

Допустим, что функция  $f(x)$  и приближение  $x_0$  к корню  $\xi$  таковы, что  $C = 1$ , а  $v = 0.1$ , т.е. абсолютные погрешности убывают по закону  $|\xi - x_k| \leq 0.1^{2^k}$ . Подставляя сюда  $k = 1, 2, 3, \dots$ , имеем:  $|\xi - x_1| \leq 10^{-2}$ ,  $|\xi - x_2| \leq 10^{-4}$ ,  $|\xi - x_3| \leq 10^{-8}$ ,  $|\xi - x_4| \leq 10^{-16}$  и т.д.

Как видим, *квадратично сходящийся процесс в идеале*, т.е. если он реализуется точно, *должен давать удвоение числа верных знаков на каждой итерации, начиная с некоторой*. Такой высокий темп установления верных цифр искомого корня не только позволяет получить корень с большой точностью сравнительно небольшим количеством вычислений, но и обеспечивает хорошую численную устойчивость метода, а также меньшую критичность к правилу окончания итерационного процесса (часто здесь применяют упрощенное правило остановки

$$|x_k - x_{k-1}| \leq \varepsilon \Rightarrow \xi \approx x_k).$$

## 5.5. ПРИМЕНЕНИЕ МЕТОДА НЬЮТОНА К ВЫЧИСЛЕНИЮ ЗНАЧЕНИЙ ФУНКЦИЙ

Элементарные функции чаще всего вычисляются с помощью приближения их подходящими многочленами. В некоторых же случаях для этих целей применяют итерационные методы, в частности, базирующиеся на методе Ньютона.

Пусть требуется найти значение заданной функции  $\varphi$  в заданной точке  $a$ . Считая  $a$  произвольной точкой из области  $D(\varphi)$  или какой-либо ее подобласти, функциональное соответствие

$$x = \varphi(a)$$

зададим неявно уравнением

$$F(a, x) = 0 \quad (5.22)$$

таким, чтобы: 1) оно было локально эквивалентным (в окрестности точки  $a$ ) данному; 2) функция  $F$  была дифференцируема по второму аргументу; 3) функции  $F$  и  $F'$  были легко вычислимы.

При каждом фиксированном  $a$  уравнение (5.22) можно считать уравнением типа (5.1) и получать приближенно его корень – требуемое значение  $x = \varphi(a)$  – методом Ньютона (5.14). Для уравнения (5.22) формула (5.14) принимает вид

$$x_{k+1} = x_k - \frac{F(a, x_k)}{F'_x(a, x_k)}, \quad (5.23)$$

где  $k = 0, 1, 2, \dots$ , а  $x_0$  – задаваемое начальное приближение к  $\varphi(a)$ .

В качестве более конкретного примера применения такого подхода выведем из формулы (5.23) *правило Ньютона вычисления арифметических корней*.

Пусть  $a$  – данное положительное число, а  $n \geq 2$  – данный натуральный показатель корня.

Очевидна связь между задачей вычисления вещественного значения

$$x = \sqrt[n]{a}$$

и задачей нахождения положительного корня уравнения

$$x^n - a = 0.$$

Приняв  $F(a, x) := x^n - a$ , находим  $F'_x(a, x) = nx^{n-1}$ , и согласно (5.23), процесс приближений к  $\sqrt[n]{a}$  определяем формулой

$$x_{k+1} = x_k - \frac{x_k^n - a}{nx_k^{n-1}}$$

или в другом виде

$$x_{k+1} = \frac{1}{n} \left[ (n-1)x_k + \frac{a}{x_k^{n-1}} \right],$$

где  $k = 0, 1, 2, \dots$ , а  $x_0 > 0$  задается.

Еще из глубокой древности известен частный случай полученного правила Ньютона – *процесс Герона*

$$x_{k+1} = \frac{1}{2} \left( x_k + \frac{a}{x_k} \right),$$

применяемый для извлечения квадратных корней. Здесь  $F(a, x) = x^2 - a$ ,  $F'(a, x) = 2x$ ,  $F''(a, x) = 2$ , т.е. при  $x > 0$  выполняются предварительные условия теоремы 6, и для монотонной сходимости  $(x_k)$  к  $\sqrt{a}$  достаточно лишь удовлетворить условию Фурье (5.21), т.е. взять  $x_0$  таким, чтобы было  $x_0^2 > a$ . Легко убедиться, что это условие будет выполнено, если взять  $x_0 = 2^{0.5m}$ , если  $m$  – четное, и  $x_0 = 2^{0.5(m+1)}$ , если  $m$  – нечетное, где  $m \in Z$  такое, что  $a = 2^m q$  и  $q \in [0.5; 1]$  ( $a$  “зажимается” между двумя соседними целыми степенями двойки:  $2^{m-1} \leq a < 2^m$ ).

Другим примером использования метода Ньютона в форме (5.23) может служить вычисление обратной величины данного числа  $a$ , иначе говоря, выполнение операции деления с помощью других арифметических операций.

---

<sup>\*)</sup> Не заботясь о монотонности последовательности приближений, обычно ограничиваются заданием  $x_0 = 2^{[0.5m]}$ .

Аналогично предыдущему зададим искомое  $x = \frac{1}{a}$  как корень уравнения

$$a - \frac{1}{x} = 0.$$

Подставляя  $F(a, x) = a - \frac{1}{x}$  и  $F'_x(a, x) = \frac{1}{x^2}$  в (5.23), после упрощения получаем итерационный процесс без делений\*

$$x_{k+1} = x_k(2 - ax_k), \quad k = 0, 1, 2, \dots \quad (5.24)$$

Здесь также считаем  $a$  и  $x$  положительными, и очевидно, что для любых  $x > 0$

$$F' > 0, \text{ а } F'' = -\frac{2}{x^3} < 0,$$

т.е. условия теоремы 5.6 будут выполнены, если взять  $x_0 \in \left(0, \frac{1}{a}\right)$ .

Непосредственное изучение процесса (5.24) показывает, что область выбора начального приближения  $x_0$  может быть в два раза расширена (правда монотонность приближений при этом уже не гарантируется).

Действительно, рассмотрим связь между поправками  $(k+1)$ -го и  $k$ -го приближений:

$$\frac{1}{a} - x_{k+1} = \frac{1}{a} - 2x_k + ax_k^2 = a\left(\frac{1}{a} - x_k\right)^2$$

(равенство типа (5.7), подтверждающее, что (5.24) – процесс второго порядка). Отсюда последовательным итерированием приходим к равенству

$$\frac{1}{a} - x_k = \frac{1}{a}(1 - ax_0)^{2^k},$$

из которого следует, что сходимость  $(x_k)$  к  $\frac{1}{a}$  имеет место в случае, когда

$$|1 - ax_0| < 1, \text{ т.е. при } x_0 \in \left(0, \frac{2}{a}\right).$$

Обычно за начальное приближение, удовлетворяющее условию  $x_0 \in \left(0, \frac{1}{a}\right)$ , берут число  $2^{-m}$ , где  $m \in \mathbb{Z}$  то же, что и в предыдущем примере.

\* Легко увидеть аналогию между этим процессом и рассмотренным ранее процессом Шульца второго порядка для обращения матриц (ср (5.24) с (3.34) при  $m = 1$ ).

## 5.6. МОДИФИКАЦИИ МЕТОДА НЬЮТОНА. МЕТОД СЕКУЩИХ

Вновь обратимся к теоремам 5.4–5.6 о сходимости метода Ньютона. Их условия предполагают неравенство нулю производной данной функции  $f(x)$  на промежутке  $[a; b]$ , где применяется метод. А это означает, что они регламентируют применение метода Ньютона только для нахождения простых нулей функции  $f(x)$ , поскольку для кратного корня  $\xi$  уравнения (5.1) имеет место равенство  $f'(\xi) = 0$ . Действительно, пусть  $\xi$  –  $m$ -кратный корень ( $m \geq 2$ ); тогда функция  $f(x)$  представима<sup>\*)</sup> в виде  $f(x) = (x - \xi)^m f_1(x)$  и ее производная  $f'(x) = (x - \xi)^{m-1} \times [m f_1(x) + (x - \xi) f_1'(x)]$  обращается в нуль при  $x = \xi$ .

Согласно (5.14), формально нужно, чтобы производная не равнялась нулю в точках  $x_k$  последовательности приближений, в предельной же точке  $\xi$  допустимо обращение производной в нуль. Как показывает пример 5.3, итерационный процесс Ньютона может сходиться и в этом случае, т.е. когда  $\xi$  является кратным корнем уравнения (5.1), но сходимостью при этом – только линейная.

Если заведомо известно число  $m$  – показатель кратности корня  $\xi$ , то для ускорения сходимости метода Ньютона в формулу (5.14) рекомендуется ввести корректирующий множитель  $m$ :

$$x_{k+1} = x_k - m \cdot \frac{f(x_k)}{f'(x_k)}. \quad (5.25)$$

Такую модификацию будем называть *методом Ньютона–Шрёдера*<sup>\*\*)</sup>. Доказательство сверхлинейной сходимости этого метода можно найти в [47], где он называется *методом Ньютона с параметром*, а также в [50] и в [43], где имеется ссылка на содержащую этот метод работу Э.Шрёдера.

**Пример 5.3.** Для функции  $f(x) = (x - 1)^2$  корень  $\xi = 1$  – двукратный. Подстановка этой  $f(x)$  и  $f'(x) = 2(x - 1)$  в формулу (5.14) определяет процесс

$$x_{k+1} = x_k - \frac{x_k - 1}{2}, \quad (5.26)$$

<sup>\*)</sup> Это представление здесь принимается за определение  $m$ -й кратности корня  $\xi$ . Часто  $\xi$  считают корнем кратности  $m$ , если

$$f(\xi) = f'(\xi) = \dots = f^{(m-1)}(\xi) = 0, \text{ а } f^{(m)}(\xi) \neq 0 \text{ (см., например, [47]).}$$

<sup>\*\*)</sup> В [23] первое введение параметра  $m$  в формулу (5.14) приписывается Е.Бодевигу (1949 г.), в связи с чем (5.25) там называют *методом Ньютона–Бодевига*.

или проще.

$$x_{k+1} = \frac{1}{2}(x_k + 1).$$

Если начать его с  $x_0 = 2$ , то на каждой последующей итерации будем получать все более близкие к  $\xi = 1$  значения:

$$x_1 = 1.5, \quad x_2 = 1.25, \quad x_3 = 1.125, \quad x_4 = 1.0625, \quad \dots$$

Вычитая 1 из обеих частей (5.26), приходим к равенству

$$x_{k+1} - 1 = \frac{1}{2}(x_k - 1),$$

означающему, что сходимость  $(x_k)$  к 1 – точно линейная. В то же время введение множителя  $m=2$  в (5.26) в соответствии с (5.25) приводит к стационарной последовательности  $x_{k+1} = 1 \quad \forall k \in N_0$ .

Не следует думать, что из (5.25) всегда будет получаться  $x_{k+1} = \xi$ . Рассмотрим менее утрированный пример.

Функция  $f(x) = x(x-1)^2$  с тем же двукратным корнем  $\xi = 1$  с помощью формул (5.14) и (5.25) порождает следующие процессы:

*метод Ньютона*

$$x_{k+1} = \frac{2x_k^2}{3x_k - 1};$$

$$x_0 = 2,$$

$$x_1 = 1.6,$$

$$x_2 = 1.347368,$$

$$x_3 = 1.193517,$$

...

*метод Ньютона–Шрёдера*

$$x_{k+1} = \frac{x_k(x_k + 1)}{3x_k - 1};$$

$$x_0 = 2,$$

$$x_1 = 1.2,$$

$$x_2 = 1.000116,$$

$$x_3 = 1.000000.$$

Налицо эффективность коррекции метода Ньютона введением в него показателя кратности корня.

Цель всех последующих видоизменений основной формулы (5.14) метода Ньютона – уменьшение вычислительных затрат, связанных с необходимостью вычисления производной на каждом итерационном шаге.

Самый простой выход на этом пути – использование на каждом шаге одного и того же шагового множителя  $\frac{1}{f'(x_0)}$ , т.е. счет по формуле

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_0)}, \quad k = 0, 1, 2, \dots \quad (5.27)$$



Такой метод называют *модифицированным* или *упрощенным методом Ньютона*<sup>\*)</sup>. Он имеет очевидную геометрическую интерпретацию: в начальной точке  $x_0$  проводится касательная к графику  $y = f(x)$  (первый шаг основного и модифицированного методов Ньютона совпадают), а во всех последующих точках  $x_1, x_2, \dots$  проводятся прямые, параллельные этой касательной (рис. 5.8).

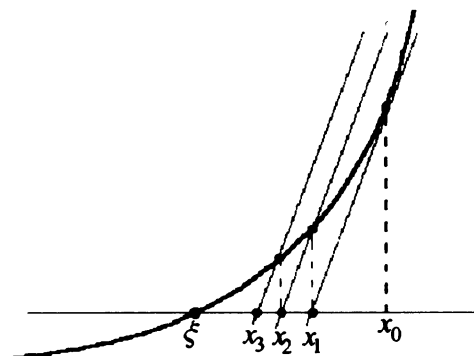


Рис. 5.8. Приближения к корню нелинейного уравнения модифицированным (упрощенным) методом Ньютона

При такой модификации метод Ньютона утрачивает высокую скорость сходимости (процесс не реагирует на изменение наклона кривой к оси абсцисс при приближении к корню) и вместо квадратичной имеет лишь скорость сходимости геометрической прогрессии, что станет очевидным несколько позже (см. п. 5.7).

На получение сверхлинейной скорости сходимости при видоизменении метода Ньютона (5.14) можно надеяться в случае, когда  $f'(x_k)$  при каждом  $k \in N$  подменяется не одним и тем же числом  $f'(x_0)$ , а некоторым близким к  $f'(x_k)$  значением, которое может быть найдено (при каждом  $k$  свое) через значения данной функции. Для таких аппроксимаций<sup>\*\*)</sup>  $f'(x_k)$  можно использовать, например, определение производной. Имеем:

$$f'(x_k) = \lim_{h \rightarrow 0} \frac{f(x_k + h) - f(x_k)}{h},$$

и при малых  $h$  (произвольного знака) получаем приближенное равенство

$$f'(x_k) \approx \frac{f(x_k + h) - f(x_k)}{h}, \quad (5.28)$$

<sup>\*)</sup> Иногда его называют еще *огрубленным методом Ньютона* [47].

<sup>\*\*)</sup> Approximare (лат.) – приближаться

позволяющее производную приближенно подменять так называемым *разностным отношением*. Подстановка (5.28) в (5.14) приводит к итерационной формуле

$$x_{k+1} = x_k - \frac{f(x_k) \cdot h}{f(x_k + h) - f(x_k)}, \quad (5.29)$$

где  $k=0, 1, 2, \dots$ , а  $h$  – малый параметр, которым должен распорядиться вычислитель.

Ясно, что при каждом  $k$  в формуле (5.28) может быть свое значение  $h$ , т.е. в формуле (5.29) вместо постоянного параметра  $h$  имеет смысл использовать связанный с номером итерации параметр  $h_k$ , т.е. вести вычисления по формуле

$$x_{k+1} = x_k - \frac{f(x_k)h_k}{f(x_k + h_k) - f(x_k)}. \quad (5.30)$$

Итерационный метод, определяемый формулами (5.29) или (5.30), назовем *разностным методом Ньютона\**.

Так как равенство (5.28) можно сделать сколь угодно точным за счет малости шага  $h$  разностного отношения (теоретически; практически это далеко не так из-за потерь точности при вычитании близких чисел), то по непрерывности можно утверждать асимптотически квадратичную скорость сходимости разностного метода Ньютона при определенных условиях.

Рассмотрим соображения, которыми следует руководствоваться при задании последовательности параметров  $h_k$  в разностном методе Ньютона (5.30). В любом случае будем исходить из постулата, что с ростом  $k$  значения  $|h_k|$  должны убывать, чтобы при приближении  $x_k$  к корню  $\xi$  производная  $f'(x_k)$  все более точно аппроксимировалась разностным отношением  $[f(x_k + h_k) - f(x_k)]/h_k$ .

Первое, что можно здесь предложить, так это задать какое-либо значение  $h_0$ , а каждое последующее значение параметра получать рекуррентным равенством  $h_{k+1} = \delta h_k$ , где  $\delta \in (0; 1)$  – некоторое фиксированное число. Например, можно положить  $h_0 = 0.1$ ,  $h_1 = 0.01$ ,  $h_2 = 0.001$  и т.д. Очевиден недостаток такого подхода – отсутствие связи между скоростью сходимости  $(x_k)$  к  $\xi$  и скоростью убывания  $|h_k|$  (может оказаться, что  $x_k$  еще не имеет достаточной близости к  $\xi$ , а значение  $|h_k|$  настолько мало, что значения  $f(x_k + h_k)$  и  $f(x_k)$  реально не различимы; противоположная ситуация чревата большой потерей скорости сходимости или, еще хуже, нарушением канонического развития итерационного процесса).

\* Другие названия – *конечно-разностный* [22] и *дискретный* [42] *метод Ньютона*. Сходимость этого метода изучается в [22], там же можно найти рекомендации по сопряжению шага дискретизации  $h$  с точностью машинных вычислений

Если учесть, что при зафиксированных в теореме 5.4 условиях (А)  $f(x_k) \rightarrow 0$  с той же скоростью, что и  $x_k \rightarrow \xi$  (см. неравенство (5.18)), есть смысл полагать в (5.30)  $h_k := f(x_k)$ . Разумеется, это можно делать на той стадии итерационного процесса, когда значения  $|f(x_k)|$  уже достаточно малы (иначе теряет силу (5.28)). При таких  $h_k$  формула (5.30) принимает вид

$$x_{k+1} = x_k - \frac{(f(x_k))^2}{f(x_k + f(x_k)) - f(x_k)} \quad (5.31)$$

и называется *методом Стеффенсена*. Подчеркнем еще раз его сугубо локальный характер сходимости, но зато *сходимость* эта *квадратичная* [43, 47].

При приблизительно равных затратах на вычисление значений данной функции и ее производной ни один из рассмотренных выше вариантов разностного метода Ньютона не дает выигрыша по сравнению с основным методом (5.14), поскольку каждый из них требует два вычисления функции на каждом итерационном шаге, не увеличивая при этом скорость сходимости. Построим такую модификацию метода Ньютона, развивая далее его разностный аналог (5.30), в которой на один шаг итерации приходилось бы только одно вычисление функции.

Опираясь на то, что необходимым условием сходимости некоторой последовательности  $x_k$  к пределу  $\xi$ , как это следует из (5.8), является сходимость к нулю последовательности разностей  $(x_k - x_{k-1})$  (причем с той же скоростью, см. (5.10)), положим в (5.30)

$$h_k := x_{k-1} - x_k, \text{ откуда } x_{k-1} = x_k + h_k.$$

В результате этого из (5.30) получаем итерационный процесс

$$x_{k+1} = x_k - \frac{f(x_k)(x_{k-1} - x_k)}{f(x_{k-1}) - f(x_k)}, \quad (5.32)$$

где  $k = 1, 2, 3, \dots$ , а  $x_0$  и  $x_1$  должны задаваться.

Формула (5.32) определяет новый метод как *двухшаговый* (результат  $(k+1)$ -го шага зависит от результатов  $k$ -го и  $(k-1)$ -го шагов) и на каждой итерации требует вычисления только одного значения функции, другое же значение, фигурирующее в этой формуле, передается с предыдущего шага. Сравнив (5.32) с формулой (5.4), полученной из геометрических соображений, легко понять, что  $x_{k+1}$  есть абсцисса точки пересечения с осью  $Ox$  прямой, проведенной через точки  $(x_{k-1}, f(x_{k-1}))$  и  $(x_k, f(x_k))$ , т.е. секущей (рис. 5.9). Отсюда название этого метода — *метод секущих*<sup>\*)</sup>.

<sup>\*)</sup> В [5] методом секущих называют метод хорд (с фиксированным концом).

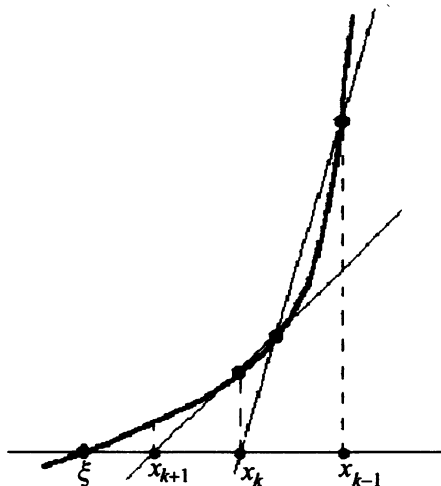


Рис. 5.9. Приближения к корню методом секущих

Можно сказать, что метод секущих и метод хорд определяются совершенно одноклассными формулами, но порождающие их идеологии различны, что сказывается на свойствах и скорости сходимости генерируемых ими последовательностей приближений.

Выясним, по какому закону убывают погрешности приближений, получаемых методом секущих (5.32) в условиях (А) теоремы 5.4.

Вычитая равенство (5.32) из равенства  $\xi = \xi$  и применяя в знаменателе формулу Лагранжа, имеем:

$$\xi - x_{k+1} = \xi - x_k + \frac{f(x_k)(x_{k-1} - x_k)}{f(x_{k-1}) - f(x_k)} = \xi - x_k + \frac{f(x_k)}{f'(\mu_k)}, \quad (5.33)$$

где  $\mu_k$  — некоторая точка между приближениями  $x_k$  и  $x_{k-1}$ . Далее воспользуемся формулой Тейлора, согласно которой можно записать

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

откуда при  $x = x_k$ ,  $\theta = \theta_k$  с учетом того, что  $\xi$  — корень уравнения (5.1), следует представление

$$f(x_k) = f'(\xi)(x_k - \xi) + \frac{1}{2}f''(\theta_k)(x_k - \xi)^2.$$

Подставим это в (5.33) и вынесем в правой части общий множитель  $\xi - x_k$ :

$$\xi - x_{k+1} = \frac{\xi - x_k}{f'(\mu_k)} \cdot \left[ f'(\mu_k) - f'(\xi) + \frac{1}{2}f''(\theta_k)(\xi - x_k) \right].$$

Вновь применяя формулу Лагранжа – теперь к разности производных – и переходя к модулям, получаем неравенство

$$|\xi - x_{k+1}| \leq \frac{|\xi - x_k|}{|f'(\mu_k)|} \cdot \left[ |f''(\nu_k)| \cdot |\xi - \mu_k| + \frac{1}{2} |f''(\theta_k)| |\xi - x_k| \right].$$

Так как для сходящейся к  $\xi$  нестационарной последовательности  $(x_k)$  справедливы неравенства

$$|\xi - \mu_k| < |\xi - x_{k-1}|, \quad |\xi - x_k| < |\xi - x_{k-1}|,$$

то в итоге можно записать искомую связь погрешностей в виде

$$|\xi - x_{k+1}| < \frac{3\beta}{2\alpha} |\xi - x_k| \cdot |\xi - x_{k-1}|, \quad (5.34)$$

где  $\alpha$  и  $\beta$  – константы из условий (А).

Ясно, что неравенство (5.34) характеризует *метод секущих* как *сверхлинейно сходящийся процесс*. Конкретный порядок метода секущих устанавливается следующим образом\*).

Обозначим для краткости

$$\varepsilon_k := |\xi - x_k|, \quad C := \frac{3\beta}{2\alpha}$$

и, используя записанное в этих обозначениях неравенство (5.34)

$$\varepsilon_{k+1} < C \varepsilon_k \varepsilon_{k-1},$$

получаем последовательно несколько первых оценок  $\varepsilon_k$  через степени  $\varepsilon_0$  (полагая по определению  $\varepsilon_1 < \varepsilon_0$ ). Имеем:

$$\text{при } k=1 \quad \varepsilon_2 < C \varepsilon_1 \varepsilon_0 < \frac{1}{C} (C \varepsilon_0)^2;$$

$$\text{при } k=2 \quad \varepsilon_3 < C \varepsilon_2 \varepsilon_1 < \frac{1}{C} (C \varepsilon_0)^3;$$

$$\text{при } k=3 \quad \varepsilon_4 < C \varepsilon_3 \varepsilon_2 < \frac{1}{C} (C \varepsilon_0)^5;$$

$$\text{при } k=4 \quad \varepsilon_5 < C \varepsilon_4 \varepsilon_3 < \frac{1}{C} (C \varepsilon_0)^8;$$

$$\text{при } k=5 \quad \varepsilon_6 < C \varepsilon_5 \varepsilon_4 < \frac{1}{C} (C \varepsilon_0)^{13} \quad \text{и т.д.}$$

Обратив внимание на то, что показатели в правых частях этих неравенств подчиняются закону "каждый последующий есть сумма двух предыдущих", нетрудно доказать в общем виде, что

$$\varepsilon_k < \frac{1}{C} (C \varepsilon_0)^{\Phi_k}, \quad (5.35)$$

\* Близкое к данному изложению можно найти в [43].

где  $(\Phi_k)$  – последовательность *чисел Фибоначчи*<sup>\*)</sup>, определяемая рекуррентным соотношением

$$\Phi_{k+1} = \Phi_{k-1} + \Phi_k, \quad k=1, 2, 3, \dots, \quad \Phi_0 = \Phi_1 = 1.$$

Для общего члена  $\Phi_k$  этой последовательности известна *формула Бинэ*:

$$\Phi_k = \frac{1}{\sqrt{5}} \left[ \left( \frac{1+\sqrt{5}}{2} \right)^{k+1} - \left( -\frac{1+\sqrt{5}}{2} \right)^{-(k+1)} \right].$$

Поскольку с ростом  $k$  роль второго члена в выражении  $\Phi_k$  становится ничтожной, для достаточно больших значений  $k$  на основе (5.35) можно считать справедливым асимптотическое неравенство

$$|\xi - x_k| < C_1 \cdot v \left( \frac{1+\sqrt{5}}{2} \right)^k,$$

где  $C_1$  и  $v$  – некоторые новые постоянные (которые легко выписать). Полученная оценка вида (5.9) позволяет утверждать, что *метод секущих имеет порядок по крайней мере*  $\frac{1+\sqrt{5}}{2} \approx 1.618$ .

**Замечание 5.4.** Для многошаговых итерационных методов иногда вводят специфическое понятие порядка сходимости. Согласно [22], *j-шаговый итерационный метод, генерирующий сходящуюся к  $\xi$  последовательность  $(x_k)$ , называется j-шагово сходящимся с порядком  $p$ , если*

$|\xi - x_{k+j}| \leq C |\xi - x_k|^p$ . Так как из неравенства (5.34) после его усиления получается

$$|\xi - x_{k+1}| < \frac{3\beta}{2\alpha} |\xi - x_{k-1}|^2,$$

то можно сказать, что *метод секущих сходится двухшагово квадратично*.

Высокий порядок скорости сходимости метода секущих в сочетании с минимальными вычислительными затратами – одно вычисление значения функции на один итерационный шаг – выводит этот метод на первое место по эффективности решения скалярных уравнений вида (5.1) среди прочих итерационных методов. Это подтверждается как теоретическими, так и практическими наблюдениями<sup>\*\*)</sup>.

При применении метода секущих возникают вопросы, связанные с началом итерационного процесса и с его окончанием. Поскольку касатель-

<sup>\*)</sup> Фибоначчи – введением знаменитого итальянского математика Леонардо Пизанского (1180–1240 гг.).

<sup>\*\*)</sup> При сравнении с методом Ньютона и другими методами, использующими производные, в [43] предлагается приравнять работу по вычислению значений функций и ее производных. Единица такой работы в [43] называется *горнером*, а в [50] – *единицей объема информационного запроса*.

ная к кривой есть предельное положение секущей, выбор начальной точки  $x_0$  в методе секущих нужно осуществлять по тому же принципу, что и в методе касательных, например, привлекая условие Фурье (5.21); вторая же из начальных точек  $x_1$ , требуемая в двухшаговом методе (5.32), может быть взята в непосредственной близости от  $x_0$  (понятие близости здесь, разумеется, условно), желательнее между точкой  $x_0$  и искомым корнем  $\xi$ .

Окончание счета по методу секущих, учитывая его быструю сходимость, можно контролировать с помощью проверок на малость модулей невязок  $|f(x_k)|$  или поправок  $|x_k - x_{k-1}|$ . Однако главное здесь — это суметь вовремя остановить процесс вычислений, не дожидаясь момента, когда погрешности вычислений начнут превосходить погрешность метода вследствие вычитания приближенно вычисляемых близких значений  $f(x_{k-1})$  и  $f(x_k)$  в знаменателе расчетной формулы (5.32). В этом плане, т.е. в численной устойчивости, метод секущих проигрывает методу Ньютона.

Если почти все рассмотренные выше методы можно отнести к классу методов линеаризации, имея ввиду, что в их основе лежит подмена исходной нелинейной модели линейной, построенной тем или иным способом, то следующим шагом должно быть построение классов методов "параболизации". Квадратичная модель (парабола) может быть получена, например, по формуле Тейлора или квадратичной интерполяцией, но в любом случае соответствующие итерационные формулы, хотя и дают более быстро сходящиеся последовательности приближений к корню, либо содержат старшие производные данной функции, либо являются слишком громоздкими и сложными как для исследования, так и для их применения.

Более важно обратить внимание на локальную сходимость таких простых и достаточно быстро сходящихся методов, как метод Ньютона и метод секущих, условия которой не так часто удается обеспечить. В связи с этим встает задача построения *гибридных алгоритмов* на базе двух или нескольких методов, соединяя быструю сходимость одних с глобальнойходимостью других. Принципы комбинирования методов в таких алгоритмах могут быть различными: можно "стартовать" с глобально сходящегося "медленного" метода и подключить быстросходящийся метод на финише для уточнения значения корня, а можно сразу начать процесс вычислений "быстрым" методом, но проводить корректировку получаемых им значений, пользуясь глобально сходящимся методом. Последний подход порождает, например, следующий простейший гибридный алгоритм.

#### **Метод Ньютона – метод половинного деления**

Шаг 0. Задать начальное приближение  $x_0$ , положить  $k := 0$ .

Шаг 1. Вычислить  $\tilde{x}_k = x_k - \frac{f(x_k)}{f'(x_k)}$ .

Шаг 2. Если  $|f(\tilde{x}_k)| < |f(x_k)|$ , то  $x_{k+1} := \tilde{x}_k$ ,

иначе  $\tilde{x}_k := \frac{1}{2}(x_k + \tilde{x}_k)$  и возвратиться к началу шага 2.

Шаг 3. Проверить на останов (работа алгоритма либо прекращается с  $\xi := x_{k+1}$ , либо продолжается переходом к шагу 1 с  $k := k+1$ ).

Приведенный алгоритм учитывает, что метод Ньютона выработывает локально правильное направление (убывания функции), но продвижение в этом направлении может оказаться чрезмерным, что и корректируется с помощью деления отрезка пополам, если не выполняется *условие релаксации*  $|f(\tilde{x}_k)| < |f(x_k)|$  в шаге 2 (рис. 5.10).

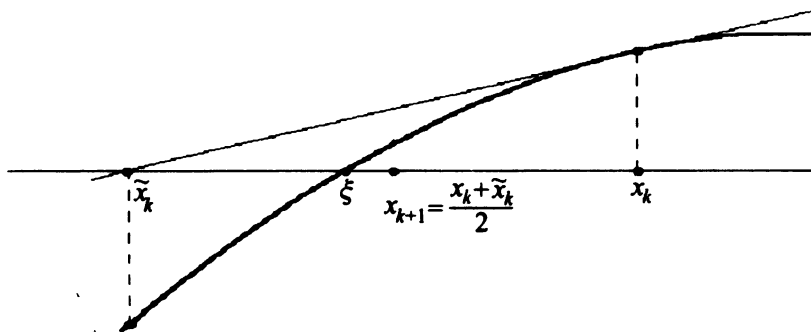


Рис. 5.10. Иллюстрация одного шага гибридного метода Ньютона – половинного деления

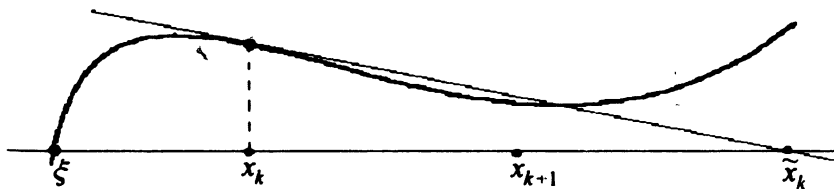


Рис. 5.11. Пример поведения гибридного метода Ньютона – половинного деления, когда  $x_{k+1}$  дальше от корня  $\xi$ , чем  $x_k$

Такой гибридный метод трудно считать глобально сходящимся (рассмотрите, например, его дальнейшее поведение в ситуации, изображенной на рис. 5.11), но он позволяет расширить границы применимости метода Ньютона и хоть как-то вести процесс поиска корня в условиях неопределенности знаков производных. Разумеется, этот алгоритм весьма схематичен.



чен и требует некоторых усилий на его детализовку. Особенно важно решить, как в тех или иных случаях выполнить шаг 3. Большую роль здесь играет правильное сопряжение задаваемой точности решения задачи с погрешностью метода (иначе, с остаточной погрешностью) и с точностью выполнения арифметических операций на используемой ЭВМ (вычислительной погрешностью), а также с точностью вычисления значений функций, что, вообще говоря, не одно и то же. Важные сведения о таких критериях окончания процессов поиска корней с привязкой их к реальным компьютерам можно почерпнуть в книге [22].

