

**В.С. РЯБЕНЬКИЙ**

**ВВЕДЕНИЕ  
В ВЫЧИСЛИТЕЛЬНУЮ  
МАТЕМАТИКУ**





ФИЗТЕХОВСКИЙ  
УЧЕБНИК

В.С. РЯБЕНЬКИЙ

ВВЕДЕНИЕ  
В ВЫЧИСЛИТЕЛЬНУЮ  
МАТЕМАТИКУ

Издание третье,  
исправленное и дополненное

*Рекомендовано Учебно-методическим объединением  
высших учебных заведений Российской Федерации  
по образованию в области прикладных математики  
и физики в качестве учебного пособия для студентов  
высших учебных заведений по направлению  
«Прикладные математика и физика»*



МОСКВА  
ФИЗМАТЛИТ<sup>®</sup>  
2008

УДК 519.6 (075.8)

ББК 22.19

Р 98

Серия «Физтеховский учебник»

*Редакционный Совет*

*Кудрявцев Н. Н. (председатель)*

*Белоусов Ю. М.,*

*Гладун А. Д.,*

*Кондранин Т. В.,*

*Петров И. Б.,*

*Половинкин Е. С.,*

*Самарский Ю. А.,*

*Сон Э. Е.,*

*Холодов А. С.*

Рябенький В. С. **Введение в вычислительную математику.** — 3-е изд., испр. и доп. — М.: ФИЗМАТЛИТ, 2008. — 288 с. — (Физтеховский учебник). — ISBN 978-5-9221-0926-0.

В книге изложены основные понятия и идеи, используемые для преобразования математических моделей к виду, удобному для вычисления с помощью компьютера. Изложение ведется на материале вычислительных задач математического анализа, алгебры и дифференциальных уравнений. Впервые в учебной литературе отражен метод разностных потенциалов для численного решения краевых задач математической физики.

Для студентов и преподавателей механико-математических и физических факультетов университетов, МФТИ, МИФИ, технических вузов.

Рекомендовано Учебно-методическим объединением высших учебных заведений Российской Федерации по образованию в области прикладных математики и физики в качестве учебного пособия для студентов высших учебных заведений по направлению «Прикладные математика и физика».

## **ПРЕДИСЛОВИЕ К ТРЕТЬЕМУ ИЗДАНИЮ**

Настоящее издание отличается от двух предыдущих (1994 и 2000) тем, что добавлено изложение идеи построения вычислительных схем метода конечных элементов на базе проекционного метода Галёркина (гл. 11, § 2). Кроме того, полностью переработана гл. 13, дающая представление о методе разностных потенциалов и совпадающая теперь в основном с введением в монографию [17].

Заметим еще, что цикл сопровождающих учебник компьютерных учебных демонстрационных работ, о котором говорилось в предисловии к первому изданию, к настоящему времени существенно развит и издан в виде книги [29].

2007 г.

*В.С. Рябенький*

## **ПРЕДИСЛОВИЕ К ПЕРВОМУ ИЗДАНИЮ**

Предлагаемая книга соответствует вводному полугодовому курсу вычислительной математики для студентов Московского физико-технического института, различные варианты которого более двадцати лет читает автор. Книга адресована студентам технических вузов различного профиля.

Цель книги — изложить достаточно стабильные и общие понятия и идеи, относящиеся к преобразованию математических моделей различных прикладных задач к виду, удобному для вычисления их решений с помощью компьютеров. Изучение книги должно давать возможность сравнительно легко доучиваться для работы в разнообразных и быстро развивающихся конкретных областях. Общая методология излагается и иллюстрируется в книге на материале численных методов математического анализа, алгебры и дифференциальных уравнений, поскольку эти методы лучше развиты, содержат достаточно совершенные алгоритмы и часто применяются в разнообразных конкретных прикладных задачах.

Наряду с традиционным для учебной литературы по вычислительной математике материалом книга дает представление о граничных уравнениях с проекторами и методе разностных потенциалов для их численного решения, а также содержит изложение локальных формул гладкого восполнения функций.

При изложении метода конечных разностей использована книга С.К. Годунова и В.С. Рябенького «Разностные схемы».

Структура книги допускает различную глубину изучения материала. Это достигается расположением материала и подбором задач.

Для более полного и глубокого изучения многих вопросов, освещенных в учебнике, можно воспользоваться книгами К.И. Бабенко [1], Н.С. Бахвалова, Н.П. Жидкова и Г.М. Кобелькова [2], С.К. Годунова и В.С. Рябенького [5, 6], Н.Н. Калиткина [9], О.В. Локуциевского и М.Б. Гаврикова [11], Г.И. Марчука [13], А.А. Самарского [19], А.А. Самарского и Е.С. Nikolaeva [21], а также другими книгами и журнальной литературой, указанными в тексте.

Отметим еще, что изучение книги можно сделать более наглядным и увлекательным, если воспользоваться предназначенным для этого циклом компьютерных учебных демонстрационно-лабораторных работ [29]. Этот цикл создан сотрудниками кафедры вычислительной математики МФТИ В.Д. Ивановым, П.Н. Коротиным, В.И. Косаревым, И.Б. Петровым (руководитель работ), В.Б. Пироговым, В.С. Рябеньким (научный руководитель), Д.С. Северовым, А.Г. Тормасовым, С.В. Устюжниковым. В программной реализации цикла участвовали студенты В.В. Бойков, Д.Л. Будько, К.Б. Бухаров, А.Ю. Езерский, А.Б. Константинов, С.А. Корытник, Ю.П. Кравченко, Ю.Д. Крикунов, Д.В. Лунев, В.А. Торгашов, А.А. Терехин, Г.Л. Химичев.

В заключение автор сердечно благодарит академика О.М. Белоцерковского, который в момент создания кафедры вычислительной математики МФТИ поручил ему чтение первого обязательного общего курса вычислительной математики. Олег Михайлович неоднократно обсуждал с автором цели курса в условиях физтеха и соответствующий этим целям характер курса. Подготовка и чтение лекций оказались началом работы автора над предлагаемой книгой.

Автор сердечно благодарит также всех своих коллег по работе в МФТИ и особенно В.В. Демченко, В.Д. Иванова, В.И. Косарева, В.Я. Митницкого, Н.П. Онуфриеву, И.Б. Петрова, В.Б. Пирогова, Л.М. Стрыгину, Р.П. Федоренко, А.С. Холодова и Л.А. Чудова за многолетнее сотрудничество, оказавшее влияние на книгу.

Автор сердечно благодарит своего коллегу К.В. Брушлинского, прочитавшего книгу в рукописи, за ряд полезных замечаний.

Автор выражает глубокую благодарность профессору К.И. Бабенко и профессору О.В. Локуциевскому, с которыми также плодотворно обсуждал вопросы преподавания.

1987 г.

*В.С. Рябенький*

## ВВЕДЕНИЕ

Современная вычислительная математика ориентирована на использование компьютеров для прикладных расчетов. Любые математические приложения начинаются с построения модели явления (изделия, действия, ситуации или другого объекта), к которому относится изучаемый вопрос. Классическими примерами математических моделей могут служить определенный интеграл, уравнение колебаний маятника, уравнение теплообмена, уравнения упругости, уравнения электромагнитного поля и другие уравнения математической физики. Назовем еще для контраста модель формальных рассуждений — алгебру Буля.

Основополагающими средствами изучения математических моделей являются аналитические методы: получение точных решений в частных случаях (например, табличные интегралы), разложения в ряды. Определенную роль издавна играли приближенные вычисления. Например, для вычисления определенного интеграла использовались квадратурные формулы.

Появление в середине XX века электронных вычислительных машин (компьютеров) радикально расширило возможности приложения математических методов в традиционных областях (механике, физике, технике) и вызвало бурное проникновение математических методов в нетрадиционные области (управление, экономику, химию, биологию, психологию, лингвистику, экологию и т. п.).

Компьютер дает возможность запоминать большие (но конечные) массивы чисел и производить над ними арифметические операции и сравнения с большой (но конечной) скоростью по заданной вычислителем программе. Поэтому с помощью компьютера можно реализовывать только те математические модели, которые описываются конечными наборами чисел, и использовать конечные последовательности арифметических действий, а также сравнений чисел по величине (для автоматического управления дальнейшими вычислениями).

В традиционных областях математическими моделями служат функции, производные, интегралы, дифференциальные уравнения. Для использования компьютеров эти исходные модели надо приближенно заменить такими, которые описываются конечными наборами чисел с указанием конечных последовательностей действий (конечных алгоритмов) для их обработки. Например, функцию следует заменить таблицей; производную

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$

заменить приближенной формулой

$$f'(x) \approx \frac{f(x+h) - f(x)}{h};$$

определенный интеграл — суммой; краевую задачу для дифференциального уравнения — задачей об отыскании упорядоченного набора чисел (таблицы) значений решения в узлах некоторой сетки, причем так, чтобы выбор шага сетки позволял достигать любой требуемой точности. Оказывается, из двух на первый взгляд равноценных способов один может оказаться принципиально непригодным из-за того, что доставляемое им приближенное решение не стремится к искомому при уменьшении шага сетки, или из-за катастрофически сильной чувствительности к погрешностям округления.

Теория таких моделей и алгоритмов составляет предмет вычислительной математики. Эта теория тесно связана с теориями приближения и интерполяции функций, уравнений с частными производными, интегральных уравнений, информационной сложности функциональных классов, алгоритмов, а также с языками программирования для расчетов на компьютере и т. п. Современные вычислительные методы позволяют, например, рассчитать характеристики обтекания газом тела заданной формы, что недоступно аналитическим методам (подобно вычислению нетабличных интегралов).

С использованием компьютеров стал возможен вычислительный эксперимент, т. е. расчет в целях проверки гипотез, а также в целях наблюдения за поведением модели, когда заранее не известно, что именно заинтересует исследователя. В процессе численного эксперимента происходит по существу уточнение исходной математической постановки задачи. В процессе расчетов на компьютере происходит накопление информации, что дает возможность в конечном счете отбирать наиболее интересные ситуации. На этом пути сделано много наблюдений и открытий, стимулирующих развитие теории и имеющих важные практические применения.

С помощью компьютеров возможно применение математических методов и в нетрадиционных областях, где не удается построить компактные математические модели вроде дифференциальных уравнений, но удается построить модели, доступные запоминанию и изучению на компьютере. Модели для компьютеров в этих случаях представляют собой цифровое кодирование схемы изучаемого объекта (например, языка) и отношений между его элементами (словами, фразами). Сама возможность изучения таких моделей на компьютере стимулирует появление этих моделей, а для создания обозримой модели необходимо выявление законов, действующих в исходных объектах. С другой стороны, получаемые на компьютере результаты (например, машинный перевод упрощенных текстов с одного языка на другой) вносят критерий практики в оценку теорий (например, лингвистических теорий), положенных в основу математической модели.

Благодаря компьютерам стало возможным рассматривать вероятностные модели, требующие большого числа пробных расчетов, имитационные модели, которые отражают моделируемые свойства объекта без упрощений (например, функциональные свойства телефонной сети).

Разнообразие задач, где могут быть использованы компьютеры, очень велико. Для решения каждой задачи нужно знать многое, связанное именно с этой задачей. Естественно, этому нельзя научиться впрок.

Целью этой книги является сообщение тех основных понятий, идей и методов, владение которыми позволяет сравнительно быстро научиться работать в конкретных областях. Это реализуется на материале вычислительных задач алгебры, математического анализа, дифференциальных уравнений, поскольку здесь методы хорошо развиты и употребляются в далеких друг от друга областях.

Назовем некоторые общие понятия и идеи, которые требуют внимания и наполняются конкретным содержанием в зависимости от задачи, которую предстоит решать с помощью компьютера. Это — дискретизация задачи; обусловленность задачи; погрешность численного метода; вычислительная устойчивость алгоритма; сравнение алгоритмов по полноте используемой ими входной информации, по объему используемой памяти и числу арифметических действий. Алгоритмы могут обладать возможностью распараллеливания для одновременного проведения вычислений на многопроцессорном компьютере. Среди основных методов вычислительной математики одним из плодотворных является комбинированное использование аналитических и компьютерных средств.

Здесь, во введении, мы предварительно познакомим читателя с перечисленными понятиями. Это даст некоторое общее представление о предмете вычислительной математики и подготовит к изучению дальнейшего материала.

## § 1. Дискретизация

Пусть требуется найти приближенное решение какой-либо задачи, в которой в качестве входных данных участвует какая-либо функция  $f(x)$ , определенная на всем (бесконечном) множестве точек отрезка  $0 \leq x \leq 1$ . Значения этой функции при каждом фиксированном  $x$  можно получить измерениями или вычислениями. Для запоминания этой функции в памяти компьютера необходимо приближенно описать ее таблицей значений на некотором конечном множестве отдельных точек  $x_0, x_1, \dots, x_n$ . Это — простейший пример дискретизации задачи: от задачи запоминания функции на отрезке  $[0, 1]$  мы перешли к задаче запоминания таблицы значений на дискретном множестве точек  $x_0, x_1, \dots, x_n$  из этого отрезка.

Пусть функция  $f(x)$  имеет достаточное число производных, а нам требуется вычислить ее производную  $f'(x)$  в данной точке  $x$ . Задачу отыскания

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h},$$

содержащую предельный переход, можно заменить приближенно задачей вычисления по одной из формул

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}, \quad (1)$$

$$f'(x) \approx \frac{f(x) - f(x-h)}{h}, \quad (2)$$

$$f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}. \quad (3)$$

Для замены производной  $f''(x)$  можно воспользоваться формулой

$$f''(x) \approx \frac{f(x+h) - 2f(x) + f(x-h)}{h^2}. \quad (4)$$

Все эти формулы уточняются при уменьшении  $h$ , а при каждом фиксированном  $h$  определены для конечных наборов значений функции и используют только арифметические операции. Эти формулы — примеры дискретизации задачи о вычислении производных  $f'(x)$ ,  $f''(x)$ .

Рассмотрим краевую задачу

$$\begin{aligned} \frac{d^2y}{dx^2} - x^2y &= \cos x, \quad 0 \leq x \leq 1, \\ y(0) &= 2, \quad y(1) = 3, \end{aligned} \quad (5)$$

об отыскании функции  $y(x)$ , определенной на отрезке  $0 \leq x \leq 1$ . Для построения приближенной дискретной модели этой задачи осуществим следующие два шага.

1) Разобьем отрезок  $0 \leq x \leq 1$  на  $N$  равных частей длины  $h = N^{-1}$ , а вместо функции  $y(x)$  будем искать набор значений  $y_0, y_1, \dots, y_N$  этой функции в точках  $x_k = kh$  ( $k = 0, 1, \dots, N$ ). В точках  $x_k$  ( $k = 1, 2, \dots, N-1$ ) заменим производную  $y''(x)$  приближенно по формуле (4) и получим

$$\frac{y_{k+1} - 2y_k + y_{k-1}}{h^2} - (kh)^2 y_k = \cos kh, \quad k = 1, 2, \dots, N-1. \quad (6)$$

2) Кроме того, в силу граничных условий (5) положим

$$y_0 = 2, \quad y_N = 3. \quad (7)$$

Система  $(N+1)$  линейных уравнений (6), (7) относительно того же числа неизвестных  $y_0, y_1, \dots, y_N$  является дискретным аналогом задачи (5).

Есть основания полагать, что с ростом  $N$  решение задачи (6), (7) есть все более точная таблица значений решения задачи (5) (в дальнейшем это будет показано).

Обозначим исходную (континуальную) краевую задачу через  $M_\infty$ , а дискретную краевую задачу (6), (7) через  $M_N$ . Тогда можно сказать, что задаче  $M_\infty$  мы сопоставили бесконечную последовательность дискретных задач  $M_N$  ( $N = 2, 3, \dots$ ).

Вычисляя решение задачи  $M_N$  при каком-либо фиксированном  $N$ , мы имеем дело с конечным набором чисел, задающих входные данные, и с конечным набором чисел  $y_0, y_1, \dots, y_N$ , подлежащих отысканию. Однако вычислительная математика обычно ставит своей целью предложить именно последовательность уточняющихся дискретных моделей  $M_N$ , так как это дает возможность выбрать такое  $N$ , которое обеспечивает выполнение требований к точности.

Переход от континуальной задачи  $M_\infty$  к последовательности  $\{M_N\}$  ее дискретных моделей возможен многими способами. Пусть  $\{M_N\}$ ,  $\{M'_N\}$  — какие-нибудь две последовательности таких моделей, причем вычисление решений дискретных задач  $M_N$ ,  $M'_N$  требует равных затрат. Тогда предпочтение надо отдать тому способу дискретизации, при котором решение дискретной задачи может служить решением исходной задачи с заданной точностью при меньшем значении  $N$ .

Бывает, что из двух, казалось бы, равноценных способов дискретизации  $M_N$  и  $M'_N$  один при возрастании  $N$  дает все более точное приближение к решению континуальной задачи  $M_\infty$ , а другой приводит к «приближенному решению» задачи, которое с ростом  $N$  теряет какое-либо сходство с искомым решением. С этой ситуацией и способами ее преодоления мы встретимся в части III книги.

### Задача

Пусть  $f(x)$  имеет ограниченные производные до требуемого порядка.

Показать, что погрешности приближенных формул (1)–(4) суть величины соответственно  $O(h)$ ,  $O(h)$ ,  $O(h^2)$ ,  $O(h^2)$ .

## § 2. Обусловленность

Во всякой задаче требуется по входным данным сделать заключение о каких-либо свойствах решения. Похожие, на первый взгляд, задачи могут резко отличаться чувствительностью интересующих нас свойств решения к возмущению входных данных. Если эта чувствительность «мала», то задача считается хорошо обусловленной; в противном случае — плохо обусловленной. Плохо обусловленные задачи обычно не только предъявляют высокие требования к точности задания входных данных, но и более трудны для вычислений.

Пример. Пусть концентрация  $y = y(t)$  некоторого вещества в момент времени  $t$  есть функция, удовлетворяющая дифференциальному уравнению

$$\frac{dy}{dt} - 10y = 0.$$

Фиксируем произвольно  $t_0$  ( $0 \leq t_0 \leq 1$ ) и делаем приближенное измерение  $y_0^*$  концентрации  $y_0 = y(t_0)$ , получив

$$y|_{t=t_0} = y_0^*.$$

Задача состоит в определении концентрации  $y = y(t)$  в произвольный момент времени  $t$  из отрезка  $0 \leq t \leq 1$ .

Если бы число  $y_0 = y(t_0)$  было известно точно, то можно было бы указать точную формулу

$$y(t) = y_0 e^{t-t_0}$$

для концентрации. Но мы знаем лишь приближенное значение  $y_0^* \approx y_0$  числа  $y_0$ . Поэтому вместо  $y(t) = y_0 e^{10(t-t_0)}$  мы можем указать лишь приближенную формулу  $y^*(t) = y_0^* e^{10(t-t_0)}$ . Очевидно, что погрешность  $y^* - y$  выражается формулой

$$y^*(t) - y(t) = (y_0^* - y_0) e^{10(t-t_0)}, \quad 0 \leq t \leq 1.$$

Допустим, нам нужно произвести замер  $y_0^*$  с такой точностью  $\delta$ ,  $|y_0^* - y_0| < \delta$ , чтобы гарантировать некоторую заданную точность  $\varepsilon > 0$  всюду на отрезке  $0 \leq t \leq 1$ , т. е. гарантировать оценку

$$|y^*(t) - y(t)| < \varepsilon, \quad 0 \leq t \leq 1.$$

Очевидно,  $\max_{0 \leq t \leq 1} |y^*(t) - y(t)| = |y^*(1) - y(1)| = |y_0^* - y_0| e^{10(1-t_0)}$ .

Отсюда получаем следующее требование к точности  $\delta$  измерения  $y_0$ :

$$\delta \leq \varepsilon e^{-10(1-t_0)}.$$

Пусть измерение  $y_0$  производится в момент  $t_0 = 0$ . Тогда требование к точности  $\delta$  измерения будет в  $e^{10}$  раз, т. е. в тысячи раз, выше, чем требуемая гарантированная точность  $\varepsilon$  результата. Ответ весьма чувствителен к погрешности задания входных данных, т. е.  $y_0$ , и задача плохо обусловлена.

Если измерение производить при  $t_0 = 1$ , то можно взять  $\delta = \varepsilon$ , т. е. достаточно измерения с гораздо меньшей точностью, чем в случае  $t_0 = 0$ , и задача хорошо обусловлена.

### Задачи

1. На каком из двух отрезков,  $x \in [1/2, 1]$  или  $x \in [-1, 0]$ , задача вычисления  $y = (1+x)/(1-x)$  по заданному  $x$  лучше обусловлена?

2. Пусть  $y = \sqrt{2} - 1$ . Можно написать также  $y = (\sqrt{2} + 1)^{-1}$ . Какая из двух формул чувствительнее к погрешности при приближенном задании  $\sqrt{2}$  в виде конечной десятичной дроби?

Указание. Сравнить модули производных функций  $(x-1)$  и  $(x+1)^{-1}$ .

### § 3. Погрешность

Во всякой вычислительной задаче по некоторым входным данным задача требуется найти ответ на поставленный вопрос. Если ответ на вопрос задачи можно дать с абсолютной точностью, то погрешность отсутствует. Но обычно удается найти ответ лишь с некоторой погрешностью. Погрешность вызывается тремя причинами.

Первая причина — некоторая неопределенность при задании входных данных, которая приведет к соответствующей неопределенности в ответе: ответ может быть указан лишь с некоторой погрешностью, которая носит название *неустранимой погрешности*.

Вторая причина: если мы ликвидируем неопределенность в задании входных данных, фиксируя какие-либо входные данные, а затем будем вычислять ответ с помощью какого-нибудь приближенного метода, то найдем не в точности тот ответ, который соответствует этим фиксированным входным данным. Возникает *погрешность, связанная с выбором приближенного метода вычислений*.

Третья причина: сам выбранный нами приближенный метод реализуется неточно из-за *погрешностей округлений* при вычислениях на реальном компьютере.

Погрешность результата складывается, таким образом, из неустранимой погрешности, погрешности метода и погрешности округлений.

Проиллюстрируем эти понятия.

1. *Неустранимая погрешность*. Пусть задача состоит в вычислении значения  $y$  некоторой функции  $y = f(x)$  при некотором  $x = t$ . Число  $t$  и функция  $f(x)$  служат входными данными задачи, а число  $y(t)$  — решением.

Пусть функция  $f(x)$  известна лишь приближенно, например,  $f(x) \approx \sin x$ , причем известно лишь, что  $f(x)$  отличается от  $\sin x$  не более, чем на некоторую величину  $\varepsilon > 0$ :

$$\sin x - \varepsilon \leqslant f(x) \leqslant \sin x + \varepsilon. \quad (1)$$

Пусть значение аргумента  $x = t$  получается приближенным измерением, в результате которого получаем некоторое  $x = t^*$ , причем известно лишь, что  $t$  лежит в пределах

$$t^* - \delta \leqslant t \leqslant t^* + \delta, \quad (2)$$

где  $\delta > 0$  — число, характеризующее точность измерения.

Из рис. 1 видно, что величиной  $y = f(t)$  может оказаться любая точка отрезка  $[a, b]$ , где  $a = \sin(t^* - \delta) - \varepsilon$ ,  $b = \sin(t^* + \delta) + \varepsilon$ . Понятно,

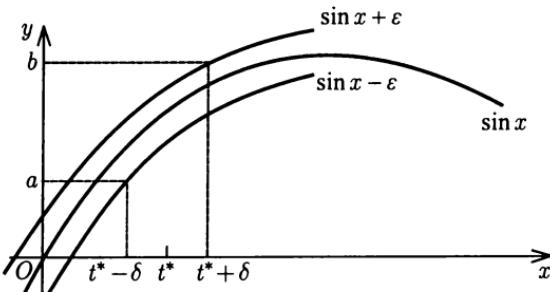


Рис. 1

что, приняв за приближенное значение числа  $y = f(t)$  любую точку  $y^*$  отрезка  $[a, b]$ , можно гарантировать оценку погрешности

$$|y - y^*| \leq |b - a|. \quad (3)$$

Эту оценку погрешности нельзя существенно уменьшить при имеющихся неполных входных данных. Самая малая погрешность, которую можно гарантировать, получается, если принять за  $y^*$  середину отрезка  $[a, b]$ , положив

$$y^* = y_{\text{opt}}^* = \frac{a+b}{2}.$$

Из рис. 1 видно, что гарантирована оценка

$$|y - y_{\text{opt}}^*| \leq \frac{|b-a|}{2}. \quad (4)$$

Это неравенство станет точным равенством, если  $y(t) = a$  или  $y(t) = b$ .

Таким образом,  $|b - a|/2$  и есть та неустранимая (неуменьшаемая) погрешность, которую можно гарантировать при имеющихся неопределенных входных данных в случае самого удачного выбора приближенного решения  $y^*$ .

Оптимальная оценка (4) ненамного лучше оценки (3). Поэтому мы не отступим от здравого смысла, если о любой точке  $y^* \in [a, b]$ , а не только о точке  $y_{\text{опт}}^*$ , условимся говорить, что она есть приближенное решение задачи вычисления числа  $y(t)$ , найденное с неустранимой погрешностью, а вместо  $|b - a|/2$  из (4) за величину неустранимой погрешности примем (условно) число  $|b - a|$ .

**2. Погрешность метода.** Положим  $y^* = \sin t^*$ . Число  $y^*$  принадлежит отрезку  $[a, b]$  и является неулучшаемым приближенным решением задачи, погрешность которого удовлетворяет оценке (3) и неустранима. Точка  $y^* = \sin t^*$  выбрана среди других точек отрезка  $[a, b]$ , потому что она задается удобной для дальнейшей работы формулой.

Для вычисления числа  $y^* = \sin t^*$  на компьютере воспользуемся разложением функции  $\sin x$  в ряд:

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots \quad (5)$$

Для вычисления  $y^*$  можно воспользоваться одним из приближенных выражений

Выбирая для приближенного вычисления  $y^*$  одну из формул (6), мы тем самым выбираем метод вычисления.

Величина  $|y^* - y_n^*|$  есть погрешность метода вычислений.

Фактически выбранный нами метод вычисления зависит от параметра  $n$  и позволяет добиваться, чтобы погрешность метода была меньше любой наперед заданной величины за счет выбора этого параметра.

Очевидно, нет смысла добиваться, чтобы погрешность метода была существенно (во много раз) меньше неустранимой погрешности. Поэтому нет смысла брать число  $n$  слишком большим. Однако в случае, если  $n$  выбрано слишком маленьким так, что погрешность метода существенно больше неустранимой погрешности, то выбранный метод не полностью использует информацию о решении, содержащуюся во входных данных, теряя часть этой информации.

3. *Погрешность округлений.* Допустим, мы зафиксировали метод вычислений, положив  $y^* \approx y_n^*$ . При вычислении  $y_n^*$  по формуле (6) на реальном компьютере в результате округлений мы получим некоторое число  $\tilde{y}_n^*$ . Погрешность  $|y_n^* - \tilde{y}_n^*|$  будем называть *погрешностью округлений*.

Эта погрешность не должна быть существенно больше погрешности метода вычислений. В противном случае произойдет потеря точности метода за счет погрешностей округлений.

### Задачи \*)

1. Пусть требуется вычислить значение  $y = f(x)$  функции  $f(x)$  по неполным входным данным  $x^*$  ( $x^* - \delta \leq x \leq x^* + \delta$ ).

Какова неустранимая погрешность, вызванная неполным значением входных данных, в зависимости от  $x^*$  и  $\delta > 0$ :

- $f(x) = \sin x$ ;
- $f(x) = \ln x$ ,  $x > 0$ ?

При каких значениях  $x^*$ , полученных приближенным измерением неопределенной величины  $x$  с погрешностью  $\delta$ , в задаче б) можно указать лишь одностороннюю оценку для  $\ln x$  сверху? Укажите эту оценку.

2\*. Пусть функция  $f(t)$  задана таблицей своих значений в точках  $t_k = kh$  ( $h = 1/N$ ,  $k = 0, \pm 1, \pm 2, \dots$ ). Пусть, кроме этой таблицы, о функции  $f(t)$  известно еще, что  $\max_x |f''(x)| \leq 1$ .

Показать, что неполные входные данные, содержащиеся в таблице, вообще говоря, не позволяют восстановить в произвольной наперед заданной точке  $t$  функцию точнее, чем с неустранимой погрешностью  $\epsilon(h) = h^2/\pi^2$ .

Указание. Показать, что наряду с функцией  $f(t) \equiv 0$ , имеющей нулевую таблицу, функция  $\varphi(t) = (h^2/\pi^2) \sin(N\pi t)$  также имеет нулевую таблицу и удовлетворяет условию  $\max_x |\varphi''(x)| \leq 1$ , при этом  $\max_t |f(t) - \varphi(t)| = h^2/\pi^2$ .

3. Задана некоторая функция  $y = f(x)$ , о которой известно, что ее вторая производная не превосходит по модулю единицы.

Показать, что погрешность приближенного равенства

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

не превосходит величины  $h$ .

\*) Символом \* будут отмечены задачи повышенной сложности.

**4\***. Пусть некоторая функция  $y = f(x)$  имеет вторую производную  $f''(x)$ , которая по модулю не превосходит единицы. При каждом  $x$  значение  $f(x)$  получается приближенным измерением величины  $f(x)$  и оказывается равным некоторому числу  $f^*(x)$ . Пусть известно лишь, что точность измерения гарантирует справедливость оценки

$$|f(x) - f^*(x)| \leq \varepsilon,$$

где  $\varepsilon > 0$  — число, характеризующее точность измерений. Пусть требуется приближенно вычислить  $f'(x)$ .

а) Как выбрать  $h$ , чтобы гарантированная погрешность приближенной формулы

$$f'(x) \approx \frac{f^*(x+h) - f^*(x)}{h}$$

была наименьшей?

б) Показать, что неустранимая погрешность задачи вычисления  $f'(x)$  при заданных (не вполне определенных) входных данных не меньше  $O(\sqrt{\varepsilon})$ .

Указание к п. б). Функции

$$f(x) \equiv 0, \quad f^*(x) = \varepsilon \sin \frac{x}{\sqrt{\varepsilon}}$$

имеют вторые производные, не превосходящие единицы по модулю, и  $\max_x |f(x) - f^*(x)| \leq \varepsilon$ . В то же время

$$\left| \frac{df^*}{dx} - \frac{df}{dx} \right| = \left| \sqrt{\varepsilon} \cos \frac{x}{\sqrt{\varepsilon}} \right| = O(\sqrt{\varepsilon}).$$

Сопоставляя результаты пп. а), б), проверить, что метод вычисления производной из п. а) дает неулучшаемый по порядку погрешности  $O(\sqrt{\varepsilon})$  результат, а также показать, что порядок неустранимой погрешности в точности совпадает с  $O(\sqrt{\varepsilon})$ .

5. Для запоминания сведений о линейной функции  $f(x) = kx + b$ ,  $\alpha \leq x \leq \beta$ , удовлетворяющей неравенствам  $0 \leq f(x) \leq 1$ , имеется шесть занумерованных клеток, в каждую из которых можно записать одну из десяти цифр: 0, 1, ..., 9.

Какова неустранимая погрешность при восстановлении функции, если эти шесть клеток заполнены одним из указанных здесь способов?

а) В первые три клетки записаны первые три цифры, стоящие после запятой в записи  $f(\alpha)$  в виде десятичной дроби, а в три оставшиеся клетки — первые три цифры, стоящие после запятой в записи  $f(\beta)$  в виде десятичной дроби.

б) Пусть  $\alpha = 0$ ,  $\beta = 10^{-2}$ . В первые три клетки записаны первые три цифры десятичной записи числа  $k$ , в четвертую клетку — 0 или 1 в зависимости от знака числа  $k$ , а в оставшиеся две клетки записаны первые две цифры после запятой из десятичной записи числа  $b$ .

в)\* Показать, что при любом способе задания функции из заданного класса с помощью шестиместной таблицы указанного вида неустранимая погрешность не меньше  $0,5 \cdot 10^{-3}$ .

Указание. Построить  $10^6$  функций из заданного класса, для любых двух из которых максимум модуля разности не меньше  $10^{-3}$ .

## § 4. О методах вычисления

Пусть для изучения некоторого объекта построена его математическая модель, которая и подлежит изучению средствами вычислительной математики.

Например, математической моделью малых колебаний маятника может оказаться следующая задача:

$$\begin{aligned} \frac{d^2y}{dt^2} + y = 0, \quad t \geq 0, \\ y(0) = 0, \quad \left. \frac{dy}{dt} \right|_{t=0} = 1, \end{aligned} \tag{1}$$

где  $y(t)$  — отклонение маятника в момент времени  $t$ .

При изучении гармонических колебаний с помощью этой математической модели, т. е. задачи Коши (1), некоторую пользу может принести знание физической сущности моделируемого физического объекта. В данной задаче, например, можно предвидеть, что решение носит колебательный (периодический) характер. Однако математическая модель (1) после ее построения становится самостоятельным объектом, для изучения которого можно применять любые математические средства, в том числе и те, которые не имеют никакого отношения к физическому происхождению задачи. Это сильно расширяет возможности исследователя.

Так, например, значение решения  $y = \sin t$  задачи (1) в заданный момент времени  $t = z$ , т. е. представление числа  $\sin z$  в виде десятичной дроби с заданным числом знаков, можно осуществить с помощью ряда

$$\sin z = z - \frac{z^3}{3!} + \frac{z^5}{5!} - \dots,$$

воспользовавшись его частичной суммой, хотя представление функции  $\sin t$  в виде степенного ряда, использованное нами, не имеет никакого физического смысла.

Для численного решения какой-либо задачи на компьютере возможны различные численные методы, или различные алгоритмы. Одни из них могут быть много лучше других.

В этой книге мы изложим хорошие алгоритмы для многих часто встречающихся классов задач, а пока объясним, чем могут различаться алгоритмы.

Пусть для вычисления решения  $y$  некоторой задачи по входным данным, совокупность которых обозначим  $X$ , имеются алгоритмы  $A_1$ ,  $A_2$ , которые дают приближения  $y_1^* = A_1(X)$ ,  $y_2^* = A_2(X)$ . При этом может оказаться следующее.

1. Алгоритм  $A_2$  может быть точнее алгоритма  $A_1$ , т. е.

$$|y - y_1^*| \gg |y - y_2^*|.$$

Например, будем вычислять  $y = \sin x|_{x=0,1}$  по формуле вида

$$y_n^* = \sum_{k=1}^n (-1)^{k-1} \frac{x^{2k-1}}{(2k-1)!}. \quad (2)$$

Примем за  $A_1$  алгоритм, соответствующий  $n = 1$ , а за  $A_2$  алгоритм, соответствующий  $n = 2$ . Очевидно, что

$$|\sin 0,1 - y_1^*| \gg |\sin 0,1 - y_2^*|.$$

2. Алгоритмы обладают одинаковой точностью, но вычисление  $y_1^* = A_1(X)$  требует много больше арифметических действий, чем вычисление  $y_2^*$ .

Пусть, например, требуется найти

$$y = 1 + x + x^2 + \dots + x^{1023} \quad \left( y = \frac{1 - x^{1024}}{1 - x} \right)$$

при  $x = 0,99$ . За  $A_1$  примем алгоритм, осуществляющий вычисления прямо по заданной формуле, возводя 0,99 поочередно в степени 1, 2, ..., 1023 и складывая. В качестве  $A_2$  примем алгоритм, осуществляющий вычисление по формуле

$$y = \frac{1 - 0,99^{1024}}{1 - 0,99}.$$

По точности алгоритмы совпадают (оба абсолютно точны), но первый требует гораздо больше арифметических операций.

А именно, вычисляя последовательно

$$x, \quad x^2 = x \cdot x, \quad \dots, \quad x^{1023} = x^{1022} \cdot x,$$

придется проделать 1022 умножения. В то же время при вычислении  $0,99^{1024}$  требуется всего 10 умножений:

$$\begin{aligned} 0,99^2 &= 0,99 \cdot 0,99, \quad 0,99^4 = (0,99^2) \cdot (0,99^2), \dots \\ &\dots, \quad 0,99^{1024} = (0,99^{512}) \cdot (0,99^{512}). \end{aligned}$$

3. Алгоритмы обладают одинаковой точностью, но  $A_1(X)$  вычислиительно устойчив, а  $A_2(X)$  неустойчив.

Например, для вычисления  $y = \sin x$  с заданной точностью  $\varepsilon = 10^{-3}$ ,  $|y - y^*| \leq 10^{-3}$  воспользуемся суммой

$$y_1^* = y_1^*(x) = \sum_{k=1}^n (-1)^{k-1} \frac{x^{2k-1}}{(2k-1)!}, \quad (3)$$

где  $n = n(\varepsilon)$  выбирается так, чтобы выполнялось неравенство

$$|y - y_1^*| \leq 10^{-3},$$

приняв вычисление этой суммы за алгоритм  $A_1$ . Если  $|x| \leq \pi/2$ , то

$$\frac{1}{(2n-1)!} \left(\frac{\pi}{2}\right)^{2n-1} \leq 10^{-3}$$

при  $n = 5$ , и

$$y_1^* = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!}.$$

Очевидно, что вычисления по этой формуле слабо реагируют на погрешности округления при вычислении каждого из слагаемых, а поэтому алгоритм  $A_1$  в данном случае вычислительно устойчив.

Пусть  $|x| \gg 1$ ; например,  $x = 100$ . Тогда для достижения требуемой точности число  $n$  должно удовлетворять неравенству

$$\frac{100^{2n+1}}{(2n+1)!} \leq 10^{-3},$$

из которого следует, что для  $n$  заведомо выполнено неравенство  $n > 48$ . Но при вычислении суммы (3) первые ее члены будут очень большими. Малая относительная погрешность при их вычислении дает большую абсолютную погрешность, и алгоритм  $A_1$  вычислительно неустойчив.

Укажем устойчивый алгоритм  $A_2$ . Записываем заданное  $x$  в виде  $x = l\pi + z$  ( $|z| \leq \pi/2$ ,  $l$  целое). Тогда

$$\begin{aligned} \sin x &= (-1)^l \sin z, \\ y^* &= A_2(x) = (-1)^l \left( z - \frac{z^3}{3!} + \frac{z^5}{5!} - \frac{z^7}{7!} \right). \end{aligned}$$

Этот алгоритм по устойчивости совпадает с алгоритмом  $A_1$  для  $|x| \leq \pi/2$ .

4. Алгоритм может быть сходящимся или расходящимся.

Пусть требуется вычислить  $y = \ln(1+x)$ . Воспользуемся рядом

$$y = \ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots \quad (4)$$

и положим

$$y^*(x) \approx y_n^* = \sum_{k=1}^n (-1)^{k+1} \frac{x^k}{k}.$$

Получим метод вычисления  $y = \ln(1+x)$ , зависящий от номера  $n$  как от параметра.

Если  $|x| = q < 1$ , то  $\lim_{n \rightarrow \infty} y_n^*(x) = y(x)$ , т. е. погрешность алгоритма с ростом  $n$  стремится к нулю. Но если  $x > 1$ , то  $\lim_{n \rightarrow \infty} y_n^*(x) = \infty$ , так как ряд (4) имеет радиус сходимости  $r = 1$ . В этом случае алгоритм не пригоден для вычислений.

В книге мы познакомимся и с некоторыми другими характеристиками алгоритмов. Встречаются алгоритмы, допускающие параллельные

вычисления или требующие последовательного выполнения операций, самонастраивающиеся на специфику входных данных или учитывающие ее лишь отчасти, алгоритмы логически простые или более сложные.

### Задачи

1. Предложить алгоритм вычисления  $y = \ln(1 + x)$ , пригодный при  $x > 1$ .
2. Рассмотрим задачу об определении последовательности  $x_0, x_1, \dots, x_N$ , удовлетворяющей уравнению

$$2x_n - x_{n+1} = 1 + (n/N)^2, \quad n = 0, 1, \dots, N - 1,$$

и дополнительному условию

$$x_0 + x_N = 1. \quad (5)$$

Предлагаются следующие два алгоритма. Полагаем

$$x_n = u_n + cv_n, \quad n = 0, 1, \dots, N.$$

В алгоритме  $A_1$  определяем  $u_n$  ( $n = 0, 1, \dots, N$ ) как решение уравнения

$$2u_n - u_{n+1} = 1 + n^2/N^2, \quad n = 0, 1, \dots, N - 1, \quad (6)$$

при условии

$$u_0 = 0. \quad (7)$$

Последовательность  $v_n$  определяем равенствами

$$2v_n - v_{n+1} = 0, \quad n = 0, 1, \dots, N - 1, \quad (8)$$

$$v_0 = 1. \quad (9)$$

Число  $c$  определяем из условия (5). При этом значения  $u_n$ ,  $v_n$  последовательно вычисляем по формулам

$$u_{n+1} = 2u_n - (1 + n^2/N^2), \quad n = 0, 1, \dots,$$

$$v_{n+1} = 2^n, \quad n = 0, 1, \dots.$$

Алгоритм  $A_2$  состоит в том, что  $u_n$  определяем как решение системы (6), но вместо условия (7) используем условие  $u_N = 0$ . Последовательность  $v_n$  определяем как решение уравнений (8), но вместо условия (9) используем  $v_N = 1$ .

- а) Проверить, что второй алгоритм устойчив, а первый очень неустойчив.
- б) Попытаться найти решение на компьютере, используя поочередно оба метода при  $N = 10$ , а затем при  $N = 100$ .

# ЧАСТЬ I

## ТАБЛИЧНОЕ ЗАДАНИЕ И ИНТЕРПОЛЯЦИЯ ФУНКЦИЙ. КВАДРАТУРЫ

---

Одним из основных в математике является понятие функции. Функция  $y = f(x)$  ( $a \leq x \leq b$ ) может быть задана некоторой формулой, например,  $y = x^2$ , которую можно хранить и использовать в компьютере в виде программы, по которой для каждого фиксированного значения  $x$  вычисляется значение  $y = x^2$ .

Однако, как правило, функция  $y = f(x)$  задается приближенно тем или иным конечным набором чисел — некоторой таблицей, обрабатывая которую, можно получить приближенное значение функции при каждом фиксированном  $x$ . Этой таблицей может служить, например, конечный набор первых коэффициентов разложения функции в степенной ряд.

Например, для функции

$$e^x, \quad 0 \leq x \leq 1, \quad e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \dots + \frac{x^n}{n!} + \dots,$$

таблицей может служить конечный набор чисел

$$1, \quad \frac{1}{1!}, \quad \dots, \quad \frac{1}{n!};$$

$n$  задано.

Чем больше натуральное  $n$ , тем точнее можно восстановить функцию по таблице значений  $n$  первых коэффициентов ее разложения в степенной ряд, пользуясь расшифровывающей эту таблицу формулой

$$e^x \approx 1 + \frac{x}{1!} + \frac{x^2}{2!} + \dots + \frac{x^n}{n!}.$$

Однако обычно таблица значений функции  $y = f(x_n)$  получается в результате измерения или вычисления ее значений на некотором наборе точек  $x_0, x_1, \dots, x_n \in [a, b]$ . Тогда возникает задача восстановления (интерполяции) функции в точках  $x$ , не совпадающих с  $x_0, x_1, \dots, x_n$ .

Наиболее употребительны и удобны алгебраическая и тригонометрическая интерполяции. Мы рассмотрим оба этих способа интерполяции. Рассмотрим здесь также задачу вычисления определенных интегралов по таблице функций, поскольку основные способы получения формул (квадратур) для приближенного вычисления определенных интегралов тесно связаны с интерполяционными формулами.

# ГЛАВА 1

## АЛГЕБРАИЧЕСКАЯ ИНТЕРПОЛЯЦИЯ

Пусть заданы точки  $x_0, x_1, \dots, x_n$  и соответственные значения  $f(x_0), f(x_1), \dots, f(x_n)$  функции  $f(x)$  в этих точках. Соответствие

$x_0$	$x_1$	$\dots$	$x_n$
$f(x_0)$	$f(x_1)$	$\dots$	$f(x_n)$

будем называть *таблицей значений функции  $f(x)$  в узлах  $x_0, x_1, \dots, x_n$* .

Это название несколько условно, так как значение  $f(x_j)$  может записываться бесконечной десятичной дробью (например,  $\sqrt{3}$ ), а для работы на компьютере все числа должны быть округлены до десятичных (или двоичных) дробей с конечным числом знаков.

*Алгебраическим интерполяционным многочленом* назовем многочлен

$$P_n(x) = c_0 + c_1x + \dots + c_nx^n$$

степени не выше  $n$ , который в узлах  $x_0, x_1, \dots, x_n$  принимает значения  $f(x_0), f(x_1), \dots, f(x_n)$  соответственно.

### § 1. Существование и единственность интерполяционного многочлена

#### 1. Интерполяционный многочлен в форме Лагранжа.

Теорема 1. Пусть заданы узлы  $x_0, x_1, \dots, x_n$ , среди которых нет совпадающих, и значения  $f(x_0), f(x_1), \dots, f(x_n)$  функции в этих узлах. Тогда существует один и только один многочлен  $P_n(x) = P_n(x, f, x_0, x_1, \dots, x_n)$  степени не выше  $n$ , принимающий в заданных узлах  $x_k$  заданные значения  $f(x_k)$ .

Доказательство. Сначала покажем, что существует не более, чем один интерполяционный многочлен  $P_n(x)$ , а затем построим его.

Если бы таких многочленов было два,  $P_n^I(x)$  и  $P_n^{II}(x)$ , то их разностью  $R_n(x) = P_n^I(x) - P_n^{II}(x)$  был бы многочлен степени не выше  $n$ , обращающийся в нуль в  $(n+1)$  точках  $x_0, x_1, \dots, x_n$ . Но каждый многочлен, отличный от тождественного нуля, имеет ровно столько корней, считая их кратности, какова его степень. Поэтому  $R_n(x) \equiv 0$ , т. е.  $P_n^I(x) \equiv P_n^{II}(x)$ . Единственность доказана.

Введем теперь вспомогательные многочлены

$$l_k(x) = \frac{(x - x_0)(x - x_1)\dots(x - x_{k-1})(x - x_{k+1})\dots(x - x_n)}{(x_k - x_0)(x_k - x_1)\dots(x_k - x_{k-1})(x_k - x_{k+1})\dots(x_k - x_n)}.$$

Очевидно, что  $l_k(x)$  есть многочлен степени  $n$  и что выполняются равенства

$$l_k(x)|_{x_j} = \begin{cases} 1, & x_j = x_k, \\ 0, & x_j \neq x_k, \end{cases} \quad j = 0, 1, \dots, n.$$

Многочлен  $P_n(x)$ , заданный равенством

$$\begin{aligned} P_n(x) &= P_n(x, f, x_0, x_1, \dots, x_n) = \\ &= f(x_0)l_0(x) + f(x_1)l_1(x) + \dots + f(x_n)l_n(x), \end{aligned} \quad (1)$$

и есть искомый интерполяционный многочлен.

Действительно, он имеет степень не выше  $n$ , так как каждое слагаемое  $f(x_j)l_j(x)$  есть многочлен степени не выше  $n$ . Кроме того, для него, очевидно, выполнены равенства  $P_n(x_j) = f(x_j)$  ( $j = 0, 1, \dots, n$ ).  $\square$

Мы не только доказали теорему, но и выписали интерполяционный многочлен  $P_n(x)$  в виде формулы (1), которая называется *записью интерполяционного многочлена в форме Лагранжа*.

Употребительны и другие записи (единственного) интерполяционного многочлена  $P_n(x, f, x_0, x_1, \dots, x_n)$ . Особенно часто используют запись в форме Ньютона.

**2. Интерполяционный многочлен в форме Ньютона. Разностные отношения.** Пусть функция  $f(x)$  в точках  $x_a, x_b, x_c, x_d$  и т. д. принимает некоторые значения  $f(x_a), f(x_b), f(x_c), f(x_d)$  и т. д. Разностное отношение нулевого порядка  $f(x_k)$  функции  $f(x)$  в точке  $x_k$  определим как значение функции в этой точке:

$$f(x_k) = f(x_k), \quad k = a, b, c, d, \dots. \quad (2)$$

Разностное отношение первого порядка  $f(x_k, x_t)$  функции  $f(x)$  для (произвольной) пары точек  $x_k, x_t$  определим через разностные отношения нулевого порядка:

$$f(x_k, x_t) = \frac{f(x_t) - f(x_k)}{x_t - x_k}.$$

Вообще разностное отношение  $f(x_0, x_1, \dots, x_n)$   $n$ -го порядка определим через разностные отношения порядка  $n - 1$ , положив

$$f(x_0, x_1, \dots, x_n) = \frac{f(x_1, x_2, \dots, x_n) - f(x_0, x_1, \dots, x_{n-1})}{x_n - x_0}. \quad (3)$$

Интерполяционный многочлен  $P_n(x, f, x_0, x_1, \dots, x_n)$  можно записать в следующей форме Ньютона:

$$\begin{aligned} P_n(x, f, x_0, x_1, \dots, x_n) &= f(x_0) + (x - x_0)f(x_0, x_1) + \dots \\ &\dots + (x - x_0)(x - x_1)\dots(x - x_{n-1})f(x_0, x_1, \dots, x_n). \end{aligned} \quad (4)$$

Несколько ниже мы докажем справедливость формулы (4), а пока установим некоторые следствия из нее.

**Следствие 1. Справедливо равенство**

$$\begin{aligned} P_n(x, f, x_0, x_1, \dots, x_{n-1}, x_n) &= P_{n-1}(x, f, x_0, x_1, \dots, x_{n-1}) + \\ &+ f(x_0, x_1, \dots, x_n)(x - x_0)(x - x_1)\dots(x - x_{n-1}). \end{aligned} \quad (5)$$

Доказательство очевидно.

**Следствие 2.** Разностное отношение  $f(x_0, x_1, \dots, x_n)$  порядка  $n$  равно коэффициенту  $c_n$  при члене  $x^n$ , входящем в интерполяционный многочлен

$$P_n(x, f, x_0, x_1, \dots, x_n) = c_n x^n + c_{n-1} x^{n-1} + \dots + c_0,$$

т. е. справедливо равенство

$$f(x_0, x_1, \dots, x_n) = c_n. \quad (6)$$

**Доказательство.** Очевидно, что в правую часть формулы (4) член  $x^n$  входит с коэффициентом  $f(x_0, x_1, \dots, x_n)$ .  $\square$

**Следствие 3.** Разностное отношение  $f(x_0, x_1, \dots, x_n)$  обращается в нуль в том и только том случае, если  $f(x_0), f(x_1), \dots, f(x_n)$  суть значения некоторого многочлена  $Q_m$ , степень которого  $m$  строго меньше  $n$ .

**Доказательство.** Если  $f(x_0, x_1, \dots, x_n) = 0$ , то из формулы (4) видно, что интерполяционный многочлен  $P_n(x, f, x_0, x_1, \dots, x_n)$ , принимающий при  $x_j$  значения  $f(x_j)$  ( $j = 0, 1, \dots, n$ ), есть многочлен степени меньше  $n$ , поскольку в силу равенства (6) коэффициент  $c_n$  при  $x^n$  равен нулю. Обратно: ввиду единственности интерполяционного многочлена степени не выше  $n$  многочлен  $Q_m(x)$  совпадает с интерполяционным многочленом  $P_n(x, f, x_0, \dots, x_n) = c_n x^n + c_{n-1} x^{n-1} + \dots + c_0$ . В силу того, что  $m < n$ , из равенства  $Q_m(x) = P_n(x, f, x_0, x_n)$  следует, что  $c_n = 0$ . В силу равенства (6) тогда  $f(x_0, x_1, \dots, x_n) = 0$ .  $\square$

**Следствие 4.** Разностное отношение  $f(x_0, x_1, \dots, x_n)$  не изменяется при произвольной перестановке его аргументов  $x_0, x_1, \dots, x_n$ .

**Доказательство.** Переставим узлы  $x_0, x_1, \dots, x_n$  так, что на месте с номером  $j$  окажется один из узлов  $x_0, x_1, \dots, x_n$ , который мы обозначим  $x'_j$  ( $j = 0, 1, \dots, n$ ). Очевидно, что интерполяционный многочлен от нумерации узлов не зависит;  $P_n(x, f, x_0, x_1, \dots, x_n) \equiv \underset{x}{\equiv} P_n(x, f, x'_0, x'_1, \dots, x'_n)$ . Поэтому наряду с записью (4) справедлива запись

$$\begin{aligned} P_n(x, f, x_0, x_1, \dots, x_n) &= f(x'_0) + f(x'_0, x'_1)(x - x'_0) + \dots \\ &\quad + f(x'_0, x'_1, \dots, x'_n)(x - x'_0)(x - x'_1)\dots(x - x'_{n-1}), \end{aligned} \quad (7)$$

так что в силу равенства (6) имеет место

$$f(x'_0, x'_1, \dots, x'_n) = c_n. \quad (8)$$

Сравнивая равенства (6) и (8), убеждаемся в справедливости доказываемого утверждения  $f(x_0, x_1, \dots, x_n) = f(x'_0, x'_1, \dots, x'_n)$ .  $\square$

**Следствие 5.** Справедливо равенство

$$f(x_0, x_1, \dots, x_n) = \frac{f(x_n) - P_{n-1}(x_n, f, x_0, x_1, \dots, x_{n-1})}{(x_n - x_0)(x_n - x_1)\dots(x_n - x_{n-1})}. \quad (9)$$

**Доказательство.** В равенстве (5) положим  $x = x_n$ . Тогда левая часть примет значение  $f(x_n)$ , и формула (9) станет очевидной.  $\square$

**Теорема 2.** Интерполяционный многочлен  $P_n(x, f, x_0, x_1, \dots, x_n)$  допускает запись в форме Ньютона, т. е. имеет место формула (4).

**Доказательство.** Воспользуемся индукцией по  $n$ . При  $n = 0$  формула (4) справедлива. Допустим, что ее справедливость уже установлена для  $n = 1, 2, \dots, k$ . Покажем, что она имеет место и для  $n = k + 1$ , т. е. докажем равенство

$$P_{k+1}(x, f, x_0, x_1, \dots, x_k, x_{k+1}) = P_k(x, f, x_0, x_1, \dots, x_k) + \\ + f(x_0, x_1, \dots, x_k, x_{k+1})(x - x_0)(x - x_1)\dots(x - x_k). \quad (10)$$

Заметим, что в силу предположения индукции о справедливости равенства (4) для  $n \leq k$  доказательства следствий 1–5, которые мы провели, опираясь на справедливость равенства (4), сохраняют силу при  $n \leq k$ .

Переходя к доказательству равенства (10), сначала покажем, что многочлен  $P_{k+1}(x, f, x_0, x_1, \dots, x_k, x_{k+1})$  можно записать в форме

$$P_{k+1}(x, f, x_0, x_1, \dots, x_k, x_{k+1}) = P_k(x, f, x_0, x_1, \dots, x_k) + \\ + \frac{f(x_{k+1}) - P_k(x_{k+1}, f, x_0, x_1, \dots, x_k)}{(x_{k+1} - x_0)\dots(x_{k+1} - x_k)}(x - x_0)\dots(x - x_k). \quad (11)$$

Очевидно, что в правой части равенства (11) стоит многочлен степени не выше  $k + 1$ , принимающий в точках  $x_j$  значения  $f(x_j)$  ( $j = 0, 1, \dots, k + 1$ ). Поэтому выражение в правой части равенства (11) есть интерполяционный многочлен

$$P_{k+1}(x, f, x_0, x_1, \dots, x_{k+1}).$$

Сравнивая равенства (10) и (11), видим, что для доказательства равенства (10) надо установить равенство

$$f(x_0, x_1, \dots, x_{k+1}) = \frac{f(x_{k+1}) - P_k(x_{k+1}, f, x_0, x_1, \dots, x_k)}{(x_{k+1} - x_0)\dots(x_{k+1} - x_{k-1})(x_{k+1} - x_k)}. \quad (12)$$

В силу следствия 4

$$P_k(x, f, x_0, x_1, \dots, x_k) = P_k(x, f, x_1, x_2, \dots, x_k, x_0) = \\ = P_{k-1}(x, f, x_1, x_2, \dots, x_k) + \\ + f(x_1, x_2, \dots, x_k, x_0)(x - x_1)(x - x_2)\dots(x - x_k). \quad (13)$$

Воспользуемся формулой (13) при  $x = x_{k+1}$  и придадим правой части равенства (12) следующий вид:

$$\frac{f(x_{k+1}) - P_k(x_{k+1}, f, x_0, x_1, \dots, x_k)}{(x_{k+1} - x_0)\dots(x_{k+1} - x_{k-1})(x_{k+1} - x_k)} = \\ = \frac{1}{x_{k+1} - x_0} \cdot \frac{f(x_{k+1}) - P_{k-1}(x_{k+1}, f, x_1, \dots, x_k)}{(x_{k+1} - x_1)\dots(x_{k+1} - x_k)} - \\ - \frac{f(x_1, x_2, \dots, x_k, x_0)}{x_{k+1} - x_0}. \quad (14)$$

В силу следствия 5 уменьшаемое в правой части равенства (14) совпадает с выражением

$$\frac{1}{x_{k+1} - x_0} f(x_1, x_2, \dots, x_{k+1}).$$

В силу следствия 4 в вычитаемом можно переставить аргументы так, что оно совпадет с  $\frac{f(x_0, x_1, \dots, x_k)}{x_{k+1} - x_0}$ .

Таким образом, правая часть равенства (14) есть

$$\frac{f(x_1, x_2, \dots, x_{k+1}) - f(x_0, x_1, \dots, x_k)}{x_{k+1} - x_0} = f(x_0, x_1, \dots, x_{k+1}),$$

так что равенство (14) совпадает с доказываемым равенством (12).  $\square$

**Теорема 3.** Пусть  $x_0 < x_1 < \dots < x_n$ , функция  $f(x)$  определена на отрезке  $x_0 \leq x \leq x_n$  и имеет на этом отрезке производную порядка  $n$ . Тогда

$$n! f(x_0, x_1, \dots, x_n) = f^{(n)}(\xi), \quad (15)$$

где  $\xi$  — некоторая точка отрезка  $[x_0, x_n]$ .

**Доказательство.** Функция

$$\varphi(x) \equiv f(x) - P_n(x, f, x_0, \dots, x_n) \quad (16)$$

обращается в нуль в  $(n + 1)$  точках  $x_0, x_1, \dots, x_n$ . По теореме Ролля ее производная обращается в нуль хотя бы в одной точке между каждыми двумя соседними нулями функции  $\varphi(x)$ . Таким образом, функция  $\varphi'(x)$  обращается в нуль не менее, чем в  $n$  точках. Аналогично  $\varphi''(x)$  обращается в нуль по крайней мере в одной точке между каждыми двумя нулями функции  $\varphi'(x)$  и имеет поэтому не менее, чем  $(n - 1)$  нулей.

Рассуждая аналогично, убедимся, что  $\varphi^{(n)}(x)$  имеет хотя бы один нуль. Обозначим его  $\xi$ , так что  $\varphi^{(n)}(\xi) = 0$ . Продифференцируем тождество (16) ровно  $n$  раз и положим после этого  $x = \xi$ :

$$0 = \varphi^{(n)}(\xi) = f^{(n)}(\xi) - \frac{d^n}{dx^n} P_n(x, f, x_0, x_1, \dots, x_n) \Big|_{x=\xi}. \quad (17)$$

Но

$$\begin{aligned} \frac{d^n}{dx^n} P_n(x, f, x_0, x_1, \dots, x_n) &= \frac{d^n}{dx^n} [P_{n-1}(x, f, x_0, x_1, \dots, x_{n-1}) + \\ &+ f(x_0, x_1, \dots, x_n)(x - x_0)(x - x_1)\dots(x - x_{n-1})] \equiv \\ &\equiv 0 + n! f(x_0, x_1, \dots, x_n). \end{aligned}$$

Поэтому из выражения (17) следует равенство (15).  $\square$

**Теорема 4.** Значения  $f(x_0), f(x_1), \dots, f(x_n)$  выражаются через разностные отношения  $f(x_0), f(x_0, x_1), \dots, f(x_0, x_1, \dots, x_n)$  формулами

$$f(x_i) = f(x_0) + (x_i - x_0)f(x_0, x_1) + (x_i - x_0)(x_i - x_1)f(x_0, x_1, x_2) + \dots + (x_i - x_0)(x_i - x_1)\dots(x_i - x_{n-1})f(x_0, x_1, \dots, x_n), \quad i = 0, 1, \dots, n,$$

т. е. формулами вида

$$f(x_i) = a_{i0}f(x_0) + a_{i1}f(x_0, x_1) + \dots + a_{in}f(x_0, \dots, x_n), \quad i = 0, 1, \dots, n. \quad (18)$$

Для доказательства можно воспользоваться равенствами  $f(x_i) = P_n(x, f, x_0, x_1, \dots, x_n)|_{x=x_i}$  и записью интерполяционного многочлена в форме (4).

**3. Сравнение записей в форме Лагранжа и Ньютона.** Для вычисления значения функции  $f(x)$  в точке  $x$ , не являющейся узлом интерполяции, можно положить  $f(x) \approx P_n(x, f, x_0, x_1, \dots, x_n)$ .

Пусть  $P_n(x, f, x_0, x_1, \dots, x_n)$  уже найден, но мы решили для уточнения привлечь еще один узел  $x_{n+1}$  и значение  $f(x_{n+1})$  в нем. Тогда для вычисления  $P_{n+1}(x, f, x_0, x_1, \dots, x_n, x_{n+1})$  с помощью формулы (1) нужно заново провести всю работу. Для вычисления же по формуле Ньютона

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

нужно досчитать только поправку

$$f(x_0, x_1, \dots, x_{n+1})(x - x_0)(x - x_1)\dots(x - x_n).$$

Кстати, сразу будет видно, насколько она велика.

**4. Обусловленность задачи построения интерполяционного многочлена.** Пусть узлы интерполяции  $x_0, x_1, \dots, x_n$  лежат на некотором отрезке  $a \leq x \leq b$ . Пусть  $f(x_0), f(x_1), \dots, f(x_n)$  — заданные числа. Соответствующий интерполяционный многочлен  $P_n(x) = P_n(x, f, x_0, x_1, \dots, x_n)$  будем для краткости обозначать  $P_n(x, f)$ .

Придадим значениям  $f(x_j)$  некоторые возмущения  $\delta f(x_j)$  ( $j = 0, 1, \dots, n$ ), и интерполяционный многочлен  $P_n(x, f)$  заменится многочленом  $P_n(x, f + \delta f)$ . Из записи (1) видно, что  $P_n(x, f + \delta f) = P_n(x, f) + P_n(x, \delta f)$ , так что возмущение, которое претерпевает интерполяционный многочлен, есть  $P_n(x, \delta f)$ . Это возмущение при заданных  $x_0, x_1, \dots, x_n$  зависит только от  $\delta f$ , но не от  $f$ . Примем за меру чувствительности интерполяционного многочлена  $P_n(x, f)$  к возмущениям  $\delta f$  в узлах наименьшее число  $L_n$ , при котором для каждого  $\delta f$  выполняется неравенство

$$\max_{a < x < b} |P_n(x, \delta f)| \leq L_n \max_j |\delta f(x_j)|.$$

Числа  $L_n = L_n(x_0, x_1, \dots, x_n, a, b)$ ,  $n = 0, 1, \dots$ , называют **константами Лебега**. Эти числа растут с ростом  $n$ . Их поведение при возрастании  $n$  существенно зависит от расположения точек  $x_j$  ( $j = 0, 1, \dots, n$ ) на отрезке  $[a, b]$ .

Если, например,  $n = 1$ ,  $x_0 = a$ ,  $x_1 = b$ , то  $L_1 = 1$ . Если  $x_0 \neq a$ ,  $x_1 \neq b$ , то  $L_1 \geq \frac{b-a}{2|x_1-x_0|}$ , т. е. чувствительность интерполяции может быть сколь угодно сильной, если  $x_1$  и  $x_0$  достаточно мало различаются. Читатель легко проверит эти утверждения об  $L_1$ .

В случае равномерно расположенных узлов

$$x_j = a + jh, \quad j = 0, 1, \dots, n, \quad h = \frac{b-a}{n},$$

можно показать, что

$$2^{n-1} > L_n > 2^{n-3} \frac{1}{n-1} \cdot \frac{1}{n-3/2}, \quad (19)$$

т. е. чувствительность результата интерполяции к погрешностям при задании  $f(x_j)$  будет резко возрастать с ростом  $n$ . Погрешности при задании  $f(x_j)$  неизбежны как при получении значений путем измерений, так и в результате округлений.

Пусть теперь  $a = -1$ ,  $b = 1$ , а узлы заданы формулой

$$x_j = -\cos \frac{(2j+1)\pi}{2(n+1)}, \quad j = 0, 1, \dots, n. \quad (20)$$

Можно показать, что в случае (20)

$$L_n \leq \frac{2}{\pi} \ln n + 1, \quad (21)$$

т. е. с ростом  $n$  константы Лебега, в отличие от (19), растут очень медленно, так что в этом случае вычислительная неустойчивость не является препятствием для использования интерполяционных многочленов высокой степени.

**5. О плохой сходимости интерполяции по равноотстоящим узлам.** Не следует думать, что для каждой непрерывной функции  $f(x)$ ,  $x \in [a, b]$ , интерполяционный многочлен  $P_n(x, f)$ , построенный по значениям  $f(x_j)$  в равноотстоящих узлах  $x_j = a + jh$ ,  $x_0 = a$ ,  $x_n = b$ , с ростом  $n$  все меньше уклоняется от функции  $f(x)$ . Можно показать, например, что для функции  $f = \frac{1}{x^2 + 0,25}$ , имеющей производные всех порядков, при  $a = -1$  и  $b = 1$  уклонение  $\max_{-1 \leq x \leq 1} |f(x) - P_n(x)|$  не стремится к нулю с ростом  $n$ .

### Задачи

1. Требуется вычислить  $f(1,14)$  с помощью линейной, квадратической и кубической интерполяций, используя следующую таблицу:

$x$	1,08	1,13	1,20	1,27	1,31
$f(x)$	1,302	1,386	1,509	1,217	1,284

Осуществить расчет с помощью интерполяционных многочленов в формах Лагранжа и Ньютона.

**2.** Пусть  $x_j = jh$  ( $j = 0, \pm 1, \pm 2, \dots$ ) — узлы, расположенные с шагом  $h$ . Проверить равенство

$$f(x_{k-1}, x_k, x_{k+1}) = \frac{f(x_{k+1}) - 2f(x_k) + f(x_{k-1})}{2! h^2}.$$

**3.** Пусть  $a = x_0$ ,  $a < x_1 < b$ ,  $x_2 = b$ .

Вычислить значение константы Лебега  $L_2$ , если  $x_1 = (a + b)/2$ . Показать, что при  $x_1 \rightarrow a$  или  $x_1 \rightarrow b$  константа Лебега  $L_2 = L_2(x_0, x_1, x_2, a, b)$  неограниченно возрастает.

**4.** Проверить утверждение об интерполяции функции  $f(x) = \frac{1}{x^2 + 0,25}$  из п. 5 экспериментально, выводя графики  $f(x)$  и  $P_n(x, f)$  на экран компьютера.

## § 2. Классическая кусочно многочленная интерполяция

Высокая чувствительность интерполяционных многочленов к погрешностям при задании таблицы и возможное отсутствие сходимости последовательности  $P_n(x, f)$  с ростом  $n$  при равноотстоящих узлах заставляет использовать кусочно многочленную интерполяцию.

**1. Определение кусочно многочленной интерполяции.** Пусть функция  $f(x)$ , определенная на  $[a, b]$ , задана в виде  $f = (f(x_0), f(x_1), \dots, f(x_n))$  ( $a = x_0 < x_1 < x_2 < \dots < x_n = b$ ). Для восстановления значений функции между узлами  $x_0, x_1, \dots, x_n$  можно воспользоваться функцией, которая между каждыми двумя соседними узлами является многочленом заданной невысокой степени, например, первой, второй, третьей и т. д. Соответствующая интерполяция называется *кусочно линейной*, *кусочно квадратичной*, *кусочно кубической* и т. д.

В случае кусочно линейной интерполяции на отрезке  $x_k \leq x \leq x_{k+1}$  для аппроксимации функции  $f(x)$  используется линейный интерполяционный многочлен  $P_1(x, f, x_k, x_{k+1})$ . В случае квадратичной интерполяции на отрезке  $x_k \leq x \leq x_{k+1}$  можно воспользоваться одним из многочленов  $P_2(x, f, x_k, x_{k+1}, x_{k+2})$  или  $P_2(x, f, x_{k-1}, x_k, x_{k+1})$ .

Аналогично строятся кусочно многочленные интерполяции произвольной степени  $s$ . Именно, при заданном  $s$  на отрезке  $[x_k, x_{k+1}]$  используем интерполяционный многочлен  $P_s(x, f, x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s})$ , где  $j$  — одно из чисел  $0, 1, \dots, s - 1$ . Желательно, чтобы отрезок  $[x_k, x_{k+1}]$  был возможно «ближе» к середине отрезка  $[x_{k-j}, x_{k-j+s}]$ . Поэтому целесообразно выбирать в качестве  $j$  то из чисел  $0, 1, \dots, s - 1$ , которое ближе к  $s/2$ , и вместо  $P_s(x, f, x_{k-j}, \dots, x_{k-j+s})$  использовать сокращенное обозначение, положив

$$P_s(x, f, x_{k-j}, \dots, x_{k-j+s}) = P_s(x, f_{kj}).$$

**2. Формула для погрешности интерполяции.** Займемся оценкой погрешности

$$R_s(x) \equiv f(x) - P_s(x, f_{kj}), \quad x_k \leq x \leq x_{k+1}, \quad (1)$$

возникающей при приближенной замене  $f(x)$  многочленом  $P_s(x, f_{kj})$ . В основе лежит следующая общая теорема о формуле для погрешности.

**Теорема 1.** Пусть  $f(t)$  — функция, определенная на некотором отрезке  $\alpha \leq t \leq \beta$  и имеющая непрерывную производную некоторого порядка  $s + 1$ . Пусть  $t_0, t_1, \dots, t_s$  — произвольный набор попарно различных точек из отрезка  $[\alpha, \beta]$ ,  $f(t_0), f(t_1), \dots, f(t_s)$  — значения функции  $f(t)$  в этих точках, а  $P_s(t)$  — интерполяционный многочлен степени не выше  $s$ , построенный по этим значениям. Тогда погрешность интерполяции  $R_s(t) = f(t) - P_s(t)$  можно представить в виде

$$R_s(t) = \frac{f^{(s+1)}(\xi)}{(s+1)!} (t - t_0)(t - t_1) \dots (t - t_s), \quad (2)$$

где  $\xi = \xi(t)$  — некоторая точка из интервала  $\alpha < t < \beta$ .

**Доказательство.** Фиксируем произвольное  $t = \bar{t} \in [\alpha, \beta]$  и докажем формулу (1) при этом  $\bar{t}$ . Если  $\bar{t}$  совпадает с одним из узлов интерполяции,  $\bar{t} = t_j$  ( $j = 0, 1, \dots, s$ ), то погрешность  $f(\bar{t}) - P_s(\bar{t})$  обращается в нуль. Формула (2) также дает  $R_s(\bar{t}) = 0$  при произвольном  $\xi$ , так что формула (2) доказана для  $t = t_j$  ( $j = 0, 1, \dots, s$ ).

Докажем теперь формулу (1) для  $t = \bar{t}$  в предположении, что фиксированное  $\bar{t} \in [\alpha, \beta]$  не совпадает ни с одним из узлов интерполяции. Введем вспомогательную функцию

$$\varphi(t) = f(t) - P_s(t) - k(t - t_0)(t - t_1) \dots (t - t_s), \quad (3)$$

выбрав число  $k$  так, чтобы при  $t = \bar{t}$  функция  $\varphi(t)$  обращалась в нуль, т. е. положив

$$k = \frac{f(\bar{t}) - P_s(\bar{t})}{(\bar{t} - t_0)(\bar{t} - t_1) \dots (\bar{t} - t_s)}. \quad (4)$$

Числитель в формуле (4) есть значение  $R_s(\bar{t})$  погрешности; поэтому из этой формулы следует, что

$$R_s(\bar{t}) = k(\bar{t} - t_0)(\bar{t} - t_1) \dots (\bar{t} - t_s). \quad (5)$$

Функция  $\varphi(t)$  обращается в нуль в  $(s+2)$  точках  $\bar{t}, t_0, t_1, \dots, t_s$ . Производная  $\varphi'(t)$  обращается в нуль не менее, чем в  $(s+1)$  точках, так как по теореме Ролля между каждыми двумя соседними точками, где функция  $\varphi(t)$  обращается в нуль, найдется точка, где обращается в нуль ее производная  $\varphi'(t)$ . Аналогично  $\varphi''(t)$  имеет не менее, чем  $s$  нулей,  $\varphi^{(3)}(t)$  — не менее, чем  $(s-1)$  нулей и т. д., наконец,  $(s+1)$ -я

производная  $\frac{d^{s+1}}{dx^{s+1}}\varphi(t)$  обращается в нуль хотя бы в одной точке  $\xi \in (\alpha, \beta)$ . Заметим, что

$$\frac{d^{s+1}}{dt^{s+1}}t^{s+1} = (s+1)!,$$

а также  $(t - t_0)(t - t_1)\dots(t - t_s) = t^{s+1} + Q_s(t)$ , где  $Q_s(t)$  — некоторый многочлен степени  $s$ . Заметим еще, что

$$\frac{d^{s+1}}{dx^{s+1}}P_s(t) \equiv \frac{d^{s+1}}{dx^{s+1}}Q_s(t) \equiv 0.$$

Возьмем производную порядка  $s+1$  от функции  $\varphi(t)$ , заданной формулой (3). Получим

$$\varphi^{(s+1)}(t) = f^{(s+1)}(t) - k(s+1)!.$$

Отсюда при  $t = \xi$  в силу  $\varphi^{(s+1)}(\xi) = 0$  следует, что

$$k = \frac{f^{(s+1)}(\xi)}{(s+1)!}.$$

Подставляя найденное  $k$  в равенство (5), получаем формулу для  $R_s(\bar{t})$ , которая в силу произвольности выбора  $\bar{t}$  совпадает с формулой (2).  $\square$

**Теорема 2.** В условиях предыдущей теоремы справедлива оценка

$$\max_{\alpha \leq t \leq \beta} |R_s(t)| \leq \frac{1}{(s+1)!} \max_{\alpha \leq t \leq \beta} |f^{(s+1)}(t)|(\beta - \alpha)^{s+1}. \quad (6)$$

**Доказательство.** Заметим, что для любого значения  $t \in [\alpha, \beta]$  каждое из выражений  $t - t_0, t - t_1, \dots, t - t_s$  не превосходит по модулю число  $\beta - \alpha$ , а затем воспользуемся формулой (2):

$$\begin{aligned} |R_s(t)| &= \frac{1}{(s+1)!} |f^{(s+1)}(\xi)(t - t_0)(t - t_1)\dots(t - t_s)| \leq \\ &\leq \frac{1}{(s+1)!} \max_{\alpha \leq t \leq \beta} |f^{(s+1)}(t)|(\beta - \alpha)^{s+1}. \end{aligned} \quad (7)$$

Поскольку  $t \in [\alpha, \beta]$  в левой части выражения (7) произвольно, то отсюда следует (6).  $\square$

Подчеркнем, что оценка (6) доказана нами для произвольного расположения узлов  $t_0, t_1, \dots, t_s$  на отрезке  $[\alpha, \beta]$ .