## 5.7. ЗАДАЧА О НЕПОДВИЖНОЙ ТОЧКЕ. МЕТОД ПРОСТЫХ ИТЕРАЦИЙ

Нельзя не заметить, что все расчетные формулы, определяющие уже изученные методы решения скалярных уравнений (5.1), такие как метод Ньютона (5.14) и его модификации (5.25), (5.27), (5.29), (5.31), имеют вид

$$x_{k+1} = \varphi(x_k), \quad k = 0, 1, 2, \dots, \quad (5.36)$$

где  $\varphi(x)$  – некоторая функция, для каждого метода своя, так или иначе связанная с исходной функцией  $f(x)$ . Попытаемся понять, каким требованиям должна удовлетворять функция  $\varphi(x)$ , чтобы последовательность  $(x_k)$ , определяемая этим самым общим одношаговым итерационным способом (5.36), называемым *методом простых итераций*<sup>\*)</sup>, была сходящейся, и как построить функцию  $\varphi(x)$  по функции  $f(x)$ , чтобы эта последовательность сходилась к корню данного уравнения (5.1).

Сразу отметим, что функцию  $\varphi(x)$  будем считать непрерывной в исследуемой области оси  $Ox$ . Поэтому, если определяемая формулой (5.36) последовательность  $(x_k)$  окажется сходящейся к некоторому числу  $\xi$ , то, переходя к пределу в равенстве (5.36), получаем

$$\xi = \varphi(\xi), \quad (5.37)$$

т.е.  $\xi = \lim_{k \rightarrow \infty} x_k$  – корень уравнения

$$x = \varphi(x). \quad (5.38)$$

Решение уравнений такого вида, наряду с (5.1), представляет самостоятельный интерес; нахождение их корней называется *задачей о неподвижной точке*. Это название связано с тем, что точка  $\xi$  при отображении с помощью  $\varphi$  из  $R_1$  в  $R_1$  остается на месте (если разумеется, таковая существует).

<sup>\*)</sup> Другие названия. *метод итераций, метод последовательных приближений*. Далее, как и в гл. 3, будем также использовать аббревиатуру МПИ

Существование и единственность корня уравнения (5.38) основывается на *принципе сжимающих отображений* или, иначе, *принципе неподвижной точки*.

**Определение 5.3.** *Непрерывная функция  $\varphi(x)$  называется сжимающей (функцией сжатия) на отрезке  $[a; b]$ , если:*

- 1)  $\varphi(x) \in [a; b] \quad \forall x \in [a; b]$ ;
- 2)  $\exists q \in (0; 1): |\varphi(x_2) - \varphi(x_1)| \leq q|x_2 - x_1| \quad \forall x_1, x_2 \in [a; b]$ .

Графическое толкование применения сжимающего отображения  $\varphi$  (как функции множества) к промежутку сжатия  $[a; b]$  предоставляет рис. 5.12.

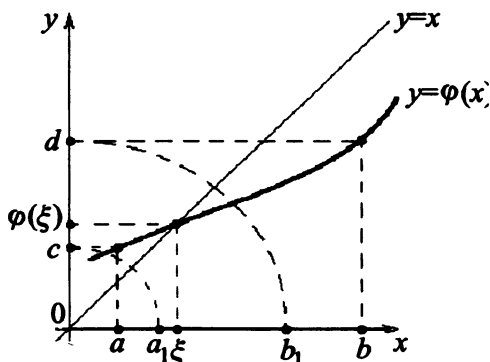


Рис. 5.12. Сжатие отрезка  $[a; b]$  возрастающей функцией  $\varphi(x)$

Как видно из этого рисунка, если  $[a; b]$  рассматривать как область определения функции сжатия  $\varphi(x)$ , то соответствующая ей область значений  $[\varphi(a); \varphi(b)]$  на оси ординат (отрезок  $[c; d]$ ), будучи перенесенным на ось абсцисс (отрезок  $[a_1; b_1]$ ), целиком содержится в  $[a; b]$ . Применяя к  $[a_1; b_1]$  те же рассуждения, что и к  $[a; b]$ , получим  $[a_2; b_2] \subset [a_1; b_1]$  и т.д. В итоге образуется бесконечная последовательность вложенных отрезков  $[a; b] \supset [a_1; b_1] \supset [a_2; b_2] \supset \dots \supset [a_k; b_k] \supset \dots$ ,

причем их длины убывают по закону

$$b_k - a_k \leq q^k (b - a) \rightarrow 0 \quad \text{как } k \rightarrow \infty.$$

Следовательно, в условиях сжатия эта последовательность имеет единственную общую точку  $(\xi)$ , которая переходит сама в себя, т.е. является неподвижной точкой отображения  $\varphi$ . При этом, очевидно, последова-

тельность  $(a_k)$  левых концов этих промежутков монотонно сходится к  $\xi$  слева, а последовательность  $(b_k)$  правых концов – справа. Так как при условии возрастания  $\varphi(x)$ , как это показано на рис. 5.12, в этом процессе

$$\begin{aligned} a_1 &= \varphi(a), & a_2 &= \varphi(a_1), & a_3 &= \varphi(a_2), & \dots, \\ b_1 &= \varphi(b), & b_2 &= \varphi(b_1), & b_3 &= \varphi(b_2), & \dots, \end{aligned}$$

то можно утверждать, что МПИ (5.36) будет давать монотонно сходящуюся к  $\xi$  последовательность, если ее начинать с  $x_0 = a$  или с  $x_0 = b$ . Так же монотонно возрастающая и монотонно убывающая последовательности приближений будут получаться по формуле (5.36) и в случаях, когда за  $x_0$  будет браться любая точка из промежутков  $[a, \xi]$  и  $(\xi, b]$  соответственно.

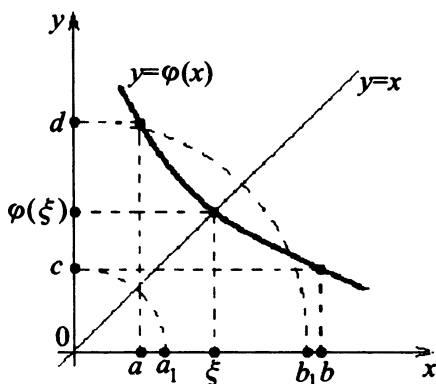


Рис. 5.13. Сжатие отрезка  $[a; b]$  убывающей функцией  $\varphi(x)$

При условии убывания сжимающей функции  $\varphi(x)$ , т.е. в случае, изображенном на рис. 5.13, начинающиеся с концов  $a$  и  $b$  промежутка сжатия последовательности выстраиваются следующим образом:

$$\begin{aligned} a, & \quad b_1 = \varphi(a), & a_2 &= \varphi(b_1), & b_3 &= \varphi(a_2), & \dots, \\ b, & \quad a_1 = \varphi(b), & b_2 &= \varphi(a_1), & a_3 &= \varphi(b_2), & \dots \end{aligned}$$

Каждая из них сходится к неподвижной точке  $\xi$ , и элементы каждой из этих последовательностей с удалением от начала дают все более хорошие приближения то с недостатком, то с избытком. Такую сходимость к  $\xi$  имеет и любая другая последовательность  $(x_k)$ , получаемая по формуле (5.36) при любом  $x_0 \in [a; b]$ . Отсюда другой термин, применяемый к неподвижной точке  $\xi$ , – **центр итерации** [43].

Более удобно иллюстрировать геометрически поведение итерационной последовательности  $(x_k)$ , определяемой МПИ (5.36), не отмечая значения  $\varphi(x_k)$  на оси ординат, а отражая их на ось абсцисс с помощью биссектрисы координатного угла  $y = x$ . Такие иллюстрации для случаев монотонного возрастания (ломаная типа "ступеньки") и монотонного убывания (ломаная типа "спираль") сжимающей функции  $\varphi(x)$  показаны на рис. 5.14 и 5.15 соответственно (обоснуйте!).

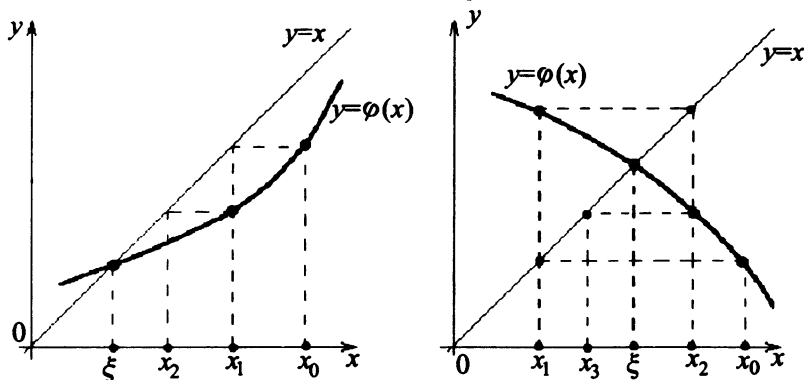


Рис. 5.14, 5.15. Монотонные и двусторонние приближения к корню методом простых итераций

Итогом проведенных выше рассуждений является следующий вывод:

если на некотором промежутке  $[a; b]$  функция  $\varphi(x)$  удовлетворяет условиям сжатия, зафиксированным определением 5.3, то: 1) уравнение (5.38) имеет и притом единственный корень  $\xi \in [a; b]$ ; 2) к этому корню со скоростью геометрической прогрессии сходится определяемая МПИ (5.36) последовательность  $(x_k)$ , начинающаяся с любого  $x_0 \in [a; b]$ , причем скорость сходимости тем выше, чем меньше коэффициент сжатия<sup>\*)</sup>  $q \in (0; 1)$ ; 3) если функция  $\varphi(x)$  монотонно возрастает на  $[a; b]$ , то приближения  $x_k$  к  $\xi$  также будут монотонными, если же  $\varphi(x)$  убывает, то процесс (5.36) порождает двусторонние приближения к корню  $\xi$ .

<sup>\*)</sup> Иначе, константа Липшица, в соответствии с применяемым к неравенству  $|\varphi(x_2) - \varphi(x_1)| \leq q|x_2 - x_1|$  термином условие Липшица (или Коши-Липшица, если известно  $q \in (0; 1)$ ).

Этот вывод имеет, в основном, качественный характер. Дальнейшие исследования будут направлены на выявление конструктивных требований к  $\varphi(x)$ , обеспечивающих для нее выполнение условий сжатия, и получение оценок близости генерируемых МПИ (5.36) приближений  $x_k$  к неподвижной точке  $\xi$ .

**Теорема 5.7.** Пусть функция  $\varphi(x)$  определена и дифференцируема на отрезке  $[a; b]$ . Тогда если выполняются условия:

- 1)  $\varphi(x) \in [a; b] \quad \forall x \in [a; b]$ ,
- 2)  $\exists q: |\varphi'(x)| \leq q < 1 \quad \forall x \in (a; b)$ ,

то уравнение (5.38) имеет и притом единственный на  $[a; b]$  корень  $\xi$ ; к этому корню сходится определяемая методом простых итераций (5.36) последовательность  $(x_k)$ , начинающаяся с любого  $x_0 \in [a; b]$ ; при этом справедливы оценки погрешности ( $\forall k \in N$ ):

$$|\xi - x_k| \leq \frac{q}{1-q} |x_k - x_{k-1}|, \quad (5.39)$$

$$|\xi - x_k| \leq \frac{q^k}{1-q} |x_1 - x_0|. \quad (5.40)$$

**Доказательство.** Взяв произвольное  $x_0$  из  $[a; b]$ , согласно условию 1) заключаем, что  $x_1 = \varphi(x_0)$  также принадлежит  $[a; b]$ . По индукции можно утверждать, что все члены последовательности  $(x_k)$ , генерируемой МПИ (5.36), благодаря условию 1), не могут покинуть отрезок  $[a; b]$ .

Вычтем из равенства (5.36) такое же равенство

$$x_k = \varphi(x_{k-1})$$

и к правой части полученного равенства

$$x_{k+1} - x_k = \varphi(x_k) - \varphi(x_{k-1})$$

применим формулу Лагранжа; согласно которой на интервале, определяемом точками  $x_{k-1}$  и  $x_k$  (а значит, на интервале  $(a; b)$ ), найдется точка  $\theta_k$  такая, что

$$\varphi(x_k) - \varphi(x_{k-1}) = \varphi'(\theta_k)(x_k - x_{k-1}).$$

Следовательно,

$$x_{k+1} - x_k = \varphi'(\theta_k)(x_k - x_{k-1})$$

и, в силу условия 2), справедливо неравенство

$$|x_{k+1} - x_k| \leq q |x_k - x_{k-1}|, \quad (5.41)$$

которое можно расценивать как выполнение главного условия сжатия на элементах итерационной последовательности  $(x_k)$  с коэффициентом сжатия  $q$ . С другой стороны, неравенство (5.41) указывает на факт и скорость сближения членов этой последовательности.

Для разностей соседних более удаленных от начала членов последовательности  $(x_k)$  на основе (5.41) получаем:

$$\begin{aligned} |x_{k+2} - x_{k+1}| &\leq q|x_{k+1} - x_k| \leq q^2|x_k - x_{k-1}|; \\ |x_{k+3} - x_{k+1}| &\leq q|x_{k+2} - x_{k+1}| \leq q^3|x_k - x_{k-1}|; \\ &\dots \end{aligned} \quad (5.42)$$

$$|x_{k+i} - x_{k+i-1}| \leq q^i|x_k - x_{k-1}| \quad \forall i \in N_0, \quad k \in N.$$

Используя эти неравенства, оценим близость между  $x_{k+m}$  (где  $m \in N$ ) и  $x_k$ , вычитая и прибавляя все промежуточные члены  $x_{k+m-1}, x_{k+m-2}, \dots, x_{k+1}$  и применяя свойство "модуль суммы не превосходит суммы модулей". Имеем:

$$\begin{aligned} |x_{k+m} - x_k| &\leq |x_{k+m} - x_{k+m-1}| + |x_{k+m-1} - x_{k+m-2}| + \dots + \\ &+ |x_{k+2} - x_{k+1}| + |x_{k+1} - x_k| \leq (q^m + q^{m+1} + \dots + q^2 + q)|x_k - x_{k-1}| = \\ &= \frac{q - q^{m+1}}{1 - q}|x_k - x_{k-1}|. \end{aligned}$$

Но, в свою очередь,

$$|x_k - x_{k-1}| \leq q^{k-1}|x_1 - x_0|, \quad (5.43)$$

что можно установить либо итерируя неравенство (5.41) при  $k=1, 2, \dots$ , либо непосредственно из (5.42), полагая  $k=1$ , а затем  $i = k - 1$ . Подставляя (5.43) в полученную выше оценку

$$|x_{k+m} - x_k| \leq \frac{q}{1 - q}(1 - q^m)|x_k - x_{k-1}|, \quad (5.44)$$

имеем

$$|x_{k+m} - x_k| \leq \frac{q^k}{1 - q}(1 - q^m)|x_1 - x_0|. \quad (5.45)$$

Поскольку правая часть неравенства (5.45) при фиксированном  $m \in N$  и  $k \rightarrow \infty$  стремится к нулю,  $(x_k)$  — последовательность Коши и, в силу замкнутости отрезка, имеет предел  $\xi \in [a; b]$ . Так как дифференцируемая функция непрерывна, то этот предел — корень уравнения (5.38). Его единственность на  $[a; b]$  доказывается от противного: предположив, что наряду с  $\xi$  есть другой корень  $\tau \in [a; b]$ , т.е. имеет место равенство  $\tau = \varphi(\tau)$ , а значит,  $\xi - \tau = \varphi(\xi) - \varphi(\tau)$ , по формуле Лагранжа получаем

$$\exists \theta \in (a; b) : \xi - \tau = \varphi'(\theta)(\xi - \tau);$$

последнее же равенство возможно лишь при  $\xi = \tau$ , поскольку по условию производная не может быть равна единице.

Перейдя к пределу в неравенствах (5.44) и (5.45) при  $m \rightarrow \infty$ , получаем оценки (5.39) и (5.40) соответственно, что и завершает доказательство теоремы.

Анализируя условия и доказательство теоремы 5.7, нельзя не заметить, что основным требованием к функции  $\varphi(x)$ , обеспечивающим сходимость МПИ (5.36) с оценками (5.39), (5.40), является условие малости модуля производной. Требование же  $\varphi(x) \in [a; b]$  нужно лишь постольку, поскольку оно должно обеспечить попадание значений  $\varphi(x)$  на промежутки, где выполняются другие требования. Очевидно, что достаточно его выполнения только на элементах итерационной последовательности  $(x_k)$  и практически это может проверяться непосредственно в процессе счета по формуле (5.36). Теоретически же, когда  $[a; b]$  — не вся числовая ось (тогда надобности в этом условии нет), имеется несколько возможностей заменить требование  $\varphi(x) \in [a; b]$  более конструктивным, т.е. априори проверяемым условием.

Будем исходить из того, что имеется некоторая точка  $x_0$ , которую можно взять в качестве начального приближения и в окрестности которой функция  $\varphi(x)$  дифференцируема и имеет малую по модулю производную. Тогда существование и единственность корня  $\xi$  уравнения (5.38) и сходимость к нему начатого с  $x_0$  процесса (5.36) можно связать с величиной  $r$  радиуса этой окрестности точки  $x_0$ .

**Теорема 5.8.** Пусть на отрезке  $[x_0 - r; x_0 + r]$  функция  $\varphi(x)$  определена, дифференцируема и ее производная удовлетворяет неравенству

$$|\varphi'(x)| \leq q < 1. \quad (5.46)$$

Тогда если величина  $r$  такова, что

$$|x_0 - \varphi(x_0)| \leq r(1 - q), \quad (5.47)$$

то на  $[x_0 - r; x_0 + r]$  имеется единственный корень  $\xi$  уравнения (5.38), и к нему сходится начатый с  $x_0$  метод простых итераций (5.36) с оценками погрешности (5.39), (5.40).

**Доказательство.** Так как  $1 - q < 1$ , то, в силу (5.47),  $x_1 = \varphi(x_0) \in [x_0 - r; x_0 + r]$ . Предположив, что  $x_k = \varphi(x_{k-1})$  принадлежит  $r$ -окрестности точки  $x_0$  при некотором  $k \in N$ , покажем, что там же будет и  $x_{k+1} = \varphi(x_k)$ . Действительно, поскольку, согласно предположению, на

отрезке  $[x_0 - x_k; x_0 + x_k] \subset [x_0 - r; x_0 + r]$  функция  $\varphi(x)$  определена и дифференцируема, то на этом отрезке найдется точка  $\mu_k$  такая, что

$$\begin{aligned} x_{k+1} - x_0 &= \varphi(x_k) - \varphi(x_0) + \varphi(x_0) - x_0 = \\ &= \varphi'(\mu_k)(x_k - x_0) + \varphi(x_0) - x_0. \end{aligned}$$

Отсюда, переходя к модулям, получаем неравенство

$$|x_{k+1} - x_0| \leq qr + r(1-q) = r,$$

означающее, что  $x_{k+1} \in [x_0 - r; x_0 + r]$ .

Таким образом, на элементах итерационной последовательности  $(x_k)$  выполняется первое требование теоремы 5.7 для отрезка  $[a; b]$  с  $a := x_0 - r$ ,  $b := x_0 + r$ , и вместе с требованием (5.46) оно обеспечивает справедливость заключения теоремы.

Полученные для МПИ простые оценки погрешности (5.39) и (5.40) можно использовать в практических вычислениях как для завершения итерационного процесса (5.36) по правилу:

$$|x_k - x_{k-1}| \leq \frac{1-q}{q} \varepsilon \Rightarrow \xi := x_k (\pm \varepsilon), \quad (5.48)$$

так и для предварительной прикидки числа итераций, достаточного для получения корня с заданной точностью  $\varepsilon$ :

$$\frac{q^k}{1-q} |x_1 - x_0| \leq \varepsilon \Leftrightarrow k \geq \frac{1}{\ln q} \ln \frac{\varepsilon(1-q)}{|x_0 - \varphi(x_0)|}$$

(или  $k \geq \frac{\ln \varepsilon - \ln r}{\ln q}$  в условиях теоремы 5.8).

**Замечание 5.5.** Легко видеть, что часто применяемый на практике простой критерий окончания процесса итераций (5.36) по выполнении неравенства  $|x_k - x_{k-1}| \leq \varepsilon$  обоснован в двух случаях: когда  $-1 < \varphi'(x_k) \leq 0$ , т.е. МПИ дает двусторонние приближения (корень  $\xi$  всегда "зажат" между любыми двумя соседними приближениями, расстояние между которыми служит эффективной оценкой погрешности, см. рис. 5.15), и когда  $0 \leq \varphi'(x_k) \leq q < 1$ , но при этом  $q \leq \frac{1}{2}$  (тогда  $\frac{1-q}{q} \geq 1$  и (5.48) будет заведомо выполнено). Использование этого упрощенного критерия окончания при значениях  $q$ , близких к концам интервала  $(0;1)$ , чревато либо лишними итерациями, либо недоитерированием.

Обратимся к связи между уравнением (5.1) и задачей о неподвижной точке – уравнением (5.38). Переписав (5.38) в виде

$$x - \varphi(x) = 0, \quad (5.49)$$



можно сказать, что это есть уравнение (5.1) с  $f(x) := x - \varphi(x)$ , и применять к (5.49) все рассмотренные в предыдущих пунктах рассуждения и методы.

Приведение уравнения (5.1) к виду (5.38) можно осуществлять множеством способов, но при этом всегда следует помнить, что это приведение нужно выполнять так, чтобы полученное уравнение соответствующего вида было не только эквивалентным (5.1), но и пригодным для проведения итераций, т.е. удовлетворяло оговоренным в теоремах 5.7, 5.8 условиям. Кроме того, попутно могут учитываться такие требования к получающемуся при этом методу итераций, как простота расчетной формулы, быстрота сходимости (малость  $q$ ), характер сходимости (монотонность или двусторонность приближений). Если уравнение (5.1) имеет несколько корней, то для нахождения каждого из них формируется своя задача о неподвижной точке.

Иногда преобразование уравнения (5.1) к виду (5.38) не вызывает больших затруднений и выполняется непосредственно. Например, уравнение

$$4 - 2x - \sin x = 0$$

достаточно записать в виде

$$x = 2 - 0.5 \sin x,$$

чтобы сказать, что оно имеет и притом единственное в  $R$  решение, к которому сходится при любом  $x_0 \in R$  последовательность

$$x_{k+1} = 2 - 0.5 \sin x_k$$

со скоростью геометрической прогрессии со знаменателем  $q \leq 0.5$ , поскольку функция  $\varphi(x) = 2 - 0.5 \sin x$  имеет производную  $\varphi'(x) = -0.5 \cos x$ , абсолютная величина которой не превосходит 0.5 при любых  $x \in R$ .

В общем случае переход от (5.1) к (5.38) осуществляют так: умножают левую и правую части уравнения (5.1) на отличный от нуля параметр  $-\lambda$  и к обеим частям прибавляют по  $x$ ; в результате получается равносильное (5.1) уравнение

$$x = x - \lambda f(x), \quad (5.50)$$

которое имеет вид (5.38) с  $\varphi(x) := x - \lambda f(x)$ . Далее параметр  $\lambda$  подбирается таким, чтобы производная  $\varphi'(x) = 1 - \lambda f'(x)$  в нужной области была малой по модулю (а если надо, то чтобы еще имела определенный знак).

Конкретные рекомендации по фиксированию  $\lambda$  в (5.50) могут быть даны в случае, когда, например, известны оценки сверху и снизу для производной исходной функции  $f(x)$ . А именно, пусть

$$0 < \alpha \leq f'(x) \leq \gamma < \infty$$

(если производная  $f'(x)$  отрицательна, можно заменить уравнение  $f(x) = 0$  на уравнение  $-f(x) = 0$ , т.е. работать с функцией  $-f(x)$ ). Тогда, соответственно,

$$1 - \lambda\gamma \leq \varphi'(x) \leq 1 - \lambda\alpha, \quad (5.51)$$

и значит,

$$|\varphi'(x)| \leq q(\lambda) := \max\{|1 - \lambda\alpha|, |1 - \lambda\gamma|\}.$$

Анализируя двойное неравенство (5.51), можно увидеть, что при любых  $\lambda \in \left(0; \frac{2}{\gamma}\right)$  будет  $q(\lambda) < 1$ . В частности, при  $\lambda = \frac{1}{\gamma}$  имеет место неравенство

$$0 \leq \varphi'(x) \leq 1 - \frac{\alpha}{\gamma} < 1,$$

обеспечивающее монотонную сходимость соответствующего МПИ со скоростью, определяемой оценками (5.39), (5.40) при  $q = 1 - \frac{\alpha}{\gamma}$ . Оптимальным

же значением  $\lambda$  является  $\lambda = \lambda_0 := \frac{2}{\alpha + \gamma}$ . При этом значении  $\lambda$  границы неравенства (5.51) таковы:

$$1 - \lambda\gamma = \frac{\alpha - \gamma}{\alpha + \gamma}, \quad 1 - \lambda\alpha = \frac{\gamma - \alpha}{\alpha + \gamma},$$

т.е. максимум  $|\varphi'(x)|$ , равный  $q(\lambda_0) = \frac{\gamma - \alpha}{\alpha + \gamma}$ , достигается на каждом из элементов двухэлементного множества  $\{|1 - \lambda\alpha|, |1 - \lambda\gamma|\}$ .

В отсутствие нужных оценок для  $f'(x)$  можно предложить следующие рассуждения. Если известно, что искомый корень находится в окрестности заданной точки  $x_0$ , где производная меняется не очень быстро, возьмем  $\lambda$  таким, чтобы  $\varphi'(x_0) = 0$ . Тогда по непрерывности производная должна остаться малой в окрестности  $x_0$ , т.е. можно рассчитывать на сходимость получающегося при этом итерационного процесса. Имеем:

$$1 - \lambda f'(x_0) = 0 \Rightarrow \lambda = \frac{1}{f'(x_0)}.$$

Подставляя это  $\lambda$  в (5.50), получаем уравнение

$$x = x - \frac{f(x)}{f'(x_0)},$$

а соответствующий ему МПИ (5.36) имеет вид

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_0)},$$

в котором узнаем *модифицированный метод Ньютона*<sup>\*)</sup> (5.27).

## 5.8. УСКОРЕНИЕ СХОДИМОСТИ ПОСЛЕДОВАТЕЛЬНЫХ ПРИБЛИЖЕНИЙ

Как явствует из предыдущего пункта, МПИ (5.36) имеет лишь линейную сходимость, причем в случаях, когда производная функции  $\varphi(x)$  близка к единице, эта сходимость может быть весьма медленной.

Один путь ускорения сходимости МПИ – построение *нестационарных процессов* на основе МПИ. На стадии приведения уравнения (5.1) к задаче о неподвижной точке (5.38) можно подменить (5.1) не однопараметрическим семейством уравнений (5.50), а последовательностью уравнений

$$x = x - \lambda_k f(x), \quad k = 0, 1, 2, \dots,$$

и в соответствующем таким уравнениям МПИ

$$x_{k+1} = x_k - \lambda_k f(x_k) \quad (5.52)$$

параметры  $\lambda_k$  подбирать так, чтобы при этом учитывалась информация, получаемая на предыдущем шаге. Если, например, выбор  $\lambda_k$  подчинить

соображениям, что  $\varphi'(x_k) = 1 - \lambda_k f'(x_k) = 0$ , т.е. взять  $\lambda_k = \frac{1}{f'(x_k)}$ ,

то (5.52) будет определять основной метод Ньютона, сходящийся, как уже известно (см. п.5.4), квадратично. Можно выбирать и другие стратегии фиксирования  $\lambda_k$ , но вряд ли здесь следует ожидать нечто принципиально новое, что нельзя более естественно получить с помощью формулы Тейлора и интерполяционных формул.\*\*)

Другой путь – это алгоритмическое построение последовательностей, так или иначе "паразитирующих" на последовательности приближений МПИ (5.36), т.е. получаемых с помощью несложных арифметических манипуляций над несколькими членами последовательности  $(x_k)$  и в ре-

---

<sup>\*)</sup> В [41, 42] метод простых итераций, примененный к уравнению (5.50), т.е. процесс вида

$$x_{k+1} = x_k - \alpha f(x_k),$$

при условии, что  $\text{sign } \alpha = \text{sign } f'(x_k)$ , называют *методом хорд* (или параллельных хорд), и модифицированный метод Ньютона считают частным случаем этого метода.

<sup>\*\*)</sup> Довольно обширная и глубокая теория итерационных методов, в основу которой положено понятие *итерационной функции*, содержится в монографии известного американского математика Дж. Трауба [50].

зультате имеющих более быструю сходимость. Для всех таких методов характерны *многошаговость*, *экономичность* (поскольку более быстрая сходимость по сравнению с базовой последовательностью достигается без дополнительного вычисления значений функций), а также *сложность исследования* условий и скорости сходимости, отсюда — *отсутствие эффективных априорных оценок погрешностей*. Возможны ситуации, когда новый метод окажется сходящимся, в то время как базовый для него МПИ расходится.

Рассмотрим два таких метода ускорения сходимости последовательности (5.36), не делая попыток их строгого обоснования, а ограничиваясь рациональными рассуждениями при их выводе, а также наглядными примерами, демонстрирующими их эффективность. *Наличие неподвижной точки  $\xi$  и дифференцируемость функции  $\varphi(x)$  далее всюду предполагается.*

**а)  $\Delta^2$ - процесс Эйткена<sup>\*)</sup>**

Пусть  $(x_k)$  — последовательность, получаемая по формуле (5.36). Вычитая (5.36) из (5.37), имеем

$$\xi - x_{k+1} = \varphi(\xi) - \varphi(x_k),$$

а уменьшив здесь индекс на единицу, получаем

$$\xi - x_k = \varphi(\xi) - \varphi(x_{k-1}).$$

К правым частям этих равенств применим формулу Лагранжа, согласно которой найдутся точки  $c_k$  и  $c_{k-1}$  такие, что

$$\varphi(\xi) - \varphi(x_k) = \varphi'(c_k)(\xi - x_k)$$

и

$$\varphi(\xi) - \varphi(x_{k-1}) = \varphi'(c_{k-1})(\xi - x_{k-1}).$$

Таким образом, имеют место следующие связи между ошибками соседних приближений:

$$\xi - x_{k+1} = \varphi'(c_k)(\xi - x_k), \quad \xi - x_k = \varphi'(c_{k-1})(\xi - x_{k-1}).$$

Предположим, что в той окрестности корня  $\xi$ , в которой находятся точки  $x_{k-1}$  и  $x_k$ , производная  $\varphi'(x)$  меняется не очень быстро. Это допущение позволяет считать, что

$$\varphi'(c_k) \approx \varphi'(c_{k-1}) \approx \eta$$

(где  $\eta$  некоторое число), и значит,

$$\xi - x_{k+1} \approx \eta(\xi - x_k), \quad \xi - x_k \approx \eta(\xi - x_{k-1}).$$

<sup>\*)</sup> Метод (читается: "дельта-два-процесс") заложен в публикациях А. Айткена, относящихся к 1926–1955 гг. Первоначально был предназначен для улучшения метода Бернулли решения алгебраических уравнений (см. далее п. 5.10)

Беря отношение этих приближенных равенств, избавляемся от  $\eta$  :

$$\frac{\xi - x_{k+1}}{\xi - x_k} \approx \frac{\xi - x_k}{\xi - x_{k-1}}, \quad (5.53)$$

и разрешаем полученное приближенное уравнение относительно  $\xi$  :

$$\xi \approx \frac{x_{k+1}x_{k-1} - x_k^2}{x_{k+1} - 2x_k + x_{k-1}}.$$

Последнее выражение можно использовать на завершающем этапе применения метода простых итераций (5.36), чтобы получить более точное приближение к корню  $\xi$  с помощью трех последних членов последовательности  $(x_k)$ . В развитие же метода обозначим правую часть этого приближенного равенства через  $\tilde{x}_{k+1}$  и придадим его выражению другой вид:

$$\tilde{x}_{k+1} := \frac{x_{k+1}x_{k-1} - x_k^2}{x_{k+1} - 2x_k + x_{k-1}} = x_{k+1} - \frac{(x_{k+1} - x_k)^2}{x_{k+1} - 2x_k + x_{k-1}}.$$

Более коротко это записывается так:

$$\tilde{x}_{k+1} = x_{k+1} - \frac{(\Delta x_k)^2}{\Delta^2 x_{k-1}}, \quad (5.54)$$

где  $\Delta x_k := x_{k+1} - x_k$ ,  $\Delta^2 x_{k-1} := \Delta x_k - \Delta x_{k-1} = x_{k+1} - 2x_k + x_{k-1}$  — так называемые *конечные разности первого и второго порядков* соответственно. Отсюда название (5.54)  $\Delta^2$ -преобразование или  $\Delta^2$ -процесс Эйткена.<sup>\*)</sup>

Организация вычислений на основе этого преобразования может быть различной. Наиболее целесообразным считается применение  $\Delta^2$ -ускорения (5.54) через два шага МПИ на третий<sup>\*\*)</sup>. Поскольку в этом комбинированном методе нахождения корня  $\xi$  уравнения (5.38) участвует хорошо изученный МПИ (5.36), для останова процесса вычислений можно использовать в подходящей его фазе вполне надежный критерий (5.48).

Примером реализации такого метода может служить следующий алгоритм.

\*) Вместо (5.54) для  $\Delta^2$ -преобразования Эйткена используют и другое, эквивалентное (5.54), представление  $\tilde{x}_{k+1} = x_{k-1} - \frac{(\Delta x_{k-1})^2}{\Delta^2 x_{k-1}}$ .