Для конкретного фиксированного расположения узлов она может быть несколько улучшена. Например, в случае линейной интерполяции и расположения узлов  $t_0, t_1$  в концах  $\alpha, \beta$  отрезка  $\alpha \leq t \leq \beta$  получим

$$\begin{aligned} |R_1(t)| &= \left| \frac{f''(\xi)}{2!}(t - \alpha)(t - \beta) \right| \leq \\ &\leq \frac{1}{2} \max_{\alpha \leq t \leq \beta} |f''(t)| \max_{\alpha \leq t \leq \beta} |(t - \alpha)(t - \beta)| = \frac{1}{8} \max_{\alpha \leq t \leq \beta} |f''(t)|(\beta - \alpha)^2 \end{aligned} \quad (8)$$

и соответственно

$$\max_{\alpha \leq t \leq \beta} |R_1(t)| \leq \frac{1}{8} \max_{\alpha \leq t \leq \beta} |f''(t)|(\beta - \alpha)^2, \quad (9)$$

в то время как оценка (6) при  $s = 1$  принимает вид

$$\max_{\alpha \leq t \leq \beta} |R_1(t)| \leq \frac{1}{2} \max_{\alpha \leq t \leq \beta} |f''(t)|(\beta - \alpha)^2.$$

Воспользуемся теоремами 1, 2 для оценки погрешности (1) кусочно многочленной интерполяции функции  $f(x)$  на отрезке  $x_k \leq x \leq x_{k+1}$ . Положим

$$\begin{aligned} \alpha &= x_{k-j}, \quad \beta = x_{k-j+s}, \\ t_0 &= \alpha = x_{k-j}, \quad t_1 = x_{k-j+1}, \dots, t_s = \beta = x_{k-j+s}. \end{aligned}$$

Далее очевидно, что

$$\max_{x_k \leq x \leq x_{k+1}} |R_s(x, f_{kj})| \leq \max_{\alpha \leq x \leq \beta} |R_s(x, f_{kj})|,$$

а в силу (6) отсюда следует неравенство

$$\begin{aligned} \max_{x_k \leq x \leq x_{k+1}} |R_s(x, f_{kj})| &\leq \\ &\leq \frac{1}{(s+1)!} \max_{x_{k-j} \leq x \leq x_{k-j+s}} |f^{(s+1)}(x)| (x_{k-j+s} - x_{k-j})^{s+1}. \quad (10) \end{aligned}$$

Если величина  $|f^{(s+1)}(x)|$  сильно меняется на отрезке  $[a, b]$ , то для получения оценки (10), гарантирующей заданную точность, шаг сетки и число  $x_{k-j+s} - x_{k-j}$  должны быть меньше там, где  $|f^{(s+1)}(x)|$  больше.

В случае равноотстоящих узлов оценка (10) влечет

$$\max_{x_{k-1} \leq x \leq x_{k+1}} |R_s(x, f_{kj})| \leq \frac{s^{s+1}}{(s+1)!} \max_{x_{k-j} \leq x \leq x_{k-j+s}} |f^{(s+1)}(x)| h^{s+1}, \quad (11)$$

где  $h$  — шаг сетки узлов интерполяции.

Отметим особо случай кусочно линейной интерполяции. В этом случае  $s = 1$ ,  $\alpha = x_k$ ,  $\beta = x_{k+1}$ . Опираясь на оценку (9), получаем

$$\max_{x_k \leq x \leq x_{k+1}} |R_1(x)| \leq \frac{1}{8} \max_{x_k \leq x \leq x_{k+1}} |f''(x)| (x_{k+1} - x_k)^2. \quad (12)$$

Будем считать в дальнейшем, что  $x_0, x_1, \dots, x_n$  — равноотстоящие узлы, так что  $x_{k+1} - x_k = h = (b - a)/n$ . Из оценки (11) следует, что

$$\max |R_s(x, f_{kj})| \leq \text{const} \cdot \max_{x_{k-j} \leq x \leq x_{k-j+s}} |f^{(s+1)}(x)| h^{s+1}, \quad (13)$$

где постоянная не зависит от шага сетки  $h$ .

### 3. Приближение производных функции, заданной своими значениями в узлах сетки.

Теорема 3. Пусть  $f(x)$  определена на отрезке  $[\alpha, \beta]$  и имеет непрерывную производную некоторого порядка  $s+1$  на этом отрезке. Пусть  $x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s}$  — узлы интерполяции, причем  $\alpha = x_{k-j} < \dots < x_{k-j+s} = \beta$ . Тогда для вычисления производных

$$\frac{d^q f(x)}{dx^q}, \quad q = 1, 2, \dots, s,$$

на отрезке  $x_k \leq x \leq x_{k+1}$  можно воспользоваться интерполяционным многочленом  $P_s(x, f_{kj})$ , положив

$$\frac{d^q f(x)}{dx^q} \approx \frac{d^q}{dx^q} P_s(x, f_{kj}), \quad x_k \leq x \leq x_{k+1}; \quad (14)$$

при этом погрешность удовлетворяет оценке

$$\begin{aligned} \max_{x_k \leq x \leq x_{k+1}} \left| \frac{d^q f(x)}{dx^q} - \frac{d^q}{dx^q} P_s(x, f_{kj}) \right| &\leq \\ &\leq \frac{\max_{x_{k-j} \leq x \leq x_{k-j+s}} |f^{(s+1)}(x)|}{(s-q+1)!} (x_{k-j+s} - x_{k-j})^{s+1-q}. \end{aligned} \quad (15)$$

Доказательство. Функция  $\varphi(x) = f(x) - P_s(x, f_{kj})$  обращается в нуль в точках  $x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s}$ . Поэтому производная функции  $\varphi(x)$  обращается в нуль по крайней мере в  $s$  точках, поскольку между каждыми двумя нулями функции  $\varphi(x)$  по теореме Ролля есть нуль функции  $\varphi'(x)$ . Аналогично  $d^q \varphi(x)/dx^q$  имеет не менее, чем  $(s-q+1)$  нулей на отрезке  $x_{k-j} \leq x \leq x_{k-j+s}$ .

Это значит, что  $\frac{d^q f(x)}{dx^q}$  и многочлен  $\frac{d^q}{dx^q} P_s(x, f_{kj})$  степени не выше  $s-q$  совпадают в некоторых  $(s-q+1)$  точках, т. е. многочлен  $P_s^{(q)}(x, f_{kj})$  является интерполяционным для функции  $f^{(q)}(x)$  на отрезке  $x_{k-j} \leq x \leq x_{k-j+s}$  степени не выше  $s-q$  по некоторым  $(s-q+1)$  узлам, лежащим на этом отрезке. Функция  $f^{(q)}(x)$  имеет производную порядка  $s-q+1$ :

$$\frac{d^{s-q+1}}{dx^{s-q+1}} f^{(q)}(x) = \frac{d^{s+1}}{dx^{s+1}} f(x).$$

Поэтому можно воспользоваться теоремой 2, приняв  $\alpha = x_{k-j}$ ,  $\beta = x_{k-j+s}$ , и в силу оценки (6) написать

$$\begin{aligned} \max_{\alpha \leq x \leq \beta} |f^{(q)} - P_s^{(q)}(x, f_{kj})| &\leq \\ &\leq \frac{1}{(s-q+1)!} \max_{\alpha \leq x \leq \beta} |f^{(s+1)}(x)| (\beta - \alpha)^{s-q+1}, \end{aligned} \quad (16)$$

откуда в силу  $\alpha \leq x_k < x_{k+1} \leq \beta$  следует оценка (15).  $\square$

**4. Оценка неустранимой погрешности при приближении функции по ее значениям в узлах интерполяции и выбор степени кусочно многочленной интерполяции.** Пусть функция  $f(x)$  определена на отрезке  $[0, \pi]$ , и пусть заданы ее значения в узлах равномерной сетки  $x_k = k\pi/n$  ( $k = 0, 1, \dots, n$ ). По таблице  $f(x_0), f(x_1), \dots, f(x_n)$  в принципе нельзя восстановить функцию  $f(x)$  точно, потому что различные функции могут совпадать в точках  $x_k$  ( $k = 0, 1, \dots, n$ ), т. е. иметь одинаковые таблицы. Если, например, о функции  $f(x)$ , кроме ее таблицы, известно лишь, что она непрерывна, то ее нельзя восстановить в точке  $x \neq x_k$  ( $k = 0, 1, \dots, n$ ) ни с какой гарантированной точностью.

Пусть о функции  $f(x)$  известно, что она имеет производную некоторого порядка  $s+1$ , причем

$$\max_x |f^{(s+1)}(x)| \leq M_s = \text{const.} \quad (17)$$

Укажем две функции из этого класса,  $M_s = 1$ :

$$f^I(x) = \frac{\sin nx}{n^{s+1}}, \quad f^{II}(x) = -\frac{\sin nx}{n^{s+1}},$$

для которых таблицы совпадают (обе чисто нулевые):

$$f^I(x_k) = f^{II}(x_k) = 0, \quad k = 0, 1, \dots, n,$$

и которые уклоняются друг от друга на величину порядка  $h^{s+1}$ :

$$\max_{0 \leq x \leq 1} |f^I(x) - f^{II}(x)| = \max_{0 \leq x \leq \pi} 2 \left| \frac{\sin nx}{n^{s+1}} \right| = 2h^{s+1}.$$

Таким образом, по таблице функции при дополнительном знании лишь оценки (17) в принципе нельзя восстановить функцию всюду на отрезке  $0 \leq x \leq \pi$  с точностью больше  $O(h^{s+1})$ . Другими словами, погрешность  $O(h^{s+1})$  неустранима при восстановлении функции  $f(x)$  ( $0 \leq x \leq \pi$ ) по ее таблице.

Очевидно, что

$$\max_x \left| \frac{d^q}{dx^q} f^I(x) - \frac{d^q}{dx^q} f^{II}(x) \right| = 2 \frac{1}{n^{s-q+1}} = 2h^{s-q+1},$$

так что неустранимая погрешность при восстановлении производной  $d^q f(x)/dx^q$  есть  $O(h^{s-q+1})$ .

Как мы видели, восстановление функции  $f(x)$  или ее производной  $d^q f(x)/dx^q$  по приближенным формулам (14) имеет погрешность, по порядку совпадающую с неустранимой. Если использовать интерполяцию степени  $q < s$ , то погрешность будет  $O(h^{q+1})$ , т. е. будет происходить потеря порядка малости погрешности, дополнительная к той  $O(h^{s+1})$ , которая обусловлена заданием функции ее таблицей и неустранима.

С другой стороны, использование интерполяции степени  $q > s$  не может повысить порядок малости (неустранимой) погрешности  $O(h^{s+1})$

и ускорить сходимость при  $h \rightarrow 0$ . В этом смысле степень  $s$  кусочно многочленной интерполяции функций, удовлетворяющих условию (17), является оптимальной.

**Замечание.** Проведенные рассуждения относятся к поведению погрешности при  $h \rightarrow 0$ . При фиксированном  $h$  интерполяция некоторой степени  $q < s$  может оказаться более точной, чем интерполяция степени  $s$ . Кроме того, если значения  $f(x_k)$  ( $k = 0, 1, \dots, n$ ) заданы приближенно с некоторым числом десятичных знаков, то потеря точности при интерполяции за счет приближенного задания  $f(x_k)$  ( $k = 0, 1, \dots, n$ ) при увеличении  $s$  возрастает (растут константы Лебега (19) из § 1). Поэтому интерполяция высокой степени (выше третьей) применяется редко.

**5. Насыщаемость (гладкостью) кусочно многочленной интерполяции.** Пусть  $f(x)$  определена на отрезке  $[a, b]$ , и пусть задана ее таблица  $f(x_k)$  в равноотстоящих узлах  $x_k$  ( $k = 0, 1, \dots, n$ );  $h = (b - a)/n$ . Мы видели, что погрешность кусочно многочленной интерполяции степени  $s$  с помощью интерполяционных многочленов  $P_s(x, f_{kj})$  на отрезке  $x_k \leq x \leq x_{k+1}$  в случае, если  $f^{(s+1)}(x)$  существует и ограничена, имеет порядок  $O(h^{s+1})$ . Если об  $f(x)$  известно лишь, что она имеет ограниченную производную некоторого порядка  $q + 1$  ( $q < s$ ), то неустранимая погрешность при ее восстановлении по таблице есть  $O(h^{q+1})$ . Можно показать, что при интерполяции с помощью  $P_s(x, f_{kj})$  порядок  $O(h^{q+1})$  достигается. Если  $f(x)$  имеет ограниченную производную порядка  $q + 1$  ( $q > s$ ), то погрешность интерполяции с помощью  $P_s(x, f_{kj})$  остается равной  $O(h^{s+1})$ , т. е. порядок погрешности не реагирует на дополнительную сверх ( $s + 1$ ) производных гладкость функции  $f(x)$ . Это свойство кусочно многочленной интерполяции называют *свойством насыщаемости (гладкостью)*.

### Задачи

1. Каков должен быть шаг  $h$  таблицы функции  $f(x) = \sin x$ , чтобы при кусочно линейной интерполяции погрешность не превосходила  $10^{-6}$ ?

2. Каков должен быть шаг  $h$  таблицы функции  $f(x) = \sin x$ , чтобы погрешность при кусочно квадратичной интерполяции не превосходила  $10^{-6}$ ?

3. Значения  $f(x)$  могут быть измерены в заданной точке с погрешностью  $|\delta f| \leq 10^{-4}$ .

С каким шагом  $h$  разумно составить таблицу функции  $f(x)$ , если восстанавливать функцию по ее таблице с помощью кусочно линейной интерполяции?

4\*. Вопрос задачи 3, но для кусочно квадратичной интерполяции.

5. Пусть

$$f'(x) \approx \begin{cases} \frac{f(x+h) - f(x)}{h}, \\ \frac{f(x+h) - f(x-h)}{2h}, \end{cases} \quad (18)$$

(19)

и пусть  $|f''(x)| \leq 1$  и  $|f'''(x)| \leq 1$ .

- а) При каких  $h$  можно гарантировать, что погрешность меньше  $10^{-3}$ ?
- б)\* Пусть  $f$  задана с погрешностью  $\delta$ . Какова наибольшая точность, достижимая по формулам (18), (19), и как выбрать  $h$ ?
- в)\* Показать, что результат, полученный для оптимального  $h$  по формуле (19), не улучшаем по порядку погрешности относительно  $\delta$ .

### § 3. Кусочно многочленная гладкая интерполяция (сплайны)

Классическая кусочно линейная, кусочно квадратичная и вообще кусочно многочленная интерполяции заданной степени  $s$  приводят к интерполирующей функции (интерполянту), которая в узлах интерполяции, вообще говоря, не имеет производной даже первого порядка.

Существуют два типа кусочно многочленных интерполянтов — локальные и нелокальные сплайны, обладающие заданным числом производных всюду, включая узлы интерполяции.

**1. Локальная интерполяция гладкости  $s$  и ее свойства.** Пусть заданы узлы  $x_l$  интерполяции и значения функции  $f(x_l)$  в них. Зададим натуральное (целое положительное) число  $s$ , фиксируем натуральное  $j$  ( $0 \leq j \leq s - 1$ ). Каждой точке  $x_l$  сопоставим интерполяционный многочлен  $P_s(x, f_{lj})$ , построенный по значениям  $f(x_{l-j}), f(x_{l-j+1}), \dots, f(x_{l-j+s})$  в узлах  $x_{l-j}, x_{l-j+1}, \dots, x_{l-j+s}$ . Кусочно многочленный локальный сплайн  $\varphi(x, s)$ , имеющий непрерывные производные порядка  $s$ , определим равенствами

$$\varphi(x, s) = Q_{2s+1}(x, k), \quad x \in [x_k, x_{k+1}], \quad k = 0, \pm 1, \dots, \quad (1)$$

где  $Q_{2s+1}(x, k)$  — многочлен степени не выше  $2s + 1$ , определяемый равенствами

$$\left. \frac{d^m Q_{2s+1}(x, k)}{dx^m} \right|_{x=x_k} = \left. \frac{d^m P_s(x, f_{kj})}{dx^m} \right|_{x=x_k}, \quad m = 0, 1, 2, \dots, s, \quad (2)$$

$$\left. \frac{d^m Q_{2s+1}(x, k)}{dx^m} \right|_{x=x_{k+1}} = \left. \frac{d^m P_s(x, f_{k+1,j})}{dx^m} \right|_{x=x_{k+1}}, \quad m = 0, 1, 2, \dots, s. \quad (3)$$

**Теорема 1.** Существует один и только один многочлен степени не выше  $2s + 1$ , удовлетворяющий (2), (3).

**Доказательство.** Система  $(2s + 2)$  линейных уравнений относительно коэффициентов многочлена

$$Q_{2s+1}(x, k) = c_{0k} + c_{1k}x + \dots + c_{2s+1,k}x^{2s+1},$$

получающегося из системы (2), (3) при замене правых частей нулями, имеет только тривиальное решение

$$c_{ik} \underset{i}{\equiv} 0.$$

Действительно, в противном случае существовал бы многочлен  $Q_{2s+1}$  степени не выше  $2s+1$ , который в силу (2), (3) имел бы корни  $x = x_k$ ,  $x = x_{k+1}$  кратности  $s+1$  каждый, т. е. имел бы  $(2s+2)$  корней, считая их кратности. В силу этого противоречия определитель системы (2), (3) отличен от нуля, и она имеет одно и только одно решение при любых правых частях.  $\square$

**Теорема 2.** Пусть  $f(x)$  есть многочлен степени не выше  $s$ . Тогда интерполянт  $\varphi(x, s)$  совпадает с этим многочленом.

**Доказательство.** Докажем тождество  $\varphi(x, s) = f(x)$  на отрезке  $x_k \leq x \leq x_{k+1}$  между соседними узлами, т. е. докажем, что  $Q_{2s+1}(x, k) = f(x)$ . В силу единственности интерполяционного многочлена  $P_s(x, f_{kj}) \equiv P_s(x, f_{k+1,j}) \equiv f(x)$ . Теперь видно, что многочлен  $Q_{2s+1}(x, k) \equiv f(x)$  дает решение системы (2), (3).  $\square$

**Теорема 3.** Кусочно многочленная интерполирующая функция  $\varphi(x, s)$ , определенная равенством (1), в узлах интерполяции  $x_l$  совпадает с заданными в них значениями  $f(x_l)$  ( $l = 0, \pm 1, \dots$ ). Кроме того,  $\varphi(x, s)$  имеет всюду в области своего определения непрерывную производную порядка  $s$ .

**Доказательство.** В произвольном узле  $x_l$  функции  $Q_{2s+1}(x, l-1)$  и  $Q_{2s+1}(x, l)$  в силу равенств (2), (3) имеют производные порядка  $m = 0, 1, \dots, s$ , совпадающие с соответствующими производными одного и того же интерполяционного многочлена  $P_s(x, f_{lj})$ . В силу (1) это доказывает теорему.  $\square$

Запишем  $Q_{2s+1}(x, k)$  в виде

$$Q_{2s+1}(x, k) = P_s(x, f_{kj}) + R_{2s+1}(x, k), \quad (4)$$

обозначив через  $R_{2s+1}(x, k)$  поправку к классическому интерполяционному многочлену  $P_s(x, f_{kj})$ .

**Теорема 4.** Поправку  $R_{2s+1}(x, k)$  можно записать в виде

$$R_{2s+1}(x, k) = (x_{k+1} - x_k)^{s+1} f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1}) q_{2s+1}\left(\frac{x - x_k}{x_{k+1} - x_k}, k\right), \quad (5)$$

где  $f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1})$  — разностное отношение порядка  $s+1$ ,

$$q_{2s+1}(X, k) = \left( \frac{x_{k+s-j+1} - x_{k-j}}{x_{k+1} - x_k} \right) \sum_{r=0}^s \left\{ \left[ \prod_{i=1}^s \left( X - \frac{x_{k-j+i} - x_k}{x_{k+1} - x_k} \right) \right]_{X=1}^{(r)} \right\} l_r(X), \quad (6)$$

$$X = \frac{x - x_k}{x_{k+1} - x_k}, \quad l_r(X) = \frac{X^{s+1} (X-1)^r}{r! s!} \sum_{m=0}^{s-r} (-1)^m \frac{(s+m)!}{m!} (X-1)^m. \quad (7)$$

(Выражение  $[...]_{X=1}^{(r)}$  означает производную по  $X$  порядка  $r$ , вычисленную при  $X = 1$ .)

**Замечание 1.** Можно считать, что формула локальной интерполяции гладкости  $s$  получена путем замены в интерполяционном многочлене  $P_{s+1}(x, f_{kj})$ , записанном в форме Ньютона

$$P_{s+1}(x, f_{kj}) = P_s(x, f_{kj}) + f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1}) \varphi_{s+1}(x, k),$$

$$\varphi_{s+1}(x, k) = (x - x_{k-j})(x - x_{k-j+1}) \dots (x - x_{k-j+s}),$$

многочлена  $\varphi_{s+1}(x, k)$  многочленом

$$(x_{k+1} - x_k)^{s+1} q_{2s+1} \left( \frac{x - x_k}{x_{k+1} - x_k}, k \right).$$

Для иллюстрации теоремы 4 заметим, что в силу формул (4)–(7) в случае  $s = 0$ ,  $j = 0$  локальный сплайн  $\varphi(x, s)$  осуществляет кусочно линейную интерполяцию. В наиболее интересном для приложений случае  $s = 2$ ,  $j = 1$  и  $x_{k+j} - x_k = h = \text{const}$

$$R_5(x, k) = P_2(x, f_{k1}) + \frac{h^3}{2!} \cdot \frac{f(x_{k+2}) - 3f(x_{k+1}) + 3f(x_k) - f(x_{k-1})}{h^3} \times \\ \times \left( \frac{x - x_k}{h} \right)^3 \frac{x - x_{k+1}}{h} \left( 3 - \frac{2(x - x_k)}{h} \right). \quad (8)$$

Доказательство теоремы 4 приведем в конце параграфа.

**Теорема 5.** Пусть функция  $f(x)$  определена всюду на отрезке  $x \in [x_{k-j}, x_{k-j+s+1}]$  и имеет на этом отрезке ограниченную производную порядка  $s+1$ . Тогда на отрезке  $x \in [x_k, x_{k+1}]$  справедливы приближенные равенства

$$f^{(m)}(x) \approx \frac{d^m \varphi(x, s)}{dx^m}, \quad m = 0, 1, \dots, s, \quad (9)$$

с оценками погрешности

$$\left| f^{(m)}(x) - \frac{d^m \varphi(x, s)}{dx^m} \right| \leq \text{const} \frac{(x_{k+s-j+1} - x_{k-j})^{s+1}}{(x_{k+1} - x_k)^m} \times \\ \times \max_{x_{k-j} \leq x \leq x_{k-j+s+1}} |f^{(s+1)}(x)|, \quad m = 0, 1, \dots, s. \quad (10)$$

**Доказательство.** В силу равенства (1) и формулы (4)

$$\left| f^{(m)}(x) - \frac{d^m \varphi(x, s)}{dx^m} \right| \leq |f^{(m)}(x) - P_s^{(m)}(x, f_{kj})| + |R_{2s+1}^{(m)}(x, k)|. \quad (11)$$

Но в теореме 3 из § 2 установлена оценка

$$|f^{(m)}(x) - P_s^{(m)}(x, f_{kj})| \leq \text{const}_1 \cdot h^{s+1-m}. \quad (12)$$

Далее, учитывая, что  $\frac{d^m}{dx^m} = \frac{1}{(x_{k+1} - x_k)^m} \cdot \frac{d^m}{dX^m}$ ,  $X = \frac{x - x_k}{x_{k+1} - x_k}$ , получаем

$$\left| \frac{d^m R_{2s+1}(x, k)}{dx^m} \right| = (x_{k+1} - x_k)^{s+1-m} |f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1})| \times \\ \times \left| \frac{x_{k+s-j+1} - x_{k-j}}{x_{k+1} - x_k} \right| \left| \sum_{r=0}^s \left\{ \left[ \prod_{i=1}^s \left( X - \frac{x_{k-j+i} - x_k}{x_{k+1} - x_k} \right) \right]_{X=1}^{(r)} \right\} l_r^{(m)}(X) \right|. \quad (13)$$

В силу теоремы 3 из § 1

$$(s+1)! f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1}) = f^{(s+1)}(\xi), \quad (14)$$

$$\xi \in [x_{k-j}, x_{k-j+s+1}].$$

Тогда из (13) очевидно, что

$$\left| \frac{d^m R_{2s+1}(x, k)}{dx^m} \right| \leqslant \\ \leqslant \text{const}_2 \cdot |x_{k+1} - x_k|^{s+1-m} \max_{x_{k-j} \leqslant x \leqslant x_{k-j+s+1}} \left| \frac{d^{s+1} f(x)}{dx^{s+1}} \right|, \quad (15)$$

где  $\text{const}_2$  зависит только от величины отношения

$$\frac{x_{k-j+s+1} - x_{k-j}}{x_{k+1} - x_k} \quad (16)$$

и остается ограниченной, если величина (16) в процессе измельчения сетки остается ограниченной.

Из выражений (11), (12), (15) следует оценка (10).  $\square$

Сделаем два замечания о неулучшаемости построенных нами локальных сплайнов (1) в некотором смысле.

**Замечание 2.** Оценки (10) в случае постоянного шага  $x_{k+1} - x_k = (b - a)/n = h$  гарантируют сходимость  $\varphi^{(m)}(x, s)$  к  $f^{(m)}(x)$  с порядком  $h^{s+1-m}$  ( $m = 0, 1, \dots, s$ ). В предположении, что  $f(x)$  имеет ограниченные производные только до порядка  $s + 1$ , более точно восстановить функцию  $f(x)$  и ее производные по значениям  $f(x_l)$  ( $l = 0, \pm 1, \dots$ ) нельзя, поскольку эти значения не содержат достаточной для этого информации, как мы показали в п. 4 § 2.

**Замечание 3.** Для вычисления  $\varphi(x, s)$  при каждом  $x$  используется не более  $(s+2)$  узлов интерполяции. Эта характеристика локальности формулы (1) неулучшаема в следующем смысле.

Если потребовать, чтобы использовалось не более  $(s+1)$  узлов интерполяции, то для сходимости  $\varphi^{(m)}(x, s)$  к  $f^{(m)}(x)$  с порядком  $h^{s+1-m}$  пришлось бы отказаться от условия непрерывности даже первых производных от  $\varphi(x, s)$  и ограничиться классической кусочно многочленной интерполяцией степени не выше  $s$ .

Если задано конечное число  $(n+1)$  узлов  $x_0, x_1, \dots, x_n$ , лежащих на отрезке  $a \leqslant x \leqslant b$ ,  $a = x_0 < x_1 < \dots < x_n = b$ , то формулу (1)

вблизи концов отрезка надо изменить, определив интерполант  $\varphi(x) = \varphi(x, s, a, b)$  равенствами

$$\varphi(x, s, a, b) = \begin{cases} P_s(x, f_{jj}), & \text{если } x_0 \leq x \leq x_j, \\ Q_{2s+1}(x, k), & \text{если } x \in [x_k, x_{k+1}] \\ & (j \leq k \leq j + n - s - 1), \\ P_s(x, f_{n+j-s,j}), & \text{если } x_{n+j-s} \leq x \leq x_n. \end{cases}$$

Напомним, что обозначение  $P_n(x, f_{kj})$  введено в конце п. 1 § 2.

Интерполант  $\varphi(x, s, a, b)$ , определенный на  $[a, b]$ , имеет непрерывные производные до порядка  $s$ . Для него имеют смысл утверждения, аналогичные теоремам 1–5.

Локальная гладкая интерполяция введена автором в 1952 г. для функций на многомерной прямоугольной сетке с постоянным шагом; приведенные здесь формулы получены автором в статье: *Ryaben'kii V.S. Local splines // Comp. Math., Mech., and Enginier.* — V. 5, № 2. Р. 211–225.

## 2. Нелокальная гладкая кусочно многочленная интерполяция.

Пусть заданы узлы  $x_0, x_1, \dots, x_n$ ,  $a = x_0 < x_1 < \dots < x_n = b$  и значения функции  $f(x_i)$  ( $i = 0, 1, \dots, n$ ) в этих узлах. Поставим задачу: на каждом отрезке  $x_k \leq x \leq x_{k+1}$  найти кубический многочлен  $P_3(x, k)$  так, чтобы возникающая при этом на отрезке  $a \leq x \leq b$  кусочно многочленная функция совпадала с заданной функцией в узлах и имела непрерывные производные до порядка  $s = 2$ . Эта функция зависит от двух произвольных постоянных и называется *кубическим сплайном Шонберга*. Сплайны минимальной для заданного  $s$  степени построены Шонбергом не только для  $s = 2$ , но и для всех натуральных  $s$ .

Пусть  $f(x)$  имеет на  $[a, b]$  ограниченную производную некоторого порядка  $s + 1$ . Тогда непрерывный с производными порядка  $s$  сплайн при подходящих постоянных сохраняет неулучшаемые приближающие свойства классической, а также локальной гладкой кусочно многочленной интерполяций в том смысле, что для него в случае, если существует  $f^{(s+1)}(x)$ , выполнены оценки типа (10). Однако сплайны Шонберга теряют свойство локальности, присущее как классической кусочно многочленной интерполяции, так и локальной гладкой интерполяции: коэффициенты многочлена, задающего интерполант на каком-либо отрезке  $x_k \leq x \leq x_{k+1}$ , зависят от значений  $f(x_0), f(x_1), \dots, f(x_n)$  функции во всех узлах  $x_0, x_1, \dots, x_n$ , число которых  $n$  растет в случае измельчения сетки.

Отметим еще, что (нежелательное) свойство насыщаемости гладкостью, неизбежно присущее локальным классической и гладкой интерполяциям, несмотря на потерю локальности, остается и у нелокальных сплайнов: для функций  $f(x)$ , имеющих более  $(s + 1)$  производных, нелокальные сплайны гладкости  $s$  при постоянном шаге  $h = (b - a)/n$  сетки имеют погрешность  $O(h^{s+1})$  того же порядка малости по отношению к  $h = 1/n$ , что и для функций  $f(x)$ , имеющих только  $(s + 1)$ -ю производную.

Нелокальные сплайны Шонберга при заданном числе непрерывных производных реализуются многочленами наименьшей возможной степени  $n$  (например,  $m = 3$ , если  $s = 2$  — кубические сплайны Шонберга).

Подробнее о сплайнах Шонберга и других сплайнах см. [1, 2, 9, 13, 17, 23] и имеющуюся там библиографию.

### 3. Доказательство теоремы 4. Коэффициенты многочлена

$$Q_{2s+1}(x, k) = c_0 + c_1 x + \dots + c_{2s+1} x^{2s+1} \quad (17)$$

определяются из системы линейных уравнений (2), (3). Правые части уравнений, составляющих систему (2), имеют вид

$$a_0^{(m)} f_{k-j} + a_1^{(m)} f_{k-j+1} + \dots + a_s^{(m)} f_{k-j+s},$$

а составляющих систему (3) — вид

$$b_0^{(m)} f_{k-j+1} + b_1^{(m)} f_{k-j+2} + \dots + b_s^{(m)} f_{k-j+s+1}, \quad m = 0, 1, \dots, s,$$

где  $a_i^{(m)}, b_i^{(m)}$  ( $i = 0, 1, \dots, s$ ) — некоторые числа, не зависящие от  $f_{k-j}, f_{k-j+1}, \dots, f_{k-j+s+1}$ .

Поэтому решение  $c_0, c_1, \dots, c_{2s+1}$  системы (2), (3) состоит из чисел, которые определяются заданными  $f_{k-j}, f_{k-j+1}, \dots, f_{k-j+s+1}$  по формулам вида

$$c_r = \alpha_0^{(r)} f_{k-j} + \alpha_1^{(r)} f_{k-j+1} + \dots + \alpha_{s+1}^{(r)} f_{k-j+s+1}, \quad r = 0, 1, \dots, 2s+1, \quad (18)$$

где  $\alpha_l^{(r)}$  — некоторые числа, не зависящие от  $f_{k-j}, f_{k-j+1}, \dots, f_{k-j+s+1}$ . Подставив выражения (18) в (17) и приведя подобные члены, содержащие  $f_{k-j+l}$  ( $l = 0, 1, \dots, s+1$ ), получим запись многочлена (17) в виде

$$Q_{2s+1}(x, k) = \sum_{l=0}^{s+1} f_{k-j+l} p(x, l), \quad (19)$$

где  $p(x, l)$ ,  $l = 0, 1, \dots, s+1$ , — некоторые многочлены от  $x$ . Заменим в (19) числа  $f_{k-j+l}$  по формулам (18) из § 1. Многочлен  $Q_{2s+1}(x, k)$  можно записать в виде

$$Q_{2s+1}(x, k) = \sum_{l=0}^s f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+l}) q_l(x, k) + \\ + (x_{k+1} - x_k)^{s+1} f(x_{k-j}, x_{k-j+2}, \dots, x_{k-j+s+1}) \tilde{q}_{2s+1}(x, k), \quad (20)$$

где  $q_l(x, k)$ ,  $l = 0, 1, \dots, s$ , и  $\tilde{q}_{2s+1}(x, k)$  — некоторые многочлены, не зависящие от  $f(x_{k-j}), f(x_{k-j+1}), \dots, f(x_{k-j+s+1})$ . Воспользуемся этой независимостью для отыскания многочленов  $q_l(x, k)$ . Зададим  $f_{k-j}, f_{k-j+1}, \dots, f_{k-j+s}$  произвольно, а  $f_{k-j+s+1}$  определим равенством

$$f_{k-j+s+1} = P_s(x, f_{k,j})|_{x=x_{k-j+s+1}}.$$

В силу теоремы 2 в этом случае

$$Q_{2s+1}(x, k) = P_s(x, f_{kj}),$$

а в силу следствия 3 из формулы Ньютона (см. § 1) разностное отношение  $f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1})$  порядка  $s + 1$  обращается в нуль. Таким образом, формула (20) примет вид

$$P_s(x, f_{kj}) = \sum_{l=0}^s f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+l}) q_l(x, k). \quad (21)$$

В силу произвольности значений  $f_{k-j}, f_{k-j+1}, \dots, f_{k-j+s}$  отсюда следует, что  $q_l(x, k) = (x - x_{k-j})(x - x_{k-j+1}) \dots (x - x_{k-j+l})$ ,  $l = 0, 1, \dots, s$ . Из (21) следует также, что равенству (20) можно придать вид

$$\begin{aligned} Q_{2s+1}(x, k) &= P_s(x, f_{kj}) + \\ &+ (x_{k+1} - x_k)^{s+1} f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1}) \tilde{q}_{2s+1}(x, k). \end{aligned} \quad (22)$$

Для вычисления многочлена  $\tilde{q}_{2s+1}(x, k)$  зададим

$$f_{k-j} = f_{k-j+1} = \dots = f_{k-j+s} = 0, \quad f_{k-j+s+1} = 1. \quad (23)$$

При условиях (23) формула (22) примет вид

$$Q_{2s+1}(x, k) = (x_{k+1} - x_k)^{s+1} f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1}) \tilde{q}_{2s+1}(x, k). \quad (24)$$

Заметим, что при условиях (23) интерполяционный многочлен  $P_{s+1}(x, f, x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1})$  в силу записи в форме Лагранжа совпадает с многочленом

$$\prod_{l=0}^s \frac{x - x_{k-j+l}}{x_{k-j+s+1} - x_{k-j+l}}, \quad (25)$$

а в силу записи в форме Ньютона (см. (4) из § 1) имеет вид

$$f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1}) \prod_{l=0}^s (x - x_{k-j+l}). \quad (26)$$

Приравнивая (25) и (26) друг другу, получаем, что при условиях (23)

$$f(x_{k-j}, x_{k-j+1}, \dots, x_{k-j+s+1}) = \left[ \prod_{l=0}^s (x_{k-j+s+1} - x_{k-j+l}) \right]^{-1}. \quad (27)$$

Таким образом, из (24) получаем

$$\tilde{q}_{2s+1}(x, k) = (x_{k+1} - x_k)^{-s-1} Q_{2s+1}(x, k) \prod_{l=0}^s (x_{k-j+s+1} - x_{k-j+l}), \quad (28)$$

где выражения в правой части вычислены при условиях (23).

Преобразуем (28) к требуемому виду (6), (7). Для этого найдем явную формулу для  $Q_{2s+1}(x, k)$  при условиях (23). Заметим, что благодаря условию (2) значение  $x = x_k$  является корнем кратности  $s + 1$  многочлена  $Q_{2s+1}(x, k)$ , что  $Q(x) = Q_{2s+1}(x)/(x - x_k)^{s+1}$  есть многочлен степени  $s$ .

Разложим  $Q(x)$  по степеням  $x - x_{k+1}$ :

$$\begin{aligned} Q_{2s+1}(x, k) &= (x - x_k)^{s+1} \frac{Q_{2s+1}(x, k)}{(x - x_k)^{s+1}} = \\ &= (x - x_k)^{s+1} \sum_{r=0}^s \frac{1}{r!} \left[ \frac{Q_{2s+1}(x, k)}{(x - x_k)^{s+1}} \right]_{x=x_{k+1}}^{(r)} (x - x_{k+1})^r. \end{aligned} \quad (29)$$

Используем равенство (3) и напишем

$$\begin{aligned} \frac{d^r}{dx^r} \left[ \frac{Q_{2s+1}(x, k)}{(x - x_k)^{s+1}} \right] \Big|_{x=x_{k+1}} &= \\ &= \sum_{l=0}^r C_r^l \left( \frac{d^{r-l}}{dx^{r-l}} P_s(x, f, x_{k-j+1}, \dots, x_{k-j+s+1}) \right) \Big|_{x=x_{k+1}} \times \\ &\quad \times [(x - x_k)^{-s-1}]_{x=x_{k+1}}^{(l)}. \end{aligned} \quad (30)$$

Но при условии (23) в силу (1)

$$P_s(x, f, x_{k-j+1}, \dots, x_{k-j+s+1}) = \prod_{i=1}^s \frac{x - x_{k-j+i}}{x_{k-j+s+1} - x_{k-j+i}}. \quad (31)$$

Подставляя (31) в (30), получаем

$$\begin{aligned} \frac{d^r}{dx^r} \left[ \frac{Q_{2s+1}(x, k)}{(x - x_k)^{s+1}} \right] \Big|_{x=x_{k+1}} &= \left[ \prod_{i=1}^s (x_{k-j+s+1} - x_{k-j+i}) \right]^{-1} \times \\ &\quad \times \sum_{l=0}^r C_r^l [(x - x_k)^{-s-1}]_{x=x_{k+1}}^{(l)} \left[ \prod_{i=1}^s (x - x_{k-j+i}) \right]_{x=x_{k+1}}^{r-l}. \end{aligned} \quad (32)$$

Подставляя (32) в (29), а (29) в (28), получаем явное выражение для  $\tilde{q}_{2s+1}(x, k)$  в виде двойной суммы по  $l$  и по  $r$ , которое зависит от  $x$  только как сложная функция  $q_{2s+1}(X, k)$  от  $X = (x - x_k)/(x_{k+1} - x_k)$ :

$$\tilde{q}_{2s+1}(x, k) = q_{2s+1}(X, k).$$

Изменяя порядок суммирования в полученной формуле для  $q_{2s+1}(X, k)$ , получаем запись в форме (6), (7).  $\square$

### Задачи

1. Пусть  $x_{k+1} - x_k = h = \text{const}$ .

Выпишите многочлены  $l_t(x, s)$ , с помощью которых локальный сплайн  $\tilde{\psi}(x, s)$  заданной гладкости  $s$ , определяемый при  $x_k \leq x \leq x_{k+1}$  равенством (4), можно записать в форме

$$\tilde{\psi}(x, s) = \sum_t f(x_t) l_t(x, s).$$

Рассмотреть случаи:  $s = 0, j = 0; s = 1, j = 0; s = 2, j = 2$ .