\*\*) Применение  $\Delta^2$ -преобразования менее чем через два шага МПИ не обосновано, так как оно должно применяться к трем последовательным членам  $l$  и  $n$  и  $n$  сходящейся последовательности.

### $\Delta^2$ -алгоритм Эйткена

Шаг 0. Ввод  $x_0$  (начального приближения),  $\varphi(x)$  (исходной функции),  $q$  (оценки модуля производной),  $\varepsilon$  (допустимой абсолютной погрешности).

Шаг 1. Вычисление значений:  $x_1 := \varphi(x_0)$ ,  $x_2 := \varphi(x_1)$ .

Шаг 2.  $\Delta^2$ -ускорение:  $\tilde{x}_2 := \frac{x_0 x_2 - x_1^2}{x_2 - 2x_1 + x_0}$ .

Шаг 3. Вычисление контрольного значения:  $x_3 := \varphi(\tilde{x}_2)$ .

Шаг 4. Проверка на точность: если  $|x_3 - \tilde{x}_2| > \frac{1-q}{q} \varepsilon$ , то положить

$x_0 := \tilde{x}_2$ ,  $x_1 := x_3$ , вычислить  $x_2 := \varphi(x_1)$  и вернуться к шагу 2.

Шаг 5. Положить  $\xi := x_3$  (с точностью  $\varepsilon$ ).

Шаг ускорения по методу Эйткена (5.54) на базе последовательности  $(x_k)$ , получаемой МПИ (5.36), имеет простую геометрическую интерпретацию (рис. 5.16 и рис. 5.17 для случаев  $-1 < \varphi'(x) < 0$  и  $0 < \varphi'(x) < 1$  соответственно).

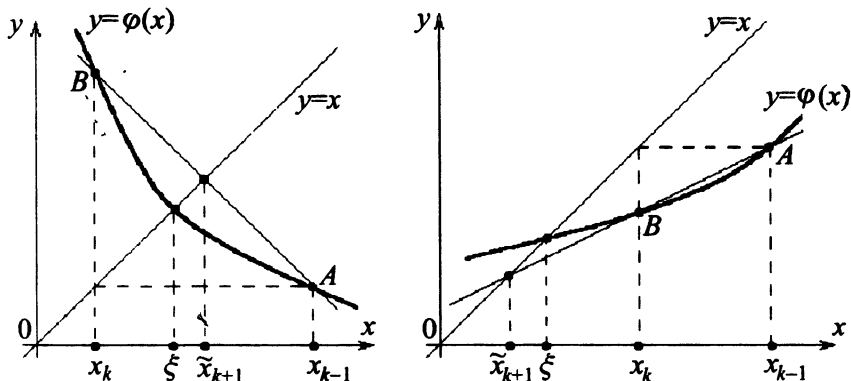


Рис. 5.16, 5.17. Геометрическая иллюстрация одного шага ускорения по методу Эйткена (случаи убывающей и возрастающей функций  $\varphi(x)$ )

Если провести прямую (хорду, секущую) через точки  $A(x_{k-1}, \varphi(x_{k-1}))$  и  $B(x_k, \varphi(x_k))$  кривой  $y = \varphi(x)$ , то абсцисса точки пересечения этой прямой с прямой  $y = x$  как раз и будет определяемая  $\Delta^2$ -преобразованием Эйткена (5.54) точка  $\tilde{x}_{k+1}$  оси  $Ox$ .

Действительно, в силу (5.36), координаты точек  $A$  и  $B$  можно записать иначе:  $A(x_{k-1}, x_k)$ ,  $B(x_k, x_{k+1})$ . Значит, уравнение прямой ( $AB$ ) имеет вид

$$\frac{y - x_k}{x_{k+1} - x_k} = \frac{x - x_{k-1}}{x_k - x_{k-1}}.$$

Рассматривая его совместно с уравнением  $y=x$ , т.е. заменяя в нем  $y$  и  $x$  на  $\tilde{x}_{k+1}$ , получаем то же выражение

$$\tilde{x}_{k+1} = \frac{x_{k+1}x_{k-1} - x_k^2}{x_{k+1} - 2x_k + x_{k-1}},$$

от которого пришли к (5.54).

Такая интерпретация  $\Delta^2$ - процесса Эйткена наталкивает на мысль о его возможной связи с изученным ранее методом секущих (5.32).

Полагая  $f(x) := x - \varphi(x)$ , применим к этой функции формулу (5.32), считая при этом, что требуемые в (5.32) значения элементов последовательности  $(x_k)$  получаются не по той же формуле (5.32), а с помощью МПИ (5.36). Имеем:

$$\begin{aligned} x_k - \frac{f(x_k)(x_{k-1} - x_k)}{f(x_{k-1}) - f(x_k)} &= x_k - \frac{(x_k - \varphi(x_k))(x_{k-1} - x_k)}{x_{k-1} - \varphi(x_{k-1}) - x_k + \varphi(x_k)} = \\ &= x_k - \frac{(x_k - x_{k+1})(x_{k-1} - x_k)}{x_{k+1} - 2x_k + x_{k-1}} = \frac{x_{k+1}x_{k-1} - x_k^2}{x_{k+1} - 2x_k + x_{k-1}} = \tilde{x}_{k+1}. \end{aligned}$$

Таким образом, один шаг  $\Delta^2$ - ускорения МПИ по методу Эйткена совпадает с одним шагом метода секущих, примененного к той же последовательности МПИ.

Сделанное наблюдение в свете известных о методе секущих сведений позволяет с большой осторожностью судить об эффекте ускорения сходимости, который может принести использование метода Эйткена. Надежность его, очевидно, выше в случае, когда  $\varphi'(x) \in (-1; 0)$ , применение более актуально, если  $q > \frac{1}{2}$ , и ускорение тем эффективней, чем меньше  $|\varphi''(x)|$  в окрестности корня  $\xi$ .

Применяя метод Эйткена, не следует забывать о проблеме своевременного прерывания счета из-за потерь точности при вычитании близких чисел. Подключение  $\Delta^2$ - ускорения на ранней стадии МПИ, когда  $x_0$  далеко от  $\xi$ , может привести к расходимости процесса, по крайней мере, в случае, когда  $\varphi'(x) > 0$  (представьте эту ситуацию, глядя на рис. 5.17). В то же время иногда с помощью метода Эйткена можно получить сходимость в условиях, когда МПИ (5.36) расходится (см., например, рис. 5.18, где  $|\xi - x_k| > |\xi - x_{k-1}|$ , но  $|\xi - \tilde{x}_{k+1}| < |\xi - x_k|$ ).

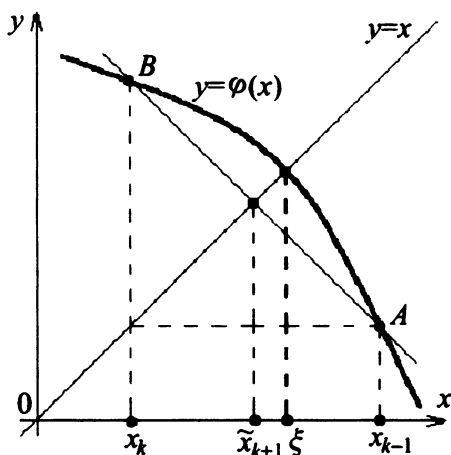


Рис. 5.18. Демонстрация возможной сходимости последовательности Эйткена ( $\tilde{x}_k$ ) при расхождении последовательности МПИ ( $x_k$ )

**Замечание 5.6.**  $\Delta^2$ - преобразование Эйткена применимо не только к последовательности приближений (5.36), но и к любым другим последовательностям, сходящимся со скоростью геометрической прогрессии. Действительно, пусть  $(x_k)$  – некоторая последовательность, линейно сходящаяся к предельной точке  $\xi$ . Тогда можно считать, что разность  $\xi - x_k$  изменяется по закону геометрической прогрессии, т.е. существуют такие постоянная  $v \in (0;1)$  и слабо изменяющаяся варианта  $C_k \approx C$ , что

$$\xi - x_{k-1} \approx Cv^{k-1}, \quad \xi - x_k \approx Cv^k, \quad \xi - x_{k+1} \approx Cv^{k+1}.$$

Отсюда получаем приближенные равенства

$$\frac{\xi - x_k}{\xi - x_{k-1}} \approx v, \quad \frac{\xi - x_{k+1}}{\xi - x_k} \approx v,$$

следствием которых является равенство (5.53), в итоге приводящее к формуле (5.54). Такой подход к выводу  $\Delta^2$ - метода Эйткена позволяет использовать его при решении других задач (см., например, замечание 4.6).

Применение  $\Delta^2$ - преобразования Эйткена к последовательностям, сходящимся квадратично, эффекта ускорения не дает [50].



## б) метод Вегстейна<sup>\*)</sup>

При выводе *метода Вегстейна* решения задачи о неподвижной точке (5.38) будем использовать как аналитические, так и геометрические соображения.

Пусть уже найдены:  $\bar{x}_k$  – элемент строящейся здесь последовательности, и  $x_{k+1} = \varphi(\bar{x}_k)$  – точка, соответствующая одному шагу МПИ, примененного к точке  $\bar{x}_k$ . Независимо от того, сходится начатый с  $\bar{x}_k$  МПИ (рис. 5.19, где  $|\xi - x_{k+1}| < |\xi - \bar{x}_k|$ ) или расходится (рис. 5.20 с  $|\xi - x_{k+1}| > |\xi - \bar{x}_k|$ ), отрезок  $AB$ , параллельный оси  $Ox$  и имеющий концами точки  $A(\bar{x}_k, \varphi(\bar{x}_k))$  и  $B(x_{k+1}, x_{k+1})$ , можно разделить точкой  $C$  так, чтобы она принадлежала вертикальной прямой  $x = \xi$  (при этом во втором случае речь идет о делении отрезка внешним образом).

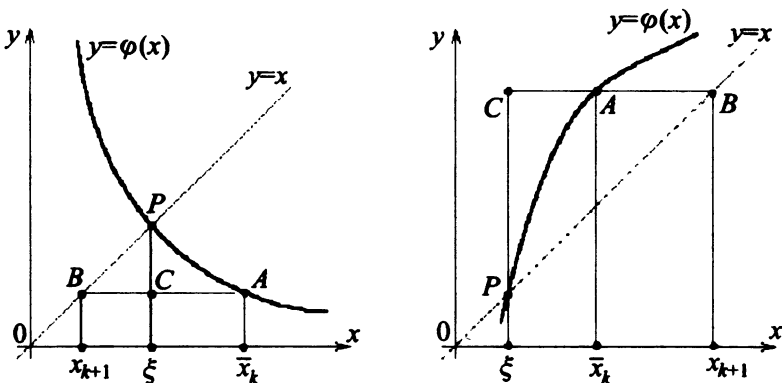


Рис. 5.19, 5.20. К построению метода Вегстейна

При любых комбинациях направлений возрастания и выпуклости графика функции  $y = \varphi(x)$  в окрестности неподвижной точки  $\xi$  имеет место равенство длин отрезков  $BC = PC$ . Различаются два случая: когда  $BC = \xi - x_{k+1}$  и когда  $BC = x_{k+1} - \xi$ . По формуле Лагранжа соответственно имеем

$$PC = \varphi(\xi) - \varphi(\bar{x}_k) = \varphi'(\theta_k)(\xi - \bar{x}_k)$$

или

<sup>\*)</sup> Ссылки на первую публикацию этого метода (1958 г.) имеются в [36, 38] (в [38] его называют *усовершенствованным методом последовательных приближений*). К сожалению, описание метода Вегстейна пока не заняло достойного места в отечественной учебной литературе по вычислительной математике. Нет упоминания о нем и в монографии [50]

$$PC = \varphi(\bar{x}_k) - \varphi(\xi) = \varphi'(\theta_k)(\bar{x}_k - \xi).$$

В любом случае можно утверждать, что существует точка  $\theta_k \in (\bar{x}_k; \xi)$  или  $\theta_k \in (\xi; \bar{x}_k)$  такая, что

$$\xi - x_{k+1} = \varphi'(\theta_k)(\xi - \bar{x}_k).$$

Разрешая это линейное уравнение относительно  $\xi$ , находим

$$\xi = x_{k+1} - \frac{x_{k+1} - \bar{x}_k}{1 - \frac{1}{\varphi'(\theta_k)}}. \quad (5.55)$$

Если бы значение  $\varphi'(\theta_k)$  было известно, то тем самым задача о неподвижной точке (5.38) была бы решена точно. Заменяем это неизвестное значение  $\varphi'(\theta_k)$  аппроксимирующим его разностным отношением:

$$\varphi'(\theta_k) \approx \frac{\varphi(\bar{x}_k) - \varphi(\bar{x}_{k-1})}{\bar{x}_k - \bar{x}_{k-1}} = \frac{x_{k+1} - x_k}{\bar{x}_k - \bar{x}_{k-1}}.$$

Подставляя приближенное значение  $\frac{1}{\varphi'(\theta_k)} \approx \frac{\bar{x}_k - \bar{x}_{k-1}}{x_{k+1} - x_k}$  в (5.55), вместо корня  $\xi$  получаем приближение к нему

$$\bar{x}_{k+1} = x_{k+1} - \frac{(x_{k+1} - x_k)(x_{k+1} - \bar{x}_k)}{(x_{k+1} - x_k) - (\bar{x}_k - \bar{x}_{k-1})}. \quad (5.56)$$

Эта итерационная формула, где  $k = 1, 2, 3, \dots$ , совместно с формулой

$$x_{k+1} = \varphi(\bar{x}_k) \quad (k = 0, 1, 2, \dots) \quad (5.57)$$

и начальными значениями  $\bar{x}_0 := x_0$ ,  $\bar{x}_1 := x_1$  полностью определяет метод Вегстейна для задачи (5.38).

Значение  $\bar{x}_2$ , получаемое по формуле Вегстейна (5.56) при заданных начальных значениях  $\bar{x}_0$  и  $\bar{x}_1$ , совпадает со значением  $\bar{x}_2$ , вычисляемым  $\Delta^2$ -процессом Эйткена. Далее, т.е. при  $k \geq 2$ , процессы (5.54) и (5.56) различаются. Учитывая, что МПИ является составной частью метода Вегстейна, в случаях, когда  $|\varphi'(x)| \leq q < 1$ , можно заканчивать процесс вычислений, как и в методе Эйткена, по выполнению критерия (5.48).

Таким образом, для реализации метода (5.56)-(5.57) может быть предложен, например, следующий алгоритм.

### алгоритм Вегстейна

- Шаг 0. Ввести  $x_0$  (начальное приближение),  $\varphi(x)$  (исходную функцию),  $q$  (оценку модуля производной),  $\varepsilon$  (допустимую абсолютную погрешность).
- Шаг 1. Вычислить  $x_1 := \varphi(x_0)$ ; положить  $\bar{x}_0 := x_0$ ,  $\bar{x}_1 := x_1$ .
- Шаг 2. Вычислить  $x_2 := \varphi(\bar{x}_1)$ .
- Шаг 3. Проверить на точность: если  $|x_2 - \bar{x}_1| > \varepsilon(1 - q)/q$ , то вычислить  $\bar{x}_2 := \frac{x_2 \bar{x}_0 - x_1 \bar{x}_1}{x_2 + \bar{x}_0 - x_1 - \bar{x}_1}$ ; переприсвоить значения  $x_0 := x_1$ ,  $\bar{x}_0 := \bar{x}_1$ ,  $x_1 := x_2$ ,  $\bar{x}_1 := \bar{x}_2$  и вернуться к шагу 2.
- Шаг 4. Положить  $\xi := x_2$  (с точностью  $\varepsilon$ ).

Разумеется, проверку на точность в подобном алгоритме можно устраивать иную (что просто необходимо, если метод Вегстейна применяется в случаях, когда  $|\varphi'(x)| > 1$ ). Если нет угрозы большой потери точности из-за вычитания близких чисел, то заканчивать работу алгоритма Вегстейна лучше выводом значения  $\xi := \bar{x}_2$ . Для вычисления значения  $\bar{x}_2$  в этом алгоритме применена равносильная (5.56) формула

$$\bar{x}_{k+1} = \frac{x_{k+1} \bar{x}_{k-1} - x_k \bar{x}_k}{x_{k+1} + \bar{x}_{k-1} - x_k - \bar{x}_k},$$

имеющая несколько отличную от (5.56) структуру.

Как показывают многочисленные эксперименты с уравнениями вида  $x = \varphi(x)$ , особый интерес среди которых вызывают случаи, когда простые итерации дают расходящиеся последовательности  $(x_k)$ , метод Вегстейна имеет определенные преимущества перед методом Эйткена по количеству обращений к вычислению значений  $\varphi(x)$  для получения корня с заданной точностью. Чаще всего метод Вегстейна еще и позволяет в более широких пределах варьировать выбор начальной точки  $x_0$ . Результаты сравнения этих двух методов на нескольких таких уравнениях представлены в табл. 5.1. В двух ее последних столбцах указано количество вычислений значений функции  $\varphi(x)$  (горнеров), потребовавшееся для достижения точки  $x^* \approx \xi$  (приближенного значения корня) такой, что  $|x^* - \varphi(x^*)| < 10^{-7}$ . (Приведенные данные получены Ковалевым П.В. на IBM PC/AT - 286).

Таблица 5.1

Функция $\varphi(x)$	Корень $\xi$ (точность $10^{-6}$ )	Значение производной $\varphi'(\xi)$ (точность $10^{-2}$ )	Начальное приближение $x_0$	Число горнеров	
				метод Эйткена	метод Вегстейна
$x^3 + 2$	-1521379	6.94	-2	18	8
			-0.5	16	9
			1	190	10
$(1-x^2)^2$	0.524889	-1.52	0	6	5
			1.3	8	5
			1.7	20	8
	1.490216	7.28	4	расходится	14
$\frac{1-x}{x}$	0.618034	-2.62	0.8	10	7
			-3	6	5
	-1618034	-0.38	-0.1	10	8
$\frac{1}{x}$	1	-1	1.2	6	5
			10	14	10
$e^x \ln x$	1	1	0.2	20	16
			1.2	20	15

Обратим внимание, что лишь для одного корня из семи в данной таблице можно говорить о сходимости МПИ.

## 5.9. НЕЛИНЕЙНЫЕ УРАВНЕНИЯ С ПАРАМЕТРОМ. БИФУРКАЦИИ<sup>\*)</sup>

Материал п.5.7 можно интерпретировать так: нелинейная непрерывная математическая модель (5.38) некоего явления изучается путем построения и исследования соответствующей дискретной модели (5.36). Связь между этими моделями на отрезке  $[a; b]$  устанавливается при выполнении двух следующих условий:

$$\varphi(x) \in [a; b] \quad \forall x \in [a; b] \quad (\text{отображение в себя})$$

и

$$|\varphi'(x)| < 1 \quad \forall x \in (a; b) \quad (\text{сжатие}).$$

(5.58)

А именно, согласно теореме 5.7, эти условия являются достаточными для существования и единственности на  $[a; b]$  решения  $\xi$  непрерывной задачи (5.38), причем оно может быть получено как предел последовательности  $(x_k)$  (т.е. как решение дискретной задачи (5.36)), начинающейся с лю-

<sup>\*)</sup> Сведения, излагаемые в этом пункте, имеют ознакомительный характер и опираются на статью [3]. Из других литературных источников, посвященных рассматриваемым здесь вопросам, отметим еще книгу [58].

бой точки  $x_0 \in [a; b]$ . Последнее можно расценить как устойчивость в данном смысле решения  $\xi$  дискретной модели (5.36).

Продолжим изучение взаимосвязи непрерывной и дискретной одномерных нелинейных моделей в следующем русле.

Во-первых, возьмем за основу и будем рассматривать некоторую конкретную дискретную модель, а для ее исследования привлечем соответствующую непрерывную модель.

Во-вторых, попытаемся выяснить, к чему может привести нарушение условий сходимости МПИ, т.е. условий (5.58), применительно к данной модели, и какие "тайны" могут скрываться за термином *нелинейность*.

Преследуя эти цели, введем в дискретное (5.36) и непрерывное (5.38) уравнения вещественный параметр  $\lambda$ , т.е. будем изучать связь между моделями вида

$$x_{k+1} = \varphi(x_k, \lambda), \quad k = 0, 1, 2, \dots; \quad x_0 \in [a; b] \quad (5.59)$$

и

$$x = \varphi(x, \lambda), \quad x \in [a; b]. \quad (5.60)$$

Заметим, что не так редки ситуации, когда в приложениях математики первичными являются именно дискретные модели, а их непрерывные аналоги нужны для того, чтобы воспользоваться хорошо развитой теорией математического анализа и плодами вычислительной математики, которая, в основном, построена по принципу "непрерывная задача  $\rightarrow$  дискретная аппроксимация".

Представим себе следующую весьма идеализированную картину. Пусть на некоторой ограниченной территории, например, на острове, может прокормиться не более  $N$  животных определенного вида, и пусть в начальный момент наблюдений за ними их количество было  $g \in (0; N)$ . Будем считать, что животные ежегодно приносят потомство, и скорость размножения характеризуется некоторым параметром  $\alpha > 0$ . Тогда, если через  $g_k$  обозначить численность животных в  $k$ -й год после начала наблюдения, то можно предположить, что закон ежегодного изменения численности популяции грубо описывается моделью

$$g_{k+1} = \alpha g_k (N - g_k), \quad k = 0, 1, 2, \dots \quad (5.61)$$

В пользу принятия такой модели говорят следующие рассуждения. Если значение  $g_0$  начальной численности мало, то второй сомножитель в начале процесса почти постоянен, и все зависит от коэффициента роста  $\alpha$ : при малых  $\alpha$ , т.е. при низкой скорости размножения численность животных будет снижаться и в конце концов популяция гибнет. Если же  $g_k \in [0; N]$  возрастает и приближается к максимально возможному значению  $N$ , то за счет близости к нулю второго сомножителя численность популяции естественно начнет снижаться.

Чтобы облегчить исследование модели (5.61), упростим ее заменой переменных. Переписав (5.61) в виде

$$\frac{g_{k+1}}{N} = N\alpha \cdot \frac{g_k}{N} \left(1 - \frac{g_k}{N}\right)$$

и положив  $x_k := \frac{g_k}{N}$ ,  $\lambda := N\alpha$ , приходим к уравнению<sup>\*)</sup>

$$x_{k+1} = \lambda x_k (1 - x_k), \quad (5.62)$$

где  $k = 0, 1, 2, \dots$ , а значения  $x_k$  в соответствии со смыслом задачи должны принадлежать отрезку  $[0; 1]$ .

На равенство (5.62) можно смотреть как на МПИ (5.59), применяемый к задаче о неподвижной точке вида (5.60), т.е. к задаче о корнях уравнения

$$x = \lambda x (1 - x) \quad (5.63)$$

в области  $x \in [0; 1]$ .

Условившись не использовать далее в обозначениях функции ее явную зависимость от параметра  $\lambda$ , положим

$$\varphi(x) := \lambda x (1 - x)$$

и преобразуем эту квадратичную функцию к виду

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

Из последнего следует, что  $\max \varphi(x) = \frac{\lambda}{4}$  при  $x = \frac{1}{2}$  и что  $\varphi(x)$  отображает отрезок  $[0; 1]$  в  $\left[0; \frac{\lambda}{4}\right]$ . Значит, при  $\lambda \in [0; 4]$  функция  $\varphi(x)$  осуществляет на отрезке  $[0; 1]$  отображение в себя, т.е. при этих  $\lambda$  элементы последовательности  $(x_k)$ , получаемой с помощью равенства (5.62), при любом  $x_0 \in [0; 1]$  не выйдут за пределы  $[0; 1]$ , другими словами, определены при любом  $k = 0, 1, 2, \dots$ <sup>\*\*)</sup>.

Для производной данной функции  $\varphi(x)$  имеем:

$$\varphi'(x) = \lambda(1 - 2x) \quad \text{и} \quad \max_{x \in [0; 1]} |\varphi'(x)| = \lambda.$$

Следовательно, при  $\lambda < 1$  отображение  $\varphi(x)$  является сжимающим на  $[0; 1]$  и имеет единственную неподвижную точку  $\xi_1 \in [0; 1]$ , а именно

<sup>\*)</sup> Это конкретное уравнение называют *логистическим* [58].

<sup>\*\*)</sup> При  $\lambda > 4$  нарушается соответствие между дискретной (5.62) и непрерывной (5.63) моделями. Например, при  $x_0 = \frac{1}{2}$  уже  $x_1 = \frac{\lambda}{4} > 1$ , т.е. не принадлежит множеству, на котором определена функция  $\varphi(x)$ .

$\xi_1 = 0$ , которая является пределом последовательности  $(x_k)$  при любом начальном значении  $x_0 \in [0; 1]$  (популяция гибнет по причине недостаточной скорости воспроизводства). Геометрическая иллюстрация этого случая, соответствующего выполнению условий теоремы 5.7, показана на рис.5.21.

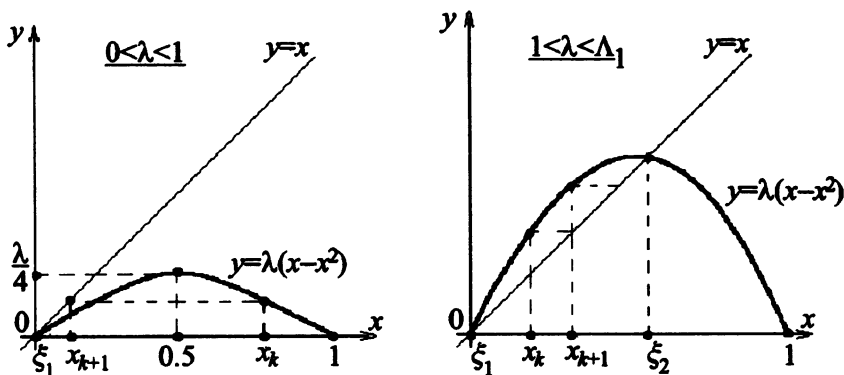


Рис. 5.21, 5.22. Сходимость МПИ (5.62) к корню  $\xi_1 = 0$  логистического уравнения (5.63) при  $\lambda \in (0; 1)$  и к корню  $\xi_2 = \frac{\lambda - 1}{\lambda}$  при  $\lambda \in (1; \Lambda_1)$

При  $1 \leq \lambda \leq 4$  нарушается одно из условий сходимости МПИ:  $\varphi(x)$  не является функцией сжатия. Что же это влечет?

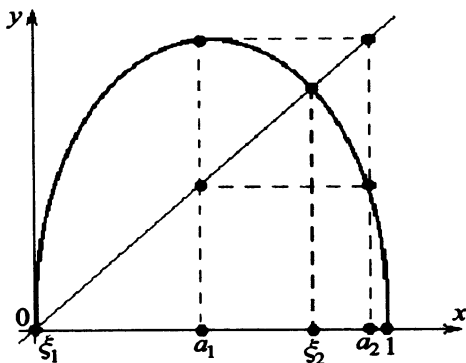


Рис. 5.23. Предельное поведение последовательности  $x_{k+1} = \varphi(x_k)$ , где  $\varphi(x) := \lambda(x - x^2)$  при  $\lambda \in (\Lambda_1; \Lambda_2)$

Очевидно, уравнение (5.63) по-прежнему сохраняет решение  $\xi_1 = 0 \in [0; 1]$ . Но при переходе  $\lambda$  через 1 это решение теряет устойчивость и появляется второе решение  $\xi_2 = \frac{\lambda - 1}{\lambda} \in [0; 1]$ , которое следует считать устойчивым, поскольку теперь именно к нему будет сходиться любая последовательность, определяемая начатым с  $x_0 \in [0; 1]$  МПИ (5.62) (рис. 5.22). Произошло явление, которое носит название **бифуркация решений**: вместо одного решения на рассматриваемом промежутке стало два решения<sup>\*)</sup>.

Сходимость  $(x_k)$  к  $\xi_2$  будет наблюдаться не при всех  $\lambda \in [1; 4]$ . Оказывается, существует число  $\Lambda_1 > 1$  такое, при переходе  $\lambda$  через которое начнет происходить **заикливание** последовательности  $(x_k)$ . А именно, какое бы ни взяли  $x_0 \in (0; 1)$ , начатая с него и продолжаемая по формуле (5.62) последовательность будет обладать тем свойством, что все ее четные члены будут иметь предел одно число, а нечетные – другое. Это означает, что найдутся числа  $a_1, a_2 \in (0; 1)$  (при каждом  $\lambda$  свои) такие, что  $a_2 = \varphi(a_1)$  и  $a_1 = \varphi(a_2)$ , причем  $a_1 \neq a_2 \neq \xi_i$  ( $i = 1, 2$ ) (см. рис. 5.23). В этом случае говорят, что дискретное отображение (5.62) имеет **устойчивый цикл периода 2** и обозначают его  $S^2$ . Относя это к исходной модельной задаче с животными, можно сказать, что при значениях  $\lambda > \Lambda_1$  численность популяции будет меняться периодически с периодом в два “года”.

Зная ситуацию качественно, нетрудно найти точно пороговое значение  $\Lambda_1$ , при котором появляется устойчивый цикл  $S^2$ .

Действительно, если известно, что на  $(0; 1)$  имеются точки  $a_1, a_2$  такие, что  $a_2 = \varphi(a_1)$ ,  $a_1 = \varphi(a_2)$ , то значит,  $a_2 = \varphi(\varphi(a_2))$ ,  $a_1 = \varphi(\varphi(a_1))$ , т.е.  $a_1$  и  $a_2$  – неподвижные точки отображения  $\varphi^{\circledast}(x) := \varphi(\varphi(x))$ , иначе, – корни уравнения  $x = \varphi(\varphi(x))$ .

Так как эта суперпозиция сохраняет старые неподвижные точки  $\xi_1 = 0$  и  $\xi_2 = (\lambda - 1)/\lambda$ , то уравнение

$$x = \lambda(\lambda x(1 - x))(1 - \lambda x(1 - x))$$

должно иметь четыре корня, из которых два известны (рис. 5.24). Исключив из этого уравнения известные корни, приходим к квадратному уравнению

$$x^2 - \frac{\lambda + 1}{\lambda} x + \frac{\lambda + 1}{\lambda^2} = 0.$$

<sup>\*)</sup> Bifurcus (лат.) – раздвоенный. Термин введен К. Якоби в 1834 году. Теория бифуркаций заложена Анри Пуанкаре в конце XIX века.



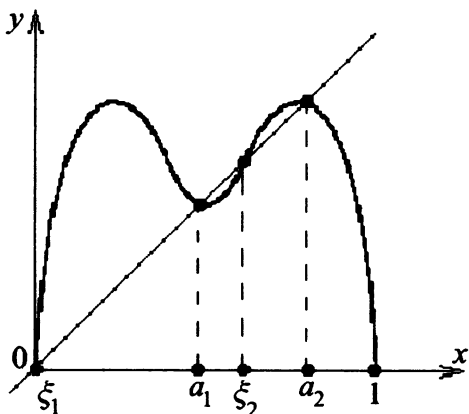


Рис. 5.24. График функции  $\varphi^{\textcircled{2}}(x) := \varphi(\varphi(x))$  при  $\varphi(x) := \lambda(x - x^2)$ ,  $\lambda \in (\Lambda_1, \Lambda_2)$  и ее неподвижные точки  $(\xi_1, \xi_2, a_1, a_2)$

Положительное значение  $\lambda$ , которое служит границей области положительности дискриминанта  $D = (\lambda^2 - 2\lambda - 3)/\lambda^2$  этого уравнения, как раз и есть искомое значение  $\Lambda_1$ , начиная с которого появляются новые устойчивые неподвижные точки  $a_1, a_2$ , т.е. цикл  $S^2$ . Очевидно, это  $\Lambda_1 = 3$ .

Устойчивость неподвижных для  $\varphi^{\textcircled{2}}(x)$  точек  $a_{1,2} = (\lambda + 1 \pm \sqrt{\lambda^2 - 2\lambda - 3})/2\lambda$  в том смысле, что они становятся точками четно-нечетного притяжения для последовательности  $(x_k)$ , устанавливает-

ся непосредственной проверкой условия  $\left| \frac{d\varphi^{\textcircled{2}}(x)}{dx} \right|_{x=a_{1,2}} < 1$ .

Дальнейшее увеличение  $\lambda$  в дискретном уравнении (5.62) ведет к тому, что начиная с некоторого  $\Lambda_2$  заикливание будет иметь более сложный характер: при каждом  $\lambda \in (\Lambda_2, \Lambda_3)$ , где  $\Lambda_2 > 3$ ,  $\Lambda_3 < 4$ , найдутся числа  $a_1, a_2, a_3, a_4$  (зависящие от  $\lambda$ ) такие, что  $a_2 = \varphi(a_1)$ ,  $a_3 = \varphi(a_2)$ ,  $a_4 = \varphi(a_3)$ ,  $a_1 = \varphi(a_4)$ , и члены последовательности  $(x_k)$  будут поочередно все сильнее притягиваться к этим числам, с какого  $x_0 \in (0; 1)$  ни начался бы процесс (5.62). Говорят, что в этом случае имеет место устойчивый цикл периода 4, т.е.  $S^4$ .

Такой процесс образования новых циклов происходит с увеличением параметра  $\lambda$  и далее. Точки  $\Lambda_1, \Lambda_2, \Lambda_3, \dots$ , в которых имеет место зарождение циклов  $S^2, S^4, S^8, \dots$ , называются *точками бифуркации удвоения периода*.