**2.** Найти точные значения констант Лебега

$$L_s = \max_{\substack{|f_m|=1, \\ k-j \leq m \leq k+s+1}} \max_{x_k \leq x \leq x_{k+1}} |Q_{2s+1}(x)|$$

для локальных гладких сплайнов при  $s = 2, j = 1$ .

**3.** Доказать, что в случае  $x_{t+1} - x_t = h = \text{const}$

$$\begin{aligned} & \max_{x_k \leq x \leq x_{k+1}} \left| \frac{d^m Q_{2s+1}(x)}{dx^m} \right| \leq \\ & \leq c \max_{0 \leq j \leq s-m+1} |f(x_{k-j+i}, x_{k-j+i+1}, \dots, x_{k-j+i+m})|, \quad m = 0, 1, \dots, s, \end{aligned}$$

где  $c = c(s, j, m)$  не зависит от  $h$ .

Вычислить  $c$  в случаях:  $s = 0, j = 0; s = 2, j = 1$ .

#### § 4. Интерполяция функций двух переменных

Пусть функция двух переменных  $f(x, y)$  задана в узлах некоторой правильной или нерегулярной сетки. Как восстановить ее приближенно в точках  $(x, y)$ , не принадлежащих множеству узлов?

**1. Случай прямоугольной сетки.** Пусть сетка образована пересечением прямых  $x = x_k$  ( $k = 0, \pm 1, \dots$ ) и прямых  $y = y_l$  ( $l = 0, \pm 1, \dots$ ). Считаем, что  $x_{k+1} > x_k, y_{l+1} > y_l$  при любых целых  $k, l$ . Значение функции в узле  $(x_k, y_l)$  обозначим через  $f_{kl}$ . Для вычисления функции в точке  $(\bar{x}, \bar{y})$  можно воспользоваться аппаратом кусочно многочленной интерполяции заданной степени  $s$  для функций одного переменного.

Для этого сначала осуществляется кусочно многочленная интерполяция заданной степени по  $x$  на каждой из прямых  $y = y_l$ . Затем при каждом интересующем нас значении  $x = \bar{x}$  осуществляется кусочно многочленная интерполяция (той же или другой степени) по  $y$  вдоль прямой  $x = \bar{x}$  по значениям функции  $f(x, y)$  в точках  $(\bar{x}, y_l)$ , полученным на первом шаге процесса. Например, в случае кусочно линейной интерполяции по обоим аргументам в прямоугольнике  $x_k \leq x \leq x_{k+1}, y_l \leq y \leq y_{l+1}$  этот процесс приводит к интерполяционному многочлену

$$\begin{aligned} P(x, y) = & f_{kl} \frac{(x - x_{k+1})(y - y_{l+1})}{(x_k - x_{k+1})(y_l - y_{l+1})} + f_{k+1,l} \frac{(x - x_k)(y - y_{l+1})}{(x_{k+1} - x_k)(y_l - y_{l+1})} + \\ & + f_{k+1,l+1} \frac{(x - x_k)(y - y_l)}{(x_{k+1} - x_k)(y_{l+1} - y_l)} + f_{k,l+1} \frac{(x - x_{k+1})(y - y_l)}{(x - x_{k+1})(y_{l+1} - y_l)}. \end{aligned}$$

**2. Треугольные сетки.** Пусть функция  $f(x, y)$  определена в некоторой криволинейной области  $D$ , причем известно, что в «горловине» функция быстро изменяется.

В таком случае для табличного задания функции прямоугольная сетка неудобна, так как она не связана с формой границы области  $D$  и, кроме того, ее нельзя сгустить только в «горловине». Можно воспользоваться треугольной сеткой (рис. 2), свободной от этих недостатков, и принять за узлы интерполяции вершины треугольников.

Если  $f_\alpha, f_\beta, f_\gamma$  — значения  $f(x, y)$  в вершинах  $\alpha, \beta, \gamma$  треугольника, то можно приближенно вычислить значение функции внутри треугольника с помощью линейной функции

$$f(x, y) \approx P(x, y) = ax + by + c,$$

подобрав  $a, b, c$  из условий

$$ax_\alpha + by_\alpha + c = f_\alpha,$$

$$ax_\beta + by_\beta + c = f_\beta,$$

$$ax_\gamma + by_\gamma + c = f_\gamma,$$

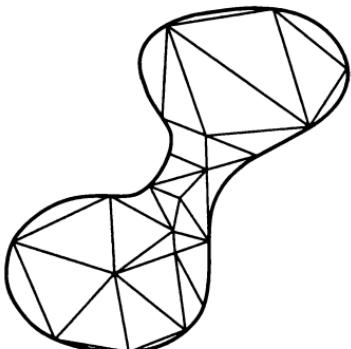


Рис. 2

где  $(x_\alpha, y_\alpha), (x_\beta, y_\beta), (x_\gamma, y_\gamma)$  — соответственно координаты вершин  $\alpha, \beta, \gamma$ . Погрешность интерполяции для функции  $f(x, y)$  с непрерывными вторыми производными будет  $O(h^2)$ , где  $h$  — длина наибольшей стороны треугольника.

Существуют способы построения кусочно многочленной интерполяции более высоких порядков, а также кусочно многочленной гладкой интерполяции. Аналогичные конструкции существуют и для интерполяции функций многих переменных.

Подчеркнем, что объем таблиц, обеспечивающих возможность восстановления функции данной гладкости с заданной точностью, быстро возрастает с ростом числа аргументов, а алгоритмы усложняются.

### Задачи

1. Пусть функция  $z = f(x, y)$  имеет вторые производные, ограниченные по модулю единицей.

Как экономно выбрать точки сетки на плоскости, чтобы по значениям функции  $z = f(x, y)$  в выбранных точках можно было восстановить функцию в любой точке квадрата  $|x| \leq 1, |y| \leq 1$  с погрешностью, не превосходящей  $\varepsilon = 10^{-3}$ ?

2. Решить предыдущую задачу при дополнительном предположении, что окружности  $x^2 + y^2 = r^2$  при каждом  $r$  являются линиями уровня функции  $z = f(x, y)$ .

3\*. Решить задачу 1 в дополнительном предположении, что функция  $z = f(x, y)$  есть функция вида  $z = \varphi(x)\psi(y)$ , где  $\varphi(x), \psi(y)$  — некоторые функции одного аргумента.

## ГЛАВА 2

## ТРИГОНОМЕТРИЧЕСКАЯ ИНТЕРПОЛЯЦИЯ

Наряду с алгебраической интерполяцией, изложенной в гл. 1, широко применяется интерполяция с помощью тригонометрических многочленов вида

$$Q\left(\cos \frac{2\pi}{L}x, \sin \frac{2\pi}{L}x\right) = \sum_{k=0}^n a_k \cos \frac{2\pi k}{L}x + \sum_{k=1}^n b_k \sin \frac{2\pi k}{L}x. \quad (1)$$

Здесь  $n$  — натуральные числа,  $L$  — положительное число,  $a_k, b_k$  — вещественные коэффициенты.

Тригонометрический многочлен  $Q\left(\cos \frac{2\pi}{L}x, \sin \frac{2\pi}{L}x, f\right)$ , совпадающий с периодической функцией  $f(x)$ ,  $f(x+L) = f(x)$ , в узлах интерполяции  $x_m = \frac{L}{N}m + x_0$  ( $m = 0, 1, \dots, N - 1$ ), можно выбрать так, чтобы он обладал определенными преимуществами перед алгебраическим интерполяционным многочленом, построенным по значениям функции в узлах

$$x_m = \frac{L}{N}m + x_0, \quad m = 0, 1, \dots, N - 1.$$

Во-первых, погрешность тригонометрической интерполяции

$$R_N(x, f) = f(x) - Q\left(\cos \frac{2\pi}{L}x, \sin \frac{2\pi}{L}x, f\right)$$

равномерно стремится к нулю при  $N \rightarrow \infty$ , если  $f(x)$  имеет хотя бы вторую производную, причем скорость убывания погрешности автоматически учитывает гладкость  $f(x)$ , т. е. возрастает с ростом числа  $(r+1)$  производных. Именно, мы докажем, что

$$\max_x |R_N(x)| = O\left(\frac{M_{r+1}}{N^{r+1}}\right), \quad M_{r+1} = \max_x \left| \frac{d^{r+1}f(x)}{dx^{r+1}} \right|.$$

Во-вторых, чувствительность тригонометрического интерполяционного многочлена к погрешностям в задании значений  $f_m$  в узлах «почти» не возрастает с ростом числа узлов.

Эти два положительных свойства тригонометрической интерполяции (возрастание точности при увеличении гладкости и вычислительную устойчивость) можно придать, как будет показано, и алгебраической интерполяции функций на отрезке за счет специального выбора узлов интерполяции и использования алгебраических многочленов Чебышёва, обладающих многими замечательными свойствами.

При первом чтении книги можно ограничиться приведенной справкой о содержании гл. 2 и сразу перейти к гл. 3.

## § 1. Интерполяция периодических функций

Пусть  $f(x)$  — периодическая с некоторым периодом  $L$  функция:

$$f(x + L) \underset{x}{\equiv} f(x), \quad (1)$$

заданная на сетке

$$x_m = \frac{L}{N} m + x_0, \quad m = 0, \pm 1, \dots, \quad (2)$$

где  $N$  — некоторое натуральное число. Для краткости будем обозначать  $f(x_m)$  через  $f_m$ :

$$f_{m+N} \underset{m}{\equiv} f_m. \quad (3)$$

### 1. Важный случай выбора узлов интерполяции и соответствующий тригонометрический интерполяционный многочлен.

Теорема 1. Пусть  $x_0 = L/(2N)$ ,  $N = 2(n+1)$ ,  $n$  натуральное. При произвольном задании значений  $f_m$  периодической с периодом  $L$  функции в узлах сетки (2) существует один и только один интерполяционный тригонометрический многочлен

$$Q_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right) = \sum_{k=0}^n a_k \cos \frac{2\pi k}{L} x + \sum_{k=1}^{n+1} b_k \sin \frac{2\pi k}{L} x, \quad (4)$$

удовлетворяющий равенствам

$$Q_n|_{x=x_m} = f_m, \quad m = 0, \pm 1, \dots. \quad (5)$$

Коэффициенты этого многочлена задаются формулами

$$a_0 = \frac{1}{N} \sum_{m=0}^{N-1} f_m, \quad (6)$$

$$a_k = \frac{2}{N} \sum_{m=0}^{N-1} f_m \cos k \left( \frac{2\pi}{N} m + \frac{\pi}{N} \right), \quad k = 1, 2, \dots, n, \quad (7)$$

$$b_k = \frac{2}{N} \sum_{m=0}^{N-1} f_m \sin k \left( \frac{2\pi}{N} m + \frac{\pi}{N} \right), \quad k = 1, 2, \dots, n, \quad (8)$$

$$b_{n+1} = \frac{1}{N} \sum_{m=0}^{N-1} f_m (-1)^m. \quad (9)$$

**Доказательство.** Рассмотрим множество всех вещественных периодических функций  $f_m$ :

$$f_{m+N} = f_m, \quad m = 0, \pm 1, \dots, \quad (10)$$

определенных в точках  $x_m = -\frac{L}{N} m + \frac{L}{2N}$ . Будем рассматривать эти функции только при  $m = 0, 1, \dots, N-1$ , поскольку при остальных  $m$  значения этих функций однозначно восстанавливаются по указанным значениям в силу периодичности (10).

Совокупность этих функций с обычными операциями сложения и умножения на вещественные числа образует линейное пространство. Размерность этого пространства, которое обозначим  $F_N$ , есть  $N$ , поскольку система функций

$$\tilde{\psi}_m^{(k)} = \begin{cases} 0, & \text{если } m \neq k - 1, \\ 1, & \text{если } m = k - 1, \end{cases} \quad k = 1, 2, \dots, N,$$

образует базис. Действительно, каждую функцию  $f \in F_N$ ,  $f = \{f_m\}$  ( $m = 0, 1, \dots, N - 1$ ), можно, и притом единственным образом, представить в виде линейной комбинации функций  $\tilde{\psi}_m^{(k)}$ . Введем в  $F_N$  скалярное умножение

$$(f, g) = \frac{1}{N} \sum_{m=0}^{N-1} f_m g_m.$$

Покажем теперь, что система функций  $\{\xi^{(k)}, \eta^{(k)}\}$ , задаваемая равенствами

$$\xi_m^{(0)} = \cos 0 \cdot x_m \equiv 1, \quad (11)$$

$$\xi_m^{(k)} = \sqrt{2} \cos \left( \frac{2k\pi}{L} x_m \right), \quad k = 1, 2, \dots, n, \quad (12)$$

$$\eta_m^{(k)} = \sqrt{2} \sin \left( \frac{2k\pi}{L} x_m \right), \quad k = 1, 2, \dots, n, \quad (13)$$

$$\eta_m^{(n+1)} = \sin \left( \frac{2(n+1)\pi}{L} x_m \right) = (-1)^m \quad (14)$$

образует ортонормированный базис пространства  $F_N$ .

Общее число функций (11)–(14) равно размерности  $N$  пространства  $F_N$ . Поэтому остается доказать равенства

$$(\xi^{(k)}, \xi^{(k)}) = 1, \quad k = 0, 1, \dots, n, \quad (15)$$

$$(\eta^{(k)}, \eta^{(k)}) = 1, \quad k = 1, 2, \dots, n+1, \quad (16)$$

$$(\xi^{(r)}, \xi^{(s)}) = 0, \quad r \neq s; \quad r, s = 0, 1, \dots, n, \quad (17)$$

$$(\eta^{(r)}, \eta^{(s)}) = 0, \quad r \neq s; \quad r, s = 1, 2, \dots, n+1, \quad (18)$$

$$(\xi^{(r)}, \eta^{(s)}) = 0, \quad r = 0, 1, \dots, n; \quad s = 1, 2, \dots, n+1. \quad (19)$$

Для доказательства равенств (15)–(19) предварительно заметим, что при любых  $N$  и  $\gamma$

$$\frac{1}{N} \sum_{m=0}^{N-1} 1 = 1, \quad (20)$$

$$\sum_{m=0}^{N-1} \cos l \left( \frac{2\pi}{N} m + \gamma \right) = 0, \quad l = 1, 2, \dots, N-1, \quad (21)$$

$$\sum_{m=0}^{N-1} \sin l \left( \frac{2\pi}{N} m + \gamma \right) = 0, \quad l = 1, 2, \dots, N-1. \quad (22)$$

Действительно, справедливость (20) очевидна. Далее проверим (21):

$$\begin{aligned} \sum_{m=0}^{N-1} \cos l \left( \frac{2m\pi}{N} + \gamma \right) &= \\ &= \frac{1}{2} \sum_{m=0}^{N-1} \left[ \exp \left\{ i \left( \frac{2lm\pi}{N} + l\gamma \right) \right\} + \exp \left\{ -i \left( \frac{2lm\pi}{N} + l\gamma \right) \right\} \right] = \\ &= \frac{1}{2} e^{il\gamma} \sum_{m=0}^{N-1} \left( \exp \left\{ i \frac{2l\pi}{N} \right\} \right)^m + \frac{1}{2} e^{-il\gamma} \sum_{m=0}^{N-1} \left( \exp \left\{ -i \frac{2l\pi}{N} \right\} \right)^m = \\ &= \frac{1}{2} e^{il\gamma} \frac{1 - e^{i \cdot 2l\pi}}{1 - \exp \{ i \cdot 2l\pi/N \}} + \frac{1}{2} e^{-il\gamma} \frac{1 - e^{-i \cdot 2l\pi}}{1 - \exp \{ -i \cdot 2l\pi/N \}} = 0 + 0 = 0, \\ l &= 1, 2, \dots, N-1. \end{aligned}$$

Равенство (22) доказывается аналогично.

Теперь установим справедливость (15)–(19). Равенство (15) при  $k = 0$  и равенство (16) при  $k = n + 1$  совпадают с (20). Если  $k = 1, 2, \dots, n$ , то (15), (16) справедливы в силу (20), (21):

$$\begin{aligned} (\xi^{(k)}, \xi^{(k)}) &= \\ &= \frac{2}{N} \sum_{m=0}^{N-1} \cos^2 \left( \frac{2k\pi}{N} m + \frac{k\pi}{N} \right) = \frac{1}{N} \sum_{m=0}^{N-1} \left[ 1 + \cos \left( \frac{4k\pi}{N} m + \frac{2k\pi}{N} \right) \right] = \\ &= \frac{1}{N} \sum_{m=0}^{N-1} 1 + \sum_{m=0}^{N-1} \cos 2k \left( \frac{2\pi}{N} m + \frac{\pi}{N} \right) = 1 + 0 = 1, \\ (\eta^{(k)}, \eta^{(k)}) &= \\ &= \frac{2}{N} \sum_{m=0}^{N-1} \sin^2 \left( \frac{2k\pi}{N} m + \frac{\pi}{N} \right) = \frac{1}{N} \sum_{m=0}^{N-1} \left[ 1 - \cos \left( \frac{4k\pi}{N} m + \frac{2\pi}{N} \right) \right] = 1. \end{aligned}$$

Докажем (17):

$$\begin{aligned} (\xi^{(r)}, \xi^{(s)}) &= \frac{2}{N} \sum_{m=0}^{N-1} \cos r \left( \frac{2\pi m}{N} + \frac{\pi}{N} \right) \cos s \left( \frac{2\pi m}{N} + \frac{\pi}{N} \right) = \\ &= \frac{1}{N} \sum_{m=0}^{N-1} \left[ \cos(r+s) \left( \frac{2\pi m}{N} + \frac{\pi}{N} \right) + \cos(r-s) \left( \frac{2\pi m}{N} + \frac{\pi}{N} \right) \right] = 0 \end{aligned}$$

(мы воспользовались (21), учитывая, что  $1 \leq |r \pm s| \leq N-1$ ).

Равенство (18) доказывается аналогично, но вместо тождества  $2 \cos \alpha \cos \beta = \cos(\alpha + \beta) - \cos(\alpha - \beta)$ , использованного при доказательстве (17), нужно использовать тождество  $2 \sin \alpha \sin \beta = \cos(\alpha - \beta) - \cos(\alpha + \beta)$ . Для доказательства (19) используется тождество  $2 \sin \alpha \cos \beta = \sin(\alpha + \beta) + \sin(\alpha - \beta)$ .

Итак, установлено, что (11)–(14) — ортонормированный базис пространства  $F_N$ . Поэтому каждая функция  $f = \{f_m\} \in F_N$  может быть представлена в виде

$$f_m = \sum_{k=0}^n a_k \cos \frac{2\pi k}{L} x_m + \sum_{k=1}^{n+1} b_k \sin \frac{2\pi k}{L} x_m,$$

т. е. в виде линейной комбинации элементов базиса (11)–(14). Умножая сеточные функции, входящие в левую и правую части этого равенства, скалярно на  $\xi^{(r)}$  или  $\eta^{(s)}$  ( $r = 0, 1, \dots, n$ ;  $s = 1, 2, \dots, n+1$ ), получаем равенства

$$\begin{aligned} a_0 &= (f, \xi^{(0)}), \\ a_k &= \sqrt{2} (f, \xi^{(k)}), \quad k = 1, 2, \dots, n, \\ b_k &= \sqrt{2} (f, \eta^{(k)}), \quad k = 1, 2, \dots, n, \\ b_{n+1} &= (f, \eta^{(n+1)}), \end{aligned}$$

которые совпадают с формулами (6)–(9).  $\square$

Сетка  $x_m = \frac{L}{N} m + \frac{L}{2N}$  ( $m = 0, \pm 1, \dots$ ), использованная для задания сеточных функций  $f \in F_N$  в теореме 1, симметрична относительно точки  $x = 0$ , так что вместе с точкой  $x = x_m$  сетка содержит точку  $x = -x_m = x_{-(m+1)}$ . Поэтому можно говорить о четных или о нечетных сеточных функциях.

Сеточная функция  $f(x)$ ,  $x = x_m$ , четная, если  $f(-x) = f(x)$ ,  $x = x_m$ , или

$$f(x_m) = f(x_{-(m+1)}), \quad m = 0, \pm 1, \dots \quad (23)$$

Сеточная функция  $f(x)$ ,  $x = x_m$ , нечетная, если  $f(-x) = -f(x)$ ,  $x = x_m$ , или

$$f(x_m) = -f(x_{-(m+1)}), \quad m = 0, \pm 1, \dots \quad (24)$$

Теорема 2. Пусть на сетке  $x_m = \frac{L}{N} m + \frac{L}{2N}$ ,  $N = 2(n+1)$ , задана четная периодическая с периодом  $N$  сеточная функция  $f_m$ . Тогда интерполяционный многочлен (4) примет вид

$$Q_n = \sum_{k=0}^n a_k \cos \frac{2\pi k}{L} x, \quad (25)$$

где

$$a_0 = \frac{1}{n+1} \sum_{m=0}^n f_m, \quad (26)$$

$$a_k = \frac{2}{n+1} \sum_{m=0}^n f_m \cos k \left( \frac{\pi}{n+1} m + \frac{\pi}{2(n+1)} \right), \quad k = 1, 2, \dots, n. \quad (27)$$

Доказательство. В силу четности (23) из (6), (7) получаем (25), (26), а из (8), (9) следует, что  $b_k \equiv 0$ .  $\square$

**Теорема 3.** Пусть на сетке  $x_m = \frac{L}{N} m + \frac{L}{2N}$  задана нечетная сеточная функция  $f_m$ . Тогда интерполяционный многочлен (4) примет вид

$$Q_n = \sum_{k=1}^{n+1} b_k \sin \frac{2\pi k}{L} x, \quad (28)$$

где

$$b_k = \frac{2}{n+1} \sum_{m=0}^n f_m \sin \frac{2\pi k}{L} x_m, \quad k = 1, 2, \dots, n, \quad (29)$$

$$b_{n+1} = \frac{1}{n+1} \sum_{m=0}^n f_m (-1)^m. \quad (30)$$

**Доказательство.** Благодаря нечетности (24) коэффициенты  $a_k$  ( $k = 0, 1, \dots, n$ ) в силу формул (6), (7) окажутся равными нулю, формулы (8), (9) совпадут с (29), (30), а многочлен (4) примет вид (28).  $\square$

**2. Чувствительность интерполяционного многочлена к погрешностям задания функции в узлах интерполяции.** Оценим, какова чувствительность интерполяционного многочлена (4) к погрешностям в задании значений  $f_m$ . Пусть вместо  $f = \{f_m\}$  задана сеточная функция  $f + \delta f = \{f_m + \delta f_m\}$ . Тогда вместо многочлена (4) получим многочлен

$$Q_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f + \delta f \right) = Q_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right) + \delta Q_n.$$

Из формул (6)–(9) видно, что возникающая погрешность есть

$$\delta Q_n = Q_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, \delta f \right).$$

Таким образом, мерой чувствительности интерполяционного многочлена (4) к возмущению  $\delta f$  входных данных могут служить числа

$$L_n = \sup_{f \in F_N} \frac{\max_x \left| Q_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right) \right|}{\max_m |f_m|}, \quad n = 1, 2, \dots, \quad (31)$$

называемые *константами Лебега*. Итак,

$$\max_x |\delta Q_n| \leq L_n \max_m |f_m|.$$

**Теорема 4.** Константы Лебега  $L_n$  тригонометрических интерполяционных многочленов (4) удовлетворяют оценкам

$$L_n \leq 2n. \quad (32)$$

**Доказательство.** Из формул (6), (9) следует, что

$$|a_k| \leq 2 \max_m |f_m|, \quad (33)$$

$$|b_k| \leq 2 \max_m |f_m|. \quad (34)$$

Из (4), (33) и (34) следует, что

$$\left| Q_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right) \right| \leq \sum_{k=0}^n |a_k| + \sum_{k=1}^{n+1} |b_k| \leq 4(n+1) \max_m |f_m|. \quad (35)$$

Оценка (35) выполнена для любой  $f \in F_N$  и любого  $N$ . Отсюда следует (32).  $\square$

### 3. Оценка погрешности интерполяции.

Теорема 5. Пусть  $f(x)$  — периодическая с некоторым периодом  $L$  функция, имеющая непрерывную производную некоторого порядка  $r+1$ :

$$\max |f^{(r+1)}(x)| = M_{r+1}. \quad (36)$$

Пусть  $f \in F_N$ ,  $f = \{f_m\}$  — таблица значений этой функции в точках сетки

$$x_m = \frac{L}{N} m + \frac{L}{2N}, \quad m = 0, \pm 1, \dots, \quad N = 2(n+1),$$

а  $Q_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right)$  — соответствующий интерполяционный многочлен (4). Тогда для погрешности интерполяции

$$R_{n+1}(x) = f(x) - Q_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right)$$

справедлива оценка

$$|R_{n+1}(x)| \leq \frac{CM_{r+1}}{n^{r+1}}, \quad C = C(L) = \text{const.} \quad (37)$$

**Доказательство.** Проведем предварительные построения. Представим  $f(x)$  в виде суммы ряда Фурье

$$f(x) = S_{n+1}(x) + \delta S_{n+1}(x),$$

где

$$S_{n+1}(x) = \sum_{k=0}^n \alpha_k \cos \frac{2\pi k}{L} x + \sum_{k=1}^{n+1} \beta_k \sin \frac{2\pi k}{L} x,$$

$$\delta S_{n+1} = \sum_{k=n+1}^{\infty} \alpha_k \cos \frac{2\pi k}{L} x + \sum_{k=n+2}^{\infty} \beta_k \sin \frac{2\pi k}{L} x,$$

а  $\alpha_k, \beta_k$  — коэффициенты разложения функции  $f(x)$  в ряд Фурье.

Ниже мы используем оценку

$$|\delta S_{n+1}(x)| \leq \left( \frac{L}{2\pi} \right)^{r+1} \frac{M_{r+1}}{n^r}.$$

Докажем ее. Из формул

$$\alpha_k = \frac{2}{L} \int_0^L f(x) \cos \frac{2\pi k}{L} x dx, \quad \beta_k = \frac{2}{L} \int_0^L f(x) \sin \frac{2\pi k}{L} x dx,$$

интегрируя их  $(r+1)$  раз по частям, получаем оценки

$$|\alpha_k| \leq \left(\frac{L}{2\pi}\right)^{r+1} \frac{2M_{r+1}}{k^{r+1}}, \quad |\beta_k| \leq \left(\frac{L}{2\pi}\right)^{r+1} \frac{2M_{r+1}}{k^{r+1}}.$$

Отсюда

$$\max_x |\delta S_{n+1}(x)| \leq \sum_{k=n+1}^{\infty} (|\alpha_k| + |\beta_k|) \leq 4 \left(\frac{L}{2\pi}\right)^{r+1} \frac{M_{r+1}}{rn^r} = A \frac{M_{r+1}}{n^r}, \quad (38).$$

$A = \text{const.}$

Заметим, что  $S_{n+1}(x)$  есть тригонометрический многочлен вида (4). В силу единственности интерполяционного тригонометрического многочлена  $Q_n\left(\cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, S_{n+1}\right)$  этот многочлен совпадает с  $S_{n+1}(x)$ :

$$Q_n\left(\cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, S_{n+1}\right) = S_{n+1}(x). \quad (39)$$

Далее, в силу оценок (32), (38)

$$\left|Q_n\left(\cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, \delta S_{n+1}\right)\right| \leq L_n A \frac{M_{r+1}}{n^r} \leq 4A \frac{M_{r+1}}{n^{r-1}}. \quad (40)$$

В силу (38), (40) получим требуемую оценку (37):

$$\begin{aligned} |R_{n+1}(x)| &= \left|f(x) - Q_n\left(\cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f\right)\right| = \\ &= \left|f(x) - Q_n\left(\cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, S_{n+1}\right) - \right. \\ &\quad \left. - Q_n\left(\cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, \delta S_{n+1}\right)\right| = \\ &= \left|(f(x) - S_{n+1}(x)) - Q_n\left(\cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, \delta S_{n+1}\right)\right| = \\ &= \left|\delta S_{n+1}(x) - Q_n\left(\cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, \delta S_{n+1}\right)\right| \leq \\ &\leq |\delta S_{n+1}(x)| + \left|Q_n\left(\cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, \delta S_{n+1}\right)\right| \leq \\ &\leq \frac{AM_{r+1}}{n^r} + \frac{2AM_{r+1}}{n^{r-1}} \leq \frac{3AM_{r+1}}{n^{r-1}}. \quad \square \end{aligned}$$

Можно показать, что оценки (32) констант Лебега, а также оценки погрешности интерполяции, доказанные нами, можно заменить более сильными.

**4. Еще один случай выбора узлов при тригонометрической интерполяции.** Отметим еще одну важную для приложений формулу тригонометрической интерполяции, посвятив ей следующую теорему.

**Теорема 6.** Пусть  $N = 2n$ . Тогда при произвольном задании значений  $f_m$  периодической с периодом  $L$  функции в узлах сетки

$$x_m = \frac{L}{N} m, \quad m = 0, \pm 1, \dots,$$

существует один и только один интерполяционный тригонометрический многочлен

$$\tilde{Q}_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right) = \sum_{k=0}^n \tilde{a}_k \cos \frac{2\pi k}{L} x + \sum_{k=1}^{n-1} \tilde{b}_k \sin \frac{2\pi k}{L} x, \quad (41)$$

удовлетворяющий равенствам

$$\tilde{Q}_n \Big|_{x=x_m} = f_m, \quad m = 0, \pm 1, \dots$$

Коэффициенты этого многочлена задаются формулами

$$\begin{aligned} \tilde{a}_0 &= \frac{1}{N} \sum_{m=0}^{N-1} f_m, \quad \tilde{a}_n = \frac{1}{N} \sum_{m=0}^{N-1} f_m (-1)^m, \\ \tilde{a}_k &= \frac{2}{N} \sum_{m=0}^{N-1} f_m \cos \frac{2\pi km}{N}, \quad k = 1, 2, \dots, n-1, \\ \tilde{b}_k &= \frac{2}{N} \sum_{m=0}^{N-1} f_m \sin \frac{2\pi km}{N}, \quad k = 1, 2, \dots, n-1. \end{aligned} \quad (42)$$

Доказательство аналогично доказательству теоремы 1, и мы его опускаем.

Отметим, что в случае четной сеточной функции  $f(x)$ ,  $f(-x_m) = f(x_m)$ , или  $f_m = f_{-m}$ , формулы (42) принимают вид

$$\begin{aligned} \tilde{a}_0 &= \frac{1}{N} (f_0 + f_n) + \frac{2}{N} \sum_{m=1}^{n-1} f_m, \quad \tilde{a}_n = \frac{1}{N} [f_0 + (-1)^n f_n], \\ \tilde{a}_k &= \frac{2}{N} [f_0 + (-1)^k f_n] + \frac{4}{N} \sum_{m=1}^{n-1} f_m \cos \frac{2\pi km}{N}, \quad k = 1, 2, \dots, n-1, \\ \tilde{b}_k &\equiv 0, \end{aligned} \quad (43)$$

а многочлен (41) принимает вид

$$\tilde{Q}_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right) = \sum_{k=0}^n \tilde{a}_k \cos \frac{2\pi km}{L}.$$

Интерполяционный многочлен  $\tilde{Q}_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right)$ , задаваемый формулами (41), (42), обладает свойством устойчивости относительно возмущения значений  $f_m$  и свойством сходимости при  $n \rightarrow \infty$  к интерполируемой функции со скоростью, реагирующей на ее гладкость. Эти свойства аналогичны установленным в теоремах 4, 5 для многочлена  $Q_n \left( \cos \frac{2\pi}{L} x, \sin \frac{2\pi}{L} x, f \right)$ , задаваемого формулой (4).

## § 2. Интерполяция функций на отрезке. Связь между алгебраической и тригонометрической интерполяциями

Пусть  $f(x)$  определена на отрезке  $-1 \leq x \leq 1$  и имеет на этом отрезке ограниченную производную некоторого порядка  $r+1$ .

Мы считаем областью определения функции  $f(x)$  отрезок  $-1 \leq x \leq 1$ , а не произвольный отрезок  $a \leq x \leq b$ , лишь для удобства. Действительно, преобразование  $x = [t(b-a) + a + b]/2$  позволяет перейти от функции  $f(x)$ , определенной на произвольном отрезке  $a \leq x \leq b$ , к функции  $F(t) \equiv f([t(b-a) + b + a]/2)$ , определенной на отрезке  $-1 \leq t \leq 1$ .

**1. Периодизация.** В силу теоремы 5 из § 1 тригонометрическая интерполяция непосредственно пригодна и эффективна для восстановления лишь гладких периодических функций по их таблицам. Поэтому

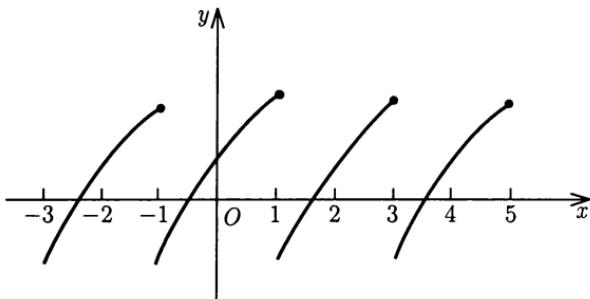


Рис. 3

использование тригонометрической интерполяции для приближенного описания функции  $f(x)$  ( $-1 \leq x \leq 1$ ) требует перехода от функции  $f(x)$  к некоторой гладкой периодической функции. Простое доопределение функции  $f(x)$  ( $-1 < x \leq 1$ ) на всей числовой оси до периодической функции с периодом  $L = 2$  приводит, вообще говоря, к разрывной функции (рис. 3). Поэтому перейдем от  $f(x)$  к функции

$$F(\varphi) \equiv f(\cos \varphi), \quad x = \cos \varphi. \quad (1)$$

Для наглядности будем считать, что функция  $F(\varphi)$  определена на единичной окружности как функция полярного угла  $\varphi$ . Значение  $F(\varphi)$  получается переносом значения  $f(x)$  из точки  $x \in [-1, 1]$  в точку  $\varphi$  на единичной окружности (рис. 4).

Очевидно, что функция  $F(\varphi)$  является четной периодической функцией с периодом  $2\pi$ . Легко видеть также, что существует и ограничена производная  $\frac{d^{r+1}}{d\varphi^{r+1}} F(\varphi)$ .

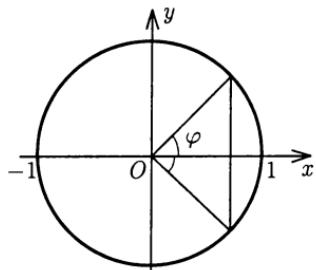


Рис. 4

**2. Тригонометрическая интерполяция.** Выберем в качестве узлов интерполяции точки

$$\varphi_m = \frac{2\pi}{N}m + \frac{\pi}{N}, \quad m = 0, \pm 1, \dots, \pm m; \quad N = 2(n+1). \quad (2)$$

Значения  $F_m = F(\varphi_m)$  функции  $F(\varphi)$  в узлах  $\varphi_m$  в силу определения функции  $F'(\varphi)$  совпадают со значениями  $f_m = f(x_m)$  исходной функции  $f(x)$  в точках  $x_m = \cos \varphi_m$ .

Для интерполяции четной функции  $F(\varphi)$  по ее значениям в узлах воспользуемся интерполяционной формулой (25) из § 1:

$$Q_n(\cos \varphi, \sin \varphi, F) = \sum_{k=0}^n a_k \cos k\varphi, \quad (3)$$

где коэффициенты  $a_k$  в силу  $F_m = f_m$  определяются формулами (26), (27) из § 1:

$$\begin{aligned} a_0 &= \frac{1}{n+1} \sum_{m=0}^n f_m, \\ a_k &= \frac{2}{n+1} \sum_{m=0}^n f_m \cos k\varphi_m, \quad k = 1, 2, \dots, n. \end{aligned} \quad (4)$$

**3. Многочлены Чебышёва. Связь между тригонометрической и алгебраической интерполяциями.** Воспользуемся равенством  $\cos \varphi = x$  и введем функции

$$T_k(x) = \cos k\varphi = \cos k \arccos x, \quad k = 0, 1, \dots \quad (5)$$

**Теорема 1.** Функции  $T_k(x)$  суть многочлены степеней  $k = 0, 1, \dots$ . При этом  $T_0(x) = 1$ ,  $T_1(x) = x$ , а  $T_2(x), T_3(x), \dots$  последовательно вычисляются по рекуррентной формуле

$$T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x) \quad k = 1, 2, \dots \quad (6)$$

**Доказательство.** Очевидно, что  $T_0(x) = \cos 0 = 1$ ,  $T_1(x) = \cos \arccos x = x$ . Воспользуемся известным тригонометрическим тождеством

$$\cos(k+1)\varphi = 2\cos \varphi \cos k\varphi - \cos(k-1)\varphi, \quad k = 1, 2, \dots,$$

которое в случае  $\varphi = \arccos x$  переходит в формулу (6).

Докажем, что  $T_k(x)$  — многочлен степени  $k$ , с помощью индукции по  $k$ .

При  $k = 0, k = 1$  это уже доказано. Фиксируем  $k \geq 1$ . Допустим, что утверждение уже доказано для  $T_j(x)$ ,  $j = 0, 1, \dots, k$ , т. е. установлено, что  $T_j(x)$  — многочлены степени  $j$  ( $j = 0, 1, \dots, k$ ). Тогда выражение в правой части формулы (6), а следовательно, и  $T_{k+1}(x)$  есть многочлен степени  $k+1$ .  $\square$

Многочлены  $T_k(x)$  были введены П.Л. Чебышёвым. Приведем несколько первых многочленов Чебышёва и их графики (рис. 5):

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_2(x) = 2x^2 - 1, \quad T_3(x) = 4x^3 - 3x.$$

Воспользуемся теперь равенством  $\varphi = \arccos x$  и перейдем в (3) от  $\varphi$  к  $x$ , положив

$$Q_n(\cos \varphi, \sin \varphi, F) \equiv P_n(x, f),$$

где

$$P_n(x, f) = \sum_{k=0}^n a_k T_k(x), \quad (7)$$

$$\begin{aligned} a_0 &= \frac{1}{n+1} \sum_{m=0}^n f_m = \frac{1}{n+1} \sum_{m=0}^n f_m T_0(x_m), \\ a_k &= \frac{2}{n+1} \sum_{m=0}^n f_m \cos k\varphi_m = \frac{2}{n+1} \sum_{m=0}^n f_m T_k(x_m). \end{aligned} \quad (8)$$

Напомним, что

$$x_m = \cos \varphi_m = \cos \frac{\pi(2m+1)}{n+1}, \quad m = 0, 1, \dots, n. \quad (9)$$

Таким образом,  $P_n(x, f)$  есть алгебраический многочлен степени не выше  $n$ , принимающий в узлах интерполяции  $x_m = \cos \varphi_m$  заданные значения  $f(x_m) = f_m$ . Узлы интерполяции  $x_m$  изображены на рис. 6.