Этому процессу бифуркаций удвоения периода можно придать наглядный вид, если отобразить на графике зависимость *элементов цикла* – значений устойчивых неподвижных точек отображений  $\varphi^{2^n}(x)$  (иначе, точек притяжения подпоследовательностей получаемой посредством МПИ (5.62) последовательности  $(x_k)$ ) – от значений параметра  $\lambda$  при  $\lambda > 1$  (рис. 5.25).

Последовательность  $(\Lambda_n)$  точек бифуркации удвоения периода обладает определенной закономерностью:

$$\frac{\Lambda_n - \Lambda_{n-1}}{\Lambda_{n+1} - \Lambda_n} \xrightarrow{n \rightarrow \infty} \delta = 4.66920\dots$$

Имеет постоянный предел, равный величине  $\alpha = 2.50290\dots$ , также отношение  $d_{n-1}/d_n$ , где через  $d_n$  обозначено расстояние от точки  $x = 0.5$  до ближайшего элемента цикла  $S^{2^{n-1}}$ , соответствующего такому значению  $\lambda$ , при котором  $x = 0.5$  является элементом того же цикла (рис. 5.25). Числа  $\delta$  и  $\alpha$  называют *постоянными Фейгенбаума* в честь открывшего эти закономерности американского математика (1978 г.).

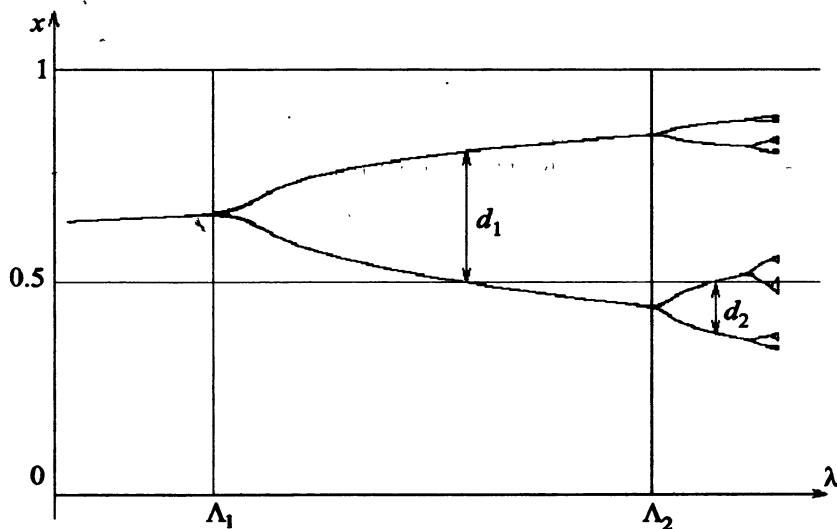


Рис. 5.25. Бифуркационная диаграмма циклов периода  $2^n$

Верхней границей значений  $\lambda$ , при которых получаемая с помощью (5.62) последовательность  $(x_k)$  ведет себя указанным образом, т.е. имеются циклы  $S^{2^n}$ , является значение  $\lambda^* \approx 3.57$ . Дальнейшее увеличение  $\lambda$  приводит к срыву цикличности. В каком-то диапазоне значений  $\lambda > \lambda^*$  будет наблюдаться бесконечное хаотическое блуждание точек последовательности  $(x_k)$  в пределах промежутка  $(0;1)$ , с какого бы  $x_0 \in (0;1)$  она не начиналась. Затем снова из хаоса возникают устойчивые циклы, происходят бифуркации удвоения периода и опять срыв в хаос. Такие чередования циклического (с разными периодами, например,  $12 \cdot 2^n$ ,  $10 \cdot 2^n$ ,  $6 \cdot 2^n$ ,  $8 \cdot 2^n$ ,  $7 \cdot 2^n$ ,  $5 \cdot 2^n$  и др.) и хаотического поведения последовательности  $(x_k)$  имеют место в процессе увеличения  $\lambda$  почти вплоть до предельного значения  $\lambda = 4$ . При этом самый большой (по  $\lambda$ ) промежуток цикличности после циклов вида  $S^{2^n}$  будет иметь цикл периода  $3 \cdot 2^n$  (при  $\lambda \geq 3.829$ ), играющий особую роль в теории бифуркаций.

Как зарождается порядок в хаосе, каковы связи между циклами разных периодов, какую роль играет последовательность обхода элементов цикла, что можно сказать об устойчивости тех или иных циклов – эти и другие вопросы возникают перед математиками, изучающими нелинейные отображения. Непростые ответы на них, достаточно наглядные в одномерном случае, позволяют понять природу многих сложных явлений (например, оценить принципиальные возможности долгосрочного прогнозирования погоды).

## 5.10. О МЕТОДАХ РЕШЕНИЯ АЛГЕБРАИЧЕСКИХ УРАВНЕНИЙ. МЕТОД БЕРНУЛЛИ

Пусть требуется найти один, несколько или все корни многочлена с действительными коэффициентами

$$P_n(x) := a_0 x^n + a_1 x^{n-1} + \dots + a_{n-1} x + a_n, \quad (5.64)$$

т.е. решить в каком-то указанном смысле уравнение

$$P_n(x) = 0. \quad (5.65)$$

Если речь идет о нахождении только действительных корней и, особенно, если нужно найти не все, а только некоторые из них, то есть резон в применении к алгебраическому уравнению (5.65) какого-нибудь из рассмотренных выше методов решения нелинейных скалярных уравнений (5.1) с  $P_n(x)$  в роли  $f(x)$ . При этом следует обратить внимание на то, что  $P_n(x)$  является определенной на всей действительной оси бесконечно дифференцируемой функцией, и здесь можно применять методы, использующие производные  $P_n(x)$ . Поскольку дифференцирование понижает

степень многочлена, вычисление значений производных потребует даже меньше арифметических действий, чем вычисление значений исходного многочлена. Отсюда – целесообразность нахождения отдельных корней многочлена (в том числе и комплексных, см. [21, 23, 38, 40, 47]) методом Ньютона или гибридным алгоритмом, его использующим.

Применение итерационных методов решения нелинейных скалярных уравнений (5.65), а также многих других методов требует многократного вычисления значений многочленов (5.64). Эта промежуточная задача для многочленов намного проще, чем задача вычисления значений трансцендентных функций: нужно всего лишь при заданном  $x = x_0$  простым перемножением находить степени  $x_0^i$  при  $i = 1, 2, \dots, n$  (это  $n - 1$  умножение), затем умножить их на коэффициенты ( $n$  умножений) и сложить результаты ( $n$  сложений). Однако в практике вычислений используют более эффективный способ вычисления значений многочленов – *схему Горнера*, позволяющую почти вдвое уменьшить количество умножений<sup>\*)</sup>. Выведем этот способ.

Согласно теореме Безу, при любом  $x_0$  найдутся числа  $b_i$  ( $i = 0, 1, 2, \dots, n$ ) такие, что

$$a_0x^n + a_1x^{n-1} + \dots + a_{n-1}x + a_n \equiv (x - x_0)(b_0x^{n-1} + b_1x^{n-2} + \dots + b_{n-2}x + b_{n-1}) + b_n,$$

причем  $b_n = P_n(x_0)$ . Выполняя умножение в правой части тождества и приравнивая коэффициенты при одинаковых степенях  $x$ , получим совокупность  $n + 1$  равенств:

$$b_0 = a_0, \quad b_1 = a_1 + x_0b_0, \quad b_2 = a_2 + x_0b_1, \quad \dots, \quad b_n = a_n + x_0b_{n-1}.$$

Таким образом, после  $n$  умножений (на одно и то же число  $x_0$ ) и  $n$  сложений приходим к искомому результату  $b_n = P_n(x_0)$ .

Однотипность вычислений, производимых в схеме Горнера, позволяяет организовать их в цикл, определяемый формулой

$$b_i = a_i + x_0b_{i-1},$$

где  $i = 1, 2, \dots, n$ ,  $b_0 := a_0$ , со значением  $P_n(x_0) = b_n$  на его выходе.

Числа  $b_0, b_1, \dots, b_n$  далее будем называть *коэффициентами схемы Горнера*. При ручном счете эти коэффициенты удобно записывать в процессе вычисления под соответствующими размещенными в один ряд коэффициентами данного многочлена (5.64) в виде следующей таблицы:

	$a_0$	$a_1$	$a_2$	$\dots$	$a_{n-1}$	$a_n$
$x_0$	$b_0$	$b_1$	$b_2$	$\dots$	$b_{n-1}$	$b_n$

<sup>\*)</sup> Такой способ был известен в Китае еще в средние века и назывался *Тянь-юань*, а затем в начале XIX века был "перестроен" в Европе англичанином Горнером и итальянцем Руффины.

**Пример 5.4.** Доказать, что число 2 является единственным рациональным корнем многочлена  $P_5(x) = x^5 - x^4 - 3x^3 + 2x + 4$ , причем простым.

Все доказательство можно отразить следующей таблицей:

	1	-1	-3	0	2	4
2	1	1	-1	-2	-2	0
1	1	2	1	-1	-3	
-1	1	0	-1	-1	-1	
2	1	3	5	8	14	
-2	1	-1	1	-4	6	

В ее первой строке находятся коэффициенты данного многочлена  $P_5(x)$ , во второй строке – коэффициенты схемы Горнера, примененной к  $P_5(x)$  при  $x_0 = 2$ , означающие, что  $P_5(2) = 0$  и  $P_5(x) = (x-2)P_4(x)$ , где  $P_4(x) = x^4 + x^3 - x^2 - 2x - 2$  (см. тождество, лежащее в основе вывода схемы Горнера). Остальные строки – результаты применения схемы Горнера к многочлену  $P_4(x)$  (что отражено в таблице подчеркиванием) в точках  $\pm 1, \pm 2$ , являющихся делителями свободного члена многочлена с целыми коэффициентами  $P_4(x)$  со старшим коэффициентом 1. Так как  $P_4(1) = -3$ ,  $P_4(-1) = -1$ ,  $P_4(2) = 14$ , и  $P_4(-2) = 6$ , то  $P_4(x)$  не имеет рациональных корней, а значит, и  $P_5(x)$  не имеет других рациональных корней, кроме 2, причем 2 не может быть двухкратным корнем.

Как видно из приведенного примера, схему Горнера удобно использовать для понижения степени алгебраического уравнения выделением линейного множителя, соответствующего известному вещественному корню. Поскольку в подавляющем большинстве случаев эти корни бывают известны лишь приближенно, при понижении степени алгебраического уравнения неизбежна потеря точности, с которой могут быть найдены следующие корни. Это обстоятельство заставляет с осторожностью относиться к выигрышу в вычислительных затратах, достигаемому таким понижением степени, когда находятся несколько корней многочлена последовательно корень за корнем.

Существует способ вычисления корня многочлена (5.64) последовательно цифра за цифрой непосредственно по схеме Горнера, применяемой к специальным образом преобразованному многочлену; в [20] такой способ называется *методом Горнера*.

Более распространен метод выделения множителей\*) [8, 23, 34, 35, 40]. Последовательное выделение линейного множителя этим методом

\*) На самом деле, имеется несколько таких методов, из которых наиболее известен *метод Лина (предпоследнего остатка)*, который и имеется здесь в виду.

базируется на применении схемы Горнера в точках последовательности приближений  $x_0, x_1, x_2, \dots$  к искомому корню  $\xi$  многочлена  $P_n(x)$ . Это оказывается равнозначным простым итерациям вида

$$x_{k+1} = \frac{P_n(0)x_k}{P_n(0) - P_n(x_k)}, \quad k = 0, 1, 2, \dots,$$

что позволяет использовать известные критерии сходимости МПИ. При нахождении пары комплексно сопряженных корней выделяют квадратичный множитель. Для этих целей нужна эффективная схема деления с остатком многочлена на квадратный трехчлен, которую нетрудно получить по аналогии с выводом схемы Горнера.

В случаях, когда вычисление корней многочлена ориентировано на применение итерационных методов решения нелинейных уравнений общего вида, встает вопрос о сужении области поиска корней. Здесь опять полезную службу может сослужить схема Горнера. Справедливо утверждение: *если все коэффициенты схемы Горнера, примененной к многочлену  $P_n(x)$  в точке  $x = \beta > 0$ , положительны, то правее  $\beta$  на оси  $Ox$  действительных корней  $P_n(x)$  нет.* Это следует из тождества

$$P_n(x) \equiv (x - \beta)(b_0x^{n-1} + b_1x^{n-2} + \dots + b_{n-2}x + b_{n-1}) + b_n,$$

показывающего, что  $P_n(x) > 0$  для любых  $x \geq \beta > 0$ , в силу положительности всех  $b_i$  ( $i = 0, 1, 2, \dots, n$ ) и разности  $x - \beta$ .

Так, в примере 5.4, глядя на приведенную там таблицу, можно сказать, что многочлен  $P_4(x)$  (а значит, и  $P_5(x)$ ) не имеет действительных корней, больших 2.

Есть более конструктивные способы указания границ всех корней многочлена.

Например, можно утверждать, что за верхнюю границу положительных корней многочлена  $P_n(x)$  с  $a_0 > 0$  можно принять число

$$R = 1 + \sqrt[t]{\frac{B}{a_0}},$$

где  $t$  — индекс первого отрицательного коэффициента в ряду  $a_1, a_2, \dots, a_n$  (иначе, разность между показателем степени многочлена и показателем степени первого отрицательного члена), а  $B$  — максимум модулей всех отрицательных коэффициентов\*).

\* В одних литературных источниках такое нахождение правой границы действительных корней называют *методом Лагранжа* [20, 21], в других — *методом Маклорена* [23].

Действительно, заменив в  $P_n(x)$  неотрицательные коэффициенты  $a_1, a_2, \dots, a_{m-1}$  нулями, а все последующие — числом  $-B$ , при  $x > 1$  имеем:

$$\begin{aligned} P_n(x) &\geq a_0 x^n - B(x^{n-m} + \dots + x + 1) = \\ &= a_0 x^n - B \frac{x^{n-m+1} - 1}{x - 1} > a_0 x^n - \frac{B x^{n-m+1}}{x - 1} = \quad (5.66) \\ &= \frac{x^{n-m+1}}{x - 1} \left[ a_0 x^{m-1} (x - 1) - B \right] > \frac{x^{n-m+1}}{x - 1} \left[ a_0 (x - 1)^m - B \right]. \end{aligned}$$

Так как последнее выражение в этой цепочке равенств и неравенств равно нулю при  $x = R$  и больше нуля при  $x > R$ , то значит,  $P_n(x) > 0$  при всех  $x \geq R$ , т.е. правее  $R$  действительных корней нет.

Любой метод нахождения верхней границы положительных корней можно приспособить для нахождения нижней (левой) границы отрицательных корней. Для этого достаточно преобразовать многочлен  $P_n(x)$  заменой  $t = -x$  и для положительных корней многочлена  $(-1)^n P_n(t)$  найти верхнюю границу  $R^*$ ; тогда число  $-R^*$  будет искомой нижней границей\*). Так же, используя известный метод нахождения верхней границы в многочленах  $P_n\left(\frac{1}{x}\right)$  и  $(-1)^n P_n\left(-\frac{1}{x}\right)$ , можно найти нижнюю границу положительных и верхнюю границу отрицательных корней многочлена  $P_n(x)$ , т.е. отделить его корни от нуля.

Зачастую более хорошие результаты показывает *метод Вестерфильда* получения симметричных границ расположения всех корней многочлена [23]. Согласно ему, *модули всех корней* (в том числе и комплексных) *приведенного многочлена*  $P_n(x)$  (т.е. при  $a_0 = 1$ ) *лежат в круге, радиус которого не превосходит суммы двух наибольших из чисел*  $\sqrt[m]{|a_m|}$ , где  $m = 1, 2, \dots, n$ .

Решение проблемы изоляции корней алгебраического уравнения не имеет особой специфики, выделяющей ее из более общего случая нелинейных скалярных уравнений. Более продвинутым здесь является решение вопроса о количестве действительных корней. Можно указать простые способы выяснения возможного количества положительных и отрицательных корней по числу перемен знаков в последовательностях коэффициентов  $P_n(x)$  и  $P_n(-x)$ .

---

\* Множитель  $(-1)^n$  поставлен ради положительности старшего коэффициента преобразованного многочлена.

Так, *теорема Декарта* говорит о том, что число положительных корней уравнения (5.65) с учетом их кратностей равно числу перемен знаков в последовательности коэффициентов  $a_0, a_1, \dots, a_n$  (без учета нулевых коэффициентов) или на четное число меньше [21].

Более точный ответ на вопрос о числе действительных корней алгебраического уравнения можно получить с помощью широко известной в алгебре *теоремы Штурма* (см. [8, 20, 21] и др.). Если при этом уже затрачены усилия на составление системы Штурма, то ее можно использовать и для нахождения промежутков изоляции действительных корней.

Имеются и другие способы нахождения границ действительных и комплексных корней алгебраических уравнений, выяснения количества положительных и отрицательных корней, а также их изоляции.

Одним из наиболее эффективных методов нахождения в с е х или почти всех корней алгебраического уравнения, как вещественных, так и комплексных, является *метод Лобачевского*, предложенный выдающимся русским математиком в 1834 году<sup>1)</sup>. Основная идея метода заключается в последовательном применении операции квадрирования корней. Суть ее такова.

С помощью обобщенной теоремы Виета легко показать, что если корни  $\xi_1, \xi_2, \dots, \xi_n$  уравнения (5.65) сильно отличаются по модулю (что, кстати, гарантирует их вещественность), а именно,  $|\xi_1| \gg |\xi_2| \gg \dots \gg |\xi_n|$ , то

$$\xi_1 \approx -\frac{a_1}{a_0}, \quad \xi_2 \approx -\frac{a_2}{a_1}, \quad \dots, \quad \xi_n \approx -\frac{a_n}{a_{n-1}}.$$

При возведении их в натуральную степень будут получаться все более удаленные друг от друга числа. Одна из макроопераций метода Лобачевского (*квадрирование корней*) состоит в том, что от данного уравнения с корнями  $\xi_1, \xi_2, \dots, \xi_n$  переходят к новому уравнению той же степени с коэффициентами

$$A_i = a_i^2 + 2 \sum_{j=1}^i (-1)^j a_{i-j} a_{i+j} \quad (i = 0, 1, 2, \dots, n) \quad (5.67)$$

и корнями  $\mu_i = -\xi_i^2$  ( $i = 1, 2, \dots, n$ ). После ее многократного применения, когда в (5.67) будет практически сведена на нет роль второго слагаемого (удвоенной суммы парных произведений), извлечением корней соответствующих степеней можно приближенно найти модули корней исходного уравнения.

<sup>1)</sup> Этот метод называют также *методом Лобачевского-Греффе* или *методом Данделена* в честь швейцарского математика Греффе и французского математика Данделена, причастных к одним из первых версий метода. Впоследствии метод Лобачевского неоднократно совершенствовался.



Более подробно о методе Лобачевского, его реализациях, разных ситуациях, с которыми можно встретиться при его применении, см. в книгах [8, 21, 23, 40]. Одна из последних версий этого не самого простого метода содержится в брошюре А.А. Беланова [7] (ориентированной, правда, на программируемые микрокалькуляторы, а не на ЭВМ).

Рассмотрим простой способ приближенного вычисления наибольшего по модулю действительного корня алгебраического уравнения (5.65), который был опубликован И. Бернулли в 1732 г. и называется *методом Бернулли*.

Запишем уравнение (5.65) в виде

$$x^n = c_1 x_1^{n-1} + c_2 x_2^{n-2} + \dots + c_{n-1} x + c_n.$$

С помощью коэффициентов  $c_i = -\frac{a_i}{a_0}$  ( $i = 1, 2, \dots, n$ ) этого уравнения будем строить последовательность  $(u_{n+k})_{k=1}^{\infty}$  по рекуррентной формуле

$$u_{n+k} = c_1 u_{n+k-1} + c_2 u_{n+k-2} + \dots + c_{n-1} u_{k+1} + c_n u_k, \quad (5.68)$$

начиная этот процесс при  $k = 1, 2, \dots, n$  со значений

$$u_1 = 0, \quad u_2 = 0, \quad \dots, \quad u_{n-1} = 0, \quad u_n = k$$

(как это рекомендовано Хильдебрандом).

На основе обобщенной теоремы Виета можно выяснить, что эта последовательность обладает замечательным свойством: если  $\xi_1, \xi_2, \dots, \xi_n$  — корни многочлена  $P_n(x)$ , то

$$\begin{aligned} \xi_1 + \xi_2 + \dots + \xi_n &= u_{n+1}, \\ \xi_1^2 + \xi_2^2 + \dots + \xi_n^2 &= u_{n+2}, \\ \dots &\dots \\ \xi_1^k + \xi_2^k + \dots + \xi_n^k &= u_{n+k} \\ \dots &\dots \end{aligned}$$

Взяв отношение последующего члена  $u_{n+k+1}$  к предыдущему  $u_{n+k}$ , выраженных через степени корней, имеем:

$$\begin{aligned} \frac{u_{n+k+1}}{u_{n+k}} &= \frac{\xi_1^{k+1} + \xi_2^{k+1} + \dots + \xi_n^{k+1}}{\xi_1^k + \xi_2^k + \dots + \xi_n^k} = \\ &= \frac{\xi_1^{k+1} \left[ 1 + \left(\frac{\xi_2}{\xi_1}\right)^{k+1} + \dots + \left(\frac{\xi_n}{\xi_1}\right)^{k+1} \right]}{\xi_1^k \left[ 1 + \left(\frac{\xi_2}{\xi_1}\right)^k + \dots + \left(\frac{\xi_n}{\xi_1}\right)^k \right]}. \end{aligned} \quad (5.69)$$

Если  $|\xi_1| > |\xi_i| \quad \forall i = 2, \dots, n$ , то, очевидно,  $\frac{u_{n+k+1}}{u_{n+k}} \xrightarrow{k \rightarrow \infty} \xi_1$ . К такому же результату приходим и в случае, когда  $\xi_1$  — не простой, а  $m$ -кратный корень (выражения в квадратных скобках в числителе и в знаменателе (5.69) имеют пределом  $m$ ).

Если сходимость последовательности отношений  $\frac{u_{n+k+1}}{u_{n+k}}$  не обнаруживается, но существует предел последовательности величин

$$\frac{u_{n+k+2}}{u_{n+k}} = \xi_1^2 \frac{1 + \left(\frac{\xi_2}{\xi_1}\right)^{k+2} + \dots + \left(\frac{\xi_n}{\xi_1}\right)^{k+2}}{1 + \left(\frac{\xi_2}{\xi_1}\right)^k + \dots + \left(\frac{\xi_n}{\xi_1}\right)^k},$$

то в этом случае данное уравнение имеет два действительных корня  $\xi_1 = -\xi_2$ , наибольших по модулю. При достаточно больших  $k$  (желательно, четных) они могут быть найдены приближенным равенством  $\xi_{1,2} \approx \pm \sqrt{\frac{u_{n+k+2}}{u_{n+k}}}$ . Хаотическое поведение последовательности  $\left(\frac{u_{n+k+2}}{u_{n+k}}\right)$

говорит о том, что преобладающими являются комплексные корни. Незначительной доработкой метод Бернулли можно приспособить и для этого случая (см. [23]).

Как видно, метод Бернулли весьма близок к степенному методу решения частичной проблемы собственных значений, более тщательно рассмотренному в п.4.2. Зная в тонкостях один из этих методов, можно многое сказать о поведении другого.

## УПРАЖНЕНИЯ

5.1. Доказать существование и единственность корня уравнения

$$(x-1)^2 e^x - 7 = 0.$$

5.2. Найти промежуток локализации (единичной длины) отрицательного корня уравнения  $x^3 - x^2 + 4 = 0$ . За сколько шагов метода половинного деления можно уточнить корень до 0.1? до 0.01? до  $10^{-6}$ ?

5.3. С точностью до 0.01 решить уравнение  $\sqrt{|x-4|} - x + 1 = 0$ :

а) методом половинного деления; б) методом хорд.

5.4. С точностью до 0.001 найти положительный корень уравнения

$$x^4 - 2x - 4 = 0:$$

а) методом Ньютона; б) методом секущих.

Уточнить корень, применив к результату б) два шага метода Стеффенсена.

5.5. Подготовить алгоритм вычисления значения функции  $y = \sqrt[3]{x}$  в точке  $x = 100$  с точностью  $\varepsilon = 10^{-6}$ , пользуясь правилом Ньютона. Сделать два приближения.

5.6. На основе метода Ньютона записать итерационный процесс, который позволял бы вычислять значения  $\frac{1}{\sqrt{a}}$  при заданных вещественных значениях  $a$ , не производя делений.

5.7. Составить гибридный алгоритм “секущих-половинного деления”. Изучить его поведение на уравнении  $x^4 + (x-2)^2 - 8 = 0$ , строя приближения к положительному корню из разных пар начальных точек  $x_0, x_1$ .

5.8. Дать обоснованное заключение о факте и скорости сходимости последовательности

$$x_{k+1} = \frac{1}{4}x_k^4 - \frac{1}{2}x_k + 1, \quad k = 0, 1, 2, \dots, \quad x_0 = 0$$

на отрезке  $[0; 1]$ . С какой точностью можно приблизиться к  $\lim_{k \rightarrow \infty} x_k$  за 10 шагов? Записать итерационный процесс Ньютона, имеющий тот же предел.

5.9. Записать сходящийся процесс простых итераций для нахождения корня уравнения  $x^3 - 2x^2 - 4x - 7 = 0$ , изолированного на промежутке  $[3; 4]$ .

5.10. Подготовить расчетные формулы для нахождения корня уравнения  $2 - \lg x - x = 0$  методом простых итераций. Выбрав начальное приближение, подсчитать, за сколько итерационных шагов можно гарантировать получение корня с точностью  $\varepsilon = 10^{-6}$ . Сделать четыре шага МПИ и по два полных шага  $\Delta^2$ -процесса Эйткена и метода Вегстейна.

**5.11.** Сравнить эффективность следующих двух подходов к вычислению значений  $P_n(x_k)$  и  $P'_n(x_k)$  при нахождении алгебраического корня уравнения  $P_n(x) = 0$  методом Ньютона  $x_{k+1} = x_k - \frac{P_n(x_k)}{P'_n(x_k)}$ :

- а) вычислив степени  $x_k^i$ , использовать их затем в  $P_n(x_k)$  и  $P'_n(x_k)$ ;
- б) применить схему Горнера сначала к многочлену  $P_n(x)$  в точке  $x_k$ , затем к его неполному частному – результату первого ее применения (правомерность этого подхода обосновать!).

**5.12.** Придумать алгоритм, позволяющий найти наибольший корень многочлена  $P_n(x)$  с заданной точностью только одними проверками на положительность коэффициентов схемы Горнера.

**5.13. а)** Доказать, что если  $R_1$  и  $R_2$  – правые границы действительных корней многочленов  $Q(x)$  и  $T(x)$  соответственно, то за верхнюю границу действительных корней многочлена  $P(x) = Q(x) + T(x)$  можно принять число  $R = \max\{R_1, R_2\}$ .

б) Для многочлена  $P_5(x) = x^5 - 2x^4 + 15x^2 - 32x + 1$  найти границы положительных корней сначала непосредственно по формуле (5.66), затем на основе а), выполняя подходящее представление  $P_5(x)$  суммой  $Q(x) + T(x)$  так, чтобы получить как можно более точную границу.

в) Найти границу модулей всех корней данного многочлена  $P_5(x)$  методом Вестерфильда.

**5.14.** Пусть  $\xi_1, \xi_2$  – вещественные корни уравнения  $x^2 - px + q = 0$ . Убедиться в справедливости равенств

$$\xi_1 + \xi_2 = u_3, \xi_1^2 + \xi_2^2 = u_4, \xi_1^3 + \xi_2^3 = u_5,$$

где  $u_i$  – элементы последовательности Бернулли-Хильдебранда, задаваемой формулой (5.68).

**5.15.** Дано уравнение  $x^4 - 3x^3 - 7x^2 + 15x + 18 = 0$ . Сделать 5–6 приближений к его старшему корню методом Бернулли.

# МЕТОДЫ РЕШЕНИЯ СИСТЕМ НЕЛИНЕЙНЫХ УРАВНЕНИЙ

Рассматривается ряд методов решения систем алгебраических и трансцендентных уравнений. Среди них метод простых итераций, метод Ньютона в разных модификациях, метод Брауна. Показывается связь между данной задачей и задачей безусловной минимизации функции нескольких переменных. Проводится сравнение методов на примере решения конкретной системы. С единых позиций изучается сходимость метода Ньютона и метода, получаемого из него применением итерационного процесса Шульца для приближенного обращения матриц Якоби.

## 6.1. ВЕКТОРНАЯ ЗАПИСЬ НЕЛИНЕЙНЫХ СИСТЕМ. МЕТОД ПРОСТЫХ ИТЕРАЦИЙ

Пусть требуется решить систему уравнений

$$\begin{cases} f_1(x_1, x_2, \dots, x_n) = 0, \\ f_2(x_1, x_2, \dots, x_n) = 0, \\ \dots \dots \dots \\ f_n(x_1, x_2, \dots, x_n) = 0, \end{cases} \quad (6.1)$$

где  $f_1, f_2, \dots, f_n$  – заданные, вообще говоря, нелинейные (среди них могут быть и линейные) вещественнозначные функции  $n$  вещественных переменных  $x_1, x_2, \dots, x_n$ .

Обозначив

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}, \quad F(\mathbf{x}) = \begin{pmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \\ \vdots \\ f_n(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} f_1(x_1, x_2, \dots, x_n) \\ f_2(x_1, x_2, \dots, x_n) \\ \vdots \\ f_n(x_1, x_2, \dots, x_n) \end{pmatrix}, \quad \mathbf{0} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix},$$

данную систему (6.1) можно записать одним уравнением

$$F(\mathbf{x}) = \mathbf{0} \quad (6.1')$$

относительно векторной функции  $F$  векторного аргумента  $\mathbf{x}$ . Таким образом, исходную задачу можно рассматривать как задачу о нулях нелинейного отображения  $F: R_n \rightarrow R_n$ . В этой постановке она является прямым обобщением основной задачи предыдущей главы – задачи построения методов нахождения нулей одномерных нелинейных отображений. Фактиче-



Если начать процесс построения последовательности  $(\mathbf{x}^{(k)})$  с некоторого вектора  $\mathbf{x}^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)})^T$  и продолжить по формуле (6.3), то при определенных условиях эта последовательность со скоростью геометрической прогрессии будет приближаться к вектору  $\mathbf{x}^*$  – неподвижной точке отображения  $\Phi(\mathbf{x})$ . А именно, справедлива следующая теорема.

**Теорема 6.1.** Пусть функция  $\Phi(\mathbf{x})$  и замкнутое множество  $M \subseteq D(\Phi) \subseteq R_n$  таковы, что:

$$1) \Phi(\mathbf{x}) \in M \quad \forall \mathbf{x} \in M;$$

$$2) \exists q < 1: \|\Phi(\mathbf{x}) - \Phi(\tilde{\mathbf{x}})\| \leq q \cdot \|\mathbf{x} - \tilde{\mathbf{x}}\| \quad \forall \mathbf{x}, \tilde{\mathbf{x}} \in M.$$

Тогда  $\Phi(\mathbf{x})$  имеет в  $M$  единственную неподвижную точку  $\mathbf{x}^*$ ; последовательность  $(\mathbf{x}^{(k)})$ , определяемая МПИ (6.3), при любом  $\mathbf{x}^{(0)} \in M$  сходится к  $\mathbf{x}^*$  и справедливы оценки

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{q}{1-q} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \frac{q^k}{1-q} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \quad \forall k \in \mathbb{N}.$$

Доказательство этой теоремы почти полностью повторяет доказательство теоремы 5.7 и даже несколько проще его. Однако и практическая ценность такой теоремы не так велика из-за неконструктивности ее условий. В случаях, когда имеется хорошее начальное приближение  $\mathbf{x}^{(0)}$  к решению  $\mathbf{x}^*$ , большой интерес для приложений может представить следующий аналог теоремы 5.8.

**Теорема 6.2.** Пусть  $\Phi(\mathbf{x})$  дифференцируема<sup>\*)</sup> в замкнутом шаре<sup>\*\*)</sup>  $S(\mathbf{x}^{(0)}, r) \subseteq D(\Phi)$ , причем  $\exists q \in (0; 1): \sup_{\mathbf{x} \in S} \|\Phi'(\mathbf{x})\| \leq q$ . Тогда

если центр  $\mathbf{x}^{(0)}$  и радиус  $r$  шара  $S$  таковы, что  $\|\mathbf{x}^{(0)} - \Phi(\mathbf{x}^{(0)})\| \leq r(1-q)$ , то справедливо заключение теоремы 6.1 с  $M = S$ .

Как видим, теорема 6.2 – это просто соответствующим образом отредактированная для многомерного случая теорема 5.8. Переложение доказательства теоремы 5.8 на этот случай менее тривиально из-за того, что

<sup>\*)</sup> Здесь и далее всюду под дифференцируемостью понимается дифференцируемость по Фреше. См приложение 2.

<sup>\*\*) Т.е. на множестве  $S$  точек  $\mathbf{x} \in R_n$  таких, что  $\|\mathbf{x} - \mathbf{x}^{(0)}\| \leq r$ .</sup>





Заметим, что как и для линейных систем, отдельные уравнения в методе (6.4) неравноправны, т.е. перемена местами уравнений системы (6.2) может изменить в каких-то пределах число итераций и вообще ситуацию со сходимостью последовательности итераций. Чтобы применить метод простых итераций (6.3) или его зейделеву модификацию (6.4) к исходной системе (6.1), нужно, как и в скалярном случае, сначала тем или иным способом привести ее к виду (6.2). Это можно сделать, например, умножив (6.1') на некоторую неособенную  $n \times n$ -матрицу  $-A$  и прибавив к обеим частям уравнения  $-AF(x) = 0$  вектор неизвестных  $x$ . Полученная система

$$x = x - AF(x)$$

эквивалентна данной и имеет вид задачи о неподвижной точке (6.2'). Проблема теперь состоит лишь в подборе матричного параметра  $A$  такого, при котором вектор-функция  $\Phi(x) := x - AF(x)$  обладала бы нужными свойствами.

## 6.2. МЕТОД НЬЮТОНА, ЕГО РЕАЛИЗАЦИИ И МОДИФИКАЦИИ

Пусть  $(A_k)$  – некоторая последовательность невырожденных вещественных  $n \times n$ -матриц. Тогда, очевидно, последовательность задач

$$x = x - A_k F(x), \quad k = 0, 1, 2, \dots,$$

имеет те же решения, что и исходное уравнение (6.1'), и для приближенного нахождения этих решений можно формально записать итерационный процесс

$$x^{(k+1)} = x^{(k)} - A_k F(x^{(k)}), \quad k = 0, 1, 2, \dots, \quad (6.5)$$

имеющий вид метода простых итераций (6.3) при  $\Phi(x) := \Phi_k(x) := x - A_k F(x)$ . В случае  $A_k \equiv A$  это, как показано в конце предыдущего пункта, – действительно МПИ с линейной сходимостью последовательности  $(x^{(k)})$ . Если же  $A_k$  различны при разных  $k$ , то формула (6.5) определяет большое семейство итерационных методов с матричными параметрами  $A_k$ . Рассмотрим некоторые из методов этого семейства.

Положим  $A_k := [F'(x^{(k)})]^{-1}$ , где

$$F'(\mathbf{x}) = J(\mathbf{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \dots & \frac{\partial f_n}{\partial x_n} \end{pmatrix}$$

— матрица Якоби вектор-функции  $F(\mathbf{x})$ . Подставив это  $A_k$  в (6.5), получаем явную формулу *метода Ньютона*

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left[ F'(\mathbf{x}^{(k)}) \right]^{-1} F(\mathbf{x}^{(k)}), \quad (6.6)$$

обобщающего на многомерный случай скалярный метод Ньютона (5.14). Эту формулу, требующую обращения матриц на каждой итерации, можно переписать в неявном виде:

$$F'(\mathbf{x}^{(k)}) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = -F(\mathbf{x}^{(k)}). \quad (6.7)$$

Применение (6.7) предполагает при каждом  $k = 0, 1, 2, \dots$  решение линейной алгебраической системы

$$F'(\mathbf{x}^{(k)}) \mathbf{p}^{(k)} = -F(\mathbf{x}^{(k)})$$

относительно векторной *поправки*  $\mathbf{p}^{(k)} = (p_1^{(k)}, p_2^{(k)}, \dots, p_n^{(k)})^T$ , а затем прибавление этой поправки к текущему приближению для получения следующего:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{p}^{(k)}.$$

К решению таких линейных систем можно привлекать самые разные методы как прямые, так и итерационные (см. гл. 2, 3) в зависимости от размерности  $n$  решаемой задачи и специфики матриц Якоби  $J(\mathbf{x}^{(k)})$  (например, можно учитывать их симметрию, разреженность и т.п.).

Сравнивая (6.7) с формальным разложением  $F(\mathbf{x})$  в ряд Тейлора

$$F(\mathbf{x}) = F(\mathbf{x}^{(k)}) + F'(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2!} F''(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)})^2 + \dots,$$

видим, что последовательность  $(\mathbf{x}^{(k)})$  в методе Ньютона получается в результате подмены при каждом  $k = 0, 1, 2, \dots$  нелинейного уравнения  $F(\mathbf{x}) = \mathbf{0}$  или, что то же (при достаточной гладкости  $F(\mathbf{x})$ ), уравнения

$$F(\mathbf{x}^{(k)}) + F'(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2!} F''(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)})^2 + \dots = \mathbf{0}$$

линейным уравнением

$$F(\mathbf{x}^{(k)}) + F'(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)}) = \mathbf{0},$$

т.е. пошаговой линеаризацией<sup>\*)</sup>. Как следствие этого факта, более тщательно изученного в  $R_1$  (см. гл.5 п.4), можно рассчитывать, что при достаточной гладкости  $F(\mathbf{x})$  и достаточно хорошем начальном приближении  $\mathbf{x}^{(0)}$  сходимость порождаемой методом Ньютона последовательности  $(\mathbf{x}^{(k)})$  к решению  $\mathbf{x}^*$  будет квадратичной и в многомерном случае. Имеется ряд теорем, устанавливающих это при тех или иных предположениях (см. [6, 21, 29, 40 и др.]). В частности, одна из таких теорем приводится ниже (теорема 6.5 в п.6).

Новым, по сравнению со скалярным случаем, фактором, осложняющим применение метода Ньютона к решению  $n$ -мерных систем, является необходимость решения  $n$ -мерных линейных задач на каждой итерации (обращения матриц в (6.6) или решения СЛАУ в (6.7)), вычислительные затраты на которые растут с ростом  $n$ , вообще говоря, непропорционально быстро. Уменьшение таких затрат – одно из направлений модификации метода Ньютона.

Если матрицу Якоби  $F'(\mathbf{x})$  вычислить и обратить лишь один раз – в начальной точке  $\mathbf{x}^{(0)}$ , то от метода Ньютона (6.6) придем к *модифицированному методу Ньютона*

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left[ F'(\mathbf{x}^{(0)}) \right]^{-1} F(\mathbf{x}^{(k)}). \quad (6.8)$$

Этот метод требует значительно меньших вычислительных затрат на один итерационный шаг, но итераций при этом может потребоваться значительно больше для достижения заданной точности по сравнению с основным методом Ньютона (6.6), поскольку, являясь частным случаем МПИ (с  $A := \left[ F'(\mathbf{x}^{(0)}) \right]^{-1}$ ), он имеет лишь скорость сходимости геометрической прогрессии<sup>\*\*)</sup>.

Компромиссный вариант – это вычисление и обращение матриц Якоби не на каждом итерационном шаге, а через несколько шагов (иногда такие методы называют *рекурсивными* [50]).

Например, простое чередование основного (6.6) и модифицированного (6.8) методов Ньютона приводит к итерационной формуле

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - A_k F(\mathbf{x}^{(k)}) - A_k F(\mathbf{x}^{(k)} - A_k F(\mathbf{x}^{(k)})), \quad (6.9)$$

<sup>\*)</sup> Обратим внимание на некоторую сознательную некорректность использования здесь термина «линейный», идущую от «школьной» привычки называть функцию  $y = ax + b$  линейной, хотя она не является линейным оператором, ибо не удовлетворяет условию однородности. В [22] в таких случаях используется термин «аффинная аппроксимация».

<sup>\*\*) Независимо от МПИ линейная сходимость модифицированного метода Ньютона (6.8) обосновывается в п.6 теоремой 6.6.</sup>

где  $A_k := [F'(x^{(k)})]^{-1}$ ,  $k = 0, 1, 2, \dots$ . За  $x^{(k)}$  здесь принимается результат последовательного применения одного шага основного, а затем одного шага модифицированного метода, т.е. *двухступенчатого процесса*

$$\begin{cases} z^{(k)} = x^{(k)} - A_k F(x^{(k)}), \\ x^{(k+1)} = z^{(k)} - A_k F(z^{(k)}). \end{cases} \quad (6.10)$$

Доказано, что такой процесс при определенных условиях порождает кубически сходящуюся последовательность  $(x^{(k)})$ .

Задачу обращения матриц Якоби на каждом  $k$ -м шаге метода Ньютона (6.6) можно попытаться решать не точно, а приближенно. Для этого можно применить, например, итерационный процесс Шульца (см. п.3.6), ограничиваясь минимумом – всего одним шагом процесса второго порядка, в котором за начальную матрицу принимается матрица, полученная в результате предыдущего  $(k-1)$ -го шага. Таким образом, приходим к *методу Ньютона с последовательной аппроксимацией обратных матриц*:

$$\begin{cases} x^{(k+1)} = x^{(k)} - A_k F(x^{(k)}), \\ \Psi_k = E - F'(x^{(k+1)})A_k, \quad A_{k+1} = A_k + A_k \Psi_k, \end{cases} \quad (6.11)$$

где  $k = 0, 1, 2, \dots$ , а  $x^{(0)}$  и  $A_0$  – начальные вектор и матрица ( $\approx [F'(x^{(0)})]^{-1}$ ) соответственно<sup>\*)</sup>. Этот метод (будем называть его более

коротко ААМН – *аппроксимационный аналог метода Ньютона*) имеет простую схему вычислений – поочередное выполнение векторных в первой строке и матричных во второй строке его записи (6.11) операций. Скорость его сходимости почти так же высока, как и у метода Ньютона. Как будет показано в п.6.6, последовательность  $(x^{(k)})$  может квадратично сходиться к решению  $x^*$  уравнения  $F(x) = 0$  (при этом матричная последовательность

$(A_k)$  также квадратично сходится к  $A^* := [F'(x^*)]^{-1}$ , т.е. в нормально развивающемся итерационном процессе (6.11) должна наблюдаться достаточно быстрая сходимость ( $\|\Psi_k\|$  к нулю).

Применение той же последовательной аппроксимации обратных матриц к простейшему рекурсивному методу Ньютона (6.9) или, что то же,

<sup>\*)</sup> Требования к степени близости  $A_0$  к  $[F'(x^{(0)})]^{-1}$ , наряду с другими требованиями, гарантирующими сходимость метода, см. далее в теореме 6.7 и следствии 6.1.

к двухступенчатому процессу (6.10) определяет его аппроксимационный аналог

$$\begin{cases} \mathbf{z}^{(k)} = \mathbf{x}^{(k)} - \mathbf{A}_k F(\mathbf{x}^{(k)}), & \mathbf{x}^{(k+1)} = \mathbf{z}^{(k)} - \mathbf{A}_k F(\mathbf{z}^{(k)}), \\ \Psi_k = \mathbf{E} - F'(\mathbf{x}^{(k+1)})\mathbf{A}_k, & \mathbf{A}_{k+1} = \mathbf{A}_k + \mathbf{A}_k \Psi_k, \end{cases} \quad (6.12)$$

который, как и (6.9), также можно отнести к методам третьего порядка. Доказательство кубической сходимости этого метода требует уже более жестких ограничений на свойства  $F(\mathbf{x})$  и близость  $\mathbf{x}^{(0)}$  к  $\mathbf{x}^*$ ,  $\mathbf{A}_0$  к  $[F'(\mathbf{x}^{(0)})]^{-1}$ , чем в предыдущем методе. Заметим, что к улучшению сходимости здесь может привести повышение порядка аппроксимации обратных матриц, например, за счет добавления еще одного слагаемого в формуле для подсчета  $\mathbf{A}_{k+1}$  (см. (3.34)):

$$\mathbf{A}_{k+1} = \mathbf{A}_k + \mathbf{A}_k \Psi_k + \mathbf{A}_k \Psi_k^2.$$

На базе метода Ньютона (6.6) можно построить близкий к нему по поведению итерационный процесс, не требующий вычисления производных. Сделаем это, заменив частные производные в матрице Якоби  $J(\mathbf{x})$  разностными отношениями, т.е. подставив в формулу (6.5) вместо  $\mathbf{A}_k$  матрицу  $[J(\mathbf{x}^{(k)}, \mathbf{h}^{(k)})]^{-1}$ , где

$$J(\mathbf{x}, \mathbf{h}) := \left( \frac{f_i(x_1, \dots, x_j + h_j, \dots, x_n) - f_i(x_1, \dots, x_j, \dots, x_n)}{h_j} \right)_{i,j=1}^n.$$

При удачном задании последовательности малых векторов  $\mathbf{h}^{(k)} = (h_1^{(k)}, \dots, h_n^{(k)})^T$  (постоянной или сходящейся к нулю) полученный таким путем *разностный* (или иначе, *дискретный*) метод Ньютона

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - [J(\mathbf{x}^{(k)}, \mathbf{h}^{(k)})]^{-1} F(\mathbf{x}^{(k)}) \quad (6.13)$$

имеет сверхлинейную, вплоть до квадратичной, скорость сходимости и обобщает на многомерный случай метод (5.29). При задании векторного параметра  $\mathbf{h}$  – шага дискретизации – следует учитывать точность машинных вычислений (*macheps*), точность вычисления значения функций  $f_i$ , средние значения получаемых приближений (см. по этому поводу [22]).

Можно связать задание последовательности  $(\mathbf{h}^{(k)})$  с какой-либо сходящейся к нулю векторной последовательностью, например, с последовательностью *невязок*  $(F(\mathbf{x}^{(k)}))$  или *поправок*  $(\mathbf{p}^{(k)})$ . Так, полагая

$h_j^{(k)} := x_j^{(k-1)} - x_j^{(k)}$ , где  $j = 1, \dots, n$ , а  $k = 1, 2, \dots$ , приходим к *простейшему методу секущих* – обобщению скалярного метода секущих (5.32):

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left[ B(\mathbf{x}^{(k)}, \mathbf{x}^{(k-1)}) \right]^{-1} F(\mathbf{x}^{(k)}), \quad (6.14)$$

где

$$B(\mathbf{x}^{(k)}, \mathbf{x}^{(k-1)}) := \left( \frac{f_i(\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_j^{(k-1)}, \dots, \mathbf{x}_n^{(k)}) - f_i(\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_j^{(k)}, \dots, \mathbf{x}_n^{(k)})}{x_j^{(k-1)} - x_j^{(k)}} \right)_{i,j=1}^n,$$

$k = 1, 2, 3, \dots$

Этот метод является двухшаговым и требует задания двух начальных точек  $\mathbf{x}^{(0)}$  и  $\mathbf{x}^{(1)}$ . Как было показано в гл.5, при  $n=1$  сходимость метода (6.14) имеет порядок  $\frac{1+\sqrt{5}}{2}$ . Можно рассчитывать на такую же скорость и в многомерном случае.

К методу секущих так же, как и к методу Ньютона, можно применить пошаговую аппроксимацию обратных матриц на основе метода Шульца. Расчетные формулы этой модификации легко выписать, заменив в совокупности формул ААН (6.11) матрицу  $F'(\mathbf{x}^{(k+1)})$  на матрицу  $B(\mathbf{x}^{(k+1)}, \mathbf{x}^{(k)})$  из (6.14).

**Замечание 6.1.** Если в одномерном случае разные подходы к линеаризации  $f(x)$  приводили к одной и той же формуле секущих (5.32), то для функции  $n$  переменных  $F(x)$  известно несколько разных обобщений этой формулы в зависимости от того, на какую основу положена линеаризация  $F(x)$  в текущей точке. Предложенный здесь простейший метод секущих (6.14) является одним из семейства методов секущих, базирующихся на аппроксимации матриц Якоби. Линейная интерполяция  $F(x)$  в  $R_n$  может привести к ряду других методов секущих (см., например, [42]). Среди множества таких методов особый интерес представляет *метод секущих Бройдена* [22]. В простейшем варианте его можно реализовать по формулам

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \mathbf{B}_k^{-1} F(\mathbf{x}^{(k)}), \\ \mathbf{B}_{k+1} &= \mathbf{B}_k + \frac{F(\mathbf{x}^{(k+1)}) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})^T}{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2^2}. \end{aligned} \quad (6.15)$$

В этом процессе, если его начать с  $\mathbf{B}_0 \approx F'(\mathbf{x}^{(0)})$ , матрицы  $\mathbf{B}_k$ , не являющиеся прямыми аппроксимациями матриц Якоби  $F'(\mathbf{x}^{(k)})$ , а получаемые в результате уточнения («пересчета») соответствующих матриц предыдущего шага, в какой-то мере близки к  $F'(\mathbf{x}^{(k)})$ , однако в общем случае нельзя утверждать сходимость  $(\mathbf{B}_k)$  к  $F'(\mathbf{x}^*)$  (что может сыграть даже полезную роль при вырождении матриц  $F'(\mathbf{x}^*)$ ).

**Замечание 6.2.** Для останова процесса вычислений в быстроходящихся методах таких, как метод Ньютона, методы секущих и т.п., часто вполне успешно применяют простой критерий:

$$\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| < \varepsilon \Rightarrow \text{stop, } \mathbf{x}^* \approx \mathbf{x}^{(k)}. \quad (6.16)$$

Это можно объяснить двумя причинами. Во-первых, оценки погрешности здесь довольно «дороги». Имеется в виду как их получение (особенно для различных модификаций базовых методов), так и их реальное применение. Во-вторых, в силу своей быстрой сходимости, к моменту достижения требуемой малости нормы поправки эти методы набирают такую инерцию, что зачастую «проскакивают» установленный порог точности, т.е. выход по критерию (6.16) дает значение  $\|\mathbf{x}^* - \mathbf{x}^{(k)}\|$  значительно (иногда на несколько порядков) меньшее, чем  $\varepsilon$ , см. численный пример в п.6.5.<sup>\*)</sup> Отслеживать факт сходимости в процессе итераций для того, чтобы реагировать на возможную расходимость в случаях, когда заранее не обеспечены условия сходимости применяемого метода, можно с помощью текущих проверок на уменьшение от шага к шагу поправок и невязок, т.е. выполнение неравенств

$$\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| < \|\mathbf{x}^{(k-1)} - \mathbf{x}^{(k-2)}\| \quad \text{и} \quad \|F(\mathbf{x}^{(k)})\| < \|F(\mathbf{x}^{(k-1)})\|.$$

<sup>\*)</sup> Как отмечается в [22], для метода Ньютона установлены неравенства

$$0.5 \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \|\mathbf{x}^* - \mathbf{x}^{(k-1)}\| \leq 2 \|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|.$$

подтверждающие высказанные соображения. В соответствии с этими неравенствами срабатывание критерия (6.16) для метода Ньютона означает, что уже вектор  $\mathbf{x}^{(k-1)}$  может быть принят за решение  $\mathbf{x}^*$ , но поскольку подсчитан вектор  $\mathbf{x}^{(k)}$ , полагаем

$$\mathbf{x}^* \approx \mathbf{x}^{(k)}$$

### 6.3. МЕТОД БРАУНА

В отличие от пошаговой линеаризации векторной функции  $F(\mathbf{x})$ , приведшей к методу Ньютона (6.6), Брауном (1966 г.)<sup>\*)</sup> предложено проводить на каждом итерационном шаге поочередную линеаризацию компонент вектор – функции  $F(\mathbf{x})$ , т.е. линеаризовать в системе (6.1) сначала функцию  $f_1$ , затем  $f_2$  и т.д., и последовательно решать получаемые таким образом уравнения. Чтобы не затенять эту идею громоздкими выкладками и лишними индексами, рассмотрим вывод расчетных формул метода Брауна в двумерном случае.

Пусть требуется найти решение системы

$$\begin{cases} f(x, y) = 0, \\ g(x, y) = 0, \end{cases} \quad (6.17)$$

и пусть уже получены приближения  $x_k, y_k$ .

Подменим первое уравнение системы (6.17) линейным, полученным по формуле Тейлора для функции двух переменных:

$$f(x, y) \approx f(x_k, y_k) + f'_x(x_k, y_k)(x - x_k) + f'_y(x_k, y_k)(y - y_k) = 0.$$

Отсюда выражаем  $x$  (обозначим этот результат через  $\tilde{x}$ ):

$$\tilde{x} = x_k - \frac{1}{f'_x(x_k, y_k)} [f(x_k, y_k) + f'_y(x_k, y_k)(y - y_k)]. \quad (6.18)$$

При  $y = y_k$  находим значение  $\tilde{x}_k$  переменной  $\tilde{x}$ :

$$\tilde{x}_k = x_k - \frac{f(x_k, y_k)}{f'_x(x_k, y_k)},$$

которое будем считать лишь промежуточным приближением (т.е. не  $x_{k+1}$ ), поскольку оно не учитывает второго уравнения системы (6.17).

Подставив в  $g(x, y)$  вместо  $x$  переменную  $\tilde{x} = \tilde{x}(y)$ , придем к некоторой функции  $G(y) = g(\tilde{x}(y), y)$  только одной переменной  $y$ . Это позволяет линеаризовать второе уравнение системы (6.17) с помощью формулы Тейлора для функции одной переменной:

$$g(\tilde{x}, y) \approx G(y_k) + G'(y_k)(y - y_k) = 0. \quad (6.19)$$

При нахождении производной  $G'(y)$  нужно учесть, что  $G(y) = g(\tilde{x}(y), y)$  есть сложная функция одной переменной  $y$ , т.е. применить формулу полной производной

$$G'(y) = g'_x(\tilde{x}, y) \cdot \tilde{x}'_y + g'_y(\tilde{x}, y). \quad (6.20)$$

Дифференцируя по  $y$  равенство (6.18), получаем

$$\tilde{x}'_y = - \frac{f'_y(x_k, y_k)}{f'_x(x_k, y_k)}.$$

<sup>\*)</sup> Ссылки на первоисточники можно найти в [42].



Подстановка последнего в (6.20) при  $y = y_k$ ,  $\tilde{x} = \tilde{x}_k$  дает

$$G'(y_k) = -g'_k(\tilde{x}_k, y_k) \cdot \frac{f'_y(x_k, y_k)}{f'_x(x_k, y_k)} + g'_y(\tilde{x}_k, y_k).$$

При известных значениях  $G(y_k) = g(\tilde{x}_k, y_k)$  и  $G'(y_k)$  теперь можно разрешить линейное уравнение (6.19) относительно  $y$  (назовем полученное значение  $y_{k+1}$ ):

$$y_{k+1} = y_k - \frac{G(y_k)}{G'(y_k)} = y_k - \frac{g(\tilde{x}_k, y_k) f'_x(x_k, y_k)}{f'_x(x_k, y_k) g'_y(\tilde{x}_k, y_k) - f'_y(x_k, y_k) g'_x(\tilde{x}_k, y_k)}.$$

Заменяя в (6.18) переменную  $y$  найденным значением  $y_{k+1}$ , приходим к значению  $x_{k+1}$ :

$$x_{k+1} = \tilde{x}(y_{k+1}) = x_k - \frac{1}{f'_x(x_k, y_k)} [f(x_k, y_k) + f'_y(x_k, y_k)(y_{k+1} - y_k)].$$

Таким образом, реализация метода Брауна решения двумерных нелинейных систем вида (6.17) сводится к следующему.

При выбранных начальных значениях  $x_0, y_0$  каждое последующее приближение по *методу Брауна* находится при  $k = 0, 1, 2, \dots$  с помощью совокупности формул

$$\begin{aligned} \tilde{x}_k &= x_k - \frac{f(x_k, y_k)}{f'_x(x_k, y_k)}, \\ q_k &= \frac{g(\tilde{x}_k, y_k) \cdot f'_x(x_k, y_k)}{f'_x(x_k, y_k) g'_y(x_k, y_k) - f'_y(x_k, y_k) g'_x(\tilde{x}_k, y_k)}, \\ p_k &= \frac{f(x_k, y_k) - q_k f'_y(x_k, y_k)}{f'_x(x_k, y_k)}, \\ x_{k+1} &= x_k - p_k, y_{k+1} = y_k - q_k, \end{aligned}$$

счет по которым должен выполняться в той очередности, в которой они записаны.

Вычисления в методе Брауна естественно заканчивать, когда выполнится неравенство  $\max\{|p_{k-1}|, |q_{k-1}|\} < \varepsilon$  (с результатом  $(x^*, y^*) \approx (x_k, y_k)$ ). В ходе вычислений следует контролировать немалость знаменателей расчетных формул. Заметим, что функции  $f$  и  $g$  в этом методе не равноправны, и перемена их ролями может изменить ситуацию со сходимостью.

Указывая на наличие *квадратичной сходимости* метода Брауна, в [42] отмечают, что рассчитывать на его большую по сравнению с методом Ньютона эффективность в смысле вычислительных затрат можно лишь в случае, когда фигурирующие в нем частные производные заменяются разностными отношениями.

## 6.4. О РЕШЕНИИ НЕЛИНЕЙНЫХ СИСТЕМ МЕТОДАМИ СПУСКА

Общий недостаток всех рассмотренных выше методов решения систем нелинейных уравнений – это сугубо локальный характер сходимости, затрудняющий их применение в случаях (довольно типичных), когда имеются проблемы с выбором хороших начальных приближений. Помощь здесь может прийти со стороны численных методов оптимизации – ветви вычислительной математики, обычно выделяемой в самостоятельную дисциплину. Для этого нужно поставить задачу нахождения решений данной нелинейной системы как оптимизационную или, иначе, экстремальную задачу. Ради геометрической интерпретации проводимых ниже рассуждений и их результатов ограничимся, как и в предыдущем пункте, рассмотрением системы, состоящей из двух уравнений с двумя неизвестными, т.е. системы (6.17).

Из функций  $f$  и  $g$  системы (6.17) образуем новую функцию

$$\Phi(x, y) = f^2(x, y) + g^2(x, y). \quad (6.21)$$

Так как эта функция неотрицательна, то найдется точка<sup>\*)</sup>  $(x^*, y^*)$  такая, что

$$\Phi(x, y) \geq \Phi(x^*, y^*) \geq 0 \quad \forall (x, y) \in R_2,$$

т.е.  $(x^*, y^*) = \arg \min_{x \in R_2} \Phi(x, y)$ . Следовательно, если тем или иным способом

удастся получить точку  $(x^*, y^*)$ , минимизирующую функцию  $\Phi(x, y)$ , и если при этом окажется, что  $\min_{(x, y) \in R_2} \Phi(x, y) = \Phi(x^*, y^*) = 0$ , то  $(x^*, y^*)$  – искомого решение системы (6.17), поскольку

$$\Phi(x^*, y^*) = 0 \Leftrightarrow \begin{cases} f(x^*, y^*) = 0, \\ g(x^*, y^*) = 0. \end{cases}$$

Последовательность точек  $(x_k, y_k)$  – приближений к точке  $(x^*, y^*)$  минимума  $\Phi(x, y)$  – обычно получают по рекуррентной формуле

$$\begin{pmatrix} x_{k+1} \\ y_{k+1} \end{pmatrix} = \begin{pmatrix} x_k \\ y_k \end{pmatrix} + \alpha_k \begin{pmatrix} p_k \\ q_k \end{pmatrix}, \quad k = 0, 1, 2, \dots, \quad (6.22)$$

где  $(p_k, q_k)^T$  – вектор, определяющий *направление минимизации*, а  $\alpha_k$  – скалярная величина, характеризующая величину шага минимизации (*шаговый множитель*). Учитывая геометрический смысл задачи минимизации функции двух переменных  $\Phi(x, y)$  – «спуск на дно» поверхности  $z = \Phi(x, y)$

<sup>\*)</sup> Вообще говоря, не единственная.

(см. рис. 6.1), итерационный метод (6.22) можно назвать *методом спуска*, если вектор  $(p_k, q_k)^T$  при каждом  $k$  является *направлением спуска* (т.е. существует такое  $\alpha > 0$ , что  $\Phi(x_k + \alpha p_k, y_k + \alpha q_k) < \Phi(x_k, y_k)$ ) и если множитель  $\alpha_k$  подбирается так, чтобы выполнялось *условие релаксации*  $\Phi(x_{k+1}, y_{k+1}) < \Phi(x_k, y_k)$ , означающее переход на каждой итерации в точку с меньшим значением минимизируемой функции.

Итак, при построении численного метода вида (6.22) минимизации функции  $\Phi(x, y)$  следует ответить на два главных вопроса: как выбирать направление спуска  $(p_k, q_k)^T$  и как регулировать длину шага в выбранном направлении с помощью скалярного параметра – шагового множителя  $\alpha_k$ . Приведем наиболее простые соображения по этому поводу.

При выборе направления спуска естественным является выбор такого направления, в котором минимизируемая функция убывает наиболее быстро. Как известно из математического анализа функций нескольких переменных, направление наибольшего возрастания функции в данной точке показывает ее градиент в этой точке. Поэтому примем за направление спуска вектор

$$\begin{pmatrix} p_k \\ q_k \end{pmatrix} := -\text{grad } \Phi(x_k, y_k) = -\begin{pmatrix} \Phi'_x(x_k, y_k) \\ \Phi'_y(x_k, y_k) \end{pmatrix}$$

– антиградиент функции  $\Phi(x, y)$ . Таким образом, из семейства методов (6.22) выделяем *градиентный метод*

$$\begin{pmatrix} x_{k+1} \\ y_{k+1} \end{pmatrix} := \begin{pmatrix} x_k \\ y_k \end{pmatrix} - \alpha_k \begin{pmatrix} \Phi'_x(x_k, y_k) \\ \Phi'_y(x_k, y_k) \end{pmatrix}. \quad (6.23)$$

Оптимальный шаг в направлении антиградиента – это такой шаг, при котором значение  $\Phi(x_{k+1}, y_{k+1})$  – наименьшее среди всех других значений  $\Phi(x, y)$  в этом фиксированном направлении, т.е. когда точка  $(x_{k+1}, y_{k+1})$  является точкой условного минимума. Следовательно, можно рассчитывать на наиболее быструю сходимость метода (6.23), если полагать в нем

$$\alpha_k = \arg \min_{\alpha > 0} \Phi(x_k - \alpha \Phi'_x(x_k, y_k), y_k - \alpha \Phi'_y(x_k, y_k)). \quad (6.24)$$

Такой выбор шагового множителя, называемый *исчерпывающим спуском*, вместе с формулой (6.23) определяет *метод наискорейшего спуска*.

Геометрическая интерпретация этого метода хорошо видна из рис. 6.1, 6.2. Характерны девяностоградусные изломы траектории наискорейшего спуска, что объясняется исчерпаемостью спуска и свойством

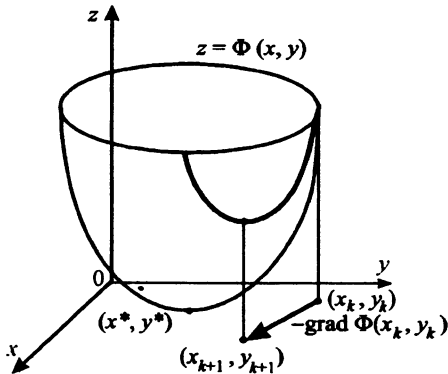


Рис. 6.1. Пространственная интерпретация метода наискорейшего спуска для функции (6.21)

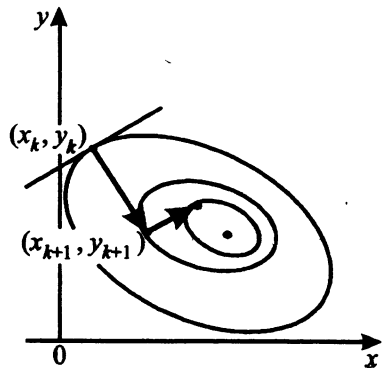


Рис. 6.2. Траектория наискорейшего спуска для функции (6.21)

градиента (а значит, и антиградиента) быть перпендикулярным к касательной к линии уровня в соответствующей точке<sup>\*)</sup>.

Наиболее типичной является ситуация, когда найти точно (аналитическими методами) оптимальное значение  $\alpha_k$  не удастся. Следовательно, приходится делать ставку на применение каких-либо численных методов одномерной минимизации и находить  $\alpha_k$  в (6.24) лишь приближенно.

Несмотря на то, что задача нахождения минимума функции одной переменной  $\varphi_k(\alpha) = \Phi(x_k - \alpha \Phi'_x(x_k, y_k), y_k - \alpha \Phi'_y(x_k, y_k))$  намного проще, чем решаемая задача, применение тех или иных численных методов нахождения значений  $\alpha_k = \text{arg min } \varphi_k(\alpha)$  с той или иной точностью требует вычисления нескольких значений минимизируемой функции. Так как это нужно делать на каждом итерационном шаге, то при большом числе шагов реализация метода наискорейшего спуска в чистом виде является

<sup>\*)</sup> Представив образно процесс движения текущей точки по траектории наискорейшего спуска на участке от положения  $(x_k, y_k)$  до  $(x_{k+1}, y_{k+1})$ , можно отметить, что спуск характеризуется пересечением линий уровня  $\Phi(x, y) = c$  все с меньшими значениями  $c$  до тех пор, пока не произойдет касание некоторой линии уровня (сплошь заполняющих  $D(\Phi)$ ); дальнейшее движение в этом направлении приведет к пересечению линий уровня  $\Phi(x, y) = c$  с увеличивающимися значениями параметра  $c$ .

достаточно высокочувствительной. Существуют эффективные схемы приближенного вычисления квазиоптимальных  $\alpha_k$ , в которых учитывается специфика минимизируемых функций (типа сумм квадратов функций) [21].

Зачастую успешной является такая стратегия градиентного метода, при которой шаговый множитель  $\alpha_k$  в (6.23) берется либо сразу достаточно малым постоянным, либо предусматривается его уменьшение, например, делением пополам для удовлетворения условию релаксации на очередном шаге. Хотя каждый отдельный шаг градиентного метода при этом, вообще говоря, далек от оптимального, такой процесс по числу вычислений функции может оказаться более эффективным, чем метод наискорейшего спуска.

Главное достоинство градиентных методов решения нелинейных систем – глобальная сходимость. Нетрудно доказать, что процесс градиентного спуска приведет к какой-либо точке минимума функции из любой начальной точки. При определенных условиях найденная точка минимума будет искомым решением исходной нелинейной системы.