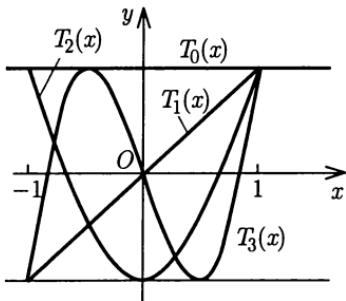


Рис. 5

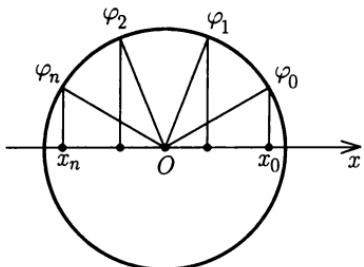


Рис. 6

Отметим, что точки

$$\varphi_m = \frac{\pi}{n+1} m + \frac{\pi}{2(n+1)}$$

( $m = 0, 1, \dots, n$ ) являются нулями функции  $\cos(n+1)\varphi$ , а  $x_m = \cos \varphi_m$  — нулями многочлена Чебышёва  $T_{n+1}(x) = \cos(n+1)\varphi$ .

Таким образом, интерполяционный многочлен  $P_n(x, f)$ , задаваемый формулой (7), реализует алгебраическую интерполяцию функции  $f(x)$  по ее значениям в точках  $x_m$ , являющихся нулями многочлена Чебышёва  $T_{n+1}(x)$ .

**4. Свойства алгебраической интерполяции с нулями многочлена Чебышёва  $T_{n+1}(x)$  в качестве узлов.** В силу равенства  $Q_n(\cos \varphi, \sin \varphi, F) = P_n(x, f)$  на алгебраический интерполяционный многочлен  $P_n(x, f)$  переносятся свойства многочлена  $Q_n(\cos \varphi, \sin \varphi, F)$ , установленные в теоремах 4, 5 из § 1. Именно, константы Лебега  $L_n$ , характеризующие чувствительность многочлена  $P_n(x, f)$  к погрешностям при задании значений  $f_m$ , удовлетворяют доказанной оценке (32) из § 1

$$L_n \leq 2n, \quad (10)$$

а погрешность интерполяции

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

равномерно стремится к нулю при  $n \rightarrow \infty$  со скоростью, зависящей от числа производных  $r + 1$ :

$$\max_{-1 \leq x \leq 1} |R_n(x)| = O\left(\frac{M_{r+1}}{n^{r+1}}\right), \quad M_{r+1} = \max_x |f^{(r+1)}(x)|. \quad (11)$$

Напомним, что при использовании вместо сетки  $x_m = \cos \varphi_m$  сетки с постоянным шагом константы Лебега быстро растут, а сходимость интерполяционного многочлена к функции  $f(x)$  может не иметь места даже для бесконечно дифференцируемых функций. Это и заставляет в случае равномерной или произвольной нерегулярной сетки использовать кусочно многочленную алгебраическую интерполяцию или сплайны (см. гл. 1).

**Замечание.** Доказанная нами в теореме 4 из § 1 оценка (32) может быть существенно усиlena: справедлива оценка

$$L_n = \frac{2}{\pi} \ln n + 1 - \theta_n, \quad 0 \leq \theta_n \leq \frac{1}{4}$$

(см., напр., [1]). Вместе с тем может быть усилено и (11).

**5. Алгоритм вычисления значений интерполяционного многочлена.** Для вычисления коэффициентов многочлена (7) по формулам (8), а также для вычисления значения этого многочлена в заданной точке  $x$  ( $-1 \leq x \leq 1$ ) нужно уметь вычислять значения многочленов  $T_k(x)$  ( $k = 0, 1, \dots$ ). Покажем, что эти значения целесообразно вычислять по формуле (6). Простота и экономичность формулы (6) очевидны. Покажем, что вычисления по ней устойчивы относительно погрешностей округления.

Рассмотрим так называемое разностное уравнение вида

$$y_{k+1} = 2xy_k - y_{k-1}, \quad (12)$$

где  $k$  — параметр. Будем искать решение этого уравнения, имеющее вид  $y_k = q^k$ , где  $q$  — некоторое число. Подставляя это выражение

в разностное уравнение, получаем следующее уравнение (характеристическое уравнение) для  $q$ :

$$q^2 - 2xq + 1 = 0, \quad q_{1,2} = x \pm \sqrt{x^2 - 1}. \quad (13)$$

В силу линейности уравнения (12) выражение

$$y_k = c_1 q_1^k + c_2 q_2^k \quad (14)$$

является его решением при любых значениях постоянных  $c_1, c_2$ . Подберем  $c_1, c_2$  из условий

$$y_0 = T_0(x) = 1, \quad y_1 = T_1(x) = x,$$

т. е.

$$c_1 + c_2 = 1, \quad c_1 q_1 + c_2 q_2 = x. \quad (15)$$

Отсюда: при  $c_1 = c_2 = 1/2$  решение (14) совпадает с  $T_k(x)$  и в силу соотношения (6) является решением уравнения (12) как функция  $k$ :

$$T_k(x) = \frac{1}{2} (x + \sqrt{x^2 - 1})^k + \frac{1}{2} (x - \sqrt{x^2 - 1})^k. \quad (16)$$

Заметим, что корни  $q_{1,2} = x \pm \sqrt{x^2 - 1}$  при  $|x| < 1$  являются комплексно-сопряженными и по модулю равны единице. Следовательно,  $q_1^k, q_2^k$  с ростом  $k$  по модулю не меняются и остаются равными единице по модулю.

Какая-либо погрешность, допущенная при некотором  $k = k_0$ , вызовет возмущение в значениях  $c_1, c_2$ , которые входят в формулу (14) при  $k > k_0$ . В силу равенств  $|q_1^k| = |q_2^k| = 1$  ( $k = 1, 2, \dots$ ) эта погрешность с ростом  $k$  не будет возрастать, что и означает вычислительную устойчивость расчетов по формуле (6) при  $|x| < 1$ .

**6. Алгебраическая интерполяция с узлами в точках экстремума многочлена Чебышёва  $T_n(x)$ .** При интерполяции функции  $F(\varphi) = f(\cos \varphi)$  воспользуемся сеткой  $\tilde{\varphi}_m = \frac{\pi}{n} m$  ( $m = 0, 1, \dots, 2n - 1$ ). В соответствии с теоремой 6 из § 1 получим тригонометрический интерполяционный многочлен

$$\tilde{Q}_n(\cos \varphi, \sin \varphi, F) = \sum_{k=0}^n \tilde{a}_k \cos k\varphi,$$

$$\tilde{a}_0 = \frac{1}{2n} (f_0 + f_n) + \frac{1}{n} \sum_{m=1}^{n-1} f_m, \quad \tilde{a}_n = \frac{1}{2n} [f_0 + (-1)^n f_n],$$

$$\tilde{a}_k = \frac{1}{n} [f_0 + (-1)^k f_n] + \frac{2}{n} \sum_{m=1}^{n-1} f_m \cos k \cdot \varphi_m,$$

$$k = 1, 2, \dots, n - 1.$$

Переходя к переменной  $x = \cos \varphi$  и обозначая  $\tilde{Q}_n(\cos \varphi, \sin \varphi, F) = \tilde{P}_n(x, f)$ , получаем

$$\tilde{P}_n(x, f) = \sum_{k=0}^n \tilde{a}_k T_k(x), \quad (17)$$

$$\begin{aligned} \tilde{a}_0 &= \frac{1}{2n}(f_0 + f_n) + \frac{1}{n} \sum_{m=1}^{n-1} f_m, \quad \tilde{a}_n = \frac{1}{2n}[f_0 + (-1)^n f_n], \\ \tilde{a}_k &= \frac{1}{n}[f_0 + (-1)^k f_n] + \frac{2}{n} \sum_{m=1}^{n-1} f_m T_k(x_m). \end{aligned} \quad (18)$$

Алгебраический интерполяционный многочлен  $\tilde{P}_n(x, f)$  с узлами интерполяции

$$\tilde{x}_m = \cos \varphi_m = \cos \frac{\pi}{n} m, \quad m = 0, 1, \dots, n,$$

наследует от тригонометрического интерполяционного многочлена  $\tilde{Q}_n(\cos \varphi, \sin \varphi, F)$  слабый рост констант Лебега с увеличением  $n$  (устойчивость относительно погрешности значений  $f_m$ ), а также увеличение порядка малости погрешности интерполяции относительно числа  $1/n$  при возрастании гладкости функции  $f(x)$ .

Обратим внимание читателя на тот очевидный факт, что в узлах интерполяции  $\tilde{x}_m$  многочлен Чебышёва  $T_n(x)$  достигает своих экстремальных на отрезке  $-1 \leq x \leq 1$  значений:  $T_n(\tilde{x}_m) = \cos \pi m = (-1)^m$ ,  $m = 0, 1, \dots, n$ , что и оправдывает заголовок этого пункта.

### Задачи

**1.** Пусть функция  $f(x)$  задана не на отрезке  $[-1, 1]$ , а на произвольном отрезке  $[a, b]$ .

Указать узлы интерполяции и выписать интерполяционные многочлены, аналогичные многочленам  $P_n(x, f)$ ,  $\tilde{P}_n(x, f)$ , построенным в этом параграфе.

**2.** Для функции  $f(x)$ ,  $-1 \leq x \leq 1$ , построить интерполяционный многочлен  $P_n(x)$ , используя в качестве узлов интерполяции нули многочлена Чебышёва  $T_{n+1}(x)$ . Вывести на экран компьютера графики  $f(x)$  и  $P_n(x)$  для  $n = 5, 10, 20, 30$ . Сделать то же самое в случае интерполяционного многочлена  $P_n(x)$ , построенного по равноотстоящим узлам  $x_k = -1 + 2k/n$ ,  $k = 0, 1, \dots, n$ . Рассмотреть случай  $f(x) = \frac{1}{x^2 + 0,25}$ . Объяснить качественное различие между обоими способами интерполяции, наблюдаемое на экране.

**3.** Введем нормированный многочлен Чебышёва  $T_n(x)$  степени  $n$ , положив  $T_n(x) = 2^{1-n} T_n(x)$ .

а) Показать, что коэффициент при  $x^n$  для  $T_n(x)$  равен единице.

б) Показать, что уклонение  $\max_{-1 \leq x \leq 1} |T_n(x)|$  многочлена  $T_n(x)$  от нуля на отрезке  $[-1, 1]$  есть  $2^{1-n}$ .

в)\* Показать, что среди всех многочленов степени  $n$  с коэффициентом 1 при  $x^n$  многочлен  $T_n(x)$  наименее уклоняется от нуля на отрезке  $[-1, 1]$ .

г) Как выбрать узлы интерполяции  $t_0, t_1, \dots, t_n$  на отрезке  $[-1, 1]$ , чтобы многочлен  $(t - t_0)(t - t_1)\dots(t - t_n)$ , входящий в формулу погрешности (2) из § 2 гл. 1, наименее уклонялся от нуля на отрезке  $[-1, 1]$ ?

4. Указать узлы интерполяции четной периодической функции  $f(\varphi)$ ,  $f(-\varphi) = f(\varphi)$ ,  $f(\varphi + 2\pi) \equiv f(\varphi)$ , для которых константы Лебега совпадают с константами Лебега для алгебраической интерполяции по равноотстоящим узлам.

## ГЛАВА 3

### ВЫЧИСЛЕНИЕ ОПРЕДЕЛЕННЫХ ИНТЕГРАЛОВ. КВАДРАТУРЫ

Определенный интеграл  $\int_a^b f(x) dx$  удается найти точно по формуле

Ньютона–Лейбница лишь в тех редких случаях, когда соответствующий неопределенный интеграл от заданной функции  $f(x)$  является табличным. Еще реже удается находить точные значения кратных интегралов  $\iint_D f(x, y) dx dy$ , где  $D$  — заданная область.

Мы укажем некоторые способы приближенного вычисления определенных и кратных интегралов, обсудим трудности, связанные с ростом размерности кратного интеграла.

Отдельный параграф посвящен комбинированному использованию аналитических и вычислительных средств при решении задач, которое широко применяется не только при вычислении интегралов.

Формулы для приближенного вычисления интеграла по таблице значений подынтегральной функции называют *квадратурными* в случае интегралов по отрезку и *кубатурными* в случае кратных интегролов.

#### § 1. Квадратурные формулы трапеций и Симпсона

Определенный интеграл  $\int_a^b f(x) dx$ , как известно, имеет геометрический смысл площади.

**1. Общая схема простейших квадратурных формул.** Для приближенного вычисления этой площади разобьем отрезок  $[a, b]$  на некоторое число  $n$  отрезков, обозначив абсциссы их концов  $x_0, x_1, \dots, x_n$ , причем

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

Затем воспользуемся кусочно многочленной интерполяцией какой-либо степени  $k \geq 1$ , заменив функцию  $y = f(x)$  возникающей при этом кусочно многочленной функцией

$$P_k(x) = P_k(x, f, x_0, x_1, \dots, x_n).$$

Функция  $P_k(x)$  на каждом отрезке  $[x_i, x_{i+1}]$  есть многочлен степени не выше  $k$ .

Для вычисления интеграла можно воспользоваться приближенным равенством

$$\int_a^b f(x) dx \approx \int_a^b P_k(x, f, x_0, x_1, \dots, x_n) dx. \quad (1)$$

Правую часть этого равенства можно записать в виде

$$\int_a^b P_k dx = \int_a^{x_1} P_k dx + \int_{x_1}^{x_2} P_k dx + \dots + \int_{x_{n-1}}^b P_k dx. \quad (2)$$

Каждое слагаемое в правой части — интеграл от многочлена степени не выше  $k$ , с которым совпадает функция  $P_k = P_k(x, f, x_0, x_1, \dots, x_n)$  на отрезке  $[x_i, x_{i+1}]$ . Его можно вычислить точно по формуле Ньютона–Лейбница, выразив ответ через числа  $x_0, x_1, \dots, x_n$  и  $f(x_0), f(x_1), \dots, f(x_n)$ . Формула (1) является, таким образом, квадратурной формулой.

Очевидно, справедлива следующая оценка погрешности квадратурной формулы:

$$\begin{aligned} \left| \int_a^b f(x) dx - \int_a^b P_k dx \right| &\leq \int_a^b |f(x) - P_k| dx \leq \\ &\leq (b-a) \max_{a \leq x \leq b} |f(x) - P_k(x, f, x_0, x_1, \dots, x_n)|. \end{aligned} \quad (3)$$

**Теорема 1.** Пусть  $f(x)$  имеет ограниченную производную порядка  $k+1$ , так что  $\max_{a \leq x \leq b} |f^{(k+1)}(x)| = M_{k+1}$ , и пусть  $h = \max_{0 \leq i \leq n-1} |x_{i+1} - x_i|$ . Тогда

$$\left| \int_a^b f(x) dx - \int_a^b P_k dx \right| \leq \text{const} \cdot |b-a|M_{k+1}h^{k+1}. \quad (4)$$

**Доказательство.** Достаточно воспользоваться следующей оценкой погрешности кусочно многочленной интерполяции:

$$\max_{a \leq x \leq b} |f(x) - P_k| \leq \text{const} \cdot M_{k+1}h^{k+1},$$

а также оценкой (3).  $\square$

Квадратурные формулы вида (1) наиболее часто употребляются при использовании кусочно линейной интерполяции  $k=1$  (форму-

ла трапеций) и при кусочно квадратичной интерполяции (формула Симпсона).

**2. Формула трапеций.** Конкретизируем формулу (1) в случае  $k = 1$  и при равноотстоящих узлах  $x_i$ ,  $x_{i+1} - x_i = (b - a)/n = h$ . В этом случае интеграл  $\int_{x_i}^{x_{i+1}} P_1 dx$  есть площадь трапеции и равен произведению полусуммы оснований этой трапеции  $(f_i + f_{i+1})/2$  на ее высоту  $h = (b - a)/n$ , т. е.

$$\int_{x_i}^{x_{i+1}} P_1(x) dx = \frac{b - a}{n} \cdot \frac{f_i + f_{i+1}}{2}. \quad (5)$$

Суммируя эти равенства почленно для  $i = 0, 1, \dots, n - 1$ , придадим формуле (1) следующий вид:

$$\int_a^b f(x) dx \approx \int_a^b P_1(x) dx = \frac{b - a}{n} \left( \frac{f_0}{2} + f_1 + \dots + f_{n-1} + \frac{f_n}{2} \right), \quad (6)$$

называемый *формулой трапеций*.

Уточним общую оценку (4) погрешности квадратурной формулы (1) для случая формулы трапеций (6). На каждом отрезке  $[x_i, x_{i+1}]$  функция  $P_1(x)$  совпадает с интерполяционным многочленом первой степени. Поэтому в силу оценки (8) из § 2 гл. 1 на каждом отрезке  $[x_i, x_{i+1}]$  справедлива оценка

$$|f(x) - P_1(x)| \leq \frac{1}{8} \max_{x_i \leq x \leq x_{i+1}} |f''(x)|h^2 \leq \frac{1}{8} \max_{a \leq x \leq b} |f''(x)|h^2.$$

Правая часть этой оценки не зависит от того, какому из отрезков  $[x_i, x_{i+1}]$  принадлежит точка  $x$ . Поэтому выписанная оценка справедлива при всех  $x \in [a, b]$ :

$$\max_{a \leq x \leq b} |f(x) - P_1(x)| \leq \frac{1}{8} \max_{a \leq x \leq b} |f''(x)|h^2. \quad (7)$$

Оценка погрешности (3) при  $k = 1$  и при равноотстоящих узлах  $x_i$ , т. е. для формулы трапеций, принимает вид

$$\left| \int_a^b f(x) dx - \int_a^b P_1(x, f, x_0, x_1, \dots, x_n) dx \right| \leq \frac{b - a}{8} \max_{a \leq x \leq b} |f''(x)|h^2. \quad (8)$$

Таким образом, эта оценка гарантирует убывание погрешности при  $h \rightarrow 0$  не медленнее, чем  $h^2$ . Она не является грубой: даже для функций, имеющих производные всех порядков, погрешность квадратурной формулы трапеций, вообще говоря, действительно убывает именно как  $h^2$ .

Это показывает следующий пример. Очевидно, что  $\int_0^1 x^2 dx = \frac{1}{3}$ . По формуле трапеций

$$\begin{aligned} \int_0^1 x^2 dx &\approx \frac{1}{n} \left( \frac{x_0^2}{2} + x_1^2 + \dots + x_{n-1}^2 + \frac{x_n^2}{2} \right) = \\ &= \frac{h^2}{n} \left[ (1^2 + 2^2 + \dots + n^2) - \frac{n^2}{2} \right] = \frac{h^2}{n} \left[ \frac{n(n+1)(2n+1)}{6} - \frac{n^2}{2} \right] = \frac{1}{3} + \frac{h^2}{6}. \end{aligned}$$

Погрешность здесь составляет  $h^2/6$ , хотя подынтегральная функция  $f(x) = x^2$  имеет производные не только второго, но и всех порядков.

Интересно напомнить тот известный факт, что для периодических функций  $f(x)$  с периодом  $b - a$  точность формулы трапеций с постоянным шагом  $h$  для  $\int_a^b f(x) dx$  самонастраивается на гладкость: порядок точности по  $h$  возрастает с ростом числа производных периодической функции; именно, имеет место

**Теорема 2.** Пусть  $f(x)$  — периодическая функция с периодом  $X = b - a$ ,  $f(x + X) \equiv f(x)$ , имеющая ограниченные производные до некоторого порядка  $k$ ,  $\max_x |f^{(k)}(x)| = M_k$ . Тогда погрешность формулы трапеций будет  $O(M_k h^k)$ .

Мы не приводим доказательства этой теоремы.

**3. Формула Симпсона.** Эта квадратурная формула получается на основе замены подынтегральной функции кусочно квадратичной. Для достаточно гладких функций ее погрешность при  $h \rightarrow 0$  убывает как  $h^4$ , т. е. быстрее, чем для формулы трапеций.

Будем считать  $n$  четным,  $n = 2k$ , а узлы равноотстоящими,  $x_{j+1} - x_j = h = (b - a)/n$ . На каждом отрезке  $[x_{2j}, x_{2j+2}]$  ( $j = 0, 1, \dots, k - 1$ ) заменим  $f(x)$  интерполяционным многочленом второй степени  $P_2(x, f_{2j}, f_{2j+1}, f_{2j+2})$ , который запишем в форме Лагранжа:

$$\begin{aligned} P_2(x, f_{2j}, f_{2j+1}, f_{2j+2}) &= f_{2j} \frac{(x - x_{2j+1})(x - x_{2j+2})}{(x_{2j} - x_{2j+1})(x_{2j} - x_{2j+2})} + \\ &+ f_{2j+1} \frac{(x - x_{2j})(x - x_{2j+2})}{(x_{2j+1} - x_{2j})(x_{2j+1} - x_{2j+2})} + f_{2j+2} \frac{(x - x_{2j})(x - x_{2j+1})}{(x_{2j+2} - x_{2j})(x_{2j+2} - x_{2j+1})}. \end{aligned}$$

Применив формулу Ньютона–Лейбница, проверим, что

$$\int_{x_{2j}}^{x_{2j+2}} P_2 dx = \frac{b - a}{3n} (f_{2j} + 4f_{2j+1} + f_{2j+2}). \quad (9)$$

В результате приближенной замены  $f(x)$  на  $P_2(x, f_{2j}, f_{2j+1}, f_{2j+2})$  получим приближенные равенства

$$\int_{x_{2j}}^{x_{2j+2}} f(x) dx \approx \int_{x_{2j}}^{x_{2j+2}} P_2 dx = \frac{b-a}{3n} (f_{2j} + 4f_{2j+1} + f_{2j+2}),$$

$$j = 0, 1, \dots, k-1.$$

Суммируя почленно приближенные равенства

$$\int_{x_{2j}}^{x_{2j+2}} f(x) dx \approx \frac{b-a}{3n} (f_{2j} + 4f_{2j+1} + f_{2j+2}), \quad (10)$$

получаем

$$\int_a^b f(x) dx \approx \frac{b-a}{3n} (f_0 + 4f_1 + 2f_2 + 4f_3 + \dots + 4f_{n-1} + f_n). \quad (11)$$

Формула (11) для приближенного вычисления определенного интеграла  $\int_a^b f(x) dx$  называется *формулой Симпсона*.

**Теорема 3.** Пусть  $f(x)$  имеет на отрезке  $[a, b]$  ограниченную третью производную:  $\max_{a \leq x \leq b} |f^{(3)}(x)| = M_3$ . Тогда погрешность формулы Симпсона, т. е. разность  $\int_a^b f(x) dx - S_n(f)$ , удовлетворяет оценке

$$\left| \int_a^b f(x) dx - S_n(f) \right| \leq \frac{(b-a)M_3}{12} h^3.$$

**Доказательство.** Формула Симпсона возникла в результате замены подынтегральной функции  $f(x)$  на отрезке  $[a, b]$  кусочно квадратичной функцией  $P_2(x)$ , которая на каждом отрезке  $[x_{2j}, x_{2j+2}]$  совпадает с интерполяционным многочленом  $P_2(x, x_{2j}, x_{2j+1}, x_{2j+2})$ , построенным по значениям  $f_{2j}, f_{2j+1}, f_{2j+2}$  в узлах  $x_{2j}, x_{2j+1}, x_{2j+2}$ .

Оценим разность  $R_2(x) \equiv f(x) - P_2(x)$  для  $x \in [a, b]$ . Каждая точка  $x$  принадлежит какому-то отрезку  $[x_{2j}, x_{2j+2}]$ . Но на отрезке  $[x_{2j}, x_{2j+2}]$  величина  $R_2(x)$  есть погрешность квадратичной интерполяции и выражается формулой

$$R_2(x) = \frac{f^{(3)}(\xi)}{3!} (x - x_{2j})(x - x_{2j+1})(x - x_{2j+2}), \quad \xi \in (x_{2j}, x_{2j+2}).$$

Отсюда

$$|R_2(x)| \leq \frac{M_3}{3!} \max_{x_{2j} \leq x \leq x_{2j+2}} |(x - x_{2j})(x - x_{2j+1})(x - x_{2j+2})| \leq \frac{M_3}{12} h^3.$$

В силу произвольности  $x \in [a, b]$

$$\max_{a \leq x \leq b} |R_2(x)| \leq \frac{M_3}{12} h^3.$$

Следовательно,

$$\begin{aligned} \left| \int_a^b f(x) dx - S_n(f) \right| &= \left| \int_a^b f(x) dx - \int_a^b P_2(x) dx \right| = \\ &= \left| \int_a^b (f(x) - P_2(x)) dx \right| = \left| \int_a^b R_2(x) dx \right| \leq \int_a^b |R_2(x)| dx \leq \\ &\leq \int_a^b \max_x |R_2(x)| dx \leq \frac{M_3 h^3}{12} \int_a^b dx = \frac{b-a}{12} h^3 M^3. \quad \square \end{aligned}$$

**Теорема 4.** Пусть  $f(x)$  имеет на отрезке  $[a, b]$  ограниченную четвертую производную:  $\max_{a \leq x \leq b} |f^{(4)}(x)| = M_4$ . Тогда погрешность формулы Симпсона, т. е. разность  $\int_a^b f(x) dx - S_n(f)$ , удовлетворяет оценке

$$\left| \int_a^b f(x) dx - S_n(f) \right| \leq \frac{(b-a) M_4}{18} h^4. \quad (12)$$

**Доказательство.** Формула Симпсона возникла из записи интеграла  $\int_a^b f(x) dx$  в виде

$$\int_a^b f(x) dx = \int_{x_0}^{x_2} f(x) dx + \int_{x_2}^{x_4} f(x) dx + \dots + \int_{x_{n-2}}^{x_n} f(x) dx \quad (13)$$

с дальнейшей заменой каждого слагаемого выражением

$$\frac{b-a}{3n} (f_{2j} + 4f_{2j+1} + f_{2j+2}).$$

Покажем, что при этом погрешность, вносимая в каждое из  $k = n/2$  слагаемых суммы (13), удовлетворяет одной и той же оценке

$$\left| \int_{x_{2j}}^{x_{2j+2}} f(x) dx - \frac{b-a}{3n} (f_{2j} + 4f_{2j+1} + f_{2j+2}) \right| \leq \frac{M_4}{9} h^5. \quad (14)$$

Из этой оценки вытекает и оценка (12):

$$\begin{aligned} \left| \int_a^b f(x) dx - S_n(f) \right| &= \\ &= \left| \sum_{j=0}^{k-1} \int_{x_{2j}}^{x_{2j+2}} f(x) dx - \frac{b-a}{3n} (f_{2j} + 4f_{2j+1} + f_{2j+2}) \right| \leqslant \\ &\leqslant \sum_{j=0}^{k-1} \left| \int_{x_{2j}}^{x_{2j+2}} f(x) dx - \frac{b-a}{3n} (f_{2j} + 4f_{2j+1} + f_{2j+2}) \right| \leqslant \\ &\leqslant \frac{n}{2} \cdot \frac{M_4}{9} h^5 = \frac{(b-a)M_4}{18} h^4. \end{aligned}$$

Для доказательства оценки (14) представим  $f(x)$  с помощью формулы Тейлора в окрестности точки  $x_{2j+1}$  в виде

$$f(x) = Q(x) + R(x),$$

где

$$\begin{aligned} Q(x) &= f(x_{2j+1}) + \frac{f'(x_{2j+1})}{1!}(x - x_{2j+1}) + \\ &\quad + \frac{f''(x_{2j+1})}{2!}(x - x_{2j+1})^2 + \frac{f'''(x_{2j+1})}{3!}(x - x_{2j+1})^3, \\ R(x) &= \frac{f^{(4)}(\xi)}{4!}(x - x_{2j+1})^4, \quad x_{2j} \leq \xi \leq x_{2j+2}. \end{aligned}$$

Далее,

$$\begin{aligned} \int_{x_{2j}}^{x_{2j+2}} f(x) dx - \frac{b-a}{3n} (f_{2j} + 4f_{2j+1} + f_{2j+2}) &= \\ &= \int_{x_{2j}}^{x_{2j+2}} Q(x) dx + \int_{x_{2j}}^{x_{2j+2}} R(x) dx - \\ &- \frac{b-a}{3n} [(Q_{2j} + R_{2j}) + 4(Q_{2j+1} + R_{2j+1}) + (Q_{2j+2} + R_{2j+2})] = \\ &= \left\{ \int_{x_{2j}}^{x_{2j+2}} Q(x) dx - \frac{b-a}{3n} (Q_{2j} + 4Q_{2j+1} + Q_{2j+2}) \right\} + \\ &\quad + \left[ \int_a^b R(x) dx - \frac{b-a}{3n} (R_{2j} + 4R_{2j+1} + R_{2j+2}) \right]. \end{aligned}$$

Можно проверить непосредственно, что для  $Q(x) \equiv 1$ ,  $Q(x) \equiv x$ ,  $Q(x) \equiv x^2$ ,  $Q(x) \equiv x^3$  слагаемое, стоящее в фигурных скобках, обращается в нуль. Следовательно, эти фигурные скобки обращаются

в нуль и для произвольного многочлена  $Q(x)$  степени не выше третьей. В силу этого

$$\left| \int_{x_{2j}}^{x_{2j+2}} f(x) dx - \frac{b-a}{3n} (f_{2j} + 4f_{2j+1} + f_{2j+2}) \right| = \\ = \left| \int_{x_{2j}}^{x_{2j+2}} R(x) dx - \frac{b-a}{3n} (R_{2j} + 4R_{2j+1} + R_{2j+2}) \right|. \quad (15)$$

Учтем, что  $R_{2j+1} = 0$ , а  $|R(x)| \leq \frac{M_4}{4!} h^4$ . Поэтому правая часть равенства (15) оценивается величиной

$$\int_{x_{2j}}^{x_{2j+2}} |R(x)| dx + \frac{b-a}{3n} \cdot 2 \max_x |R(x)| \leq \\ \leq \max_x |R(x)| \left( \int_{x_{2j}}^{x_{2j+2}} dx + \frac{2}{3} \cdot \frac{b-a}{n} \right) = \max_x |R(x)| \left( 2 \frac{b-a}{n} + \frac{2}{3} \cdot \frac{b-a}{n} \right) \leq \\ \leq \frac{M_4}{4!} \cdot \frac{8}{3} (b-a) \frac{1}{n} h^4 = \frac{M_4}{9} h^5. \quad \square$$

**Следствие.** Формула Симпсона точна для функций  $f(x)$ , являющихся многочленами степени не выше третьей.

**Доказательство.** Для таких функций  $M_4 = \max |f^{(4)}(x)| = 0$  и в правой части оценки (14) стоит нуль.  $\square$

Заметим, что на практике выбор числа  $n$  и соответственно шага сетки перепоручают компьютеру. Ведут расчет с некоторым  $n$ , затем с удвоенным  $n$  и т. д., пока результат не перестанет меняться в требуемом десятичном знаке.

### Задачи

1. Показать, что формула трапеций точна, если  $f(x)$  — многочлен степени не выше первой.

2\*. Доказать теорему 2 в случае  $b-a=1$ .

Указание. Воспользоваться представлением  $f(x)$  в виде ряда Фурье

$$f(x) = \sum_{k=-\infty}^{\infty} c_k e^{2k\pi xi}.$$

3. Показать, что в случае переменного шага сетки, удовлетворяющего условию  $x_{j+1} - x_j < h = \text{const}$ , использование кусочно квадратической интерполяции приведет к обобщению формулы Симпсона, порядок погрешности которой возрастает и будет  $O(M_3 h^3)$ , где  $M_3 = \max |f'''(x)|$ .

4. Показать, что для сколь угодно гладких подынтегральных функций погрешность формулы Симпсона может фактически достигать величины  $\text{const} \cdot h^4$ .

## § 2. Сочетание численных и аналитических методов при вычислении интегралов с особенностями

Даже для простейших несобственных интегралов непосредственное применение квадратурных формул наталкивается на трудности. Например, формула трапеций

$$\int_a^b f(x) dx \approx \frac{b-a}{n} \left( \frac{f_0}{2} + f_1 + \dots + f_{n-1} + \frac{f_n}{2} \right)$$

в случае  $a = 0$ ,  $b = 10$ ,  $f(x) = \cos x / \sqrt{x}$  неприменима, так как  $f_0$  не определено.

Для преодоления этих трудностей естественно воспользоваться аналитическими методами, чтобы свести задачу к вычислению определенного интеграла от гладкой ограниченной функции, а затем применить какую-либо квадратурную формулу. В рассматриваемом примере с этой целью можно сначала воспользоваться формулой интегрирования по частям, приняв  $u = \cos x$ ,  $\frac{1}{\sqrt{x}} dx = dv$ . Тогда

$$\int_0^{10} \frac{\cos x}{\sqrt{x}} dx = \cos x (2\sqrt{x}) \Big|_0^{10} + \int_0^{10} 2\sqrt{x} \sin x dx,$$

и осталось вычислить следующий определенный интеграл:

$$\int_0^{10} (\sin x) \sqrt{x} dx.$$

Возможен другой прием. А именно, разобьем интеграл  $\int_0^{10} \frac{\cos x}{\sqrt{x}} dx$  на два:

$$\int_0^{10} \frac{\cos x}{\sqrt{x}} dx = \int_0^c \frac{\cos x}{\sqrt{x}} dx + \int_c^{10} \frac{\cos x}{\sqrt{x}} dx, \quad (1)$$

где  $c > 0$  — некоторое, пока произвольное число. Для вычисления первого интеграла, входящего в правую часть равенства (1), воспользуемся разложением

$$\cos x = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots$$

Тогда

$$\int_0^c \frac{\cos x}{\sqrt{x}} dx = \int_0^c \frac{dx}{\sqrt{x}} - \int_0^c \frac{x^{3/2}}{2!} dx + \int_0^c \frac{x^{7/2}}{4!} dx - \dots \quad (2)$$

Каждый из интегралов в правой части равенства (2) вычисляется по формуле Ньютона–Лейбница. Пусть требуется вычислить искомый интеграл с погрешностью, не превосходящей заданное значение  $\varepsilon > 0$ .

Для этого достаточно ограничиться суммой нескольких первых членов ряда в правой части равенства (2). Число этих членов зависит от  $\varepsilon > 0$  и от  $c$ . Не стоит брать  $c$  большим, например,  $c = 10$ . Тогда даже при не очень малом  $\varepsilon$ , например, при  $\varepsilon = 10^{-3}$ , потребуется много членов ряда. Не стоит брать  $c \ll 1$ , так как в этом случае интеграл (2) будет вычисляться экономно, но зато подынтегральная

функция в интеграле  $\int_c^{10} \frac{\cos x}{\sqrt{x}} dx$  будет иметь большие по абсолютной

величине производные. Это потребует больших значений  $n$  в формуле трапеций (или Симпсона) для достижения заданной точности  $\varepsilon = 10^{-3}$  при вычислении этого интеграла.

Для преодоления трудности, вызываемой особенностью, иногда удается воспользоваться заменой переменной. Например, несобственный

интеграл  $\int_0^1 \frac{\cos x}{\sqrt{x}} dx$  в результате замены  $x = t^2$  переходит в определенный интеграл без особенностей

$$\int_0^1 \frac{\cos t^2}{t} \cdot 2t dt = 2 \int_0^1 \cos t^2 dt.$$

Залог успешного решения задачи — в разумном распределении трудностей вычисления между аналитическими преобразованиями и приемами, которые должен выполнить исследователь, и рутинными вычислениями, которые должен выполнить компьютер. Это вообще характерно для решения любых задач с помощью быстродействующего компьютера.

### Задача

Предложить алгоритмы для вычисления следующих несобственных интегралов:

$$\int_0^1 \frac{\sin x}{x^{3/2}} dx, \quad \int_0^{10} \frac{\ln(1+x)}{x^{3/2}} dx, \quad \int_0^\infty e^{-x^2} dx, \quad \int_0^{\pi/2} \frac{\cos x}{\sqrt{\pi/2 - x}} dx.$$

## § 3. Кратные интегралы

Вычисление кратных интегралов вида

$$I^{(m)} = I^{(m)}(f, D) = \int_D \dots \int f(x_1, \dots, x_m) dx_1 \dots dx_m \quad (1)$$

по некоторой области  $D$  пространства  $x_1, \dots, x_m$  можно свести к использованию квадратурных формул для вычисления определенных

интегралов вида  $\int\limits_a^b f(x) dx$ . Однако сложность соответствующих алгоритмов и число арифметических операций для получения ответа с заданной точностью  $\varepsilon > 0$ , вообще говоря, быстро возрастают с ростом размерности  $m$ , даже если область интегрирования — единичный куб  $D = \{|x_i| < 1, i = 1, 2, \dots, m\}$ .

Большое дополнительное усложнение возникает в случае области интегрирования  $D$  общего вида.

**1. Переход от кратного интеграла к повторному и использование квадратурных формул.** Для выяснения сути дела предположим, что  $m = 2$ , и рассмотрим интеграл вида

$$I^{(2)} = \iint_D f(x, y) dx dy \quad (2)$$

по области  $D$  (рис. 7), заключенной между кривыми

$$y = \varphi_1(x), \quad y = \varphi_2(x), \quad (3)$$

при значениях  $x \in [a, b]$ .