Главный недостаток – медленная сходимость. Доказано, что сходимость таких методов – лишь линейная, причем, если для многих методов, таких как метод Ньютона, характерно ускорение сходимости при приближении к решению, то здесь имеет место скорее обратное. Поэтому есть резон в построении гибридных алгоритмов, которые начинали бы поиск искомой точки – решения данной нелинейной системы – глобально сходящимся градиентным методом, а затем производили уточнение каким-то быстросходящимся методом, например, методом Ньютона (разумеется, если данные функции обладают нужными свойствами).

Разработан ряд методов решения экстремальных задач, которые соединяют в себе низкую требовательность к выбору начальной точки и высокую скорость сходимости. К таким методам, называемым *квазиньютоновскими*, можно отнести, например, *метод переменной метрики (Дэвидона-Флетчера-Пауэлла)*, *симметричный* и *положительно определенный методы секущих* (на основе формулы пересчета Бройдена, см. замечание 6.1), а также уже упоминавшийся ранее применительно к СЛАУ (см. гл. 3) *метод сопряженных градиентов*.

При наличии негладких функций в решаемой задаче следует отказаться от использования производных или их аппроксимаций и прибегнуть к так называемым *методам прямого поиска (циклического покоординатного спуска, Хука и Дживса, Розенброка и т.п.)*. Описание упомянутых и многих других методов такого типа можно найти в учебной и в специальной литературе, посвященной решению экстремальных задач (см., например, [12, 22, 37]).

**Замечание 6.3.** Для разных семейств численных методов минимизации могут быть рекомендованы свои критерии останова итерационного процесса. Например, учитывая, что в точке минимума дифференцируемой

функции должно выполняться необходимое условие экстремума, на конец счета градиентным методом можно выходить, когда станет достаточно малой норма градиента. Если же принять во внимание, что минимизация применяется к решению нелинейной системы, то целесообразно отслеживать близость к нулю значений минимизируемой неотрицательной функции, т.е. судить о точности получаемого приближения по квадрату его евклидовой нормы невязки.

**Замечание 6.4.** Как отмечалось в начале этого пункта, ограничение размерности решаемой системы здесь делалось сугубо из иллюстративных соображений. Ничто не помешает развить рассматриваемый подход на случай  $n$ -мерной системы (6.1), сводя ее решение к экстремальной задаче

$$\Phi(x) := \sum_{i=1}^n f_i^2(x) \rightarrow \min.$$

## 6.5. ЧИСЛЕННЫЙ ПРИМЕР

Типичное поведение рассмотренных методов решения систем нелинейных уравнений отражает следующий численный пример.

Пусть ищется решение системы

$$\begin{cases} 20 \ln(x - y) - x - y - 6 = 0, \\ 20 \sin(0.7x - 0.7y) + 7x + 7y = 0 \end{cases} \quad (6.25)$$

в окрестности точки  $x_0 = 0$ ,  $y_0 = -1$ .

Результаты применения к этой системе разных описанных выше итерационных процессов, начинающихся с данной точки  $(x_0, y_0)$  и заканчивающихся, как только выполнится неравенство

$$\max\{|x_k - x_{k-1}|, |y_k - y_{k-1}|\} < \varepsilon, \quad (6.26)$$

где  $k$  – номер итерации, представлены таблицами 6.1<sup>\*</sup> и 6.2<sup>\*</sup>). В них приведены: приближенные решения, полученные на  $k$ -й итерации с помощью семи перечисленных там методов, значения  $k$ , при которых сработал критерий (6.26) с  $\varepsilon = 0.0001$  для табл. 6.1 и с  $\varepsilon = 0.000001$  для табл. 6.2, а также векторы невязок, характеризующие некоторую меру близости указанного приближения к точному решению системы (6.25). При этом всюду шаг аппроксимации производных (начальный шаг в методах секущих и Брауна) брался равным  $\varepsilon$  в каждой компоненте.

<sup>\*</sup> Расчеты проведены Кислухиным Д.А.

Таблица 6.1

№ п/п	Метод	Приближенное решение $(x_k, y_k)^T$	Число итераций $k$	Невязка $(f(x_k, y_k), g(x_k, y_k))^T$
1	Метод Ньютона	-0.46584782 -1.67846885	3	-0.000000164961 -0.000000089180
2	Разностный метод Ньютона	-0.46584782 -1.67846886	4	-0.000000000504 -0.000000000268
3	Метод секущих	-0.46584782 -1.67846885	4	-0.000000293536 -0.000000158314
4	Метод Ньютона с аппроксимацией обратных матриц	-0.46584829 -1.67846812	3	-0.000020367752 -0.000009436432
5	Метод Брауна (с аппроксимацией производных*)	-0.46584558 -1.67846688	4	0.000000052534 0.000031958146
6	Модифицированный (упрощенный) метод Ньютона	-0.46585337 -1.67845758	5	-0.000276258352 -0.000114675484
7	Метод наискорейшего спуска	-0.46570007 -1.67823855	30	-0.001739733248 0.001882572819

Таблица 6.2

№ п/п	Метод	Приближенное решение $(x_k, y_k)^T$	Число итераций $k$	Невязка $(f(x_k, y_k), g(x_k, y_k))^T$
1	Метод Ньютона	-0.46584782 -1.67846886	4	0.000000000001 -0.000000000002
2	Разностный метод Ньютона	-0.46584782 -1.67846886	4	-0.000000000497 -0.000000000261
3	Метод секущих	-0.46584782 -1.67846886	5	0.000000000001 -0.000000000002
4	Метод Ньютона с аппроксимацией обратных матриц	-0.46584782 -1.67846886	4	0.000000000098 -0.000000000048
5	Метод Брауна (с аппроксимацией производных)	-0.46584782 -1.67846886	5	-0.000000000000 -0.000000037367
6	Модифицированный (упрощенный) метод Ньютона	-0.46584784 -1.67846880	8	-0.000001418985 -0.000000589060
7	Метод наискорейшего спуска	-0.46574667 -1.67846711	36	-0.000012735642 0.000014735135

\*) Аппроксимация производных на  $k$ -й итерации осуществлялась с шагом

$$h_k = \min\{|p_k|, |q_k|\}.$$

Число итераций и, соответственно, число нулей перед первыми значащими цифрами в невязках подкрепляет сделанные выше заявления о быстрой сходимости методов 1-5 и справедливость высказанных в замечании 6.2 соображений о надежности и даже некоторой грубости простого критерия останова (6.16) (для системы (6.25) реализованного в виде (6.26)). Сравнивая четвертые знаки после запятой у приближенных решений, полученных методом наискорейшего спуска и, например, методом Ньютона, в таблице 6.1, а также, аналогично, шестые знаки в таблице 6.2, видим, что для медленно сходящихся методов примененный критерий уже не является столь надежным (из процесса итерирования в градиентных методах, как уже упоминалось в замечании 6.3, обычно выходят по своим критериям, например, по малости нормы градиента).

## 6.6. СХОДИМОСТЬ МЕТОДА НЬЮТОНА И НЕКОТОРЫХ ЕГО МОДИФИКАЦИЙ

Многие утверждения о сходимости метода Ньютона восходят к известным результатам Л.В. Канторовича<sup>\*)</sup>, перенесшего метод Ньютона на нелинейные операторные уравнения в банаховых пространствах, в связи с чем и метод в таком общем случае часто называют *методом Ньютона-Канторовича*. В основе этих результатов лежит принцип мажорирования операторного уравнения некоторым скалярным уравнением, через корни которого оценивается абсолютная погрешность приближений, получаемых в итерационном процессе (см. [29], а также [40] и др.). Таким путем получают условия квадратичной сходимости основного (6.6) и линейной сходимости модифицированного (упрощенного) (6.8) методов Ньютона. Непосредственное применение подобной методики к анализу сходимости методов более общего вида, в частности, ААМН (6.11), затруднительно. Поэтому изберем другой путь исследования сходимости итерационных последовательностей. Вообще говоря, здесь будет эксплуатироваться та же идея скалярного мажорирования, только в несколько ином виде.

Прежде всего, попытаемся выяснить, какие условия нужно наложить на скалярную последовательность  $(p_k)$ , мажорирующую последовательности норм поправок  $x^{(k+1)} - x^{(k)}$  и невязок  $F(x^{(k)})$ , чтобы последовательность  $n$ -мерных векторов  $x^{(k)}$  сходилась к нулю вектор-функции  $F(x)$  с заданным порядком  $\mu$  ( $\geq 1$ ) скорости сходимости независимо от способа получения элементов последовательности  $(x^{(k)})$ .

---

<sup>\*)</sup> Канторович Леонид Витальевич (1912-1986) – российский академик, приобретший мировую известность своими основополагающими работами в области линейного программирования и применения функционального анализа в вычислительной математике.



**Теорема 6.3.** I. Пусть непрерывная векторная функция  $F: M \subseteq R_n \rightarrow R_n$  и последовательность векторов  $\mathbf{x}^{(k)} \in M$  таковы, что при всех  $k \in N_0$  выполняются условия:

$$1) \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq \lambda p_k;$$

$$2) \|F(\mathbf{x}^{(k)})\| \leq p_k,$$

где числа  $p_k$  определяются рекуррентным равенством

$$p_{k+1} = G_0 p_k^\mu, \quad k=0, 1, 2, \dots, \quad (6.27)$$

а  $\lambda > 0$ ,  $G_0 > 0$ ,  $p_0 > 0$  и  $\mu > 1$  — некоторые числовые параметры.

II. Тогда если  $v := G_0 p_0^{\mu-1} < 1$  и замкнутый шар

$S \left( \mathbf{x}^{(0)}, r := \lambda p_0 \sum_{i=0}^{\infty} v^{\frac{\mu^i-1}{\mu-1}} \right)$  содержится в  $M$ , то все члены последовательности  $(\mathbf{x}^{(k)})$  принадлежат  $S$ , последовательность  $(\mathbf{x}^{(k)})$  имеет предел  $\mathbf{x}^* \in S$  такой, что  $F(\mathbf{x}^*) = \mathbf{0}$ ; при этом быстрота сходимости  $(\mathbf{x}^{(k)})$  к  $\mathbf{x}^*$  характеризуется неравенством ( $\forall k \in N$ )

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{\lambda p_0}{1 - v^{\mu^k}} v^{\frac{\mu^k-1}{\mu-1}}. \quad (6.28)$$

**Доказательство.** Пользуясь равенством (6.27), выразим элементы последовательности  $(p_i)$  через ее начальный член  $p_0$  и определенную в теореме величину  $v$ :

$$\begin{aligned} p_i &= G_0 p_{i-1}^\mu = G_0 (G_0 p_{i-2}^\mu)^\mu = G_0^{1+\mu} p_{i-2}^{\mu^2} = G_0^{1+\mu+\mu^2} \cdot p_{i-3}^{\mu^3} = \dots = \\ &= G_0^{1+\mu+\mu^2+\dots+\mu^{i-1}} p_0^{\mu^i} = G_0^{\mu^i-1} \cdot p_0^{\mu^i-1} \cdot p_0^{1-\mu^i} \cdot p_0^{\mu^i} = p_0 \cdot v^{\frac{\mu^i-1}{\mu-1}}. \end{aligned}$$

Следовательно, условие 1) доказываемой теоремы можно переписать в виде

$$\|\mathbf{x}^{(i+1)} - \mathbf{x}^{(i)}\| \leq \lambda p_0 \cdot v^{\frac{\mu^i-1}{\mu-1}}. \quad (6.29)$$

\* Без всяких прочих изменений можно заменить здесь  $R_n \rightarrow R_n$  на более общий случай  $R_n \rightarrow R_m$ .

Посредством (6.29) теперь устанавливаем, что

$$\|x^{(k+1)} - x^{(0)}\| \leq \sum_{i=0}^k \|x^{(i+1)} - x^{(i)}\| \leq \sum_{i=0}^k \lambda p_0 \cdot v^{\frac{\mu'-1}{\mu-1}} \leq r,$$

т.е. все члены заданной последовательности  $(x^{(k)})$  принадлежат  $S \subseteq M$ . Покажем, что она удовлетворяет критерию Коши. С помощью того же неравенства (6.29) имеем

$$\begin{aligned} \|x^{(k+m)} - x^{(k)}\| &\leq \sum_{i=k}^{k+m-1} \|x^{(i+1)} - x^{(i)}\| \leq \lambda p_0 \sum_{i=k}^{k+m-1} v^{\frac{\mu'-1}{\mu-1}} = \\ &= \lambda p_0 \cdot v^{\frac{\mu'-1}{\mu-1}} \left( 1 + v^{\mu^k} + v^{\mu^k + \mu^{k+1}} + \dots + v^{\mu^k + \mu^{k+1} + \dots + \mu^{k+m-1}} \right) < \\ &< \lambda p_0 \cdot v^{\frac{\mu'-1}{\mu-1}} \left( 1 + v^{\mu^k} + v^{2\mu^k} + \dots + v^{(m-1)\mu^k} \right) = \lambda p_0 \cdot v^{\frac{\mu'-1}{\mu-1}} \cdot \frac{1 - v^{m\mu^k}}{1 - v^{\mu^k}}. \end{aligned}$$

Полученное неравенство, рассматриваемое при фиксированном  $m \in N$  и  $k \rightarrow \infty$ , говорит о фундаментальности  $(x^{(k)})$  и существовании предельного вектора  $x^*$  в шаре  $S$  (в силу предполагаемой замкнутости  $S$ ). С другой стороны, если в нем зафиксировать  $k$  и перейти к пределу при  $m \rightarrow \infty$ , сразу получается утверждаемая оценка (6.28). Подстановка выражения

$\mu^{k-1}$   
 $p_k = p_0 \cdot v^{\mu^{-1}}$  в условие 2) показывает, что  $\|F(x^{(k+1)})\| \rightarrow 0$  при  $k \rightarrow \infty$  (т.е. при  $x^{(k)} \rightarrow x^*$ ), а это, в силу предполагаемой непрерывности  $F(x)$ , означает, что  $x^* := \lim_{k \rightarrow \infty} x^{(k)}$  есть решение уравнения  $F(x) = 0$ . Теорема доказана.

**Замечание 6.5.** Нетрудно убедиться, что теорема 6.3 (а также следующая теорема 6.4) справедлива и при  $\mu = 1$ . При этом

$$v := G_0, \quad r := \frac{\lambda p_0}{1-v}, \quad \|x^* - x^{(k)}\| \leq \frac{\lambda p_0}{1-v^k} \cdot v^k$$

(т.е. сходимость  $(x^{(k)})$  к  $x^*$  в  $S(x^{(0)}, r)$  в этом случае – линейная).

Изменим требования к  $(x^{(k)})$  и  $F(x)$ , фигурирующие в части I теоремы 6.3, так, чтобы осталось неизменной ее констатирующая часть II.

**Теорема 6.4.** Пусть существуют такие последовательности положительных чисел  $H_k$  и  $G_k$ , удовлетворяющих условиям  $H_k \leq H_0$ ,  $G_k \leq G_0$ , и число  $\mu > 1$ , что в предположении, что  $\mathbf{x}^{(k)} \in M$ , при всех  $k \in N_0$  выполняются неравенства:

$$1) \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq H_k \cdot \|F(\mathbf{x}^{(k)})\|;$$

$$2) \|F(\mathbf{x}^{(k+1)})\| \leq G_k \cdot \|F(\mathbf{x}^{(k)})\|^\mu.$$

Тогда справедлива часть II теоремы 6.3 с  $p_0 \geq \|F(\mathbf{x}^{(0)})\|$ ,  $\lambda = H_0$ .

**Доказательство.** Определим последовательность  $(p_k)$  равенством  $p_{k+1} = G_0 p_k^\mu$  и по индукции покажем, что эта последовательность мажорирует  $\|F(\mathbf{x}^{(k)})\|$  одновременно с доказательством принадлежности векторов  $\mathbf{x}^{(k)}$  шару  $S(\mathbf{x}^{(0)}, r)$ .

По условию  $\mathbf{x}^{(0)} \in S$  и  $\|F(\mathbf{x}^{(0)})\| \leq p_0$ . Сделаем индукционное предположение, что

$$\mathbf{x}^{(i)} \in S \quad \text{и} \quad \|F(\mathbf{x}^{(i)})\| \leq p_i \quad \forall i \in \{0, 1, \dots, k\}.$$

Тогда

$$\|\mathbf{x}^{(i+1)} - \mathbf{x}^{(i)}\| \leq H_i \cdot \|F(\mathbf{x}^{(i)})\| \leq H_0 p_i = H_0 p_0 \cdot \nu^{\frac{\mu'-1}{\mu-1}}.$$

Из этого следует (см. доказательство теоремы 6.3) неравенство  $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(0)}\| \leq r$ , означающее, что  $\mathbf{x}^{(k+1)} \in S \subseteq M$ . Таким образом, значение  $F(\mathbf{x}^{(k+1)})$  существует и

$$\|F(\mathbf{x}^{(k+1)})\| \leq G_k \|F(\mathbf{x}^{(k)})\|^\mu \leq G_0 p_k^\mu = p_{k+1}.$$

Теперь можно сказать, что условия теоремы 6.3 полностью выполнены, значит справедливо и ее заключение.

Прежде чем применить доказанные выше теоремы к конкретным методам типа (6.5), выведем из формулы Тейлора<sup>\*</sup>

$$F(\mathbf{x}) \equiv F(\mathbf{x}^{(0)} + \mathbf{h}) = F(\mathbf{x}^{(0)}) + F'(\mathbf{x}^{(0)})\mathbf{h} + \frac{1}{2!}F''(\mathbf{x}^{(0)})\mathbf{h}^2 + \dots + \frac{1}{(l-1)!}F^{(l-1)}(\mathbf{x}^{(0)})\mathbf{h}^{l-1} + \bar{\omega}(\mathbf{x}^{(0)}, \mathbf{h}) \quad (6.30)$$

с остаточным членом

$$\bar{\omega}(\mathbf{x}^{(0)}, \mathbf{h}) = \frac{1}{(l-1)!} \int_{\mathbf{x}^{(0)}}^{\mathbf{x}} F^{(l)}(\mathbf{z})(\mathbf{x} - \mathbf{z})^{l-1} d\mathbf{z} \quad (6.31)$$

простое неравенство для оценивания  $\|F(\mathbf{x})\|$  в произвольной точке  $\mathbf{x} \in M$  через значения  $F$  и  $F'$  в близкой к  $\mathbf{x}$  точке  $\mathbf{x}^{(0)} \in M$ .

**Лемма 6.2.** Пусть векторная функция  $F(\mathbf{x})$  в области  $M \in R_n$  дифференцируема по Фреше и ее производная удовлетворяет условию Липшица:

$$\|F'(\mathbf{x}) - F'(\mathbf{x}^{(0)})\| \leq L \cdot \|\mathbf{x} - \mathbf{x}^{(0)}\| \quad \forall \mathbf{x}, \mathbf{x}^{(0)} \in M.$$

Тогда при любых  $\mathbf{x}, \mathbf{x}^{(0)}$  из  $M$

$$\|F(\mathbf{x})\| \leq \|F(\mathbf{x}^{(0)}) + F'(\mathbf{x}^{(0)})\| \|\mathbf{x} - \mathbf{x}^{(0)}\| + \frac{L}{2} \|\mathbf{x} - \mathbf{x}^{(0)}\|^2. \quad (6.32)$$

**Доказательство.** Запишем формулу Тейлора (6.30) с остаточным членом (6.31) для случая  $l = 1$ :

$$F(\mathbf{x}) = F(\mathbf{x}^{(0)}) + \int_{\mathbf{x}^{(0)}}^{\mathbf{x}} F'(\mathbf{z}) d\mathbf{z}.$$

Учитывая, что  $\int_{\mathbf{x}^{(0)}}^{\mathbf{x}} F'(\mathbf{x}^{(0)}) d\mathbf{z} \equiv F'(\mathbf{x}^{(0)})\|\mathbf{x} - \mathbf{x}^{(0)}\|$ , ее можно преобразовать к

виду

$$F(\mathbf{x}) = F(\mathbf{x}^{(0)}) + F'(\mathbf{x}^{(0)})\|\mathbf{x} - \mathbf{x}^{(0)}\| + \int_{\mathbf{x}^{(0)}}^{\mathbf{x}} [F'(\mathbf{z}) - F'(\mathbf{x}^{(0)})] d\mathbf{z}.$$

В последнем представлении  $F(\mathbf{x})$  интеграл по отрезку  $[\mathbf{x}^{(0)}; \mathbf{x}]$  заменой

<sup>\*</sup> Выражения типа  $F^{(i)}\mathbf{h}^i$  следует понимать как результат применения  $i$ -линейного оператора  $i$ -кратного дифференцирования  $F^{(i)}$  к вектору  $\mathbf{h}$  (см. приложение 2).

$z = x^{(0)} + \tau(x - x^{(0)})$  сведем к интегралу по абстрактной переменной  $\tau \in [0; 1]$ . Имеем:

$$F(x) = F(x^{(0)}) + F'(x^{(0)})(x - x^{(0)}) + \int_0^1 [F'(x^{(0)} + \tau(x - x^{(0)})) - F'(x^{(0)})](x - x^{(0)}) d\tau.$$

Отсюда, переходя к нормам, получаем доказываемое неравенство (6.32), предварительно оценив норму интеграла следующим образом:

$$\begin{aligned} & \left\| \int_0^1 [F'(x^{(0)} + \tau(x - x^{(0)})) - F'(x^{(0)})](x - x^{(0)}) d\tau \right\| \leq \\ & \leq \int_0^1 \|F'(x^{(0)} + \tau(x - x^{(0)})) - F'(x^{(0)})\| \cdot \|x - x^{(0)}\| d\tau \leq \\ & \leq \int_0^1 L \|x - x^{(0)}\|^2 \tau d\tau = \frac{L}{2} \|x - x^{(0)}\|^2. \end{aligned}$$

**Теорема 6.5.** Пусть функция  $F(x)$  определена и дифференцируема по Фреше в некоторой открытой области  $M \subseteq R_n$ , причем:

- 1)  $\exists L > 0: \|F'(x) - F'(\tilde{x})\| \leq L \|x - \tilde{x}\| \quad \forall x, \tilde{x} \in M;$ <sup>\*)</sup>
- 2)  $\exists [F'(x)]^{-1}$  и  $\exists C > 0: \|[F'(x)]^{-1}\| \leq C \quad \forall x \in M.$

Тогда если

$$v := 0.5LC^2 p_0 < 1, \quad \text{где } p_0 \geq \|F'(x^{(0)})\|,$$

и замкнутый шар  $S(x^{(0)}, r := Cp_0 \sum_{i=0}^{\infty} v^{2^i - 1})$  целиком содержится

в  $M$ , то все члены последовательности  $(x^{(k)})$ , определяемые методом Ньютона (6.6), начинающимися с заданного  $x^{(0)}$ , лежат в  $S \subseteq M$ ; последовательность  $(x^{(k)})$  имеет предел  $x^* \in S$ , служащий решением уравнения  $F(x) = 0$ ; справедлива оценка погрешности

$$\|x^* - x^{(k)}\| \leq \frac{Cp_0}{1 - v^{2^k}} \cdot v^{2^k - 1}.$$

<sup>\*)</sup> Согласно лемме 6.1, для дважды непрерывно дифференцируемой функции  $F(x)$  в качестве  $L$  можно брать оценку сверху величины  $\|F''(x)\|$ .

Доказательство. Непосредственно из равенства (6.6) получаем

$$\|x^{(k+1)} - x^{(k)}\| \leq \| [F'(x^{(k)})]^{-1} \| \cdot \| F(x^{(k)}) \| \leq C \cdot \| F(x^{(k)}) \|,$$

т.е. требование 1) теоремы 6.4 с  $H_k \equiv C$ .

Далее обратимся к неравенству (6.32), установленному леммой 6.2. Положив в нем  $x = x^{(k+1)}$ ,  $x^{(0)} = x^{(k)}$ , приведем (6.32) к виду

$$\| F(x^{(k+1)}) \| \leq \| F(x^{(k)}) + F'(x^{(k)})(x^{(k+1)} - x^{(k)}) \| + \frac{L}{2} \| x^{(k+1)} - x^{(k)} \|^2. \quad (6.33)$$

Но в данном случае, т.е. для метода Ньютона,

$$F(x^{(k+1)}) + F'(x^{(k)})(x^{(k+1)} - x^{(k)}) = 0,$$

поэтому

$$\| F(x^{(k+1)}) \| \leq \frac{L}{2} \| x^{(k+1)} - x^{(k)} \|^2 \leq \frac{1}{2} LC^2 \| F(x^{(k)}) \|^2.$$

Таким образом, выполнено и требование 2) теоремы 6.4 с  $G_k = \frac{1}{2} LC^2$ ,  $\mu = 2$ .

Завершает доказательство подстановка постоянных  $\lambda = H_0 = C$ ,  $G_0 = \frac{1}{2} LC^2$ ,  $\mu = 2$  в часть II теоремы 6.3.

Для модифицированного метода Ньютона требование обратимости матрицы Якоби в любой точке  $M$  заменим менее ограничительным требованием ее обратимости лишь в начальной точке  $x^{(0)}$ .

**Теорема 6.6.** Пусть для  $F: (M \subseteq R_n) \rightarrow R_n$ :

$$1) \exists F'(x): (\exists L > 0: \| F'(x) - F'(\bar{x}) \| \leq L \| x - \bar{x} \| \quad \forall \bar{x} \in M) \quad \forall x \in M;$$

$$2) \exists [F'(x^{(0)})]^{-1}, \exists C_0 > 0: \| [F'(x^{(0)})]^{-1} \| \leq C_0.$$

Тогда если при  $p_0 \geq \| F(x^{(0)}) \|$  величина

$$t := LC_0^2 p_0 \leq 0.125 \quad (6.34)$$

и замкнутый шар  $S(x^{(0)}, r := \frac{2C_0 p_0}{1 \pm \sqrt{1 - 8t}})$  содержится в  $M$ , то начатый с  $x^{(0)}$  модифицированный метод Ньютона (6.8) сходится

в  $S$  к решению  $\mathbf{x}^*$  уравнения (6.1') с оценкой погрешности

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{C_0 p_0}{1 - \nu^k} \cdot \nu^k,$$

$$\text{где } \nu := \frac{1}{2} \mp \sqrt{\frac{1}{4} - 2t}.$$

**Доказательство.** Как и в предыдущей теореме, требование 1) теоремы 6.4 получается сразу же из формулы, определяющей метод:

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq \left\| \left[ F'(\mathbf{x}^{(0)}) \right]^{-1} F(\mathbf{x}^{(k)}) \right\| \leq C_0 \|F(\mathbf{x}^{(k)})\|.$$

Для оценки  $\|F(\mathbf{x}^{(k+1)})\|$  преобразуем неравенство (6.33) так, чтобы воспользоваться равенством нулю выражения  $F(\mathbf{x}^{(k)}) + F'(\mathbf{x}^{(0)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$  в соответствии с (6.8). Имеем

$$\begin{aligned} \|F(\mathbf{x}^{(k+1)})\| &\leq \|F(\mathbf{x}^{(k)}) + F'(\mathbf{x}^{(k)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) + \\ &+ [F'(\mathbf{x}^{(k)}) - F'(\mathbf{x}^{(0)})](\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})\| + \frac{L}{2} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|^2 \leq \\ &\leq \left( \|F(\mathbf{x}^{(k)}) - F'(\mathbf{x}^{(0)})(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})\| + \frac{L}{2} \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \right) \cdot \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \leq \\ &\leq L \left( \|\mathbf{x}^{(k)} - \mathbf{x}^{(0)}\| + \frac{1}{2} (\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(0)}\| + \|\mathbf{x}^{(0)} - \mathbf{x}^{(k)}\|) \right) \cdot C_0 \cdot \|F(\mathbf{x}^{(k)})\| \leq \\ &\leq LC_0 \left( r + \frac{1}{2}(r+r) \right) \cdot \|F(\mathbf{x}^{(k)})\| = 2LC_0 r \cdot \|F(\mathbf{x}^{(k)})\|. \end{aligned}$$

Отсюда видно, что можно считать выполненным требование 2) теоремы 6.4 с  $G_k = G_0 = 2LC_0 r$  и  $\mu = 1$ .

Подстановка постоянных  $\mu = 1$ ,  $\lambda = C_0$  в заключительную часть II теоремы 6.3 показывает (с учетом замечания 6.5), что в данном случае должно быть  $r = \frac{C_0 p_0}{1 - \nu}$ , а  $\nu = 2LC_0 r < 1$ . Исключая из последних двух равенств параметр  $r$ , получаем квадратное относительно  $\nu$  уравнение

$$\nu^2 - \nu + 2t = 0,$$

оба корня которого  $\nu_{1,2} = 0,5 \mp \sqrt{0,25 - 2t}$  положительны и меньше 1, если только неотрицателен дискриминант  $0,25 - 2t$ , что обеспечивается условием (6.34). Подставляя эти значения  $\nu$ , находим связанные с ним значения радиуса  $r$  шара  $S$ . Теперь справедливость заключения данной теоремы очевидна.

<sup>\*)</sup> В выражениях  $r$  и  $\nu$  одновременно берутся либо только верхние, либо только нижние знаки.

Обратимся, наконец, к обоснованию квадратичной сходимости метода Ньютона с последовательной аппроксимацией обратных матриц, т.е. ААМН (6.11).

**Теорема 6.7.** Пусть функция  $F(\mathbf{x})$  определена и дифференцируема в  $M \subseteq R_n$ , причем

$$\exists L > 0: \|F'(\mathbf{x}) - F'(\bar{\mathbf{x}})\| \leq L\|\mathbf{x} - \bar{\mathbf{x}}\| \quad \forall \mathbf{x}, \bar{\mathbf{x}} \in M.$$

Тогда если вектор  $\mathbf{x}^{(0)}$  и матрица  $\mathbf{A}_0$  таковы, что при некотором  $\lambda > 0$  выполняются неравенства

$$\|\mathbf{E} - F'(\mathbf{x}^{(0)})\mathbf{A}_0\| \leq L\lambda^2 \|F(\mathbf{x}^{(0)})\|, \quad (6.35)$$

$$v := 4L\lambda^2 \|F(\mathbf{x}^{(0)})\| \leq 1 - \frac{\|\mathbf{A}_0\|}{2\lambda - \|\mathbf{A}_0\|} \quad (6.36)$$

и

$$S := \left\{ \mathbf{x} \in R_n \mid \|\mathbf{x} - \mathbf{x}^{(0)}\| \leq \lambda \|F(\mathbf{x}^{(0)})\| \cdot \sum_{i=0}^{\infty} v^i \right\} \subseteq M,$$

то начатый с данных  $\mathbf{x}^{(0)}$ ,  $\mathbf{A}_0$  ААМН (6.11) сходится в  $S$  к решению  $\mathbf{x}^*$  уравнения  $F(\mathbf{x}) = \mathbf{0}$  и имеет место оценка погрешности

$$\|\mathbf{x}^* - \mathbf{x}^{(k)}\| \leq \frac{\lambda \|F(\mathbf{x}^{(0)})\|}{1 - v^{2^k}} \cdot v^{2^k}. \quad (6.37)$$

**Доказательство.** Введем в рассмотрение последовательности положительных величин  $p_k, \beta_k, b_k$ , определяемых при  $k = 0, 1, 2, \dots$  равенствами

$$p_{k+1} = 4L\lambda^2 p_k, \quad p_0 := \|F(\mathbf{x}^{(0)})\|; \quad (6.38)$$

$$\beta_k = 2L\lambda^2 p_k; \quad (6.39)$$

$$b_{k+1} = \beta_k^2, \quad b_0 := L\lambda^2 p_0. \quad (6.40)$$

Очевидно невозрастание этих последовательностей. Легко также видеть, что

$$b_k = L\lambda^2 p_k \quad \forall k \in N_0. \quad (6.41)$$

Действительно, предположив равенство (6.41) верным при некотором  $k$ , имеем

$$b_{k+1} = (2L\lambda^2 p_k)^2 = L\lambda^2 \cdot 4L\lambda^2 p_k^2 = L\lambda^2 p_{k+1},$$

т.е. то, что получили бы формальной заменой в (6.41)  $k$  на  $k+1$ .



Обозначим  $\mathbf{B}_k := \mathbf{E} - F'(\mathbf{x}^{(k)})\mathbf{A}_k$  (— невязка для  $\mathbf{A}_k$  относитель-

но  $[F'(\mathbf{x}^{(k)})]^{-1}$ ).