Очевидно, что

$$I^{(2)} = \iint_D f(x, y) dx dy = \int_a^b F(x) dx, \quad (4)$$

где

$$F(x) = \int_{\varphi_1(x)}^{\varphi_2(x)} f(x, y) dy, \quad a \leq x \leq b. \quad (5)$$

Интеграл (5) можно приближенно заменить квадратурной формулой. Возьмем для определенности формулу трапеций вида

$$I^{(2)} = \int_a^b F(x) dx \approx \frac{b-a}{n} \left( \frac{F(x_0)}{2} + F(x_1) + \dots + F(x_{n-1}) + \frac{F(x_n)}{2} \right), \quad (6)$$

$$x_k = a + k \frac{b-a}{n}, \quad k = 0, 1, \dots, n.$$

Выражения

$$F(x_j) = \int_{\varphi_1(x_j)}^{\varphi_2(x_j)} f(x_j, y) dy,$$

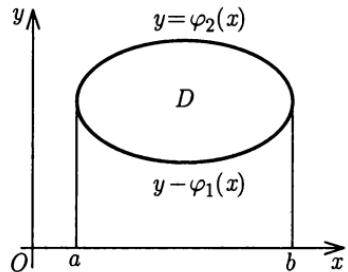


Рис. 7

входящие в формулу (6), в свою очередь можно вычислять по квадратурной формуле трапеций:

$$\begin{aligned} F(x_j) &\approx \frac{\varphi_2(x_j) - \varphi_1(x_j)}{n} \times \\ &\quad \times \left( \frac{f(x_j, y_0)}{2} + f(x_j, y_1) + \dots + f(x_j, y_{n-1}) + \frac{f(x_j, y_n)}{2} \right), \quad (7) \\ y_k &= \varphi_1(x_j) + \frac{\varphi_2(x_j) - \varphi_1(x_j)}{n} k, \quad k = 0, 1, \dots, n. \end{aligned}$$

При приближенном вычислении интеграла  $I^{(2)}$  по формулам (6), (7) с ростом  $n$  погрешность  $\varepsilon = \varepsilon(n)$  уменьшается, а число арифметических операций  $\nu(n)$  возрастает. Очевидно, что  $\nu(n) = O(n^2)$ . Если функция  $f(x, y)$  имеет ограниченные производные второго порядка, а функции  $\varphi_1(x)$ ,  $\varphi_2(x)$  достаточно гладкие, то в силу свойств формулы трапеций  $\varepsilon(n) = O(n^{-2})$ .

Прием перехода от кратного интеграла к повторному, а затем вычисления последнего с помощью многократного использования квадратурных формул может быть перенесен и на общий случай кратных интегралов  $I^{(m)}$  ( $m > 2$ ), если область  $D$  задана достаточно удобно. При этом число арифметических операций при использовании квадратурных формул с заданным  $n$  есть

$$\nu(n) = O(n^m).$$

Если  $f(x_1, x_2, \dots, x_m)$  имеет ограниченные производные второго порядка, граница  $\Gamma$  области  $D$  достаточно гладкая и в качестве квадратурных формул многократно используется формула трапеций, то погрешность  $\varepsilon$  при вычислении кратного интеграла есть  $\varepsilon = \varepsilon(n) = O(n^{-2})$ .

При заданном  $\varepsilon$  выбор  $n$  производится обычно путем экспериментальных расчетов последовательно с некоторым  $n = n_0$ , затем  $n = 2n_0, \dots$  до тех пор, пока результат не перестанет изменяться в пределах заданной точности. Существуют приемы ускорения расчетов.

**2. Эффективные методы для вычисления кратных интегралов при больших значениях размерностей  $m$ .** Приведем простейший пример, где упрощение достигается за счет учета специфики задачи.

Пусть требуется вычислить интеграл  $\iint_D f(x, y) dx dy$ , причем известно, что  $D$  — круг:  $x^2 + y^2 \leq R^2$ , а  $f(x, y)$  в действительности зависит не от  $x$  и  $y$  в отдельности, а от величины  $x^2 + y^2 = r^2$ , так что  $f(x, y) \equiv \tilde{f}(r)$ . Тогда

$$I^{(2)} = \iint_D f(x, y) dx dy = 2\pi \int_0^R \tilde{f}(r)r dr,$$

и дело свелось к вычислению обычного определенного интеграла.

О вычислении кратных интегралов см. [10, 22].

**3. Понятие о методе Монте-Карло.** Указанный выше прием редукции к обычным определенным интегралам окажется затруднительным даже в случае  $m = 2$ , если область  $D$  ограничена достаточно сложной кривой.

Будем считать, что область  $D$  лежит внутри квадрата ( $0 < x < 1$ ,  $0 < y < 1$ ). Доопределим функцию  $f(x, y)$  вне  $D$  нулем. Пусть имеется датчик случайных чисел, который выдает числа, равномерно распределенные на отрезке  $[0, 1]$ , с плотностью вероятности, равной единице. С помощью этого датчика строим последовательность точек (пар чисел)  $P_k = (x_k, y_k)$  ( $k = 1, 2, \dots$ ) и составляем числовую последовательность

$$S_k = \frac{1}{k} \sum_{j=1}^k f(x_j, y_j).$$

Легко представить, что при произвольном  $\varepsilon > 0$  вероятность того, что погрешность  $I^{(2)} - S_k$  меньше  $\varepsilon$ :

$$|I^{(2)} - S_k| < \varepsilon,$$

с ростом числа  $k$  испытаний стремится к единице. Мы не приводим здесь имеющихся точных теорем.

Смысл метода Монте-Карло особенно нагляден, если  $f(x, y) = 1$ . Тогда  $I^{(2)}$  есть площадь области  $D$ , а  $S_k$  — отношение числа тех из точек  $P_1, P_2, \dots, P_k$ , которые попали в область  $D$ , к числу  $k$  всех этих точек. Очевидно, что это отношение приблизительно равно площади области  $D$ , поскольку площадь всего квадрата, куда попадают точки, равна единице.

Привлекательность метода Монте-Карло — в логической простоте соответствующего алгоритма, который не усложняется с ростом размерности  $m$  и с усложнением формы области  $D$ . Однако в случае простых областей и узких классов функций он проигрывает методам, учитывающим эту специфику.

Метод Монте-Карло имеет и многие другие приложения [23].

## ЧАСТЬ II

# СИСТЕМЫ СКАЛЯРНЫХ УРАВНЕНИЙ

---

Числа называют еще скалярами, а поэтому уравнения и системы уравнений относительно одного или нескольких неизвестных чисел называют скалярными. Существуют еще функциональные уравнения, в частности, дифференциальные уравнения, в которых неизвестными являются не числа, а функции, но в этой части книги мы ими заниматься не будем.

Среди систем скалярных уравнений можно выделить важный класс систем линейных алгебраических уравнений, которому уделено много внимания (гл. 4, 5).

В гл. 6 рассматриваются нелинейные скалярные уравнения.

## ГЛАВА 4

### СИСТЕМЫ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ УРАВНЕНИЙ. МЕТОДЫ ОТЫСКАНИЯ ТОЧНОГО РЕШЕНИЯ

Системы линейных алгебраических уравнений (СЛАУ) высокого порядка возникают во многих приложениях. Рассмотрим координатную форму записи СЛАУ порядка  $n$ :

$$\begin{aligned} a_{11}x_1 + \dots + a_{1n}x_n &= f_1, \\ \dots &\dots \\ a_{n1}x_1 + \dots + a_{nn}x_n &= f_n, \end{aligned}$$

где  $a_{ij}$ ,  $f_i$  ( $i, j = 1, 2, \dots, n$ ) — заданные числа,  $x_j$  ( $j = 1, 2, \dots, n$ ) — искомые числа. Более общая и часто более удобная форма задания СЛАУ — операторная.

Мы обсудим формы записи совместных СЛАУ (§ 1), определим понятие обусловленности оператора и соответствующей СЛАУ (§ 2, 3), рассмотрим некоторые точные (§ 4, 6, 7) и итерационные (§ 4) методы вычисления решений СЛАУ, указав классы СЛАУ, для которых один алгоритм вычисления решений предпочтительнее другого. В § 5 устанавливается связь между задачей на минимум квадратичной функции и СЛАУ.

Отметим сразу же, что для вычисления решения СЛАУ известные формулы Крамера при сколько-нибудь значительном  $n$  ( $n > 5$ ) не применяются из-за большого числа арифметических операций и вычислительной неустойчивости.

## § 1. Формы записи совместных СЛАУ

*Линейным уравнением* относительно неизвестных  $z, u, \dots, w$  называется уравнение вида

$$\alpha z + \beta u + \dots + \gamma w = t, \quad (1)$$

где  $\alpha, \beta, \dots, \gamma$ , а также  $t$  — какие-нибудь заданные числа.

**1. Каноническая форма записи СЛАУ.** Пусть задана система  $n$  линейных уравнений относительно того же числа  $n$  неизвестных. Эту систему, очевидно, можно записать в следующей (канонической) форме:

$$\begin{aligned} a_{11}x_1 + \dots + a_{1n}x_n &= f_1, \\ \dots &\dots \\ a_{n1}x_1 + \dots + a_{nn}x_n &= f_n. \end{aligned} \quad (2)$$

Для этого надо занумеровать заданные уравнения вида (1) и неизвестные номерами  $1, 2, \dots, n$ , а затем неизвестные переобозначить через  $x_j$ , коэффициент при неизвестном  $x_j$  в уравнении, получившем номер  $i$ , — через  $a_{ij}$ , а правые части — через  $f_i$  ( $i, j = 1, 2, \dots, n$ ). Из линейной алгебры известно, что система (2) имеет решение (совместна) при любых правых частях в том и только том случае, если соответствующая однородная система имеет только тривиальное решение  $x_1 = x_2 = \dots = x_n = 0$ . Для этого необходимо и достаточно, чтобы матрица  $A$  коэффициентов системы

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \quad (3)$$

имела определитель, отличный от нуля:  $\det A \neq 0$ . Условие  $\det A \neq 0$  совместности системы (2) при любых правых частях одновременно обеспечивает и единственность решения для любого конкретного набора правых частей.

С помощью матрицы  $A$  систему (2) можно переписать в виде

$$\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix} = \begin{bmatrix} f_1 \\ \dots \\ f_n \end{bmatrix}, \quad (4)$$

или, короче,

$$Ax = f. \quad (5)$$

**2. Операторная форма записи СЛАУ.** Пусть заданы линейное пространство  $R^n$  размерности  $n$ , линейный оператор  $A : R^n \rightarrow R^n$  и некоторый элемент  $f \in R^n$ . Рассмотрим уравнение

$$Ax = f \quad (6)$$

относительно  $x \in R^n$ . Известно, что это уравнение имеет решение при произвольном  $f \in R^n$  в том и только том случае, если соответствующее однородное уравнение  $Ax = 0$  имеет только тривиальное решение  $x = 0 \in R^n$ .

Пусть элементы пространства  $R^n$  — это всевозможные упорядоченные системы чисел вида

$$z = \begin{bmatrix} z_1 \\ \dots \\ z_n \end{bmatrix}. \quad (7)$$

В случае задания элементов пространства  $R^n$  в форме (7) будем говорить, что *элементы заданы своими координатами*.

Из линейной алгебры известно, что в этом случае каждому линейному оператору  $A$  соответствует матрица вида (3), такая, что заданному  $x \in R^n$  сопоставляется  $y = Ax$  по формулам

$$\begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix}. \quad (8)$$

Уравнение (6) в этом случае совпадает с уравнением (4) или с канонической записью (2). Оно является, таким образом, лишь операторной интерпретацией системы (2). Однако элементы пространства  $R^n$  не обязательно задаются как упорядоченные системы чисел (7), а оператор  $A$  не обязательно исходно задается в форме (8).

Таким образом, операторная форма (6) задания системы линейных уравнений может отличаться от канонической формы (2).

Приведем пример прикладной задачи, которая естественно приводит к системе линейных уравнений высокого порядка, заданной в операторной форме. Этот пример будет нами многократно использоваться в книге.

**3. Разностный аналог задачи Дирихле для уравнения Пуассона.** Пусть в квадратной области  $D = \{0 < x, y < 1\}$  требуется численно решить следующую задачу:

$$\begin{aligned} -\Delta u &= -\left(\frac{\partial^2 u}{\partial x_1^2} + \frac{\partial^2 u}{\partial x_2^2}\right) = f(x, y), \quad (x, y) \in D, \\ u|_{\partial D} &= 0, \end{aligned} \quad (9)$$

где  $\partial D$  — граница области  $D$ . Для численного решения этой задачи зададим натуральное  $M$ , введем  $h = M^{-1}$  и построим сетку с узлами

$$(x_{m_1}, y_{m_2}) = (m_1 h, m_2 h), \quad m_1, m_2 = 0, 1, \dots, M. \quad (10)$$

Будем искать таблицу приближенных значений  $u_{m_1 m_2}$  решения  $u(x, y)$  задачи (9) в узлах сетки (10) (рис. 8). Для этого в каждом узле сетки  $(x_{m_1}, y_{m_2}) = (m_1 h, m_2 h)$ , лежащем внутри  $D$ , заменим

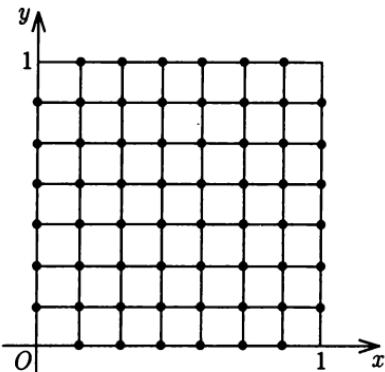


Рис. 8

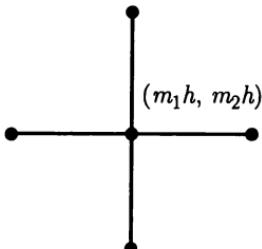


Рис. 9

дифференциальное уравнение (9) разностным, заменяя производные  $\partial^2 u / \partial x^2|_{(x_{m_1}, y_{m_2})}$ ,  $\partial^2 u / \partial y^2|_{(x_{m_1}, y_{m_2})}$  разностными отношениями

$$\begin{aligned} -\Delta^{(h)} u|_{m_1 m_2} = & -\left(\frac{u_{m_1+1, m_2} - 2u_{m_1 m_2} + u_{m_1-1, m_2}}{h^2} + \right. \\ & \left. + \frac{u_{m_1, m_2+1} - 2u_{m_1 m_2} + u_{m_1, m_2-1}}{h^2}\right) = f_{m_1 m_2}, \quad (11) \\ m_1, m_2 = 1, \dots, M-1, \end{aligned}$$

где  $f_{m_1 m_2} = f(x_{m_1}, y_{m_2})$ .

Это уравнение, составленное в точке  $(m_1 h, m_2 h)$ , связывает значение искомой сеточной функции в пяти точках (рис. 9).

Если точка сетки  $(m_1 h, m_2 h)$  лежит на расстоянии шага сетки  $h$  от границы квадрата  $D$ , то в уравнение (11) входят значения в одной или двух точках границы. Значения  $u_{m_1, m_2 \pm 1}$ ,  $u_{m_1 \pm 1, m_2}$  в таких точках положим равными нулю. Таким образом, система (11) есть система  $n = (M-1)^2$  линейных уравнений относительно такого же числа неизвестных  $u_{m_1 m_2}$  ( $m_1, m_2 = 1, 2, \dots, M-1$ ). Для получения решения задачи Дирихле с большой точностью нужно взять достаточно мелкую сетку, или большое число  $M$ . Система (11) будет тогда системой высокого порядка.

Придадим системе (11) смысл операторной формы записи (6). Для этого введем пространство  $U^{(h)} = R^n$ ,  $n = (M-1)^2$ , состоящее из всевозможных вещественных функций  $z_{m_1 m_2}$ , заданных на сетке  $(m_1 h, m_2 h)$  ( $m_1, m_2 = 1, 2, \dots, M-1$ ). Введем оператор  $A = -\Delta^{(h)}$ ,

который действует из  $U^{(h)}$  в  $U^{(h)}$ , сопоставляя заданной функции  $u^{(h)} \in U^{(h)}$  некоторую функцию  $v^{(h)} \in U^{(h)}$  по формулам

$$\begin{aligned} v^{(h)}|_{m_1 m_2} &\equiv v_{m_1 m_2} = -\Delta^{(h)} u^{(h)}|_{m_1 m_2} = \\ &= -\frac{u_{m_1+1, m_2} - 2u_{m_1, m_2} + u_{m_1-1, m_2}}{h^2} - \frac{u_{m_1, m_2+1} - 2u_{m_1, m_2} + u_{m_1, m_2-1}}{h^2}, \end{aligned} \quad (12)$$

где  $u_{n_1 n_2} = u^{(h)}|_{n_1 n_2}$  или  $u_{n_1 n_2} = 0$  в зависимости от того, лежит ли точка  $(n_1 h, n_2 h)$  внутри или на границе квадрата  $0 \leq x, y \leq 1$ .

Система линейных уравнений (11) приобретает тогда вид (6):

$$-\Delta^{(h)} u^{(h)} = f^{(h)}, \quad (13)$$

где  $f^{(h)} = \{f_{m_1 m_2}\} \in U^{(h)}$  — заданная функция, а  $u^{(h)} \in U^{(h)}$  — искомая функция.

**Замечание 1.** Систему (11), как всякую систему  $n$  линейных уравнений относительно  $n$  неизвестных, можно было бы записать в каноническом виде (2). Однако при любой нумерации уравнений и неизвестных матрица возникающей системы (2) уступала бы по простоте и обозримости заданию системы в исходной форме (11) и вполне адекватной ей операторной форме (13).

**Замечание 2.** Абстрактное уравнение в операторной форме (6) всегда можно привести к каноническому виду (2). Для этого надо выбрать какой-либо базис  $e_1, e_2, \dots, e_n$  в пространстве  $R^n$ , записать элементы  $x, f$  из  $R^n$  в виде

$$x = x_1 e_1 + \dots + x_n e_n, \quad f = f_1 e_1 + \dots + f_n e_n$$

и отождествить эти элементы с наборами их координат так, что

$$x = \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix}, \quad f = \begin{bmatrix} f_1 \\ \dots \\ f_n \end{bmatrix}.$$

Оператор  $A$ , входящий в (6), запишется соответствующей этому оператору и сделанному выбору базиса матрицей (3) по формулам вида (8). Уравнение (6) приобретет канонический вид (4), или (2). Однако преобразование (6) к виду (2) с целью вычисления решения может оказаться трудоемким и нецелесообразным.

### Задача

Отображение  $A: R^2 \rightarrow R^2$  элементов  $x = (x_1, x_2)$  пространства  $R^2$  в себя задано следующим образом. Рассматривается следующая задача Коши для системы дифференциальных уравнений:

$$dz_1/dt = z_1, \quad dz_2/dt = z_2, \quad 0 \leq t \leq 1,$$

$$z_1|_{t=0} = x_1, \quad z_2|_{t=0} = x_2.$$

Принимаем за  $y = Ax$  элемент  $y = (y_1, y_2)$ :  $y_1 = z_1(t)|_{t=1}$ ,  $y_2 = z_2(t)|_{t=1}$ .

Доказать, что оператор  $A$  линейный.

а) Выписать матрицу оператора  $A$ , если за базис в пространстве  $R^2$  принять

$$e_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad e_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

б) То же, что и в п. а), но в базисе

$$e_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad e_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

в) Вычислить  $x = A^{-1}y$ ,  $y = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$ .

## § 2. Нормы

Решение системы линейных уравнений порядка  $n$ , а также погрешность приближенного решения можно трактовать как элементы линейного пространства  $R^n$ . Для количественных суждений о погрешности нужно ввести понятие длины вектора (элемента) из пространства  $R^n$ . Вместо термина «длина вектора» в линейной алгебре употребляется термин «норма вектора». Напомним определения линейного нормированного пространства  $R^n$  и нормы линейного оператора  $A: R^n \rightarrow R^n$ .

**1. Нормированные пространства.** Будем говорить, что пространство  $R^n$  *нормировано*, если каждому элементу  $x \in R^n$  сопоставлено неотрицательное число  $\|x\|$ , причем выполнены следующие условия (аксиомы).

1°.  $\|x\| > 0$ , если  $x \neq 0$ .

2°. Для любого числа  $\lambda$  и любого  $x \in R^n$

$$\|\lambda x\| = |\lambda| \|x\|.$$

3°. Для любых  $x \in R^n$ ,  $y \in R^n$  выполнено неравенство треугольника

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

Приведем примеры норм. Пусть  $R^n$  состоит из элементов вида  $x = (x_1, x_2, \dots, x_n)$ , где  $x_j$  — числа. Можно показать, что функции  $\|x\|_1$ ,  $\|x\|_2$ , определенные равенствами

$$\|x\|_1 = \max_j |x_j|, \tag{1}$$

$$\|x\|_2 = \sqrt{\sum_j x_j^2}, \tag{2}$$

удовлетворяют аксиомам 1°–3°. Они называются *первой* и *второй нормами*.

Введем в  $R^n$  скалярное умножение  $(x, y)$ , положив

$$(x, y) = x_1 y_1 + \dots + x_n y_n \tag{3}$$

в случае вещественного пространства или

$$(x, y) = x_1 \bar{y}_1 + \dots + x_n \bar{y}_n \quad (4)$$

в случае комплексного пространства.

Напомним, что вещественное линейное пространство со скалярным умножением называется *евклидовым*, а комплексное — *унитарным*.

Можно проверить, что функция

$$\|x\|_3 = (x, x)^{1/2} \quad (5)$$

является нормой (удовлетворяет условиям  $1^\circ - 3^\circ$ ). В случае вещественного пространства эту норму называют *евклидовой*, а в случае комплексного — *эрмитовой*.

Мы привели примеры норм в случае, если элементы пространства  $R^n$  записаны как векторы  $x = \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix}$ , т. е. заданы своими координатами. Возможно задание норм и без использования координатной записи элементов. Так, в пространстве  $U^{(h)} = R^n$ ,  $n = (M - 1)^2$ , функций  $u_{m_1 m_2}$  ( $m_1, m_2 = 1, 2, \dots, M - 1$ ), которое мы ввели и использовали при рассмотрении разностного аналога уравнения Пуассона в § 1, можно

ввести нормы, положив

$$\|u^{(h)}\| = \max_{m_1, m_2} |u_{m_1 m_2}|, \quad (6)$$

$$\|u^{(h)}\| = \sum_{m_1, m_2} |u_{m_1 m_2}|. \quad (7)$$

Можно проверить, что в вещественном пространстве  $U^{(h)}$  можно ввести скалярное умножение  $(u^{(h)}, v^{(h)})$ , положив

$$(u^{(h)}, v^{(h)}) = h^2 \sum_{m_1, m_2} u_{m_1 m_2} v_{m_1 m_2}, \quad (8)$$

а затем определить соответствующую евклидову норму, положив

$$\|u^{(h)}\| = (u^{(h)}, u^{(h)})^{1/2}. \quad (9)$$

Вообще, если в вещественном или комплексном линейном пространстве  $R^n$  введено какое-либо скалярное умножение  $(x, y)$ , то возникает соответствующая евклидова или эрмитова норма

$$\|x\| = (x, x)^{1/2}. \quad (10)$$

Можно описать все скалярные умножения и соответствующие евклидовые (эрмитовые) нормы. Для этого фиксируем какое-либо одно скалярное умножение  $(x, y)$ .

Напомним, что оператором  $B^*$ , *сопряженным* (в смысле выбранного скалярного умножения) какому-либо заданному линейному оператору

ру  $B: R^n \rightarrow R^n$ , называется такой линейный оператор  $B^*: R^n \rightarrow R^n$ , для которого выполнено тождество

$$(Bx, y)_{\frac{x}{y}} = (x, B^*y).$$

Известно, что для всякого оператора  $B$  существует один и только один сопряженный ему оператор  $B^*$ .

Оператор  $B$  называют *самосопряженным*, если  $B^* = B$ .

Оператор  $B: R^n \rightarrow R^n$  называется *положительно определенным* ( $B > 0$ ), если  $(Bx, x) > 0$  для всех  $x \neq 0$ . Известно, что если  $B = B^* > 0$ , то выражение

$$[x, y]_B \equiv (Bx, y) \quad (11)$$

удовлетворяет аксиомам скалярного умножения, причем каждое скалярное умножение в  $R^n$  может быть задано формулой (11) при подходящем подборе  $B = B^* > 0$ . В соответствии с этим каждый положительно определенный и самосопряженный оператор  $B = B^* > 0$  порождает евклидову (эрмитову) норму

$$\|x\|_B = ([x, x]_B)^{1/2}. \quad (12)$$

В частности, евклидову норму (10), порожденную первоначально выбранным скалярным умножением, можно записать в форме (12), положив  $B = E$ , где  $E$  — тождественный оператор:

$$\|x\|_E = ([x, x]_E)^{1/2} = (x, x)^{1/2}.$$

**2. Норма линейного оператора.** Введем норму  $\|A\|$  оператора  $A: R^n \rightarrow R^n$ , согласованную с нормой, выбранной в пространстве  $R^n$ , положив

$$\|A\| = \max_{\substack{x \in R^n \\ x \neq 0}} \frac{\|Ax\|}{\|x\|}, \quad (13)$$

где  $\|x\|$ ,  $\|Ax\|$  — нормы элементов  $x$ ,  $Ax$  из нормированного пространства  $R^n$ .

Таким образом,  $\|A\|$  — это коэффициент растяжения того вектора  $x' \in R^n$ , который растягивается не слабее любого другого  $x \in R^n$ .

Условимся каждую матрицу

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \quad (14)$$

отождествлять с оператором  $A: R^n \rightarrow R^n$ , который действует в пространстве  $R^n$ , состоящем из элементов  $x$  вида  $x = \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix}$ , где  $x_j$  —

числа, и сопоставляет заданному  $x$  некоторый элемент  $y = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix}$  по формулам

$$y_j = \sum_i a_{ij} x_i, \quad i, j = 1, 2, \dots, n. \quad (15)$$

Благодаря этому условию приобретает смысл понятие нормы матрицы  $A$  как нормы оператора, задаваемого формулами (15).

**Теорема 1.** Для норм матрицы  $A = \{a_{ij}\}$ , согласованных с нормами  $\|x\|_1 = \max |x_j|$ ,  $\|x\|_2 = \sqrt{\sum |x_j|^2}$  в пространстве  $R^n$  векторов  $x = \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix}$ , справедливы формулы

$$\|A\|_1 = \max_i \sum_j |a_{ij}|, \quad (16)$$

$$\|A\|_2 = \max_j \sum_i |a_{ij}|. \quad (17)$$

Доказательство предоставляем читателю.

**Теорема 2.** Пусть  $R^n$  — евклидово пространство, скалярное умножение в котором будем обозначать  $(x, y)$ , и пусть  $\|x\| = (x, x)^{1/2}$  — соответствующая евклидова норма. Пусть  $A$  — самосопряженный оператор:  $(Ax, y) \equiv (x, Ay)$ . Тогда

$$\|A\| = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_j |\lambda_j|, \quad (18)$$

где  $\lambda_j$  ( $j = 1, 2, \dots, n$ ) — собственные числа оператора  $A$ .

Доказательство. Из линейной алгебры известно, что в случае  $A = A^*$  существует ортонормированный базис

$$e_1, e_2, \dots, e_n, \quad (19)$$

состоящий из собственных векторов оператора  $A$ :

$$Ae_j = \lambda_j e_j, \quad j = 1, 2, \dots, n.$$

Запишем произвольный  $x \in R^n$ , а также  $Ax$  в виде линейных комбинаций векторов базиса

$$x = x_1 e_1 + x_2 e_2 + \dots + x_n e_n, \quad (20)$$

$$Ax = \lambda_1 x_1 e_1 + \lambda_2 x_2 e_2 + \dots + \lambda_n x_n e_n. \quad (21)$$

В силу ортонормированности базиса

$$\|x\| = (x_1^2 + x_2^2 + \dots + x_n^2)^{1/2},$$

$$\|Ax\| = [(\lambda_1 x_1)^2 + (\lambda_2 x_2)^2 + \dots + (\lambda_n x_n)^2]^{1/2}.$$

Очевидно,

$$\|Ax\| \leq \max_j |\lambda_j| (x_1^2 + x_2^2 + \dots + x_n^2)^{1/2} = \max_j |\lambda_j| \|x\|.$$

Поэтому для любого  $x \in R^n$  ( $x \neq 0$ )

$$\frac{\|Ax\|}{\|x\|} \leq \max_j |\lambda_j|,$$

а для  $x = e_k$ ,  $\max |\lambda_j| = |\lambda_k|$ , выполнено точное равенство

$$\frac{\|Ae_k\|}{\|e_k\|} = |\lambda_k| = \max_j |\lambda_j|.$$

Следовательно, справедливо (18).  $\square$

### Задачи

1. В пространстве  $R^2$ , состоящем из элементов  $x = (x_1, x_2)$ , которые будем изображать точками плоскости  $Ox_1x_2$ , нарисовать совокупность точек, для которых  $\|x\| = 1$  (т. е. единичную окружность), в случае, если норма понимается в каждом из следующих трех смыслов:

$$\begin{aligned}\|x\| &= \|x\|_1 = \max |x_j|, \\ \|x\| &= \|x\|_2 = |x_1| + |x_2|, \\ \|x\| &= \|x\|_3 = (x_1^2 + x_2^2)^{1/2}.\end{aligned}$$

2. Пусть  $A: R^n \rightarrow R^n$  — произвольный линейный оператор, а  $\lambda_j$  — какое-нибудь его собственное число.

Доказать, что тогда при произвольном выборе нормы в  $R^n$  соответствующая норма  $\|A\|$  оператора  $A$  удовлетворяет неравенству  $\|A\| \geq |\lambda_j|$ .

3. Привести пример линейного оператора  $A: R^2 \rightarrow R^2$ , заданного матрицей с собственными числами  $\lambda_1 = \lambda_2 = 1$  и такого, что  $\|A\|_k > 1000$  ( $k = 1, 2, 3$ ).

4. Пусть  $R^n$  — евклидово пространство,  $A, B$  — произвольные линейные операторы.

Доказать, что тогда  $\|AB\| \leq \|A\| \|B\|$ .

5. Пусть  $R^n$  — евклидово пространство,  $A$  — произвольный оператор,  $B$  — ортогональный линейный оператор, т. е.  $(Bx, Bx) \equiv (x, x)$ .

Доказать, что тогда  $\|AB\| = \|BA\| = \|A\|$ .

6. Пусть  $R^n$  — евклидово пространство и  $A = A^*$  — самосопряженный линейный оператор  $A: R^n \rightarrow R^n$ . Доказать, что тогда  $\|A^2\| = \|A\|^2$ .

Привести пример, показывающий, что для  $A \neq A^*$  возможно  $\|A^2\| < \|A\|^2$ .

7\*. Пусть  $R^n$  — евклидово пространство и  $A: R^n \rightarrow R^n$  — произвольный линейный оператор. Доказать, что тогда

$$\|A\|^2 = \|AA^*\| = \lambda_{\max}(A^*A),$$

где  $\lambda_{\max}(A^*A)$  — наибольшее собственное число оператора  $C = A^*A$ .

8. Пусть  $R^2$  — двумерное пространство вещественных векторов  $x = (x_1, x_2)$  со скалярным умножением  $(x, y) = x_1y_1 + x_2y_2$  и нормой  $\|x\| = \sqrt{(x, x)} = \sqrt{x_1^2 + x_2^2}$ .

Вычислить соответствующую норму следующих матриц:

$$A = \begin{bmatrix} 1 & \sqrt{6} \\ \sqrt{6} & 2 \end{bmatrix}, \quad A = \begin{bmatrix} 1 & 1 \\ 5 & 2 \end{bmatrix}.$$

9. Пусть  $R^2$  — двумерное пространство вещественных векторов  $x = (x_1, x_2)$  со скалярным произведением  $(x, y) = x_1y_1 + x_2y_2$ , и пусть  $B = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$ .

Проверить, что  $B = B^* > 0$ . Вычислить нормы  $\|A\|_B$  матриц  $A$  из предыдущей задачи, согласованные с  $\|x\|_B^2 = [x, x]_B = (Bx, x)$ .

10\*. Пусть  $R^n$  — евклидово пространство,  $\|x\|^2 = (x, x)$  — соответствующая норма в нем и  $A: R^n \rightarrow R^n$  — произвольный линейный оператор.

Доказать, что  $\|A\| = \|A^*\|$ . Если  $A$  — невырожденный оператор, т. е. если  $A$  имеет обратный оператор  $A^{-1}$ , то  $(A^{-1})^* = (A^*)^{-1}$ ,  $\|(A^*)^{-1}\| = \|(A^{-1})^*\|$ .

### § 3. Обусловленность СЛАУ

Две на первый взгляд похожие системы линейных уравнений могут обладать существенно различными чувствительностями своих решений к погрешностям задания входных данных. Это видно уже для систем второго порядка  $Ax = f$ :

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 &= f_1, \\ a_{21}x_1 + a_{22}x_2 &= f_2. \end{aligned} \tag{1}$$

Будем считать, что  $a_{i1}^2 + a_{i2}^2 = 1$  ( $i = 1, 2$ ). Решение такой системы геометрически интерпретируется как точка пересечения двух прямых на плоскости  $(x_1, x_2)$ .

Пусть системе (1) соответствует пара прямых: в одном случае, как на рис. 10, а, и в другом, как на рис. 10, б. Если немного изменить  $f_1$

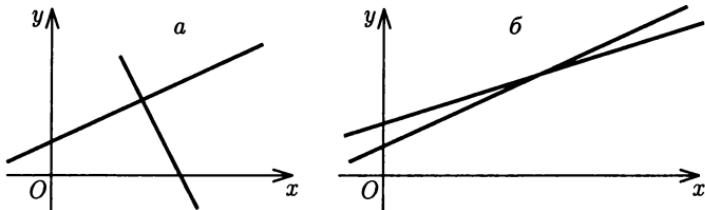


Рис. 10

или коэффициенты  $a_{11}, a_{22}$ , то соответствующая прямая на рис. 10, а или 10, б сдвинется или повернется. При этом в случае рис. 10, а точка пересечения прямых (решение) сдвинется слабо, а в случае

рис. 10, б — значительно сильнее. Чувствительность решения к возмущению входных данных можно охарактеризовать с помощью так называемого числа обусловленности  $\mu(A)$ . В дальнейшем мы увидим, что число обусловленности существенно влияет не только на чувствительность решения к заданию входных данных, но и на число арифметических операций для приближенного вычисления (методами последовательных приближений) решения уравнения  $Ax = f$  с заданной точностью.

**1. Число обусловленности.** Числом обусловленности линейного оператора  $A$ , действующего в нормированном пространстве  $R^n$ , а также числом обусловленности СЛАУ  $Ax = f$  назовем число

$$\mu(A) = \|A\| \cdot \|A^{-1}\|. \quad (2)$$

Если  $A$  — вырожденный оператор, т. е.  $A^{-1}$  не существует, то полагаем  $\mu(A) = \infty$ .

Мы условились отождествлять каждую матрицу  $A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{bmatrix}$

с линейным оператором, действующим в пространстве  $R^n$  элементов  $x$  вида  $x = (x_1, x_2, \dots, x_n)$ :  $y = Ax$ , где  $y = (y_1, y_2, \dots, y_n)$  вычисляется по формулам

$$y_i = \sum_j a_{ij} x_j.$$

Поэтому приведенное определение числа  $\mu(A)$  имеет смысл и для матриц  $A$ , так что можно говорить о числе обусловленности матрицы  $A$  и о числе обусловленности системы линейных уравнений, заданных не только в операторном, но и в каноническом виде.

Нормы  $\|A\|$ ,  $\|A^{-1}\|$  согласованы с нормами, выбранными в  $R^n$ . Поэтому число обусловленности  $\mu(A)$  также согласовано с выбором нормы в пространстве  $R^n$ . Если  $A$  — матрица, а для  $x \in R^n$  используется норма  $\|x\|_1 = \max |x_j|$ , то будем писать  $\mu_1(A)$ ; если  $\|x\|_2 = \sqrt{\sum |x_j|^2}$ , то будем писать  $\mu_2(A)$ .

Если пространство  $R^n$  — евклидово со скалярным произведением  $(x, y)$  и нормой  $\|x\| = (x, x)^{1/2}$ , то будем писать  $\mu_B(A)$ . Если в пространстве  $R^n$  на основе исходного фиксированного скалярного умножения  $(x, y)$  введены новое скалярное умножение  $[x, y]_B = (Bx, y)$  и соответствующая норма  $\|x\|_B = ([x, x]_B)^{1/2}$ , то будем обозначать соответствующее число обусловленности через  $\mu_B(A)$ .

Выясним, в чем состоит геометрический смысл числа обусловленности. Для этого рассмотрим совокупность  $S$  векторов, норма которых равна единице, т. е. единичную сферу. Среди них отметим по одному вектору  $x_{\max}, x_{\min}$ , для которых верны равенства

$$\|Ax_{\max}\| = \max_{x \in S} \|Ax\|, \quad \|Ax_{\min}\| = \min_{x \in S} \|Ax\|.$$

Читатель легко установит, что

$$\|A\| = \|Ax_{\max}\|, \quad \|A^{-1}\| = 1/\|Ax_{\min}\|.$$

Отсюда непосредственно следует, что

$$\mu(A) = \max_{x \in S} \|Ax\| / \min_{x \in S} \|Ax\|. \quad (3)$$

Теперь видно, что всегда

$$\mu(A) \geq 1. \quad (4)$$

Геометрический смысл числа  $\mu(A)$  особенно нагляден, если используется евклидова норма  $\|x\| = (x, x)^{1/2}$ , а размерность пространства  $R^n$  есть 2, т. е.  $R^n$  — плоскость. В этом случае  $S$  — единичная окружность:  $x_1^2 + x_2^2 = 1$ . При линейном преобразовании эта окружность переходит в эллипс. Число  $\mu(A)$  в соответствии с (3) есть отношение большой полуоси этого эллипса к его малой полуоси.

**Теорема 1.** Пусть оператор  $A = A^*$  самосопряжен в смысле скалярного произведения  $[x, y]_B$ . Тогда

$$\mu_B(A) = \frac{|\lambda_{\max}|}{|\lambda_{\min}|}, \quad (5)$$

где  $\lambda_{\max}$  и  $\lambda_{\min}$  — соответственно наибольшее и наименьшее по абсолютной величине собственные числа оператора  $A$ .