Докажем, что при любом  $k \in N_0$  скалярные последовательности  $(p_k)$ ,  $(\beta_k)$ ,  $(b_k)$  мажорируют последовательности норм векторов  $F(\mathbf{x}^{(k)})$  и матриц  $\Psi_k$ ,  $\mathbf{B}_k$  соответственно и, вместе с тем,  $\lambda$  ограничивает сверху  $\|\mathbf{A}_k\|$ . С этой целью сделаем индукционное предположение, что одновременно выполняются неравенства:

$$\|F(\mathbf{x}^{(k)})\| \leq p_k, \quad \|\Psi_{k-1}\| \leq \beta_{k-1}, \quad \|\mathbf{B}_k\| \leq b_k \quad \text{и} \quad \|\mathbf{A}_k\| \leq \lambda. \quad (6.42)$$

Тогда можно проделать следующие выкладки:

$$\begin{aligned} \|\Psi_k\| &= \|\mathbf{E} - F'(\mathbf{x}^{(k+1)})\mathbf{A}_k\| = \|\mathbf{B}_k + F'(\mathbf{x}^{(k)})\mathbf{A}_k - F'(\mathbf{x}^{(k+1)})\mathbf{A}_k\| \leq \\ &\leq \|\mathbf{B}_k\| + L\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| \cdot \|\mathbf{A}_k\| \leq b_k + L\lambda\|\mathbf{A}_k F(\mathbf{x}^{(k)})\| \leq \\ &\leq L\lambda^2 p_k + L\lambda^2 p_k = \beta_k \end{aligned}$$

(см. (6.11), (6.42), (6.41), (6.39));

$$\begin{aligned} \|\mathbf{B}_{k+1}\| &= \|\mathbf{E} - F'(\mathbf{x}^{(k+1)})\mathbf{A}_{k+1}\| = \|\mathbf{E} - F'(\mathbf{x}^{(k+1)})\mathbf{A}_k - F'(\mathbf{x}^{(k+1)})\mathbf{A}_k\Psi_k\| = \\ &= \|\Psi_k^2\| \leq \|\Psi_k\|^2 \leq \beta_k^2 = b_{k+1} \end{aligned}$$

(см. (6.11), (6.40) и предыдущее неравенство);

$$\begin{aligned} \|F(\mathbf{x}^{(k+1)})\| &\leq \|F(\mathbf{x}^{(k)}) - F'(\mathbf{x}^{(k)})\mathbf{A}_k F(\mathbf{x}^{(k)})\| + \frac{L}{2}\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|^2 \leq \\ &\leq \|\mathbf{B}_k\| \cdot \|F(\mathbf{x}^{(k)})\| + \frac{L}{2}\|\mathbf{A}_k\|^2 \cdot \|F(\mathbf{x}^{(k)})\|^2 \leq b_k p_k + \frac{L}{2}\lambda^2 p_k^2 = \\ &= L\lambda^2 p_k^2 + \frac{L}{2}\lambda^2 p_k^2 < 4L\lambda^2 p_k^2 = p_{k+1} \end{aligned}$$

(см. (6.33), (6.11), (6.42), (6.41), (6.38)); также из (6.11) следует, что

$$\begin{aligned} \|\mathbf{A}_{k+1}\| &\leq \|\mathbf{A}_k\|(1 + \|\Psi_k\|) \leq \|\mathbf{A}_{k-1}\|(1 + \|\Psi_{k-1}\|)(1 + \|\Psi_k\|) \leq \dots \leq \\ &\leq \|\mathbf{A}_0\| \cdot \prod_{j=0}^k (1 + \|\Psi_j\|) \leq \|\mathbf{A}_0\| \cdot \prod_{j=0}^k (1 + \beta_j); \end{aligned}$$

но  $\beta_k = \frac{1}{2}v^{2^k}$  (действительно, из предположения, что  $\beta_{k-1} = \frac{1}{2}v^{2^{k-1}}$ ,

получаем  $\beta_k = 2L\lambda^2 p_k = \frac{1}{2}(4L\lambda^2 p_{k-1})^2 = \frac{1}{2}(2\beta_{k-1})^2 = \frac{1}{2}(v^{2^{k-1}})^2 = \frac{1}{2}v^{2^k}$ ),

поэтому далее можно продолжить оценивание

$$\begin{aligned} \|A_{k+1}\| &\leq \|A_0\| \cdot \prod_{j=0}^k \left(1 + \frac{1}{2} v^{2^j}\right) \leq \|A_0\| \left(1 + \frac{1}{2} \sum_{i=1}^{2^{k+1}-1} v^i\right) = \\ &= \|A_0\| \left(1 + \frac{1}{2} \cdot \frac{v - v^{2^{k+1}}}{1 - v}\right) \leq \|A_0\| \cdot \frac{2 - v}{2 - 2v} \leq \lambda \end{aligned}$$

(последнее, в силу наложенного на  $v$  условия (6.36)).

Так как  $\|F(x^{(0)})\| = p_0$ , согласно заданию (6.38), а  $\|B_0\| \leq b_0$  по условию (6.35) и, кроме того,  $\|A_0\| \leq \lambda$  (ибо в противном случае должно быть  $v < 0$ , что противоречило бы определению  $v$  в (6.36)), то в соответствии с принципом математической индукции доказываемое мажорирование действительно имеет место при любом  $k \in N_0$ .

Итак, отберем только нужную для применения теоремы 6.3 информацию. Имеем:

$$\begin{aligned} \|x^{(k+1)} - x^{(k)}\| &\leq \|A_k\| \cdot \|F(x^{(k)})\| \leq \lambda p_k, \\ \|F(x^{(k)})\| &\leq p_k, \\ p_{k+1} &= G_0 p_k^2, \quad \text{где } G_0 = 4L\lambda^2, \quad p_0 = \|F(x^{(0)})\| \\ &\quad (\text{т.е. } \mu = 2, v = G_0 p_0 = 4L\lambda^2 p_0). \end{aligned}$$

Поскольку здесь на величину  $v$  наложено более сильное, чем в теореме 6.3, ограничение (6.36), с этим  $v$  заключительная часть теоремы 6.3 будет тем более верна. Теорема доказана.

Условие (6.36) теоремы 6.7 можно трактовать так: параметр  $\lambda > 0$  должен удовлетворять кубическому неравенству

$$4Lp_0\lambda^3 - 2L\|A_0\|p_0\lambda^2 - \lambda + \|A_0\| \leq 0.$$

Исследование этого неравенства приводит к ряду следствий [13]; наиболее простым из них (в некотором смысле) является

**Следствие 6.1.** *Квадратичная сходимость начатого с  $x^{(0)}$ ,  $A_0$  метода (6.11) обеспечивается выполнением условий*

$$\begin{aligned} \|E - F'(x^{(0)})A_0\| &\leq 0.109, \quad L\|A_0\|^2 \|F(x^{(0)})\| \leq 0.0567, \\ S(x^{(0)}, r := 2.46 \cdot \|A_0\| \cdot \|F(x^{(0)})\|) &\subseteq M. \end{aligned}$$

<sup>\*)</sup> Как показывают более тонкие исследования, фигурирующие здесь постоянные 0.109 и 0.0567 сильно занижены (см., например, [60], где в этой роли выступает одна постоянная 0.25).

**Замечание 6.6.** Приведенный в этом пункте цикл теорем можно значительно расширить и обобщить.

Во-первых, без каких-либо существенных изменений все эти утверждения могут быть доказаны для нелинейных операторных уравнений в банаховых пространствах.

Во-вторых, вместо условия Липшица на  $F'(x)$  (которое, кстати, легко можно заменить требованием ограниченности  $\|F''(x)\|$  в  $M$  постоянной  $L$ ) можно вставить более слабое *условие Гельдера*

$$\|F'(x) - F'(\bar{x})\| \leq L \|x - \bar{x}\|^\alpha, \quad \alpha \in (0; 1],$$

частным случаем которого является условие Липшица. При этом порядок  $\mu$  метода в теоремах типа теорем 6.5, 6.7 может быть установлен «плавающим» от 1 до  $1 + \alpha$  в зависимости от жесткости требований, накладываемых на исходные данные.

В-третьих, общие теоремы 6.3 и 6.4 позволяют исследовать сходимость методов более высоких порядков как с аппроксимацией обратных к матрицам Якоби матриц, так и без нее, например, методов третьего порядка (6.9) и (6.12).

Использование техники оценочных функций, ключ к которой можно найти в книге [31] Л. Коллатца, а также обобщение леммы 6.2 на случай, когда условие Липшица или Гельдера накладывается на производные более высоких порядков, позволяют указывать условия сходимости семейства методов вида

$$x^{(k+1)} = x^{(k)} - Q(x^{(k)}, A_k),$$

содержащего изученные здесь методы и ряд других методов.

В заключение отметим, что применение теорем сходимости типа теорем (6.5) - (6.7) и других подобных утверждений для конкретных нелинейных систем упирается, как правило, в проблему нахождения постоянных Липшица  $L$  для производных векторных функций. Даже в простейшем случае, когда  $L$  находится из условия  $L \geq \|F''(x)\|$ , требуется сделать оценку

величины нормы функциональной матрицы Гессе  $\left( \frac{\partial^2 f_m}{\partial x_i \partial y_j} \right)_{i,j,m=1}^n$  в не-

которой  $n$ -мерной области  $M$ , что весьма непросто.

## УПРАЖНЕНИЯ

**6.1.** Провести доказательство теоремы 6.1 по аналогии с доказательством теоремы 5.7.

**6.2.** Провести доказательство теоремы 6.2 по аналогии с доказательством теоремы 5.8, используя лемму 6.1.

**6.3.** Можно ли утверждать, что система

$$\begin{cases} x = 0.1 \sin x + 0.3 \cos y - 0.4, \\ y = 0.2 \cos x - 0.1 \sin y - 0.3 \end{cases}$$

имеет и притом единственное вещественное решение? Почему? Сделать 5 итераций МПИ (6.3), начиная процесс с нулевого начального вектора, и оценить погрешность с помощью априорной оценки теоремы 6.1. Найти решение с точностью  $\varepsilon = 10^{-6}$ , останавливая процесс вычислений на основе апостериорной оценки погрешности. Сделать 5 шагов метода по координатных итераций (6.4). Сравнить результат с предыдущим.

**6.4.** Дана система

$$\begin{cases} x^3 - y^3 + 0.1 = 0, \\ xy - 0.95 = 0 \end{cases}$$

и точка  $(x_0, y_0) = (1, 1)$ . Записать для этой системы основной (6.6) и модифицированный (6.8) методы Ньютона и сделать по 2-3 итерационных шага. Попытаться применить здесь какие-нибудь подходящие теоремы сходимости (см. теоремы 6.1, 6.2, 6.5, 6.6).

**6.5.** Сравнить по числу арифметических операций, приходящихся на реализацию одного итерационного шага при решении  $n$ -мерной нелинейной системы, следующие методы:

- а) метод Ньютона в явной форме (6.6);
- б) метод Ньютона в неявной форме (6.7);
- в) аппроксимационный аналог метода Ньютона (6.11);
- г) простейший метод секущих (6.14);
- д) разностный метод Ньютона (6.13).

Подсчет вычислительных затрат вести, предполагая, что решение линейных систем в (6.7) и обращение матриц в (6.6), (6.14) и (6.13) производится методом Гаусса и что вычисление значений функций и производных при этом во внимание не принимается, т.е. эти значения считаются уже найденными.

6.6. Записать алгоритм решения нелинейных систем такой модификацией метода Ньютона, при которой частные производные в матрицах Якоби аппроксимируются симметричными разностными отношениями:

$$\frac{\partial f_i}{\partial x_j} \approx \frac{f_i(x_1, \dots, x_j + h_j, \dots, x_n) - f_i(x_1, \dots, x_j - h_j, \dots, x_n)}{2h_j}.$$

Увеличатся ли вычислительные затраты при переходе от (6.13) к такому симметричному разностному методу?

6.7. Дополнить рассмотрение численного примера п.6.5 (табл.6.1 и 6.2) применением методов третьего порядка (6.9) (или, что то же, (6.10)) и (6.12). Провести сравнительный анализ полученных результатов.

6.8. Записать аппроксимационные аналоги (типа ААМН (6.11)) для методов секущих (6.14) и (6.15). Протестировать новые методы на системе (6.25).

6.9. Вывести расчетные формулы метода Брауна (см. п. 6.3) для трехмерного случая. Опробовать полученные формулы на системе

$$\begin{cases} x^2 + y^2 - z = 0, \\ x^2 + y^2 - z^2 = 0, \\ \ln x - \sqrt{y} + 0.8 = 0, \end{cases}$$

взяв начальное приближение  $x_0 = 0.5$ ,  $y_0 = 0.5$ ,  $z_0 = 0.5$ . Проверить результаты тем же методом Брауна, но примененным к двумерной системе, к которой легко перейти от данной исключением  $z$ .

6.10. Составить гибридный алгоритм, осуществляющий поиск решения нелинейной системы сначала методом градиентного спуска, а затем методом Ньютона или какой-либо его быстросходящейся модификацией. Применить построенный алгоритм к системе из предыдущего упражнения при различных начальных точках (в частности, при  $x_0 = y_0 = z_0 = 0$ ).

## КРАТКИЕ СВЕДЕНИЯ О НОРМАХ ВЕКТОРОВ И МАТРИЦ. СХОДИМОСТЬ В КОНЕЧНОМЕРНЫХ ПРОСТРАНСТВАХ \*)

Норма – это одна из важнейших скалярных характеристик векторов и матриц, удобная для применения, в частности, в вычислительной математике.

Будем сначала рассматривать множество  $R_n$  всех  $n$ -мерных вещественных векторов, образующее, как известно, линейное пространство. Это множество неупорядоченное, т.е. нельзя сказать какой вектор больше, а какой меньше. Чтобы получить возможность каким-то образом сравнивать векторы по величине, и вводится понятие нормы вектора. Норма обобщает понятие длины вектора, являющееся естественным и хорошо изученным в реальных пространствах  $R_2$ ,  $R_3$ . Как и многие другие математические понятия, норма определяется аксиоматически.

**Определение 1.** *Нормой вектора  $x \in R_n$  называется такое действительное число, обозначаемое  $\|x\|$ , что:*

- 1)  $\|x\| \geq 0$ , причем  $\|x\| = 0 \Leftrightarrow x = 0$ ;
- 2)  $\|\lambda x\| = |\lambda| \|x\| \quad \forall \lambda \in R_1$ ;
- 3)  $\|x + y\| \leq \|x\| + \|y\| \quad \forall y \in R_n$ .

Линейное пространство  $R_n$  с введенной в нем нормой называют **нормированным пространством** (точнее, вещественным нормированным пространством).

Легко видеть, что обычное понятие длины (модуля) вектора удовлетворяет всем аксиомам, определяющим норму, т.е. длина есть норма. Обратное неверно: определение 1 норму вектора однозначно не задает. Можно указать бесчисленное множество конструкций  $\|x\|$ , удовлетворяющих всем трем аксиомам нормы вектора. Например, нетрудно проверить, что нормой вектора имеет право называться выражение

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

---

\*) Более полное изложение см., например, в книгах [16, 21, 29, 31, 42, 45, 51, 54].

при любом  $p \in [1, \infty)$ , где  $x_i$  при  $i = 1, 2, \dots, n$  – компоненты вектора  $x$ . Это так называемая  $l_p$ -норма или норма Гельдера [16]. Фиксированием значений параметра  $p$  отсюда можно получать конкретные привлекательные по тем или иным соображениям нормы.

Так, при  $p = 2$  из (1) получаем *евклидову норму вектора*

$$\|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2},$$

являющуюся наиболее естественным расширением понятия длины вектора на  $n$ -мерный случай.<sup>\*)</sup> Множество  $R_n$  всех  $n$ -мерных векторов с введенной в нем евклидовой нормой называют *евклидовым пространством* и часто обозначают через  $E_n$ . Это пространство характерно тем, что в нем определено скалярное произведение векторов  $x = (x_1, \dots, x_n)$  и  $y = (y_1, \dots, y_n)$  равенством

$$(x, y) := \sum_{i=1}^n x_i y_i$$

и при этом, очевидно, имеет место связь

$$\|x\|_2 = \sqrt{(x, x)}.$$

При  $p = 1$  из  $l_p$ -нормы получается *норма-сумма*

$$\|x\|_1 := \sum_{i=1}^n |x_i|,$$

а переход в (1) к пределу при  $p \rightarrow \infty$  дает так называемую *норму-максимум*

$$\|x\|_\infty := \max_{i=1, n} |x_i|.$$

Указанные наиболее распространенные на практике три нормы  $\|x\|_2$ ,  $\|x\|_1$ ,  $\|x\|_\infty$  называют еще соответственно *сферической*, *октаэдрической* и *кубической* по названию поверхности, определяемой уравнением  $\|x\| = \text{const}$  (в частности,  $\|x\| = 1$ ), когда  $x$  считается переменным трехмерным радиус-вектором.

Пусть теперь рассматривается множество всех  $n \times n$ -матриц, также образующих линейное пространство.

---

<sup>\*)</sup> Разумеется, знак  $|\cdot|$  здесь может быть опущен (чего никак нельзя делать в следующих двух частных случаях). Его сохранение требуется лишь тогда, когда вместо  $R_n$  рассматривается пространство  $C_n$  векторов с комплексными компонентами.

**Определение 2.** *Нормой матрицы  $A$  называется действительное число  $\|A\|$ , удовлетворяющее условиям:*

- 1)  $\|A\| \geq 0$ , причем  $\|A\| = 0 \Leftrightarrow A = 0$ ;
- 2)  $\|\lambda A\| = |\lambda| \cdot \|A\| \quad \forall \lambda \in R$ ;
- 3)  $\|A + B\| \leq \|A\| + \|B\|$ ;
- 4)  $\|AB\| \leq \|A\| \cdot \|B\|$

( $B$  – произвольная  $n \times n$ -матрица).

Иногда норму матрицы определяют только с помощью первых трех условий (аксиом); в таком случае говорят, что определение 2 задает в пространстве матриц *мультипликативную норму*.

Как и норму вектора, норму матрицы, отвечающую определению 2, можно вводить далеко не единственным способом. Из множества всевозможных норм матрицы наибольший интерес представляют такие, которые определенным образом соотносятся с векторными нормами, поскольку, чаще всего, матрицы и векторы рассматриваются в комплексе. Так, при умножении матрицы  $A$  на вектор  $x$  получается вектор  $Ax$ , и естественно потребовать, чтобы матричная норма удовлетворяла *условию согласованности*

$$\|Ax\| \leq \|A\| \cdot \|x\|.$$

Например, если для матрицы  $A = (a_{ij})_{i,j=1}^n$  ввести норму равенством<sup>\*)</sup>

$$\|A\|_F := \sqrt{\sum_{i,j=1}^n |a_{ij}|^2},$$

то она будет удовлетворять всем четырем аксиомам определения 2 и согласована с евклидовой нормой вектора. Такую норму называют *нормой Фробениуса* [44], *евклидовой* или *шуровской нормой* [53], а также *нормой Э. Шмидта* [31].

Более сильным требованием к норме матрицы, чем условие согласованности, является условие подчиненности. А именно, норма  $n \times n$ -матрицы  $A$  называется *подчиненной нормой*  $n$ -мерного вектора  $x$ , если она задается равенством

$$\|A\| := \max_{\|x\|=1} \|Ax\|.$$

Таким образом, при заданной векторной норме за подчиненную ей норму матрицы  $A$  принимается максимум норм векторов  $Ax$ , когда  $x$  пробегает множество всех векторов, норма которых равна единице.

<sup>\*)</sup> См. предыдущую сноску



Очевидно, при всякой подчиненной норме для единичной матрицы  $E$  должно быть

$$\|E\| = \max_{\|x\|=1} \|Ex\| = \max_{\|x\|=1} \|x\| = 1.$$

Так как

$$\|E\|_F = \sqrt{\sum_{i=1}^n 1^2} = \sqrt{n},$$

то введенная выше норма Фробениуса, будучи согласованной с евклидовой нормой вектора, не является подчиненной ей. Нормой вещественной матрицы  $A$ , подчиненной евклидовой норме вектора, служит *спектральная норма*:

$$\|A\|_2 := \sqrt{\Lambda},$$

где  $\Lambda$  – наибольшее собственное число матрицы  $A^T A$  (положительное, в силу симметричности  $A^T A$ ).

Нормами матрицы  $A = (a_{ij})_{i,j=1}^n$ , подчиненными другим введенным выше нормам векторов, являются:

$$\|A\|_1 := \max_j \sum_{i=1}^n |a_{ij}| \quad \text{для } \|x\|_1,$$

$$\|A\|_\infty := \max_i \sum_{j=1}^n |a_{ij}| \quad \text{для } \|x\|_\infty.$$

В тех случаях, когда безразлично, какую из норм следует употребить, пишут просто  $\|x\|$  или  $\|A\|$ , понимая под этим любые нормы (или любые подчиненные, или любые согласованные, что бывает ясно из рассмотрения объектов, с которыми оперируют). Так, например, для любой мультипликативной нормы матрицы справедливо неравенство

$$\|A^k\| \leq \|A\|^k \quad \forall k \in N,$$

но под символом  $\|\cdot\|$  в левой и правой частях должна пониматься одна и та же норма.

Пусть имеется последовательность  $x^{(1)}, x^{(2)}, \dots, x^{(k)}, \dots$ , членами которой являются  $n$ -мерные векторы

$$x^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)}).$$

Если все последовательности  $(x_i^{(k)})$  соответствующих компонент этих векторов имеют пределы, т.е. при всех  $i=1, 2, \dots, n$  существуют

$\lim_{k \rightarrow \infty} x_i^{(k)} = x_i^*$ , то вектор  $x^* = (x_1^*, x_2^*, \dots, x_n^*)$  считают *пределом векторной*

*последовательности*  $(x^{(k)})_{k=1}^{\infty}$ . Этот факт сходимости оформляют обычным образом:

$$\lim_{k \rightarrow \infty} x^{(k)} = x^* \quad \text{или} \quad x^{(k)} \xrightarrow{k \rightarrow \infty} x^*.$$

Так же поэлементно (поэлементно) понимается и сходимость последовательностей матриц. А именно, *пределом последовательности матриц*  $A^{(1)}, A^{(2)}, \dots, A^{(k)}, \dots$  с элементами  $a_{ij}^{(1)}, a_{ij}^{(2)}, \dots, a_{ij}^{(k)}, \dots$  называется такая матрица  $A^*$ , элементами которой являются числа  $a_{ij}^* := \lim_{k \rightarrow \infty} a_{ij}^{(k)}$ .

Необходимым и достаточным условием поэлементной сходимости последовательности векторов  $x^{(k)}$  к вектору  $x^*$ , последовательности матриц  $A^{(k)}$  к матрице  $A^*$ , является условие *сходимости по норме*:

$$\|x^{(k)} - x^*\| \rightarrow 0, \quad \|A^{(k)} - A^*\| \rightarrow 0$$

соответственно. При этом важно отметить, что все нормы в конечномерных пространствах (каковыми являются векторные и матричные пространства) эквивалентны. Это надо понимать так: если доказана сходимость векторной или матричной последовательности в одной нормировке, то сходимость к тому же пределу обеспечена и при любой другой нормировке.

## ПРОИЗВОДНЫЕ ВЕКТОРНЫХ ФУНКЦИЙ <sup>\*)</sup>

Пусть  $R_n$  и  $R_m$  – два вещественных нормированных пространства с одинаковой нормой, и пусть  $F$  – некоторое заданное, в общем случае нелинейное отображение из  $R_n$  в  $R_m$ , т.е.  $F$  – векторная функция, аргументами которой являются  $n$ -мерные векторы  $x$ , а значениями – векторы  $F(x)$  размерности  $m$ . Пусть, далее,  $x = (x_1, x_2, \dots, x_n)$  – фиксированная точка некоторой открытой области  $D$ , содержащейся в области определения отображения  $F$ ,  $h$  – произвольный (текущий)  $n$ -мерный вектор такой, что  $x + h \in D$ , и наконец, пусть <sup>\*\*)</sup>  $F(x) = (f_1(x), f_2(x), \dots, f_m(x))$ .

**Определение 1.** Если существует такое линейное преобразование (отображение, оператор)  $A$ , что приращение  $F$  представимо в виде

$$F(x+h) - F(x) = Ah + \omega(x, h), \quad (1)$$

где  $\frac{\|\omega(x, h)\|}{\|h\|} \rightarrow 0$  при  $\|h\| \rightarrow 0$ , то отображение  $F$  называется

**дифференцируемым по Фреше**, линейное преобразование  $A = F'(x)$  – **производной Фреше** или **сильной производной**<sup>\*\*\*)</sup>, а его значение  $Ah = F'(x)h$  – **дифференциалом Фреше**.

Очевидно, при  $n = m = 1$  это определение совпадает с обычным определением дифференцируемости функций одной вещественной переменной, и производная Фреше в этом случае есть обычная производная.

Выясним структуру производной Фреше при  $n, m \geq 1$ . Для этого нужно найти выражения элементов  $a_{ij}$  ( $i = \overline{1, m}; j = \overline{1, n}$ ) матрицы линейного преобразования  $A$ , зависящие, вообще говоря, от  $x$ . Преследуя эту цель, будем рассматривать определяющее производную равенство (1) на координатных векторах, что правомерно, в силу произвольности  $h$ .

Положим сначала  $h := h_1 := (\delta_1, 0, \dots, 0)$ , где  $\delta_1$  – произвольное действительное число. Тогда, учитывая, что

<sup>\*)</sup> См [29, 42, 45, 51 и др.]

<sup>\*\*)</sup> В предлах этого приложения вектор-столбцы фигурируют, в основном, в виде векторов-строк без символа транспонирования

<sup>\*\*\*)</sup> Наряду с производной Фреше определяют еще **производную Гато** или **слабую производную**. Всякая функция, дифференцируемая по Фреше, дифференцируема и по Гато, но не наоборот.

$$\mathbf{x} + \mathbf{h} = \begin{pmatrix} x_1 + \delta_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}, \quad \mathbf{A}\mathbf{h} = \mathbf{A}\mathbf{h}_1 = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \begin{pmatrix} \delta_1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} a_{11}\delta_1 \\ a_{21}\delta_1 \\ \vdots \\ a_{m1}\delta_1 \end{pmatrix},$$

и предполагая отображение  $F$  дифференцируемым, т.е. равенство (1) выполненным, получим векторное равенство

$$\begin{pmatrix} f_1(x_1 + \delta_1, x_2, \dots, x_n) - f_1(x_1, x_2, \dots, x_n) \\ f_2(x_1 + \delta_1, x_2, \dots, x_n) - f_2(x_1, x_2, \dots, x_n) \\ \dots \\ f_m(x_1 + \delta_1, x_2, \dots, x_n) - f_m(x_1, x_2, \dots, x_n) \end{pmatrix} = \begin{pmatrix} a_{11}\delta_1 \\ a_{21}\delta_1 \\ \dots \\ a_{m1}\delta_1 \end{pmatrix} + \begin{pmatrix} \omega_1(\mathbf{x}, \delta_1) \\ \omega_2(\mathbf{x}, \delta_1) \\ \dots \\ \omega_m(\mathbf{x}, \delta_1) \end{pmatrix}, \quad (2)$$

где  $\omega_i(\mathbf{x}, \delta_1)$  —  $i$ -я компонента остаточного члена  $\omega(\mathbf{x}, \mathbf{h})$  дифференциала Фреше на векторе  $\mathbf{h}_1$ . Так как  $\|\omega(\mathbf{x}, \mathbf{h}_1)\| \geq |\omega_i(\mathbf{x}, \mathbf{h}_1)|$ , а  $\|\mathbf{h}\| = \|\mathbf{h}_1\| = |\delta_1|$ , то, очевидно,

$$\lim_{\|\mathbf{h}_1\| \rightarrow 0} \frac{\|\omega(\mathbf{x}, \mathbf{h}_1)\|}{\|\mathbf{h}_1\|} = 0 \Rightarrow \lim_{\delta_1 \rightarrow 0} \frac{|\omega_i(\mathbf{x}, \delta_1)|}{|\delta_1|} = 0 \Rightarrow \lim_{\delta_1 \rightarrow 0} \frac{\omega_i(\mathbf{x}, \delta_1)}{\delta_1} = 0.$$

Поэтому, деля  $i$ -е равенство, получающееся из (2) переходом к координатам, на число  $\delta_1$  и устремляя  $\delta_1$  к нулю, имеем

$$\lim_{\delta_1 \rightarrow 0} \frac{f_i(x_1 + \delta_1, x_2, \dots, x_n) - f_i(x_1, x_2, \dots, x_n)}{\delta_1} = \lim_{\delta_1 \rightarrow 0} \frac{\omega_i(\mathbf{x}, \delta_1)}{\delta_1} = a_{i1},$$

т.е. при  $i = 1, 2, \dots, m$  справедливо представление

$$a_{i1} = \frac{\partial f_i(\mathbf{x})}{\partial x_1}.$$

Рассматривая равенство (1) на других координатных векторах  $\mathbf{h} = \mathbf{h}_j = (0, \dots, \delta_j, \dots, 0)$ , аналогично предыдущему находим выражения всех элементов матрицы линейного преобразования  $\mathbf{A}$ , фигурирующей в определении 1:

$$a_{ij} = \lim_{\delta_j \rightarrow 0} \frac{f_i(x_1, \dots, x_j + \delta_j, \dots, x_n) - f_i(x_1, \dots, x_j, \dots, x_n)}{\delta_j} = \frac{\partial f_i(\mathbf{x})}{\partial x_j} \\ (i = 1, 2, \dots, m; \quad j = 1, 2, \dots, n)$$

Таким образом, производная Фреше отображения  $F: R_n \rightarrow R_m$  представляет собой матрицу частных производных от каждой компоненты  $f_i$  векторной функции  $F = (f_1, f_2, \dots, f_m)$  по каждой компоненте  $x_j$  векторного аргумента  $\mathbf{x}$ . Эта  $m \times n$ -матрица называется *матрицей Якоби* и имеет вид

$$F'(\mathbf{x}) := J(\mathbf{x}) := \begin{pmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f_2(\mathbf{x})}{\partial x_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \frac{\partial f_m(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{pmatrix}.$$

В частности, если  $F$  осуществляет отображение из  $R_n$  в  $R_1$ , т.е. если  $F(\mathbf{x}) = f(x_1, x_2, \dots, x_n)$  – функция  $n$  переменных, то

$$F'(\mathbf{x}) = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right) = \text{grad } F.$$

Следовательно, *градиент* – это матрица-строка Якоби.

Приведем некоторые свойства производной Фреше.

1. Если производная Фреше в точке  $\mathbf{x}$  существует, то она единственна.

2. Наличие всех частных производных первого порядка у функции  $F(\mathbf{x})$  еще не означает существования у  $F(\mathbf{x})$  производной Фреше.

3. Отображение  $F$ , дифференцируемое по Фреше в точке  $\mathbf{x}$ , непрерывно в этой точке (чего нельзя утверждать в случае слабой дифференцируемости).

4. Справедливо цепное правило дифференцирования композиции отображений: если отображение  $F: R_n \rightarrow R_m$  имеет производную в точке  $\mathbf{x}$ , а  $G: R_m \rightarrow R_l$  – в точке  $F(\mathbf{x})$ , то их композиция  $H := G \circ F$ , переводящая элементы пространства  $R_n$  в элементы пространства  $R_l$ , имеет производную в точке  $\mathbf{x}$ , причем  $H'(\mathbf{x}) = G'(F(\mathbf{x})) \cdot F'(\mathbf{x})$ .

Множество  $m \times n$ -матриц образует линейное пространство  $L(R_n, R_m)$  размерности  $s = n \cdot m$ , которое можно естественным образом нормировать. Поэтому к  $F'$ , как к отображению из  $R_n$  в  $L(R_n, R_m)$ , можно применить данное выше определение 1 производной Фреше, в результате чего приходим к понятию второй производной Фреше нелинейного многомерного отображения  $F$ .

**Определение 2.** Если отображение  $F: R_n \rightarrow R_m$  имеет в точке  $\mathbf{x} \in D \subseteq R_n$  дифференцируемую по Фреше производную  $F': R_n \rightarrow L(R_n, R_m)$ , то производная производной, т.е.  $(F'(\mathbf{x}))'$ , на-

зывается *второй производной Фреше* отображения  $F$  и обозначается  $F''(\mathbf{x})$ .

Ясно, что  $F''(\mathbf{x}) \in L(R_n, L(R_n, R_m))$ , т.е.  $F''(\mathbf{x})\mathbf{g}$  является линейным отображением ( $m \times n$ -матрицей) при каждом  $\mathbf{g}$  из  $R_n$ , в то время как  $F'(\mathbf{x})\mathbf{h}$  есть вектор из  $R_m$  при каждом  $\mathbf{h}$  из  $R_n$ . Следовательно,

$$F''(\mathbf{x})\mathbf{g}\mathbf{h} \in R_m \quad \forall \mathbf{g} \in R_n, \mathbf{h} \in R_n,$$

и значит, вторую производную можно интерпретировать как отображение из  $R_n \times R_n$  в  $R_m$ . При этом заметим, что  $F''(\mathbf{x})$  является билинейным отображением<sup>\*)</sup> и обладает свойством симметрии

$$F''(\mathbf{x})\mathbf{h}\mathbf{g} = F''(\mathbf{x})\mathbf{g}\mathbf{h} \quad \forall \mathbf{g}, \mathbf{h} \in R_n.$$

Для того чтобы связать  $F''(\mathbf{x})$  с частными производными функций  $f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x})$  – компонент вектора  $F$ , нужно так же, как и при выявлении структуры первой производной, рассмотреть равенство (1), примененное к  $F'(\mathbf{x})$ , на координатных векторах  $\mathbf{h} = \mathbf{h}_i, \mathbf{g} = \mathbf{g}_j$ . В результате получается, что вектор  $F''(\mathbf{x})\mathbf{h}\mathbf{g}$  при произвольных  $\mathbf{h}$  и  $\mathbf{g}$  из  $R_n$  имеет представление

$$F''(\mathbf{x})\mathbf{h}\mathbf{g} = (\mathbf{g}^T H_1(\mathbf{x})\mathbf{h}, \mathbf{g}^T H_2(\mathbf{x})\mathbf{h}, \dots, \mathbf{g}^T H_m(\mathbf{x})\mathbf{h}),$$

где  $H_j(\mathbf{x})$  – так называемая *матрица Гессе*, составленная из вторых частных производных, а именно:

$$H_j(\mathbf{x}) := \begin{pmatrix} \frac{\partial^2 f_j(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f_j(\mathbf{x})}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f_j(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f_j(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f_j(\mathbf{x})}{\partial x_2^2} & \dots & \frac{\partial^2 f_j(\mathbf{x})}{\partial x_2 \partial x_n} \\ \dots & \dots & \dots & \dots \\ \frac{\partial^2 f_j(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f_j(\mathbf{x})}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f_j(\mathbf{x})}{\partial x_n^2} \end{pmatrix}$$

Аналогично могут быть введены и изучены производные высших порядков.

<sup>\*)</sup> Отображение  $B: R_n \times R_n \rightarrow R_m$  называется *билинейным*, если оно линейно по каждому из своих аргументов, т.е. если

$$B(\alpha_1 \mathbf{g}_1 + \alpha_2 \mathbf{g}_2)\mathbf{h} = \alpha_1 B\mathbf{g}_1\mathbf{h} + \alpha_2 B\mathbf{g}_2\mathbf{h} \quad \text{и} \quad B\mathbf{g}(\beta_1 \mathbf{h}_1 + \beta_2 \mathbf{h}_2) = \beta_1 B\mathbf{g}\mathbf{h}_1 + \beta_2 B\mathbf{g}\mathbf{h}_2.$$

## ОБРАЗЦЫ ПОСТАНОВОК ЛАБОРАТОРНЫХ ЗАДАНИЙ

Лабораторная работа 1. «Метод Гаусса и  $LU$ -разложение матриц»

Дана система  $Ax = b$ , где:

вариант 1

$$A = \begin{pmatrix} 14 & -8 & -21 & 12 \\ 10 & -6 & -15 & 9 \\ 35 & -20 & -56 & 32 \\ 25 & -15 & -40 & 24 \end{pmatrix}, \quad b = \begin{pmatrix} 19 \\ 14 \\ 53 \\ 39 \end{pmatrix};$$

вариант 2

$$A = \begin{pmatrix} 1 & 0 & -3 & -9 \\ 0 & 1 & -7 & -21 \\ 3 & 12 & -92 & -279 \\ 1 & 4 & -31 & -94 \end{pmatrix}, \quad b = \begin{pmatrix} 11 \\ 22 \\ 297 \\ 100 \end{pmatrix};$$

и т.д.

1. Решить систему методом Гаусса. Предусмотреть постолбцовый выбор главного элемента и итерационное уточнение решения до достижения точности  $\varepsilon = 10^{-12}$  по евклидовой норме невязки в рамках применяемой схемы реализации метода.

2. Выполнить  $LU$ -разложение матрицы  $A$  и с его помощью получить  $\det A$  и решение  $x$  данной системы.

3. Найти матрицу  $X = A^{-1}$  двумя способами:

- решая подсистемы  $Ax^j = e^j$  системы  $AX = E$  (используя при этом фрагменты выполнения п.1);
- применяя готовые формулы, полученные на основе  $LU$ -разложения.

4. Вычислить  $\text{cond } A$  в различных простых нормах и охарактеризовать чувствительность данной системы к погрешностям исходных данных.

### Лабораторная работа 2. «Метод прогонки»

Методом прогонки найти вектор  $(u_1, u_2, \dots, u_{10})$ , являющийся решением уравнения (системы)

$$a_k u_{k-1} + b_k u_k + c_k u_{k+1} = f_k,$$

где  $k = 1, 2, \dots, 10$ ;  $a_1 = c_{10} = 0$ ; коэффициенты  $a_k$  (при  $k = 2, 3, \dots, 10$ ),  $b_k$  (при  $k = 1, 2, \dots, 10$ ),  $c_k$  (при  $k = 1, 2, \dots, 9$ ) и  $f_k$  (при  $k = 1, 2, \dots, 10$ ) задаются следующей таблицей:

	$a_k$	$b_k$	$c_k$	$f_k$
Вариант 1	$k$	$3.1k$	$-2k$	$\frac{2.1k^2 + 7.2k + 2}{k^2 + 3k + 2}$
Вариант 2	$\frac{3}{k}$	$\frac{11}{10k}$	$\frac{2}{k}$	$30.5 - \frac{41.6}{k}$
...	...	...	...	...

Что можно сказать об устойчивости прогонки в данном конкретном случае?

### Лабораторная работа 3. «Прямое и итерационное решение симметричной линейной алгебраической системы»

С точностью  $\epsilon = 10^{-6}$  найти решение системы

$$\sum_{j=1}^6 a_{ij} x_j = b_j \quad (i = 1, 2, \dots, 6)$$

с матрицей коэффициентов вида

$$A = (a_{ij}) = \begin{pmatrix} p_1 & 0.1p_1 & 0 & 0 & q & 0 \\ 0.1p_1 & p_2 & 0.1p_2 & 0 & 0 & q \\ 0 & 0.1p_2 & p_3 & 0.1p_3 & 0 & 0 \\ 0 & 0 & 0.1p_3 & p_4 & 0.1p_4 & 0 \\ q & 0 & 0 & 0.1p_4 & p_5 & 0.1p_5 \\ 0 & q & 0 & 0 & 0.1p_5 & p_6 \end{pmatrix},$$

если:

	$p_i \ (i = 1, 2, \dots, 6)$	$b_i \ (i = 1, 2, \dots, 6)$	$q$
Вариант 1	$i$	1	-0.5
Вариант 2	$10 - i$	$25 - 9i$	2
...	...	...	...



Применить:

- а) метод квадратных корней;
- б) метод Якоби (сделав предварительно подсчет числа итераций, гарантирующего получение решения с заданной точностью);
- в) метод Зейделя.

Провести сравнительный анализ примененных методов.

Попытаться уменьшить число итераций метода Зейделя, вводя релаксационный параметр и оптимизируя его значение экспериментальным путем.

#### Лабораторная работа 4. «Численное решение алгебраических проблем собственных значений»

Дана матрица:

вариант 1

$$\begin{pmatrix} 7 & -1 & -2 & 3 \\ -1 & 6 & 0 & 2 \\ -2 & 0 & 5 & 1 \\ 3 & 2 & 1 & 7 \end{pmatrix};$$

вариант 2

$$\begin{pmatrix} 5 & 2 & 0 & -1 \\ 2 & 7 & -3 & 1 \\ 0 & -3 & 9 & 4 \\ -1 & 1 & 4 & 8 \end{pmatrix};$$

и т.д.

1. Найти наибольшее по модулю собственное число и соответствующий ему собственный вектор

- а) степенным методом;
- б) методом скалярных произведений.

(В качестве начального взять вектор  $(1; 1; 1; 1)$ ).

2. Найти наименьшее по модулю собственное число и соответствующий ему собственный вектор

- а) методом обратных итераций;
- б) методом обратных итераций с отношениями Рэлея.

3. Решить полную проблему собственных значений методом вращения Якоби.

Точность  $\varepsilon = 10^{-6}$  (в евклидовой норме).

**Лабораторная работа 5. «Методы решения скалярных уравнений»**

С точностью  $\varepsilon = 10^{-6}$  найти каждый из корней уравнения:  
вариант 1

$$4x \ln^2 x - 4\sqrt{1+x} + 5 = 0;$$

вариант 2

$$x^4 e^x + \sqrt[3]{x-1} - 2 = 0;$$

и т.д.

каждым из следующих четырех способов:

- 1) методом половинного деления;
- 2) методом хорд;
- 3) методом Ньютона;
- 4) методом секущих.

Сравнить методы по числу итераций и по вычислительным затратам.  
Что можно сказать об эффективности примененных методов?

**Лабораторная работа 6. «Скалярная задача о неподвижной точке»**

С точностью  $\varepsilon = 10^{-6}$  решить уравнение:  
вариант 1

$$e^{-0.45x} - \sqrt{x-3} = 0;$$

вариант 2

$$(x-4)^3 + \ln x = 0;$$

и т.д.

- а) методом простых итераций;
- б)  $\Delta^2$ -процессом Эйткена;
- в) методом Вегстейна.

Предварительно привести уравнение к виду, пригодному для проведения итераций; доказать существование и единственность корня; выбрав начальное приближение, сделать априорную оценку количества шагов метода простых итераций.

Результаты представить по следующей форме:

Метод	Начальное приближение	Априорное число итераций	Фактическое число итераций	Полученный корень	Невязка
МПИ					
Эйткена					
Вегстейна					

Сравнить методы по требуемому количеству вычислений функций для получения решения с заданной точностью.

**Лабораторная работа 7. «Метод Бернулли вычисления корней многочлена»**

Дано уравнение:

вариант 1

$$x^4 + 1.1x^3 - 11.51x^2 - 2.331x - 0.117 = 0;$$

вариант 2

$$x^4 - 10.3x^3 - 15.75x^2 - 286.875x - 84.375 = 0;$$

и т.д.

1. Непосредственным применением метода Бернулли найти наибольший и наименьший по модулю корни.

2. Используя найденные корни, понизить по схеме Горнера степень многочлена и найти остальные корни.

**Лабораторная работа 8. «Решение систем нелинейных уравнений»**

Дана система и начальная точка:

вариант 1

$$\begin{cases} (x - y)^3 - 8(x + y) = 0, & x_0 = 2, \\ 2(x - y) + 15 \ln(x + y) - 5 = 0, & y_0 = -0.5; \end{cases}$$

вариант 2

$$\begin{cases} 0.8x^2 + 2xy + 1.3y^2 + 20x - 15y = 0, & x_0 = 0.5, \\ e^{0.6y - 0.8x} - 1.14x - 1.52y = 0, & y_0 = 1; \end{cases}$$

и т.д.

Найти решение данной системы, исходя из данной начальной точки, следующими методами:

- 1) основным методом Ньютона (явным или неявным);
- 2) разностным методом Ньютона (с разными шагами дискретизации производной);
- 3) модифицированным (упрощенным) методом Ньютона;
- 4) методом Ньютона с аппроксимацией обратных матриц;
- 5) методом Брауна;
- 6) методом градиентного спуска.

Провести сравнение всех указанных методов решения нелинейных систем на основе конкретного вычислительного материала, полученного при задании точности  $\varepsilon = 10^{-2}, 10^{-4}, 10^{-6}, 10^{-8}$ .

## Литература

1. Амосов А.А., Дубинский Ю.А., Копчёнова Н.В. Вычислительные методы для инженеров. – М.: Высшая школа, 1994.
2. Арушанян О.Б., Залёткин С.Ф. Численное решение обыкновенных дифференциальных уравнений на Фортране. – М.: Изд.-во МГУ, 1990.
3. Ахромеева Т.С., Курдюмов С.П., Малинецкий Г.Г. Парадоксы мира нестационарных структур. В кн. "Компьютеры и нелинейные явления". – М.: Наука, 1988.
4. Бабушка И., Витасек Э., Прагер М. Численные процессы решения дифференциальных уравнений. – М.: Мир, 1969.
5. Бахвалов Н.С. Численные методы. – М.: Наука, 1973.
6. Бахвалов Н.С., Жидков Н.П., Кобельков Г.М. Численные методы. – М.: Наука, 1987.
7. Беланов А.А. Решение алгебраических уравнений методом Лобачевского. – М.: Наука, 1989.
8. Березин И.С., Жидков Н.П. Методы вычислений. Т.2. – М.: Физматгиз, 1962.
9. Боглаев Ю.П. Вычислительная математика и программирование. – М.: Высшая школа, 1990.
10. Бродис В.М. Вычислительная работа в курсе математики средней школы. – М.: Изд. АПН РСФСР, 1962.
11. Ватях Е. Последовательно-параллельные вычисления. – М.: Мир, 1985.
12. Васильев Ф.П. Численные методы решения экстремальных задач. – М.: Наука, 1988.
13. Вержбицкий В.М. Выбор параметров в теоремах сходимости одного аппроксимационного аналога метода Ньютона. «Ж. Вычисл. мат. и мат. физ.», 1975, 15, № 6, с.1594-1597.
14. Вержбицкий В.М. Обращение матриц и решение нелинейных систем. – Ижевск: Изд. ИМИ, 1980.
15. Воеводин В.В. Вычислительные основы линейной алгебры. – М.: Наука, 1977.

16. *Воеводин В.В., Кузнецов Ю.А.* Матрицы и вычисления. – М.: Наука, 1984.
17. *Волков Б.А.* Численные методы. – М: Наука, 1979.
18. *Гавурин М.К.* Нелинейные функциональные уравнения и непрерывные аналоги итерационных методов. "Изв. Вузов, Матем.", 1958, № 5(6), с.18-31.
19. *Гутер Р.С., Овчинский Б.В.* Элементы численного анализа и математической обработки результатов опыта. – М.: Наука, 1970.
20. *Данилина Н.И., Дубровская Н.С. и др.* Вычислительная математика. – М.: Наука, 1985.
21. *Демидович Б.П., Марон И.А.* Основы вычислительной математики. – М.: Наука, 1970.
22. *Дэннис Дж., Шнабель Р.* Численные методы безусловной оптимизации и решения нелинейных уравнений. – М.: Мир, 1988.
23. *Загускин В.Л.* Справочник по численным методам решения уравнений. – М.: Физматгиз, 1960.
24. *Иванов В.В.* Методы вычислений на ЭВМ: Справочное пособие. – Киев: Наукова думка, 1986.
25. *Икрамов Х.Д.* Численные методы линейной алгебры (решение линейных уравнений). "Математика, кибернетика", №4. – М.: Знание, 1987.
26. *Икрамов Х.Д.* Несимметричная проблема собственных значений. – М.: Наука, 1991.
27. *Ильин В.П., Кузнецов Ю.А.* Трехдиагональные матрицы и их приложения. – М.: Наука, 1985.
28. *Калиткин Н.Н.* Численные методы. – М: Наука, 1978.
29. *Канторович Л.В., Акилов Г.П.* Функциональный анализ. – М.: Наука, 1977.
30. *Коллатц Л.* Задачи на собственные значения с техническими приложениями. – М.: Наука, 1968.

31. *Коллатц Л.* Функциональный анализ и вычислительная математика. – М.: Мир, 1969.
32. *Коллатц Л., Альбрехт Ю.* Задачи по прикладной математике. – М.: Мир, 1978.
33. *Косарев В.И.* 12 лекций по вычислительной математике. – М.: Изд-во МФТИ, 1995.
34. *Крылов В.И., Бобков В.В., Монастырный П.И.* Вычислительные методы. Т.1. – М.: Наука, 1976.
35. *Крылов В.И., Бобков В.В., Монастырный П.И.* Начала теории вычислительных методов. Линейная алгебра и нелинейные уравнения. – Минск: Наука и техника, 1985.
36. *Ланс Дж.Н.* Численные методы для быстродействующих вычислительных машин. – М.: ИЛ, 1962.
37. *Лесин В.В., Лисовец Ю.П.* Основы методов оптимизации. – М.: Изд. МАИ, 1995.
38. *Мак-Кракен Д., Дорн У.* Численные методы и программирование на Фортране. – М.: Мир, 1977.
39. Под ред. *Монастырного П.И.* Сборник задач по методам вычислений. – М.: Наука, 1994.
40. *Мысовских И.П.* Лекции по методам вычислений. – М.: Физматгиз, 1962.
41. *Ортега Дж., Пул У.* Введение в численные методы решения дифференциальных уравнений. – М.: Наука, 1986.
42. *Ортега Дж., Рейнболдт В.* Итерационные методы решения нелинейных систем уравнений со многими неизвестными. – М.: Мир, 1975.
43. *Островский А.М.* Решение уравнений и систем уравнений. – М.: ИЛ, 1963.
44. *Парлетт Б.* Симметричная проблема собственных значений. – М.: Мир, 1983.
45. *Пузачев В.С.* Лекции по функциональному анализу. – М.: Изд. МАИ, 1996.

46. *Рябенкий В.С.* Введение в вычислительную математику. – М.: Наука, 1994.
47. *Самарский А.А., Гулин А.В.* Численные методы. – М.: Наука, 1989.
48. *Самарский А.А., Николаев Е.С.* Методы решения сеточных уравнений. – М.: Наука, 1978.
49. *Тихонов А.Н., Арсенин В.Я.* Методы решения некорректных задач. – М.: Наука, 1985.
50. *Трауб Дж.* Итерационные методы решения уравнений. – М.: Мир, 1985.
51. *Треногин В.А.* Функциональный анализ. – М.: Наука, 1980.
52. *Турчак Л.И.* Основы численных методов. – М.: Наука, 1987.
53. *Уилкинсон Дж.Х.* Алгебраическая проблема собственных значений. – М.: Наука, 1970.
54. *Фаддеев Д.К., Фаддеева В.Н.* Вычислительные методы линейной алгебры. – М.: Физматгиз, 1960.
55. *Форсайт Дж., Молер К.* Численное решение систем линейных алгебраических уравнений. – М.: Мир, 1969.
56. *Хейгеман Л., Янг Д.* Прикладные итерационные методы. – М.: Мир, 1986.
57. *Хемминг Р.В.* Численные методы. – М.: Наука, 1968.
58. *Шарковский А.Н., Майстренко Ю.Л., Романенко Е.Ю.* Разностные уравнения и их приложения. – Киев: Наукова думка, 1986.
59. *Altman M.* An optimum cubically convergent iterative method of inverting a linear bounded operator in hilbert space. «Pacific J.Math.», 10, 1960, № 4, 1107-1113.
60. *Diaconu A.* Sur quelques méthodes itératives combinées. «Matematica» (RSR), 22(45), 1980, № 2, с.247-261.
61. *Schulz G.* Iterative Berechnung der reziproken Matrix. ZAMM, 13, 1933, с.57-59.

## Предметный указатель

- Алгоритм Вегстейна 194  
— гибридный 174  
— численно устойчивый 21  
Аппроксимационный аналог метода Ньютона 219  
Билинейное отображение 253  
Бифуркация решений 199  
Ведущий элемент 37  
Вековой определитель 97  
Вторая производная Фреше 253  
Главный элемент 37  
Горнер 173  
Градиент 252  
Граница погрешности 8  
Декомпозиция 42  
Диагональное преобладание 70  
Дифференциал Фреше 250  
Дифференцируемость по Фреше 250  
Евклидово пространство 246  
Единица объема информационного запроса 173  
Задача на собственные значения 96  
— неустойчивая 21  
— о неподвижной точке 176, 213  
Значащая цифра 12  
— — верная 13  
Запятая фиксированная 13  
— плавающая 14  
Зацикливание 199  
Исчерпывающий спуск 226  
Итерационная функция 186  
Итерационный процесс двухступенчатый 219  
— — квадратично сходящийся 152  
— — нестационарный 77, 186  
— — первого порядка 152  
— — стационарный 77  
Итерированный вектор 102  
Квадрирование корней 207  
Ключевой элемент 121  
Компактная схема Гаусса 44  
Конечные разности 188  
Константа Липшица 179  
Коэффициент сжатия 179  
— схемы Горнера 203  
— чувствительности 28  
Краевые условия 51  
Лемма Неймана 64  
Мантисса 14  
Масштабирование 38  
Матрица Гессе 253  
— Гильберта 24  
— итерирования (перехода) 70  
— отражения (Хаусхолдера) 130  
— плоских вращений 119  
— простой структуры 101  
— с диагональным преобладанием 43  
— сопровождающая 98  
— Хессенберга (правая почти треугольная) 130  
— Якоби 251  
Машинная бесконечность 15  
Машинное слово 15  
Машинный ноль 15  
— эpsilon 15  
Машинных чисел диапазон 14  
— — точность представления абсолютная 14  
— — — — относительная 14  
Метод Бернулли 208  
— бисекций 137, 148  
— Брауна 224  
— вариационного типа 83  
— Вегстейна 192  
— Вестерфильда 206  
— вращений 56  
— — Якоби 119  
— второго порядка 157  
— Гаусса 37  
— Гаусса-Зейделя 74  
— главных элементов 38  
— Горнера 204



**Метод градиентный** 226  
 — двухшаговый 83, 170  
 — дихотомии (вилки, проб) 148  
 — Зейделя 72  
 — итерационный 33  
 — — двухслойный 82  
 — — трехслойный 83  
 — — фон Мизеса 103  
 — касательных 154  
 — квадратных корней 50  
 — квазиньютоновский 228  
 — Лагранжа (Маклорена) 205  
 — Лина (предпоследнего остатка) 204  
 — линеаризации 154  
 — Лобачевского (Лобачевского-Греффе, Данделена) 207  
 — минимальных невязок 85  
 — наискорейшего спуска 226  
 — Некрасова 74  
 — нестационарный 82  
 — неустойчивый 21  
 — неясный 82  
 — нижней релаксации 80  
 — Ньютона 154, 217  
 — — модифицированный (упрощенный) 168, 186, 218  
 — — огрубленный 168  
 — — разностный (конечно-разностный, дискретный) 169, 220  
 — — с параметром 166  
 — — с последовательной аппроксимацией обратных матриц 219  
 — Ньютона-Канторовича 231  
 — Ньютона-Рафсона 154  
 — Ньютона-Шрёдера 166  
 — обратных итераций 112  
 — одновременных смещений 72  
 — отражений 133  
 — переменной метрики (ДФП) 228  
 — переменных направлений 83  
 — покоординатных итераций 215  
 — полной релаксации 78  
 — половинного деления 148  
 — попеременно-треугольный 83  
 — последовательной верхней релаксации 80  
 — последовательных приближений 176, 213

**Метод последовательных смещений** 72  
 — прогонки 52  
 — простых итераций 64, 176, 213  
 — прямого поиска 228  
 — прямой 32  
 — расщепления 83  
 — регуляризации 21  
 — рекурсивный 218  
 — релаксации 78  
 — Ричардсона 82  
 — секущих 170, 221  
 — — Бройдена 221, 228  
 — скалярных произведений 107  
 — сопряженных градиентов 84, 228  
 — спуска 226  
 — стационарный 82  
 — степенной 103  
 — Стеффенсена 170  
 — точный 33  
 — установления 81  
 — хорд (пропорциональных частей, линейной интерполяции) 150  
 — частных Рэлея 108  
 — Шульца 89  
 — — зейделя модификация 89  
 — явный 82  
 — — итерационный с чебышевским набором параметров 82  
 — Якоби 69  
 — — циклический с барьерами 123

**Многочлен возмущенный** 19

**Направление минимизации** 225  
 — спуска 226

**Невязка** 26, 59, 86, 157, 220

**Непрерывный аналог итерационного метода** 81

**Норма вектора** 245

- Гёльдера 246
- евклидова 246, 247
- матрицы 247
- мультипликативная 247
- подчиненная 247
- спектральная 248
- Фробениуса 247

**Норма-максимум** 246

**Норма-сумма** 246

**Нормированное пространство** 245

- Обратная задача теории погрешностей 10
- Обратные итерации 112
  - со сдвигами 113
  - с отношениями Рэлея 115
  - с переменными сдвигами 115
- Обратный ход 36
- Обреченный элемент 122
- Округление правильное 14
- Основание вещественного числа 14
- Отношение Рэлея 100
- Оценка погрешности 8
  - апостериорная 66
  - априорная 66
- Ошибка 157
- Параметр релаксации 78
- Погрешность абсолютная 8
  - безусловная 27
  - задачи 7
  - метода 7
  - неустраняемая 7
  - округлений 7
  - относительная 8
  - полная 8
  - условная 27
  - устраняемая 7
- Полная проблема собственных значений 96
- Порядок вещественного числа 14
- Последовательность Рэлея 116
- Постоянная Фейгенбаума 201
- Поправка 59, 157, 217, 220
- Правило ложного положения 150
  - Ньютона 28
  - вычисления арифметических корней 163
  - Чеботарёва 12
- Предел векторной последовательности 248
  - последовательности матриц 249
- Преобразование отражения (Хаусхолдера) 130
  - плоских вращений Гивенса 133
  - подобия 117
- Прием Гарвика 110
- Принцип А.Н.Крылова 12
  - равных влияний 10
  - неподвижной точки 177
  - сжимающих отображений 177
- Пробная точка 147
- Прогонка корректная 52
  - обратная 52
  - прямая 52
  - устойчивая 52
- Прогоночные коэффициенты 52
- Производная Гато (слабая) 250
  - Фреше (сильная) 250
- Промежуток неопределенности корня 147
  - существования — 147
- Процедура исчерпывания 138
- Процесс Герона 164
- Прямой ход 36
- Разностные отношения 169
- Реальный метод простых итераций 91
- Сверхрелаксация 80
- Сдвиги 137
- Сжимающая функция 177
- Симметризация Гаусса 77
- Система возмущенная 20
  - ленточная 51
  - нормальная 76
- Собственная пара матрицы 96
- Собственное число 96
- Собственный вектор 97
  - элемент 97
- Спектр матрицы 99
- Спектральный радиус 24
- Способ перебора 146
- Средняя скорость сходимости 154
- Схема Горнера 203
  - единственного деления 44
  - Холецкого 44, 50
- Сходимость асимптотически линейная 152
  - глобальная 154
  - квадратичная 152
  - кубическая 152
  - линейная 151
  - локальная 154
  - по норме 249
  - сверхлинейная 151
  - со скоростью геометрической прогрессии 152
  - с  $p$ -м порядком 151
  - $j$ -шаговая с порядком  $p$  173
- Счет на установление 103

- Теорема Банаха 64
- Больцано-Коши 144
  - Декарта 207
  - Островского-Рейча 79
  - Штурма 207
- Точка бифуркации удвоения периода 201
- Трехточечное разностное уравнение второго порядка 51
- Уравнение логистическое 197
- Некрасова 74
  - скалярное (числовое, конечное) 141
  - характеристическое 97
- Условие Гёльдера 242
- Коши-Липшица 179
  - Липшица 179, 235
  - релаксации 175, 226
  - согласованности норм 247
  - Фурье 161
- Усовершенствованный метод последовательных приближений 192
- Устойчивый цикл 199
- Факторизация 42
- Формула Бинэ 173
- Формулы Крамера 34
- Функция сжатия 177
- Центр итерации 178
- Частичная проблема собственных значений 96
- Частичное упорядочивание 37
- Числа Фибоначчи 173
- Число обусловленности 23
- — ненулевого простого корня 30
  - — Тогда 24
- Шаговый множитель 225
- Элемент цикла 201
- ПВР-метод 80
- ADI-метод 83
- INVIT-алгоритм 112
- LU-алгоритм 125
- LU-разложение 40
- PM-алгоритм 104
- QR-алгоритм 129
- RQI-алгоритм 115
- SP-алгоритм 107
- SOR-метод 80
- $U^T U$ -разложение 48
- regula falsi-метод 150
- $q$ -сходимость 153
- $r$ -сходимость 153
- $l_p$ -норма 246
- $\Delta^2$ -алгоритм Эйткена 189
- $\Delta^2$ -преобразование 188
- $\Delta^2$ -процесс Эйткена 188

*Учебное издание*

**Верещагин Валентин Михайлович**

**ЧИСЛЕННЫЕ МЕТОДЫ**  
(линейная алгебра и нелинейные уравнения)

Редактор *Ж.И. Яковлева*  
Художественный редактор *Ю.Э. Иванова*  
Технический редактор *Л.А. Овчинникова*  
Оригинал-макет выполнен *С.В. Высоцким*

ЛР № 010146 от 25.12.96. Изд. № ФМ-192. Сдано в набор 28.09.99.  
Подп. в печать 22.10.99. Формат 60x88<sup>1</sup>/<sub>16</sub>. Бум. газетн. Гарнитура «Таймс»  
Печать офсетная. Объем 16,66 усл. печ. л. 16,91 усл. кр.-отт.  
Тираж 8000 экз. Зак. № 9759.

ГУП издательство «Высшая школа», 101430, Москва, ГСП-4,  
Неглинная ул., д. 29/14.

Отпечатано с оригинал-макета на Государственном унитарном  
предприятии Смоленский полиграфический комбинат  
Министерства Российской Федерации по делам печати,  
телерадиовещания и средств массовых коммуникаций.  
214020, Смоленск, ул. Смольянинова, 1.

**Вержбицкий В.М.**  
В31 Численные методы (линейная алгебра и нелинейные уравнения): Учеб. пособие для вузов. — М.: Высш. шк., 2000. — 266 с.: ил.

ISBN 5-06-003654-5

В книге последовательно излагаются численные методы решения линейных алгебраических систем, обращения матриц, вычисления собственных чисел и собственных векторов матриц, а также методы решения нелинейных скалярных уравнений и систем таких уравнений. Показываются идеи, выводы и взаимосвязь методов, обсуждается их сравнительная эффективность, обосновывается сходимость, приводятся алгоритмы, численные примеры, задания для упражнений и лабораторных работ. Наряду с методическими погрешностями изучаемых процессов, уделяется внимание и погрешностям, связанным с их компьютерными реализациями.

Пособие предназначено для студентов математических и инженерных специальностей вузов и может быть полезно всем, кто связан с изучением и применением вычислительной математики.

УДК 519.6  
ББК 22.19

**В издательстве**  
**«ВЫСШАЯ ШКОЛА»**  
**ИМЕЮТСЯ В НАЛИЧИИ**  
**КНИГИ ИЗ СЕРИИ**  
**«ВЫСШАЯ МАТЕМАТИКА»**

*Архипов Г. И., Садовничий В. А., Чубариков В. Н.* Лекции по математическому анализу. Учебник. — 695 с.

Книга является учебником по курсу математического анализа и посвящена дифференциальному и интегральному исчислениям функций одной и нескольких переменных. В ее основу положены лекции, прочитанные авторами на механико-математическом факультете МГУ им. М. В. Ломоносова. В учебнике предложен новый подход к изложению ряда основных понятий и теорем анализа, а также и к самому содержанию курса. Она доступна широкому кругу читателей, а первая ее часть может быть использована при изучении ряда тем по алгебре и началам математического анализа в математических школах.

*Виноградов И. М.* Элементы высшей математики. (Аналитическая геометрия. Дифференциальное исчисление. Основы теории чисел). Учебник. — 511 с.

Книга написана на основе лекций по высшей математике, которые автор читал в Ленинградском политехническом институте, а также его широко известного учебника по теории чисел. Математика воспринимается автором этой книги как мощное орудие для решения проблем естествознания. Поэтому он стремится возможно быстрее дать в руки студентов это орудие, на большом количестве простых, но принципиально важных примеров научить пользоваться ими. При этом он не делает различия, будет ли впоследствии оно применяться в инженерной работе или в абстрактных математических исследованиях. Часть, посвященная теории чисел, является классическим учебником, который стал настольной книгой специалистов в этой области.

*Нечаев В. И.* Элементы криптографии (Основы теории защиты информации). Учебное пособие. — 109 с.

Книга является первым учебным пособием по теории защиты информации, фундаментом которой является прикладная теория чисел. В основу книги положены лекции, читавшиеся автором на математическом факультете Московского педагогического государственного университета. В пособии рассматриваются современные методы шифрования по открытому ключу и электронная подпись. Книга отличается широтой охвата материала в области защиты информации. В ней также содержится интересный исторический очерк развития криптографии. Она доступна широкому кругу читателей, начиная с учащихся старших классов математических школ.

*Привалов И. И.* Введение в теорию функций комплексного переменного. Учебник. — Изд. 14-е. — 432 с.

Книга написана известным автором и является одним из наиболее апробированных и хорошо себя зарекомендовавших учебников по теории функций комплексного переменного. Она отличается строгостью выводов и простотой изложения материала.

*Садовничий В. А.* Теория операторов. Учебник. — Изд. 3-е. — 368 с.

Учебник соответствует программе курсов «Функциональный анализ», «Теория операторов», «Анализ III», которые читаются в университетах и педагогических вузах. В книге приведены основные теоретико-множественные понятия, представлена общая теория метрических, топологических, линейных топологических и нормированных пространств, общая теория меры, измеримых функций и интеграла Лебега. Подробно рассмотрены теория операторов в гильбертовом пространстве, спектральная теория самосопряженных операторов, применения методов теории аналитических функций в спектральной теории несамосопряженных операторов, теория преобразования Фурье и обобщенные функции.

## КНИГИ, НАХОДЯЩИЕСЯ В ПЕЧАТИ

*Бахвалов Н. С., Лапин А. В., Чижонков Е. В.* Численные методы в задачах и упражнениях. Учебное пособие. — 26 л.

Учебное пособие содержит элементы теории, примеры решений задач и упражнения для самостоятельной работы. Представленные задачи разбиты по рекомендуемым темам семинарских занятий, а их подбор призван способствовать закреплению материала, излагаемого в теоретическом курсе. Типовые задачи снабжены решениями, которые могут быть использованы студентами для самостоятельного изучения предмета и овладения общими принципами применения вычислительных методов. Ответы и указания помогут преподавателям в выборе полезных и интересных задач в соответствии со спецификой вуза.

*Виноградова И. А., Олехник С. Н., Садовничий В. А.* Задачи и упражнения по математическому анализу. В 2-х ч. Учебное пособие. — 80 л.

Учебное пособие соответствует программе курса математического анализа для студентов механико-математических и математических факультетов университетов, пединститутов и технических вузов. Задачник отражает современное развитие математики и должен заменить известное пособие Б. Д. Демидовича. В отличие от задачника Б. Д. Демидовича большая часть задач в данном пособии приводится с решениями, в связи с чем оно может быть полезно при самостоятельном изучении предмета.



*Гаилов Г. И., Чубариков В. Н.* Арифметика. Алгоритмы. Сложность вычислений. Учебное пособие. — 25 л.

Книга состоит из 17 параграфов, тесно связанных целью, поставленной авторами — ознакомить читателя с материалом, который скорейшим путем приводит к современным проблемам теории чисел и теории сложности арифметических алгоритмов. Задачи в каждом параграфе сгруппированы по единству идей и содержания, иногда их объединяют связанные между собой сходные методы решения или просто схожесть формулировки. Часто группы задач заканчиваются красивой и очень трудной задачей. Но если прорешать все задачи подряд, то и она не покажется трудной. Каждый параграф сопровождается указаниями и решениями. В книге имеется достаточное количество и несложных задач.

ISBN 5-06-003654-5



9 785060 036541