**Доказательство.** Пусть  $e_1, e_2, \dots, e_n$  — ортонормированный в смысле скалярного умножения  $[x, y]_B$  базис пространства  $R^n$ , состоящий из собственных векторов оператора  $A$ , а вещественные числа  $\lambda_j$  ( $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n| > 0$ ) — соответствующие собственные числа,  $Ae_j = \lambda_j e_j$ . Тогда каждый вектор  $x$  можно записать в виде

$$x = x_1 e_1 + x_2 e_2 + \dots + x_n e_n,$$

причем

$$\begin{aligned} Ax &= \lambda_1 x_1 e_1 + \lambda_2 x_2 e_2 + \dots + \lambda_n x_n e_n, \\ \|Ax\|_B &= (\lambda_1^2 |x_1|^2 + \lambda_2^2 |x_2|^2 + \dots + \lambda_n^2 |x_n|^2)^{1/2}. \end{aligned}$$

Очевидно, что при условии  $x \in S$ , т. е. при  $\|x\|_B = 1$ ,

$$\max_{x \in S} \|Ax\|_B = |\lambda_1| = |\lambda_{\max}|, \quad \min_{x \in S} \|Ax\|_B = |\lambda_n| = |\lambda_{\min}|,$$

так что в силу (3) справедливо (5).  $\square$

**2. Число  $\mu(A)$  как характеристика системы  $Ax = f$ .**

**Теорема 2.** Пусть правая часть линейного уравнения

$$Ax = f, \quad x \in R^n, \quad f \in R^n, \quad (6)$$

где  $A$  — невырожденный линейный оператор, получила возмущение  $\Delta f$ . При этом решение  $x$  уравнения получит некоторое приращение  $\Delta x$ , так что

$$A(x + \Delta x) = f + \Delta f. \quad (7)$$

Тогда относительная погрешность  $\|\Delta x\|/\|x\|$  решения удовлетворяет неравенству

$$\frac{\|\Delta x\|}{\|x\|} \leq \mu(A) \frac{\|\Delta f\|}{\|f\|}, \quad (8)$$

причем существуют такие  $f$  и  $\Delta f$ , при которых в (8) достигается строгое равенство.

**Доказательство.** Из выражений (7), (6) следует  $A(\Delta x) = \Delta f$ ,  $\Delta x = A^{-1}(\Delta f)$ . Используем еще  $Ax = f$ . Тогда

$$\begin{aligned} \frac{\|\Delta x\|}{\|x\|} &= \frac{\|A^{-1}(\Delta f)\|}{\|x\|} = \frac{\|Ax\|}{\|x\|} \cdot \frac{\|A^{-1}(\Delta f)\|}{\|\Delta f\|} \cdot \frac{\|\Delta f\|}{\|Ax\|} = \\ &= \frac{\|Ax\|}{\|x\|} \cdot \frac{\|A^{-1}(\Delta f)\|}{\|\Delta f\|} \cdot \frac{\|\Delta f\|}{\|f\|}. \end{aligned} \quad (9)$$

Но

$$\frac{\|Ax\|}{\|x\|} \leq \|A\|, \quad \frac{\|A^{-1}(\Delta f)\|}{\|\Delta f\|} \leq \|A^{-1}\|, \quad (10)$$

так что в силу (9), (10) при любых  $f$ ,  $\Delta f$  справедливо неравенство

$$\frac{\|\Delta x\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|\Delta f\|}{\|f\|} = \mu(A) \frac{\|\Delta f\|}{\|f\|}, \quad (11)$$

т. е. справедливо (8).

Если  $\Delta f$  — тот элемент пространства  $R^n$ , для которого

$$\frac{\|A^{-1}(\Delta f)\|}{\|\Delta f\|} = \|A^{-1}\|,$$

а  $f = Ax$  — тот элемент  $f \in R^n$ , для которого

$$\frac{\|Ax\|}{\|x\|} = \|A\|,$$

то правая часть выражения (9) совпадет с правыми частями (11) и с (8), которые при этих  $f$  и  $\Delta f$  превращаются в строгие равенства.  $\square$

Подчеркнем, что решение уравнения  $Ax = f$  не при всех  $f$  одинаково чувствительно к возмущению  $\Delta f$  правой части. При заданном фиксированном  $f$  может оказаться, что  $\|Ax\|/\|x\| \ll \|A\|$ , так что оценка (9) в этом случае обеспечивает более слабую чувствительность относительной погрешности  $\|\Delta x\|/\|x\|$  решения к погрешности  $\|\Delta f\|/\|f\|$ , чем неравенство (8).

Для отыскания точного значения числа обусловленности нужно уметь найти нормы операторов  $A$ ,  $A^{-1}$ . Это обычно очень трудоемко. Если, например, оператор  $A$  задан своей матрицей и нас интересует  $\mu_1(A)$  или  $\mu_2(A)$ , то нужно найти обратную матрицу  $A^{-1}$ , после чего  $\|A\|_1$ ,  $\|A^{-1}\|_1$  или  $\|A\|_2$ ,  $\|A^{-1}\|_2$  вычисляются согласно формулам (1), (2) из § 2. Еще труднее находить число обусловленности  $\mu_B(A)$  в евклидовой норме, задаваемой каким-либо оператором  $B = B^* > 0$ . Поэтому часто ограничиваются получением оценок для  $\mu(A)$  сверху,

используя ту или иную специфику оператора  $A$ . В дальнейшем нам встретятся примеры таких оценок для  $\mu_B(A)$ .

Здесь укажем класс матриц, для которого удается получить оценку для  $\mu_1(A)$ , не отыскивая матрицу  $A^{-1}$ .

Матрица

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{bmatrix}$$

называется *матрицей с диагональным преобладанием величины  $\delta > 0$* , если

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}| + \delta, \quad i = 1, 2, \dots, n. \quad (12)$$

**Теорема 3.** Пусть  $A$  — матрица с диагональным преобладанием величины  $\delta > 0$ . Тогда существует обратная матрица  $A^{-1}$ , причем ее норма, согласованная с нормой  $\|x\|_1 = \max_j |x_j|$ , удовлетворяет оценке

$$\|A^{-1}\|_1 \leq \frac{1}{\delta}. \quad (13)$$

**Доказательство.** Зададим произвольно  $f = \begin{bmatrix} f_1 \\ \dots \\ f_n \end{bmatrix}$  и предположим, что при этом  $f$  система уравнений  $Ax = f$  имеет решение  $x = \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix}$ , причем  $\max_i |x_i| = |x_k|$ . Выпишем скалярное уравнение с номером  $k$ , входящее в систему  $Ax = f$ :

$$a_{k1}x_1 + a_{k2}x_2 + \dots + a_{kn}x_n = f_k.$$

Учтем, что  $|x_k| \geq |x_i|$  ( $i = 1, 2, \dots, n$ ), и выпишем оценку

$$\begin{aligned} |f_k| &= \left| \sum_j a_{kj}x_j \right| \geq |a_{kk}| |x_k| - \sum_{j \neq k} |a_{kj}| |x_j| \geq \\ &\geq |a_{kk}| |x_k| - \left( \sum_{j \neq k} |a_{kj}| \right) |x_k| = \left( |a_{kk}| - \sum_{j \neq k} |a_{kj}| \right) |x_k| \geq \delta |x_k|. \end{aligned}$$

Отсюда  $|x_k| \leq \frac{1}{\delta} |f_k|$ . Но  $|x_k| = \max_j |x_j| = \|x\|_1$ ,  $|f_k| \leq \max_j |f_j| = \|f\|_1$ . Поэтому

$$\|x\|_1 \leq \frac{1}{\delta} \|f\|_1. \quad (14)$$

В частности, в случае  $f = 0 \in R^n$  отсюда следует, что система  $Ax = 0$  имеет только тривиальное решение  $x = 0$ , а значит, система  $Ax = f$  имеет одно и только одно решение при любом  $f \in R^n$  и существует  $A^{-1}$ .

Оценка (14) означает, что для  $x = A^{-1}f$  при любом  $f$  имеет место оценка

$$\|A^{-1}f\|_1 \leq \frac{1}{\delta} \|f\|_1, \quad \frac{\|A^{-1}f\|_1}{\|f\|_1} \leq \frac{1}{\delta},$$

так что

$$\|A^{-1}\|_1 = \max_f \frac{\|A^{-1}f\|_1}{\|f\|_1} \leq \frac{1}{\delta}. \quad \square$$

**Следствие.** Пусть  $A$  — матрица с диагональным преобладанием величины  $\delta$ . Тогда

$$\mu_1(A) = \|A\|_1 \|A^{-1}\|_1 \leq \frac{1}{\delta} \|A\|_1. \quad (15)$$

### Задачи

1. Доказать, что числа обусловленности  $\mu_1(A)$ ,  $\mu_2(A)$  матрицы  $A$  не изменяются, если в матрице  $A$  поменять местами строки или столбцы.

2. Доказать, что для матрицы  $A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{bmatrix}$  и соответствующей транспонированной матрицы  $A^T = \begin{bmatrix} a_{11} & a_{21} & \dots & a_{n1} \\ \dots & \dots & \dots & \dots \\ a_{1n} & a_{2n} & \dots & a_{nn} \end{bmatrix}$  выполнены равенства  $\mu_1(A) = \mu_2(A^T)$ ,  $\mu_2(A) = \mu_1(A^T)$ .

3. Показать, что число обусловленности оператора  $A$  не меняется при умножении этого оператора на произвольное вещественное число  $k$  ( $k \neq 0$ ).

4. Показать, что  $\mu_B(A) = 1$  в том и только том случае, если выполнено одно из условий:

а)  $A$  — оператор подобия, т. е.  $Ax \equiv_k x$  ( $k \neq 0$ );

б)  $A$  — ортогональный оператор, т. е.  $[Ax, Ax]_B \equiv [x, x]_B$ ;

в)  $A$  — произведение оператора подобия и ортогонального оператора.

5\*. Показать, что  $\mu_B(A) = \mu_B(A_B^*)$ , где  $A_B^*$  — оператор, сопряженный оператору  $A$  в смысле скалярного умножения  $[x, y]_B$ .

6. Пусть  $A$  — матрица, причем  $\det A \neq 0$ . Умножим одну из строк на некоторое число  $k$  и обозначим результат через  $A_k$ .

Показать, что  $\mu(A_k) \rightarrow \infty$  при  $k \rightarrow \infty$ .

7\*. Доказать, что для любого линейного оператора  $A$

$$\mu_B(A_B^* A) = (\mu_B(A))^2,$$

где  $A_B^*$  — оператор, сопряженный оператору  $A$  в смысле скалярного умножения  $[x, y]_B$ .

8\*. Пусть  $A = A^*$ ,  $B = B^* > 0$  в смысле некоторого скалярного умножения  $(x, y)$ . Пусть для всех  $x \in R^n$  выполняются неравенства

$$\gamma_1(Bx, x) \leq (Ax, x) \leq \gamma_2(Bx, x),$$

$\gamma_1 > 0$ ,  $\gamma_2 > 0$  — некоторые числа. Рассмотрим оператор  $C = B^{-1}A$ .

Доказать, что справедливо неравенство

$$\mu_B(C) \leq \gamma_2/\gamma_1.$$

Примечание. Решение этой задачи мы приведем в гл. 5 в связи с ее приложениями.

#### § 4. Методы исключения Гаусса

Опишем простой метод исключения Гаусса и метод Гаусса с выбором главного элемента для точного решения систем линейных уравнений, заданных в каноническом виде

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= f_1, \\ \dots &\dots \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n &= f_n. \end{aligned} \quad (1)$$

**1. Простой метод Гаусса.** Из первого уравнения системы (1) выразим  $x_1$  через остальные. Получим

$$x_1 = a'_{12}x_2 + \dots + a'_{1n}x_n + f'_1. \quad (2)$$

Подставим это выражение вместо  $x_1$  в остальные ( $n - 1$ ) уравнений и получим систему ( $n - 1$ ) уравнений относительно ( $n - 1$ ) неизвестных  $x_2, x_3, \dots, x_n$ . Из первого уравнения этой системы выразим  $x_2$  через остальные:

$$x_2 = a'_{23}x_3 + \dots + a'_{2n}x_n + f'_2. \quad (3)$$

Действуя аналогично, при  $k = 3, 4, \dots, n - 1$  получаем

$$x_k = a'_{k,k+1}x_{k+1} + \dots + a'_{kn}x_n + f'_k, \quad (4)$$

а при  $k = n$

$$x_n = f'_n. \quad (5)$$

Равенство (5) дает значение  $x_n$ , а затем по формулам (4) последовательно найдем  $x_{n-1}, x_{n-2}, \dots, x_1$ .

Приведенный алгоритм может оказаться нереализуемым из-за деления на нуль или дать грубую ошибку в результате округлений. Объясним это.

Если  $x_1$  входит в первое уравнение системы (1) с нулевым коэффициентом  $a_{11}$ , то уже запись (2) невозможна, так как  $a'_{12} = -a_{12}/a_{11}$ . Необходимость деления на нуль может встретиться на любом шаге процесса. Если деление на нуль не встречалось и формулы (4), (5) при  $k = 1, 2, \dots, n$  получены, то может возникнуть вычислительная неустойчивость при вычислении  $x_n, x_{n-1}, \dots, x_1$ . Если, например,  $a'_{k,k+1} = 2$  ( $k = n - 1, n - 2, \dots, 1$ ), а остальные  $a'_{ij}$  обращаются в нуль, то погрешность  $\epsilon$ , допущенная при вычислении  $x_n$ , возрастет вдвое при вычислении  $x_{n-1}$ , еще вдвое возрастет при вычислении  $x_{n-2}$  и в  $2^{n-1}$  раз при вычислении  $x_1$ . Уже при небольшом  $n = 11$  погрешность возрастает в тысячу раз.

Нужно иметь в виду и опасность того, что числа  $f'_k$  могут быстро возрастать с ростом  $k$ . В этом случае малые относительные погрешности, допущенные в  $f'_k$  при вычислении по формуле (4), вносят большую по абсолютной величине погрешность в  $x_k$ .

Сформулируем достаточное условие, гарантирующее вычислительную устойчивость метода Гаусса.

**Теорема 1.** Пусть матрица  $A$  системы (1) является матрицей с диагональным преобладанием величины  $\delta > 0$ . Тогда в алгоритме простого метода Гаусса не встретится деление на нуль. Кроме того, выполнены неравенства

$$\sum_{j=1}^{n-k} |a'_{k,k+j}| < 1, \quad b_{k+1,0} = -1, \quad b_{k+1,1} = 1. \quad (7)$$

$$j=1 \quad \quad \quad k = 1, 2, \dots, n-1. \quad (7)$$

$$|f'_k| \leq \frac{2}{\delta} \max_j |f_j|, \quad (8)$$

Доказательству теоремы предпошлем следующую лемму.

**Лемма 1.** Пусть

$$\begin{aligned} b_{11}t_1 + b_{12}t_2 + \dots + b_{1r}t_r &= \varphi_1, \\ \dots \dots \dots \dots \dots \dots \dots & \\ b_{r1}t_1 + b_{r2}t_2 + \dots + b_{rr}t_r &= \varphi_r, \end{aligned} \tag{8}$$

где  $r > 1$ , есть некоторая система линейных уравнений с диагональным преобладанием величины  $\delta$ :

$$|b_{ll}| \geq \sum_{i \neq l} |b_{lj}| + \delta. \quad (9)$$

Тогда при получении из первого уравнения системы (8) выражения для  $t_1$ :

$$t_1 = b'_{12} t_2 + \dots + b'_{1r} t_r + \varphi'_1 \quad (10)$$

имеют место следующие три факта:

не встретится деление на нуль;

выполняется неравенство

$$\sum_{i=2}^r |b'_{1j}| < 1; \quad (11)$$

система  $(r-1)$ -го порядка относительно  $t_2, t_3, \dots, t_r$ , полученная из (8) исключением  $t_1$  по формуле (10), окажется системой с диагональным преобладанием (той же, что и в (9), величины  $\delta$ ).

**Доказательство.** В силу (9) для  $l = 1$

$$|b_{11}| \geq |b_{12}| + |b_{13}| + \dots + |b_{1r}| + \delta.$$

Отсюда  $b_{11} \neq 0$ , выражение (10) со значениями

$$b'_{ij} = -b_{1j}/b_{11} \quad (j = 2, 3, \dots, r), \quad \varphi'_1 = \varphi_1/b_{11}$$

имеет смысл, а также выполняется (11).

Докажем последнее, третье утверждение леммы. Действительно, подставляя (10) в уравнение с номером  $j$  ( $j > 1$ ) системы (8), получаем после приведения подобных уравнение

$$(b_{j2} + b_{j1}b'_{12})t_2 + (b_{j3} + b_{j1}b'_{13})t_3 + \dots + (b_{jr} + b_{j1}b'_{1r})t_r = \varphi'_j, \quad j = 2, 3, \dots, r. \quad (12)$$

В этой системе  $(r - 1)$ -го порядка уравнением с номером  $s$  ( $s = 1, 2, \dots, r - 1$ ) окажется уравнение

$$(b_{s+1,2} + b_{s+1,1}b'_{12})t_2 + (b_{s+1,3} + b_{s+1,1}b'_{13})t_3 + \dots + (b_{s+1,r} + b_{s+1,1}b'_{1r})t_r = \varphi'_{s+1}, \quad s = 1, 2, \dots, r - 1. \quad (13)$$

В матрице этой системы в строке с номером  $s$  стоят числа

$$(b_{s+1,2} + b_{s+1,1}b'_{12}), \quad (b_{s+1,3} + b_{s+1,1}b'_{13}), \quad \dots, \quad (b_{s+1,r} + b_{s+1,1}b'_{1r}),$$

причем на диагонали окажется элемент  $b_{s+1,s+1} + b_{s+1,1}b'_{1,s+1}$ .

Покажем, что имеет место диагональное преобладание (величины  $\delta$ ), т. е. что справедливо неравенство

$$|b_{s+1,s+1} + b_{s+1,1}b'_{1,s+1}| \geq \sum_{\substack{j=2 \\ j \neq s+1}}^r |b_{s+1,j} + b_{s+1,1}b'_{1j}| + \delta. \quad (14)$$

Докажем более сильное неравенство

$$|b_{s+1,s+1}| - |b_{s+1,1}b'_{1,s+1}| \geq \sum_{\substack{j=2 \\ j \neq s+1}}^r (|b_{s+1,j}| + |b_{s+1,1}b'_{1j}|) + \delta.$$

Оно равносильно неравенству

$$|b_{s+1,s+1}| \geq \sum_{\substack{j=2 \\ j \neq s+1}}^r |b_{s+1,j}| + |b_{s+1,1}| \sum_{j=2}^r |b'_{1j}| + \delta. \quad (15)$$

В силу неравенства (11) неравенство (15) заведомо будет выполнено, если выполняется неравенство, получающееся из (15) заменой  $\sum |b'_{1j}|$  числом 1:

$$|b_{s+1,s+1}| \geq \sum_{\substack{j=2 \\ j \neq s+1}}^r |b_{s+1,j}| + |b_{s+1,1}| + \delta = \sum_{\substack{j=1, \\ j \neq s+1}}^r |b_{s+1,j}| + \delta. \quad (16)$$

Но это неравенство справедливо по условию (9), так что неравенство (14) доказано.  $\square$

Доказательство теоремы 1. Сначала докажем формулу (4) и неравенство (6). Воспользуемся индукцией по  $k$ . В случае  $k = 1, n > 1$  справедливость формулы (4) и неравенства (6) установлена в первых

двух утверждениях леммы. При этом в силу третьего утверждения леммы система  $(n - 1)$ -го порядка относительно  $x_2, x_3, \dots, x_n$ , полученная из (1) исключением  $x_1$  по формуле (2), оказывается системой с диагональным преобладанием величины  $\delta$ . Пусть формула (4) и неравенство (6) уже доказаны для  $k = 1, 2, \dots, s$  ( $s < n - 1$ ), причем доказано также, что система порядка  $n - s$  относительно  $x_{s+1}, x_{s+2}, \dots, x_n$ , полученная из (1) исключением  $x_1, x_2, \dots, x_s$  по формуле (4), обладает диагональным преобладанием величины  $\delta$ . Рассмотрев эту систему в качестве системы (8) и использовав лемму, установим справедливость предположения индукции для  $k = s + 1$ . Доказательство (4), (6) по индукции окончено.

Докажем теперь (7). В силу теоремы 3 из § 3 для решения системы (1) справедливо неравенство

$$\max_j |x_j| \leq \frac{1}{\delta} \max_j |f_j|. \quad (17)$$

Отсюда и из (4) при любом  $k = 1, 2, \dots, n$  следует неравенство

$$\begin{aligned} |f'_k| &= \left| x_k - \sum_j a'_{kj} x_j \right| \leq |x_k| + \sum_{j=k+1}^n |a'_{kj}| |x_j| \leq \\ &\leq \max_{1 \leq j \leq n} |x_j| \left( 1 + \sum_{j=k+1}^n |a'_{kj}| \right) \leq 2 \max_j |x_j| \leq \frac{2}{\delta} \max_j |f_j|, \end{aligned}$$

которое совпадает с (7).  $\square$

Подчеркнем, что условие теоремы достаточно, но не необходимо для применимости простого метода Гаусса. Если для какой-либо системы (1) простой алгоритм метода Гаусса удалось применить на компьютере, то найденное при этом решение не будет точным только в силу конечной разрядности компьютера и связанных с этим неизбежных погрешностей округления.

Подставляя найденное приближенное решение системы (1) в левую часть и вычисляя невязки, можно судить о погрешности решения, воспользовавшись оценкой

$$\frac{\|\Delta x\|}{\|x\|} \leq \mu(A) \frac{\|\Delta f\|}{\|f\|},$$

если только число обусловленности  $\mu(A)$  или оценка для него известны.

**Теорема 2.** Порядок числа  $\nu$  арифметических операций для реализации алгоритма Гаусса есть

$$\nu = O(n^3). \quad (18)$$

Доказательство состоит в прямом подсчете.

**2. Прогонка.** Изложенный метод Гаусса особенно эффективен для системы (1) следующего специального вида:

$$\begin{aligned} b_1x_1 + c_1x_2 &= f_1, \\ a_2x_1 + b_2x_2 + c_2x_3 &= f_2, \\ a_3x_2 + b_3x_3 + c_3x_4 &= f_3, \\ \dots & \\ a_{n-1}x_{n-2} + b_{n-1}x_{n-1} + c_{n-1}x_n &= f_{n-1}, \\ a_nx_{n-1} + b_nx_n &= f_n, \end{aligned} \quad (19)$$

матрица которой трехдиагональна. Условия диагонального преобладания в этом случае принимают вид

$$\begin{aligned} |b_1| &\geq |c_1| + \delta, \\ |b_k| &\geq |a_k| + |c_k| + \delta, \quad k = 2, 3, \dots, n-1, \\ |b_n| &\geq |a_n| + \delta. \end{aligned} \quad (20)$$

Уравнения (4) принимают вид

$$x_k = A_k x_{k+1} + F_k, \quad k = 1, 2, \dots, n-1, \quad (21)$$

$$x_n = F_n, \quad (22)$$

где  $A_k, F_k$  — некоторые коэффициенты.

Положим по определению  $A_n = 0$  и условимся записывать (21), (22) единой формулой

$$x_k = A_k x_{k+1} + F_k, \quad k = 1, 2, \dots, n. \quad (23)$$

Очевидно, что в (23)

$$A_1 = -\frac{c_1}{b_1}, \quad F_1 = \frac{f_1}{b_1}. \quad (24)$$

Допустим, что уже вычислены  $A_k, F_k$  при некотором  $k$  ( $1 \leq k \leq n-1$ ). Подставляя выражение  $x_k = A_k x_{k+1} + F_k$  в  $(k+1)$ -е уравнение системы (19), получаем

$$x_{k+1} = -\frac{c_{k+1}}{b_{k+1} + a_{k+1}A_k} x_{k+2} + \frac{f_{k+1} - a_{k+1}F_k}{b_{k+1} + a_{k+1}A_k},$$

причем полагаем  $c_n = 0$ . Отсюда получаются рекуррентные соотношения

$$\begin{aligned} A_{k+1} &= \frac{-c_{k+1}}{b_{k+1} + a_{k+1}A_k}, \\ F_{k+1} &= \frac{f_{k+1} - a_{k+1}F_k}{b_{k+1} + a_{k+1}A_k}, \quad k = 1, 2, \dots, n-1. \end{aligned} \quad (25)$$

Процесс вычисления решения системы (19) распадается на вычисление прогоночных коэффициентов  $A_j, F_j$  ( $j = 1, 2, \dots, n$ ) по формулам (24), (25), а затем последовательное вычисление значений  $x_n, x_{n-1}, \dots, x_1$  по формулам (23) при  $k = n, n-1, \dots, 1$ .

Описанный алгоритм метода Гаусса для системы с трехдиагональной матрицей часто называют *прогонкой*. Коэффициенты  $A_k$ ,  $F_k$  называют *прогоночными коэффициентами*, процесс их вычисления — *прямой прогонкой*. Процесс вычисления компонент решения  $x_n, x_{n-1}, \dots, x_1$  с помощью прогоночных соотношений (23) называют *обратной прогонкой*.

Подсчитаем порядок числа арифметических действий для реализации прогонки. Для вычисления прогоночных коэффициентов требуется  $O(n)$  арифметических операций. Для обратной прогонки требуется  $n$ -кратное использование формулы (23), т. е. также  $O(n)$  арифметических операций. Общее число арифметических операций есть  $O(n)$ . Очевидно, для решения системы (19) не существует алгоритма более экономного по порядку числа арифметических действий, поскольку число  $n$  неизвестных также есть  $O(n)$ .

*Прогонкой* называют простой алгоритм метода Гаусса и в том случае, когда матрица  $A$  системы (1) содержит ненулевые элементы на  $m$  ( $3 < m \ll n$ ) соседних диагоналях, среди которых есть главная. Если  $m$  фиксировано, а  $n$  может быть любым, то число арифметических операций также есть  $O(n)$ .

Системы вида (19) высоких порядков привлекли особое внимание в начале 50-х годов XX века в связи с тем, что они возникли при использовании так называемых неявных разностных схем для уравнения теплопроводности. О таких схемах и их роли будет рассказано в части III.

И.М. Гельфанд и О.В. Локуциевский предложили алгоритм прогонки для численного решения этих систем, установив его полную адекватность задаче. Они построили континуальное замыкание алгоритма прогонки и установили вычислительную устойчивость этого алгоритма (см. И.М. Гельфанд и О.В. Локуциевский, дополнение к [5]).

Указанная работа — одна из первых, где был четко поставлен и решен вопрос об устойчивости вычислительного алгоритма. Этот вопрос приобрел особую остроту при вычислениях на компьютере, когда миллионы и миллиарды операций осуществляются в автоматическом режиме без контроля человека.

Теорема о достаточных условиях применимости метода Гаусса, доказанная в этом параграфе, является обобщением результата И.М. Гельфанда и О.В. Локуциевского об устойчивости прогонки на случай систем с заполненной матрицей.

Отметим, что идею переноса (“прогонки”) условия  $b_1x_1 + c_1x_2 = f_1$ , заданного первым уравнением трехдиагональной системы (19), реализуемую при получении прогоночных коэффициентов и соотношений (23), стали использовать и в других методах исключения, которые также называют методами прогонки [21].

**3. Метод Гаусса с выбором главного элемента.** Как мы видели, простой метод Гаусса может натолкнуться на препятствия: в процессе его реализации может встретиться деление на нуль или произойти

потеря точности из-за вычислительной неустойчивости. Приведем модификацию метода исключения Гаусса — метод Гаусса с выбором главного элемента, который в случае  $\det A \neq 0$  гарантирует от деления на нуль и повышает вычислительную устойчивость по сравнению с простым методом Гаусса.

Сначала отыскивается самый большой (один из самых больших, если их несколько) по модулю элемент матрицы  $A$ . Допустим, это элемент  $a_{ij}$  — коэффициент при  $x_j$  в уравнении с номером  $i$ . Очевидно, что  $a_{ij}$  отличен от нуля. Осуществляем перенумерацию уравнений и неизвестных так, чтобы уравнение с номером  $i$  стало первым и неизвестное  $x_j$  стало первым:  $x'_1 = x_j$ . После этого осуществляем первый шаг описанного выше простого метода Гаусса (без выбора главного элемента). При рассмотрении полученной системы порядка  $n - 1$  вновь производим выбор главного элемента, перенумерацию уравнений и неизвестных и т. д. Получаем систему равенств вида (4), (5), которые позволяют вычислить все компоненты решения, начиная с последней. На компьютерах обычно имеются стандартные программы методов Гаусса.

Простой подсчет показывает, что метод Гаусса с выбором главного элемента, как и простой метод Гаусса, требует  $O(n^3)$  арифметических операций.

Отметим еще, что метод Гаусса позволяет находить обратную матрицу  $A^{-1}$ . Столбец с номером  $j$  матрицы  $A^{-1}$  обозначаем  $\begin{bmatrix} x_{1j} \\ \dots \\ x_{nj} \end{bmatrix}$ .

Из равенства  $AA^{-1} = E$  следует, что этот столбец является решением системы

$$A \begin{bmatrix} x_{1j} \\ \dots \\ x_{nj} \end{bmatrix} = \begin{bmatrix} \delta_{1j} \\ \dots \\ \delta_{nj} \end{bmatrix}, \quad j = 1, 2, \dots, n, \quad (26)$$

где  $\delta_{kj} = 0$ , если  $k \neq j$ ,  $\delta_{jj} = 1$ .

Системы (26), которых имеется  $n$  экземпляров, различаются только обозначениями неизвестных и правыми частями. Поэтому большая часть вычислений по методу Гаусса для них совпадает. Порядок числа арифметических операций для отыскания  $A^{-1}$  при правильной организации расчета остается  $O(n^3)$ , как для решения одной системы.

**4. Замечание о других универсальных методах точного решения.** Существует ряд других методов, в которых тем или иным способом исходная система (1) приводится к более простому виду (в методе Гаусса система (1) приводится к треугольному виду (4)), а затем вычисляется решение этой «простой» системы.

Приведение к этому более простому виду можно осуществить с помощью умножения системы  $Ax = f$  на подходящую ортогональную матрицу  $C$  (т. е. матрицу, для которой  $(Cx, Cx) = (x, x)$ ). Тогда число обусловленности  $\mu_3(CA)$  новой матрицы системы  $CAx = Cf$

совпадает с числом обусловленности  $\mu_3(A)$  (докажите равенство  $\mu_3(CA) = \mu_3(A)$ ), так что вместо  $\mu_3(A)$  можно искать  $\mu_3(CA)$ , что может оказаться проще. Таким образом, основываясь на теореме 1 из § 3, одновременно с вычислением решения можно оценить погрешность, вызванную неточным заданием входных данных  $f_i$ .

**5. Об алгоритме с гарантированной оценкой погрешности.** При расчете на реальном компьютере с заданным числом разрядов наряду с влиянием неточного задания входных данных на каждой арифметической операции возникают погрешности округления. Влияние этих погрешностей округления на результат зависит не только от разрядности компьютера, но и от числа обусловленности матрицы системы, и от выбранного алгоритма. В [3] построен алгоритм, который учитывает влияние погрешностей округления на данном компьютере и выдает результат с гарантированной точностью либо в процессе вычислений устанавливает, что данная система обусловлена настолько плохо, что при расчете на компьютере с заданной разрядностью какая-либо точность не может быть гарантирована.

### Задачи

1. Вычислить решение системы

$$10^{-3}x + y = 5, \quad x - y = 6$$

простым методом Гаусса и с выбором главного элемента. Вычисления вести с двумя значащими цифрами. Сравнить и объяснить результаты.

2. Для численного решения краевой задачи

$$\frac{d^2x}{dt^2} - p^2(t)x = f(t), \quad 0 < t < 1, \quad p(t) \neq 0,$$

$$x(0) = \varphi, \quad x(1) = \psi,$$

разобьем отрезок  $0 \leq x \leq 1$  на  $N$  равных частей и будем искать приближенно таблицу значений решения  $x_0, x_1, \dots, x_N$  в точках разбиения  $t_n = nh$  ( $n = 0, 1, \dots, N$ ;  $h = N^{-1}$ ). В точках  $t_n$  ( $n = 1, 2, \dots, N - 1$ ) заменим производную разностным отношением

$$\left. \frac{d^2x}{dt^2} \right|_{t=nh} \approx \frac{x_{n+1} - 2x_n + x_{n-1}}{h^2}.$$

Получим вместо исходной задачи ее разностный аналог

$$x_0 = \varphi,$$

$$\frac{x_{n+1} - 2x_n + x_{n-1}}{h^2} - p^2(t_n)x_n = f(t_n), \quad n = 1, 2, \dots, N - 1, \quad (25)$$

$$x_N = \psi.$$

а) Показать, что для решения этой системы линейных уравнений можно и целесообразно воспользоваться прогонкой.

б) Составить программу вычисления решения системы (25) на компьютере, если  $p(t) = 1 + t$ ,  $f(t) = e^t$ .

## § 5. Связь между задачей на минимум квадратичной функции и СЛАУ

Пусть  $R^n$  есть  $n$ -мерное евклидово пространство элементов  $x$ . Зададим какой-либо линейный оператор  $A$ ,  $f \in R^n$ , число  $c$  и рассмотрим функцию от  $x$  вида

$$F(x) = (Ax, x) - 2(f, x) + c, \quad (1)$$

называемую *квадратичной функцией*.

Заметим, что эта квадратичная функция в силу  $(Ax, x) \equiv (x, A^*x) \equiv (A^*x, x)$  совпадает с квадратичной функцией

$$F(x) = (A^*x, x) - 2(f, x) + c,$$

а значит, и с функцией

$$F(x) = \left( \frac{A + A^*}{2} x, x \right) - 2(f, x) + c.$$

Без ограничения общности будем считать, что в выражении (1) оператор  $A$  самосопряженный. В противном случае мы заменили бы его самосопряженным оператором  $(A + A^*)/2$ .

Предположим, что  $A = A^*$ ,  $A > 0$ , т. е.  $(Ax, x) > 0$  в случае  $x \neq 0$ . Поставим задачу об отыскании элемента  $z \in R^n$ , придающего наименьшее значение функции  $F(x)$ :

$$F(z) = \min_{x \in R^n} F(x). \quad (2)$$

Задача (2) о минимуме квадратичной функции и задача об отыскании решения системы линейных уравнений

$$Ax = f, \quad A = A^* > 0, \quad (3)$$

равносильны. Сформулируем соответствующее утверждение в виде теоремы.

**Теорема 1.** Пусть  $A = A^* > 0$ . Существует один и только один элемент  $z \in R^n$ , придающий наименьшее значение квадратичной функции (1). Этот элемент есть решение уравнения (3).

**Доказательство.** В силу положительной определенности оператора  $A$  он невырожден, а следовательно, уравнение (3) имеет одно и только одно решение  $z \in R^n$ .

Покажем, что при любом  $\delta \in R^n$  ( $\delta \neq 0$ ) будет  $F(z + \delta) > F(z)$ :

$$\begin{aligned} F(z + \delta) &= A(z + \delta), z + \delta) - 2(f, z + \delta) + c = \\ &= [(Az, z) - 2(f, z) + c] + 2(Az, \delta) - 2(f, \delta) + (A\delta, \delta) \equiv \\ &\equiv F(z) + 2(Az - f, \delta) + (A\delta, \delta) = F(z) + (A\delta, \delta) > F(z). \square \end{aligned}$$

Установленная равносильность задач (1) и (3) позволяет сводить решение любой из них к решению другой.

Линейные уравнения вида (3) с самосопряженным и положительно определенным оператором  $A$  представляют собой важный класс СЛАУ по двум причинам.

Во-первых, СЛАУ

$$Cx = \varphi \quad (4)$$

с произвольным невырожденным линейным оператором  $C: R^n \rightarrow R^n$  сводится к СЛАУ вида (3). Достаточно положить  $A = C * C$ ,  $f = C * \varphi$ .

Во-вторых, многие краевые задачи для эллиптических уравнений являются задачами Лагранжа-Эйлера для некоторых задач о минимуме квадратичных функционалов. Поэтому естественно, что при «правильной» дискретизации этих вариационных задач возникает задача о минимуме квадратичной функции в конечномерном пространстве, которая в силу доказанной теоремы приводит к СЛАУ (3).

### Задачи

1. Задана квадратичная функция скалярных аргументов  $x_1, x_2$ :

$$F(x_1, x_2) = x_1^2 + 2x_1x_2 + 4x_2^2 - 2x_1 + 3x_2 + 5.$$

Рассмотреть эту функцию как функцию вектора  $x = (x_1, x_2)$  евклидова пространства  $R^2$  со скалярным умножением  $(x, y) = x_1y_1 + x_2y_2$  и записать ее в форме (1). Проверить, что  $A > 0$ , и найти решение соответствующей задачи (3).

- 2\*. Записать функцию  $F(x_1, x_2)$  из предыдущей задачи в виде

$$F(x) = [Ax, x]_B - 2[f, x]_B + C,$$

где

$$[x, y]_B = (Bx, y), \quad B = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}, \quad A = A_B^*.$$

## § 6. Метод сопряженных градиентов как метод точного решения СЛАУ

Рассмотрим уравнение

$$Ax = f, \quad A = A^* > 0, \quad f \in R^n, \quad x \in R^n, \quad (1)$$

где  $R^n$  — пространство со скалярным произведением  $(x, y)$  и  $A: R^n \rightarrow R^n$  — самосопряженный в смысле этого скалярного умножения положительно определенный оператор.

**1. Метод сопряженных градиентов вычисления решения.** Зададим произвольно  $x_0 \in R^n$  и построим последовательность

$$\begin{aligned} x_1 &= (E - \tau_1 A)x_0 + \tau_1 f, \\ x_{k+1} &= \alpha_{k+1}(E - \tau_{k+1} A)x_k + (1 - \alpha_{k+1})x_{k-1} + \alpha_{k+1}\tau_{k+1}f, \end{aligned} \quad (2)$$

где

$$\begin{aligned} \tau_{k+1} &= \frac{(r_k, r_k)}{(Ar_k, r_k)}, \quad r_k = Ax_k - f, \quad k = 0, 1, \dots, \\ \alpha_1 &= 1, \quad \alpha_{k+1} = \left(1 - \frac{\tau_{k+1}}{\tau_k} \cdot \frac{(r_k, r_k)}{(r_{k-1}, r_{k-1})} \cdot \frac{1}{\alpha_k}\right)^{-1}, \quad k = 1, 2, \dots. \end{aligned} \quad (3)$$

Оказывается, что существует номер  $k_0$  ( $k_0 \leq n$ ), такой, что член  $x_{k_0}$  последовательности (2) совпадает с точным решением  $x$  СЛАУ (1):

$$x = x_{k_0}, \quad k_0 \leq n. \quad (4)$$

Отметим также, что  $x_1, x_2, \dots, x_k, \dots$  являются уточняющимися с ростом номера  $k$  последовательными приближениями к решению; для заданного малого  $\epsilon > 0$  погрешность  $\|x - x_k\|$  приближения  $x_k$  при условиях, которые будут указаны в § 2 гл. 5, может стать меньше  $\epsilon$  уже при  $k \ll n$ .

В настоящее время метод сопряженных градиентов при «умеренном» числе обусловленности  $\mu(A)$  и больших  $n$  на практике используется обычно именно как метод последовательных приближений.

Вычисление точного решения при большом числе обусловленности  $\mu(A)$  и большом  $n$  с помощью последовательности (2), (3) и равенства (4) может натолкнуться на препятствие, состоящее в возможной потере вычислительной устойчивости при нахождении членов  $x_k$  последовательности (2). Однако для хорошо обусловленных систем и умеренных  $\mu(A)$  метод (2)–(4) обладает достоинствами по сравнению с методами исключения, которые мы опишем в пп. 2, 3.

**2. Произвольность формы задания оператора  $A$ .** Очевидно, что в (2) используется только возможность по заданному элементу  $y \in R^n$  находить элемент  $z = Ay$  ( $z \in R^n$ ), а также по заданным произвольным  $y \in R^n$ ,  $z \in R^n$  вычислять их скалярное произведение, т. е. число  $(y, z)$ . Поэтому не обязательно, чтобы система (1) была задана в каноническом виде

$$\sum a_{ij}x_j = f_i, \quad i = 1, 2, \dots, n, \quad (5)$$

или была приведена к такому виду. К тому же нет необходимости хранить матрицу  $A$ , которая имела бы  $n^2$  элементов, в то время как векторы  $y \in R^n$ ,  $z = Ay \in R^n$ , записанные в координатной форме, задаются лишь  $n$  числами каждый.

**3. Простота использования многопроцессорных компьютеров.** Пусть система (1) задана в канонической форме (5), так что в записи (1) оператор  $A$  — матрица, а элементы  $x, f \in R^n$  пространства

$R^n$  — наборы чисел  $x = \begin{bmatrix} x_1 \\ \dots \\ x_n \end{bmatrix}$ ,  $y = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix}$ . Вычисление по задан-

ному  $y = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix}$  каждой из  $n$  компонент вектора  $Ay = z = \begin{bmatrix} z_1 \\ \dots \\ z_n \end{bmatrix}$

производится независимо. На многопроцессорном компьютере, допускающем одновременное осуществление многих арифметических операций, вычисление всех  $n$  компонент вектора  $Ay$  по заданному  $y \in R^n$  может быть осуществлено одновременно и совершенно единообразно. Ускорение при использовании многопроцессорных компьютеров может быть получено не только в случае задания системы (1) в канонической форме, но и в случае, если  $R^n$  состоит из сеточных функций,

а значения  $Ay \in R^n$  сеточной функции в каждой точке сетки вычисляются независимо (см., например, оператор  $-\Delta^{(h)} : U^{(h)} \rightarrow U^{(h)}$  в формуле (11) из § 1).

Заметим, что общее число  $\nu$  арифметических операций при решении системы (5) методом сопряженных градиентов имеет порядок  $n^3$ , т. е.  $\nu = O(n^3)$ , как и в методе исключения Гаусса. Действительно, вычисление каждого члена  $x_{n+1}$  последовательности (2) требует применения матрицы  $A$ , т. е.  $O(n^2)$  арифметических операций, а количество  $k_0 \leq n$  таких членов может совпасть с  $n$ . Но вычисление  $Ay$  на многопроцессорном компьютере требует того же времени, что и вычисление одной компоненты при использовании однопроцессорного, так что время может быть сокращено в  $n$  раз.

Заметим, что при использовании метода Гаусса также возможно совмещение во времени многих арифметических операций.

## § 7. Конечные ряды Фурье и запись точного решения разностного аналога задачи Дирихле для уравнения Пуассона

Существуют методы, приспособленные для точного решения лишь узких классов СЛАУ, но зато более эффективные, чем универсальные. Среди них — представление решений СЛАУ в виде конечных рядов Фурье.

Пусть СЛАУ задана в виде

$$Ax = f, \quad x \in R^n, \quad f \in R^n, \quad A: R^n \rightarrow R^n, \quad (1)$$

причем  $A = A^*$ ,  $e_1, e_2, \dots, e_n$  — собственные векторы оператора  $A$ , образующие ортонормированный базис, а

$$\lambda_1, \lambda_2, \dots, \lambda_n, \quad \lambda_j \neq 0, \quad (2)$$

суть соответствующие собственные числа,  $Ae_j = \lambda_j e_j$ .

Каждый самосопряженный невырожденный оператор  $A$  имеет ортонормированный базис, состоящий из собственных векторов оператора  $A$ , а все его собственные значения вещественны и отличны от нуля. Однако здесь мы предполагаем гораздо большее: считаем, что нам фактически известны этот ортонормированный базис, состоящий из собственных векторов, и соответствующие собственные значения. Это предположение весьма сужает класс допустимых самосопряженных операторов.

Итак, пусть выполнены предположения, сформулированные выше. Для вычисления решения системы (1) запишем заданное  $f$  и искомое  $x$  в виде конечных рядов Фурье:

$$f = F_1 e_1 + F_2 e_2 + \dots + F_n e_n, \quad F_j = (f, e_j), \quad (3)$$

$$x = X_1 e_1 + X_2 e_2 + \dots + X_n e_n. \quad (4)$$

Подставляя (3), (4) в уравнение (1), получаем

$$\sum_{j=1}^n (\lambda_j X_j) e_j = \sum_{j=1}^n F_j e_j. \quad (5)$$

Приравнивая коэффициенты при  $e_j$  в левой и правой частях (5), получаем значения коэффициентов  $X_j$  представления искомого решения в виде конечного ряда Фурье (4):

$$X_j = \frac{F_j}{\lambda_j}, \quad j = 1, 2, \dots, n. \quad (6)$$

Конкретизируем эту абстрактную схему применительно к вычислению решения  $u^{(h)}$  разностного аналога задачи Дирихле для уравнения Пуассона в квадратной области (см. п. 3 § 1):

$$-\Delta^{(h)} u^{(h)} = f^{(h)}, \quad u^{(h)} \in U^{(h)}, \quad f^{(h)} \in F^{(h)}. \quad (7)$$

В этом случае удается указать собственные функции и собственные числа оператора  $-\Delta^{(h)}$ .

**1. Ряды Фурье для сеточных функций.** Рассмотрим множество всех вещественных функций  $v = \{v_m\}$ , определенных в точках  $x_m = mh$  ( $m = 0, 1, \dots, M$ ;  $h = 1/M$ ), обращающихся в нуль при  $m = 0, M$ . Совокупность этих функций с обычными операциями сложения и умножения их на вещественные числа образует линейное пространство. Размерность этого пространства есть  $M - 1$ , поскольку система функций

$$\tilde{\psi}_m^{(k)} = \begin{cases} 0, & m \neq k, \\ 1, & m = k, \end{cases} \quad k = 1, 2, \dots, M - 1,$$

очевидным образом образует базис. Действительно, каждую функцию  $v = (v_0, v_1, \dots, v_M)$ , где  $v_0 = v_M = 0$ , можно единственным образом представить в виде линейной комбинации функций  $\tilde{\psi}^{(1)}, \tilde{\psi}^{(2)}, \dots, \tilde{\psi}^{(M-1)}$ :

$$v = v_1 \tilde{\psi}^{(1)} + v_2 \tilde{\psi}^{(2)} + \dots + v_{M-1} \tilde{\psi}^{(M-1)}.$$

Введем в рассматриваемом пространстве скалярное умножение, положив

$$(v, w) = h \sum_{m=0}^M v_m w_m. \quad (8)$$

Покажем, что система функций

$$\psi^{(k)} = \left\{ \sqrt{2} \sin \frac{k\pi m}{M} \right\}, \quad k = 1, 2, \dots, M - 1, \quad (9)$$

образует ортонормированный базис в рассматриваемом пространстве, т. е.

$$(\psi^{(k)}, \psi^{(r)}) = \begin{cases} 0, & k \neq r, \\ 1, & k = r, \end{cases} \quad k, r = 1, 2, \dots, M - 1. \quad (10)$$

Для доказательства заметим, что

$$\sum_{m=0}^{M-1} \cos \frac{l\pi m}{M} = \frac{1}{2} \sum_{m=0}^{M-1} (e^{il\pi m/M} + e^{-il\pi m/M}) = \\ = \frac{1}{2} \cdot \frac{1 - e^{il\pi}}{1 - e^{i\pi/M}} + \frac{1}{2} \cdot \frac{1 - e^{-il\pi}}{1 - e^{-i\pi/M}} = \begin{cases} 0, & l \text{ четно}, \\ 1, & l \text{ нечетно}, \end{cases} \quad 0 < l < 2M.$$

Отсюда при  $k \neq r$  получаем

$$(\psi^{(k)}, \psi^{(r)}) = 2h \sum_{m=0}^M \sin \frac{k\pi m}{M} \sin \frac{r\pi m}{M} = 2h \sum_{m=0}^{M-1} \sin \frac{k\pi m}{M} \sin \frac{r\pi m}{M} = \\ = h \sum_{m=0}^{M-1} \cos \frac{(k-r)\pi m}{M} - h \sum_{m=0}^{M-1} \cos \frac{2k\pi m}{M} = 0,$$

а при  $k = r$

$$(\psi^{(k)}, \psi^{(r)}) = h \sum_{m=0}^{M-1} \cos 0 - h \sum_{m=0}^{M-1} \cos \frac{2k\pi m}{M} = hM - h \cdot 0 = 1.$$

Любая сеточная функция  $v = (v_0, v_1, \dots, v_M)$  разлагается по ортогональному базису (2) в сумму

$$v = c_1 \psi^{(1)} + c_2 \psi^{(2)} + \dots + c_{M-1} \psi^{(M-1)},$$

или

$$v_m = \sqrt{2} \sum_{k=1}^{M-1} c_k \sin \frac{k\pi m}{M}, \quad (11)$$

где

$$c_k = (v, \psi^{(k)}) = \sqrt{2} h \sum_{m=0}^M v_m \sin \frac{k\pi m}{M}.$$

Ясно, что благодаря ортонормированности базиса (2) имеем

$$(v, v) = c_1^2 + c_2^2 + \dots + c_{M-1}^2. \quad (12)$$

Сумма (11) и есть разложение сеточной функции  $v = \{v_m\}$  в конечный ряд Фурье, а равенство (12) — точный аналог равенства Парсеваля в обычной теории рядов Фурье.

Совершенно аналогично можно рассмотреть конечные ряды Фурье для функций на сеточном квадрате. Рассмотрим сетку

$$x_m = mh, \quad y_n = nh, \quad 0 \leq mh \leq 1, \quad 0 \leq nh \leq 1,$$

причем  $h = 1/M$ ,  $M$  — натуральное число. Совокупность вещественных функций  $v = \{v_{mn}\}$ , определенных в точках сетки и

обращающихся в нуль в точках, лежащих на границе квадрата, образует линейное пространство. Введем в нем скалярное умножение

$$(v, w) = h^2 \sum_{n,m=0}^M v_{mn} w_{mn}.$$

В рассматриваемом линейном пространстве размерности  $(M - 1)^2$  система функций

$$\begin{aligned} \psi^{(k,l)} &= 2 \sin \frac{k\pi m}{M} \sin \frac{l\pi n}{M}, \\ k &= 1, 2, \dots, M - 1, \quad l = 1, 2, \dots, M - 1, \end{aligned} \quad (13)$$

образует ортонормированный базис

$$(\psi^{(k,l)}, \psi^{(r,s)}) = \begin{cases} 0, & \text{если } k \neq r \text{ или } l \neq s, \\ 1, & \text{если } k = r \text{ и } l = s. \end{cases}$$

Это следует из (13), если заметить, что

$$\begin{aligned} (\psi^{(k,l)}, \psi^{(r,s)}) &= \left( 2 \sum_{m=0}^M \sin \frac{k\pi m}{M} \sin \frac{r\pi m}{M} \right) \left( 2 \sum_{n=0}^M \sin \frac{l\pi n}{M} \sin \frac{s\pi n}{M} \right) = \\ &= (\psi^{(k)}, \psi^{(r)}) (\psi^{(l)}, \psi^{(s)}). \end{aligned}$$

Любая функция  $v = \{v_{mn}\}$ , обращающаяся в нуль на границе квадрата, разлагается в конечный двумерный ряд Фурье:

$$v_{mn} = 2 \sum_{k,l=1}^{M-1} c_{kl} \sin \frac{k\pi m}{M} \sin \frac{l\pi n}{M}, \quad (14)$$

где  $c_{kl} = (v, \psi^{(k,l)})$ . Справедливо равенство Парсеваля

$$(v, v) = \sum_{k=l=1}^{M-1} c_{kl}^2. \quad (15)$$

**2. Представление решения в виде конечного ряда Фурье.** Можно проверить прямым вычислением справедливость равенств

$$-\Delta^{(h)} \psi^{(r,s)} = \frac{4}{h^2} \left( \sin^2 \frac{r\pi}{2M} + \sin^2 \frac{s\pi}{2M} \right) \psi^{(r,s)}, \quad r, s = 1, 2, \dots, M - 1, \quad (16)$$

где оператор  $-\Delta^{(h)}: U^{(h)} \rightarrow U^{(h)}$  определен в гл. 4, § 1, п. 3.

Равенства (16) означают, что функции  $\psi^{(r,s)}$ , образующие ортонормированный базис в пространстве  $U^{(h)}$ , являются собственными функциями оператора  $-\Delta^{(h)}: U^{(h)} \rightarrow U^{(h)}$ , а числа

$$\lambda_{r,s} = \frac{4}{h^2} \left( \sin^2 \frac{r\pi}{2M} + \sin^2 \frac{s\pi}{2M} \right), \quad r, s = 1, 2, \dots, M - 1, \quad (17)$$

суть отвечающие им собственные числа. В частности, отсюда следует, что оператор  $-\Delta^{(h)}$  самосопряжен.

В соответствии с общей схемой решение задачи (7) запишется в виде

$$u_{m_1 m_2} = \sum_{r,s=1}^{M-1} \frac{F_{rs}}{\lambda_{rs}} \cdot 2 \sin \frac{r\pi m_1}{M} \sin \frac{s\pi m_2}{M}, \quad (18)$$

где

$$F_{rs} = 2h^2 \sum_{m_1, m_2=1}^M f_{m_1 m_2} \sin \frac{r\pi m_1}{M} \sin \frac{s\pi m_2}{M}. \quad (19)$$

**Замечание.** В случае  $N = 2^p$ , где  $p$  — натуральное число, существует способ вычисления коэффициентов  $F^{(r,s)}$  по формуле (19) и вычисления решения  $u^{(h)}$  по формуле (18) за  $O(M^2 \ln M)$  арифметических операций. Этот замечательный алгоритм называется *быстрым преобразованием Фурье* (см., напр., [7, 12]).

Отметим, что

$$\lambda_{11} = \lambda_{11}(h) = \frac{8}{h^2} \sin^2 \frac{\pi}{2M}.$$

Очевидно, что

$$\pi^2 \leq \lambda_{11} \leq 2\pi^2. \quad (20)$$

Далее,

$$\lambda_{M-1, M-1} \sim O\left(\frac{1}{h^2}\right). \quad (21)$$

Поэтому число обусловленности оператора  $-\Delta^{(h)}$

$$\mu(-\Delta^{(h)}) = \frac{\lambda_{\max}}{\lambda_{\min}} = \frac{\lambda_{M-1, M-1}}{\lambda_{11}} \quad (22)$$

растет при  $h \rightarrow 0$  как  $O(h^{-2})$ .

Заметим еще, что наибольшим собственным значением самосопряженного оператора  $(-\Delta^{(h)})^{-1}$ , которое совпадает с нормой этого оператора, является число  $\lambda_{11}^{-1}$ , которое в силу (20) превосходит величину  $\pi^2$ . Поэтому для решения задачи (7) при любом  $f^{(h)}$  выполнено неравенство

$$\|u^{(h)}\| \leq \frac{1}{\pi^2} \|f^{(h)}\|. \quad (23)$$

### Задача

Выписать с помощью конечного ряда Фурье решение разностного аналога задачи Дирихле для уравнения Пуассона

$$-\left(\frac{u_{m_1+1, m_2} - 2u_{m_1 m_2} + u_{m_1-1, m_2}}{h^2} + \frac{u_{m_1, m_2+1} - 2u_{m_1 m_2} + u_{m_1, m_2-1}}{h^2}\right) = f_{m_1 m_2},$$

$$m_1 = 1, 2, \dots, M-1, \quad m_2 = 1, 2, \dots, M-1,$$

со следующими условиями на границе сеточного квадрата:

$$\begin{aligned} u_{0m_2} &= \varphi_{m_2}, \quad u_{Mm_2} = \psi_{m_2}, \quad m_2 = 1, 2, \dots, M - 1, \\ u_{m_10} &= \xi_{m_1}, \quad u_{m_1M} = \eta_{m_1}, \quad m_1 = 1, 2, \dots, M - 1, \end{aligned}$$

где  $f_{m_1m_2}, \varphi_{m_2}, \psi_{m_2}, \xi_{m_1}, \eta_{m_1}$  — заданные функции своих аргументов.

**Указание.** Заметим, что решение поставленной задачи во внутренних точках совпадает с решением задачи, которая возникает при замене функций  $\varphi_{m_2}, \psi_{m_2}, \xi_{m_1}, \eta_{m_1}$  тождественно обращающимися в нуль и при одновременной замене значений  $f_{m_1m_2}$  в приграничных точках сетки другими значениями  $\tilde{f}_{m_1m_2}$  (какими именно?). После этого решение внутри сеточной области записать в виде конечного ряда Фурье вида (18).

## ГЛАВА 5

# МЕТОДЫ ПОСЛЕДОВАТЕЛЬНЫХ ПРИБЛИЖЕНИЙ (ИТЕРАЦИОННЫЕ МЕТОДЫ) РЕШЕНИЯ СИСТЕМ ЛИНЕЙНЫХ АЛГЕБРАИЧЕСКИХ УРАВНЕНИЙ

Задача отыскания точного решения уравнения

$$Ax = f, \quad f \in R^n, \quad x \in R^n, \quad (1)$$

где  $A: R^n \rightarrow R^n$  — некоторый линейный оператор, не диктуется, как правило, запросами приложений. В приложениях обычно допустимо использование приближенного решения, известного с достаточной для каждого данного приложения точностью. К тому же найти точное решение, как правило, принципиально невозможно, так как входные данные (правая часть и сам оператор  $A$ ) бывают известны не вполне точно, что приводит к неустранимой погрешности в результате. Кроме того, в силу конечной разрядности компьютера неизбежны погрешности округления в процессе вычислений.

Поэтому во многих случаях для вычисления решения  $x$  уравнения (1) точным методам (методу Гаусса или другим) целесообразно предпочесть тот или иной метод последовательных приближений (итерационный метод). Каждый итерационный метод состоит в указании рекуррентного соотношения, которое по заданному произвольно нулевому приближению  $x^{(0)}$  решения  $x$  позволяет вычислить первое, второе, вообще  $p$ -е ( $p = 1, 2, \dots$ ) приближения  $x^{(p)} \in R^n$  решения  $x$ .

Итерационный процесс должен быть построен так, чтобы последовательные приближения  $x^{(p)}$  с ростом номера  $p$  стремились к решению  $x$  уравнения. Тогда для каждого  $\varepsilon > 0$  существует такой номер  $p = p(\varepsilon)$ , что выполняется неравенство

$$\|x - x^{(p)}\| < \varepsilon. \quad (2)$$

Задавая  $\varepsilon > 0$  достаточно малым, можно воспользоваться  $p$ -м приближением  $x^{(p)}$  в качестве приближенного решения с требуемой в данной задаче точностью.

Изложим некоторые итерационные методы и укажем условия, при которых эти методы целесообразно предпочесть точным или один другому.

## § 1. Методы простых итераций

Заметим, что СЛАУ

$$Ax = f \quad (1)$$

можно преобразовать к форме

$$x = (E - \tau A)x + \tau f, \quad (2)$$

причем новое уравнение (2) равносильно исходному при любом значении параметра  $\tau$  ( $\tau > 0$ ). Вообще  $Ax = f$  многими способами можно заменить равносильной системой вида

$$x = Bx + \varphi, \quad x \in R^n, \quad \varphi \in R^n, \quad (3)$$

частным случаем которой является (2).

**1. Общая схема метода простых итераций** состоит в вычислении последовательности

$$x^{(p+1)} = Bx^{(p)} + \varphi, \quad p = 0, 1, \dots, \quad (4)$$

при заданном произвольно значении  $x^{(0)}$ . Ниже мы укажем условия, при которых эта последовательность сходится к решению СЛАУ  $Ax = f$ .

Предварительно становимся на частном случае итерационной формулы (4):

$$x^{(p+1)} = (E - \tau A)x^{(p)} + \tau f, \quad p = 0, 1, \dots. \quad (5)$$

Заметим, что для вычислений по этой формуле достаточно уметь по заданному  $x^{(p)} \in R^n$  находить элемент  $Ax^{(p)}$ , получающийся в результате действия оператора  $A$ .

Таким образом, итерационный процесс вычисления решения СЛАУ  $Ax = f$ , в отличие от метода Гаусса, можно реализовать и в случае операторной формы задания СЛАУ, не выделяя какой-либо базис в  $R^n$  и не приводя систему к каноническому виду

$$\sum_{j=1}^n a_{ij} x_j = f_i, \quad i = 1, 2, \dots, n. \quad (6)$$

В процессе вычислений по формуле (5) надо держать в памяти компьютера не  $n^2$  чисел, составляющих матрицу оператора  $A$  в каком-либо базисе, а лишь  $n$  чисел, описывающих вектор  $Ax^{(p)} \in R^n$ . К тому же мы увидим, что для некоторых классов СЛАУ число арифметических

действий для получения решения с разумной точностью итерационными методами имеет порядок многое меньше, чем  $O(n^3)$ .

**Теорема 1.** Пусть в  $R^n$  фиксирована некоторая норма, причем соответствующая норма оператора равносильной системы (3) оказалась меньше единицы:

$$\|B\| = q < 1. \quad (7)$$

Тогда система (1) имеет одно и только одно решение  $x$ ; при любом  $x^{(0)} \in R^n$  последовательность (4) сходится к решению  $x$ , причем погрешность  $p$ -го приближения (или  $p$ -й итерации)

$$\varepsilon^{(p)} \equiv x - x^{(p)}$$

удовлетворяет оценке

$$\|\varepsilon^{(p)}\| = \|x - x^{(p)}\| \leq q^p \|x - x^{(0)}\| = q^p \|\varepsilon^{(0)}\|. \quad (8)$$

Тем самым норма погрешности  $\|\varepsilon^{(p)}\|$  с ростом  $p$  стремится к нулю не медленнее геометрической прогрессии  $q^p$ .

**Доказательство.** В случае  $\varphi = 0$  уравнение (3) не имеет решений, отличных от  $x = 0$ . В противном случае для решения  $x \neq 0$ ,  $\varphi = 0$  из (3) следовало бы, что

$$\|x\| = \|Bx\| \leq \|B\| \|x\| = q \|x\| < \|x\|,$$

т. е.  $\|x\| < \|x\|$ . В силу доказанного утверждения справедливо также утверждение о существовании и единственности решения СЛАУ (3), а следовательно, и (1) при произвольном  $\varphi$ .

Пусть  $x$  — решение системы (3). Зададим  $x^{(0)} \in R^n$  произвольно. Вычитая из (3) равенство (4) почленно, получаем

$$\varepsilon^{(p+1)} = B\varepsilon^{(p)}, \quad p = 0, 1, \dots$$

Отсюда имеем

$$\begin{aligned} \|x - x^{(p)}\| &= \|\varepsilon^{(p)}\| = \|B\varepsilon^{(p-1)}\| \leq q \|\varepsilon^{(p-1)}\| \leq \\ &\leq q^2 \|\varepsilon^{(p-2)}\| \leq \dots \leq q^p \|\varepsilon^{(0)}\| = q^p \|x - x^{(0)}\|. \end{aligned} \quad \square$$

**Замечание.** Условие (7) может нарушиться при каком-нибудь другом выборе нормы  $\|x\|'$ . Однако сходимость сохранится, причем оценка (8) заменится оценкой

$$\|\varepsilon^{(p)}\|' \leq cq^p \|\varepsilon^{(0)}\|', \quad (8')$$

где  $c$  — некоторая постоянная, зависящая от новой нормы, а значение  $q$  ( $q < 1$ ) прежнее.

Справедливость этого замечания следует из факта эквивалентности любых двух норм в конечномерном пространстве, состоящего в следующем. Если  $\|x\|$ ,  $\|x\|'$  — какие-нибудь нормы в  $R^n$ , то существуют такие числа  $c_1 > 0$ ,  $c_2 > 0$ , что  $c_1 \|x\|' \leq \|x\| \leq c_2 \|x\|'$  для всех  $x \in R^n$ , а  $c_1$ ,  $c_2$  от  $x$  не зависят. Легко видеть, что из (8) следует (8') при  $c = c_2/c_1$ .

Пример 1. Пусть система (1) задана в каноническом виде (6), причем  $A = \{a_{ij}\}$  — матрица с диагональным преобладанием, т. е. выполнены неравенства

$$|a_{ii}| > \sum_{j,j \neq i} |a_{ij}|, \quad i = 1, 2, \dots, n.$$

В уравнении с номером  $i$  перенесем все слагаемые  $a_{ij}x_j$  в правую часть и разделим это уравнение на  $a_{ii}$ . Возникнет система вида (3) с матрицей

$$B = \begin{bmatrix} 0 & b_{12} & b_{13} & \dots & b_{1,n-1} & b_{1n} \\ b_{21} & 0 & b_{23} & \dots & b_{2,n-1} & b_{2n} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ b_{n1} & b_{n2} & b_{n3} & \dots & b_{n,n-1} & 0 \end{bmatrix}.$$

Очевидно, существует такое число  $q$ , что

$$\sum_{j=1}^n |b_{ij}| \leq q < 1,$$

а следовательно, первая норма  $\|B\|_1 = \max_i \sum_j |b_{ij}|$  удовлетворяет оценке вида (7).

Зададим  $x^{(0)}$  произвольно, но так, чтобы погрешность  $\|\varepsilon^{(0)}\|_1 = \|x - x^{(0)}\|_1$  была, скажем, заведомо не больше единицы. И пусть нас устраивает ответ с погрешностью  $10^{-3}$ . Тогда достаточно выбрать  $p$  так, чтобы выполнялось условие  $q^p \leq 10^{-3}$ . В случае, например,  $q = 1/2$  достаточно взять  $p = 10$  независимо от  $n$ .

Общее число арифметических операций будет  $O(10n^2) = O(n^2)$ , а не  $O(n^3)$ , как в методе Гаусса. Действительно, каждое применение матрицы  $B$  требует  $n^2$  действий.

## 2. Необходимое и достаточное условие сходимости простых итераций.

**Теорема 2.** Пусть  $B: R^n \rightarrow R^n$ , где  $R^n$  —  $n$ -мерное комплексное пространство. Последовательность простых итераций

$$x^{(p+1)} = Bx^{(p)} + \varphi, \quad p = 0, 1, \dots, \tag{9}$$

сходится (в произвольной норме) при произвольном  $x^{(0)}$  в том и только том случае, если модуль каждого собственного числа  $\lambda_j$  оператора  $B$  строго меньше единицы:

$$|\lambda_j| < \rho < 1, \quad j = 1, 2, \dots, n.$$

**Доказательство.** Обозначим  $x - x^{(p)} = \varepsilon^{(p)}$ . Сначала докажем достаточность условия (9) для сходимости. Пусть условие (9) выполнено. Заметим прежде всего, что тогда число  $\lambda = 1$  не является собственным числом оператора  $B$ , и поэтому однородная система  $By = 1 \cdot y$  имеет только тривиальное решение  $y = 0$ . Следовательно,

система  $x = Bx + \varphi$  всегда имеет решение  $x$ . Сходимость равносильна тому, что последовательность

$$\varepsilon^{(p+1)} = B\varepsilon^{(p)}, \quad p = 0, 1, \dots,$$

сходится при произвольном  $\varepsilon^{(0)}$  к нулю:  $\|\varepsilon^{(p)}\| \rightarrow 0$ . Фиксируем  $\varepsilon^{(0)}$ . Очевидно, что  $\|\varepsilon^{(p)}\| \leq \|B\|^p \|\varepsilon^{(0)}\|$ . Обозначим через  $U(\lambda)$  сумму ряда векторных величин

$$U(\lambda) = \sum_{p=0}^{\infty} \frac{\varepsilon^{(p)}}{\lambda^p}.$$

Этот ряд заведомо сходится вне круга  $|\lambda| > \|B\|$  на комплексной плоскости  $\lambda$ . Очевидно, что  $\lambda U(\lambda) - \lambda \varepsilon^{(0)} = BU(\lambda)$ , или  $U(\lambda) = -\lambda(B - \lambda E)^{-1}\varepsilon^{(0)}$ . Из определения  $U(\lambda)$  видно, что  $\varepsilon^{(p)}$  является вычетом вектор-функции  $\lambda^{p-1}U(\lambda)$ :

$$\varepsilon^{(p)} = \frac{1}{2\pi i} \int_{|\lambda|=r} \lambda^{p-1} U(\lambda) d\lambda = -\frac{1}{2\pi i} \int_{|\lambda|=r} \lambda^p (B - \lambda E)^{-1} \varepsilon^{(0)} d\lambda,$$

где  $r$  — любое число, которое больше  $\|B\|$ .

В силу (9) подынтегральная вектор-функция — аналитическая функция вне круга  $|\lambda| > \rho$ , так как оператор  $(B - \lambda E)^{-1}$  существует при всех  $\lambda: |\lambda| > \rho$ . Поэтому можно деформировать контур интегрирования, выбирая  $r = \rho + \varepsilon$ , где  $\varepsilon > 0$  произвольно, не изменяя интеграла. Отсюда

$$\begin{aligned} \|\varepsilon^{(p)}\| &= \frac{1}{2\pi} \left| \int_{|\lambda|=\rho+\varepsilon} \lambda^p (B - \lambda E)^{-1} \varepsilon^{(0)} d\lambda \right| \leq \\ &\leq (\rho + \varepsilon)^p \max_{\lambda: |\lambda|=\rho+\varepsilon} \|(B - \lambda E)^{-1}\| \|\varepsilon^{(0)}\|. \end{aligned}$$

Выберем  $\varepsilon > 0$  настолько малым, чтобы было  $\rho + \varepsilon < 1$ . Тогда правая часть последнего неравенства с ростом  $p$  стремится к нулю. Достаточность доказана.

Пусть условие (9) не выполнено, так что некоторое  $\lambda_k$  удовлетворяет неравенству  $|\lambda_k| \geq 1$ . Предположим, что вопреки утверждению теоремы сходимость  $x^{(p)} \rightarrow x$  имеет место при любом выборе  $x^{(0)}$ . Выберем  $x^{(0)}$  так, чтобы  $\varepsilon^{(0)} = e_k$ , где  $e_k$  — собственный вектор оператора  $A$ , соответствующий собственному числу  $\lambda_k$ . Тогда  $\varepsilon^{(p)} = B^p \varepsilon^{(0)} = B^p e_k = \lambda_k^p e_k$ . В силу  $|\lambda_k| \geq 1$  эта последовательность не стремится к нулю. Получаем противоречие.  $\square$

Сделаем очевидное замечание, что если последовательность  $x^{(p)}$  сходится,  $x = \lim_{p \rightarrow \infty} x^{(p)}$ , то ее предел  $x$  является решением системы  $x = Bx + \varphi$ .

Отметим весьма интересное обстоятельство, которое встречается нам впервые. Задача вычисления  $x = \lim_{p \rightarrow \infty} x^{(p)}$  предельно хорошо обусловлена: результат вообще не зависит от входных данных, т. е. от

выбора  $x^{(0)}$ . В то же время сходящийся в силу теоремы 2 алгоритм вычисления последовательности  $x^{(p)}$  может быть вычислительно очень неустойчивым, т. е. очень чувствительным к ошибкам округления в процессе счета. Неустойчивость может иметь место, если наряду с выполнением  $\max |\lambda_j| = \rho < 1$  выполняется  $\|B\| > 1$ .

В случае  $\|B\| \leq 1$  норма  $\|\varepsilon^{(p)}\| = \|B^p \varepsilon^{(0)}\|$  монотонно убывает. В случае же  $\|B\| > 1$  существует такое  $\varepsilon^{(0)}$ , что  $\|\varepsilon^{(p)}\|$  сначала возрастает, а потом убывает. При этом высота «горба» может быть сколь угодно большой. Малая относительная погрешность округлений, допущенная при том  $p$ , где расположен максимум «горба», будет возрастать по норме: эта погрешность будет с ростом  $p$ , как и норма  $\|\varepsilon_p\|$ , развиваться, проходя через максимум, и т. д. Вычислительная неустойчивость может оказаться уже при небольшом превышении  $\|B\|$  над единицей и не очень больших  $n$  столь сильной, что расчет станет бессмысленным.

Строгое определение устойчивости метода простых итераций, классификацию видов возможной неустойчивости, примеры и теоремы см. в: Рябенький В.С. Об устойчивости итерационных процессов // ДАН СССР. — 1970. — Т. 193, № 3, или [6, § 47].

**3. Метод простых итераций в случае  $A = A^* > 0$ .** Рассмотрим уравнение (3), т. е.  $x = Bx + \varphi$ ,  $x \in R^n$ , в предположении, что  $R^n$  — евклидово пространство со скалярным произведением  $(x, y)$  и нормой  $\|x\| = (x, x)^{1/2}$ , а  $B = B^*$  — самосопряженный оператор. Пусть  $\nu_j$  ( $j = 1, 2, \dots, n$ ) — собственные числа оператора  $B$ . Введем число

$$q = \max_j |\nu_j|.$$

Зададим произвольно начальное приближение  $x^{(0)} \in R^n$  и построим последовательность простых итераций (4):

$$x^{(p+1)} = Bx^{(p)} + \varphi, \quad p = 0, 1, \dots. \quad (10)$$

**Лемма 1.** 1°. Если  $q < 1$ , то уравнение  $x = Bx + \varphi$  имеет решение  $\tilde{x}$ , а последовательные приближения  $x^{(p)}$  имеют погрешность  $\varepsilon^{(p)} = \tilde{x} - x^{(p)}$ , которая в норме  $\|x\| = (x, x)^{1/2}$  с ростом  $p$  стремится к нулю и удовлетворяет оценке

$$\|\varepsilon^{(p)}\| \leq q^p \|\varepsilon^{(0)}\|, \quad p = 0, 1, \dots, \quad (11)$$

причем существует такое  $x^{(0)}$ , что в (11) достигается точное равенство.

2°. Пусть система  $x = Bx + \varphi$  при данном  $\varphi$  имеет решение  $x \in R^n$ , но  $q \geq 1$ . Тогда существует начальное приближение  $x^{(0)} \in R^n$ , такое, что соответствующая последовательность простых итераций (10) не сходится к решению  $x$ .

**Доказательство.** В случае  $B = B^*$  евклидова норма оператора  $B$  совпадает с числом  $q = \max |\nu_j|$ . Поэтому утверждение 1° справедливо в силу теоремы 1.

Докажем 2°. Вычитая из тождества  $\tilde{x} = B\tilde{x} + \varphi$  почленно равенства (10), получаем равенства  $\varepsilon^{(p+1)} = B\varepsilon^{(p)}$ .

Пусть  $q \geq 1$ , причем  $q = |\nu_k|$ . Зададим  $x^{(0)}$  так, чтобы получить  $\varepsilon^{(0)} = e_k$ , где  $e_k$  — собственный вектор оператора  $B$ , соответствующий собственному числу  $\nu_k$ . Очевидно, что в таком случае погрешность  $\varepsilon^{(p)}$  не стремится к нулю с ростом  $p$ :

$$\varepsilon^{(p)} = B\varepsilon^{(p-1)} = \dots = B^p\varepsilon^{(0)} = B^p e_k = \nu_k^p e_k. \square \quad (12)$$

Рассмотрим уравнение

$$Ax = f, \quad A = A^* > 0, \quad (13)$$

относительно  $x \in R^n$ , где  $R^n$  — евклидово пространство. Зададим  $\tau > 0$  и приведем уравнение (13) к виду

$$x = (E - \tau A)x + \tau f. \quad (14)$$

Зададим произвольное нулевое приближение  $x^{(0)}$  и рассмотрим последовательность простых итераций

$$x^{(p+1)} = (E - \tau A)x^{(p)} + \tau f, \quad p = 0, 1, \dots. \quad (15)$$

Поскольку  $A = A^*$ , то все собственные значения оператора  $A$  положительны. Пусть  $\lambda_{\min}$  и  $\lambda_{\max}$  — соответственно наименьшее и наибольшее из них.

**Теорема 3. 1°.** Если  $\tau > 0$  достаточно мало, а именно, удовлетворяет неравенствам

$$0 < \tau < \frac{2}{\lambda_{\max}}, \quad (16)$$

то последовательность  $x^{(p)}$  сходится к решению  $x$  уравнения (13), причем гарантировано убывание нормы погрешности  $\|x - x^{(p)}\|$  при возрастании  $p$  в соответствии с оценкой

$$\|x - x^{(p)}\| \leq q^p \|x - x^{(0)}\|, \quad p = 0, 1, \dots, \quad (17)$$

где  $q < 1$ ,

$$q = q(\tau) = \max(|1 - \tau\lambda_{\min}|, |1 - \tau\lambda_{\max}|). \quad (18)$$

2°. Пусть  $\tau$  — произвольное число, удовлетворяющее (16). Существует начальное приближение  $x^{(0)}$ , при котором оценку (17) при выбранном  $\tau$  улучшить нельзя, так как при этом  $x^{(0)}$  соотношение (17) превращается в точное равенство.

3°. Если условие (16) нарушено, так что  $\tau \geq 2/\lambda_{\max}$ , то существует  $x^{(0)}$ , при котором с ростом  $p$  последовательность  $x^{(p)}$  не сходится к решению  $x$ .

4°. Число  $q = q(\tau)$ , задаваемое формулой (18), принимает наименьшее значение  $q_{\text{опт}} = q(\tau_{\text{опт}})$ , если  $\tau = \tau_{\text{опт}} = 2/(\lambda_{\min} + \lambda_{\max})$ . В этом случае число  $q$  принимает значение

$$q = q_{\text{опт}} = \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} = \frac{\mu(A) - 1}{\mu(A) + 1}, \quad (19)$$

где  $\mu(A) = \lambda_{\max}/\lambda_{\min}$  — число обусловленности оператора  $A$ .

Доказательство опирается на лемму 1. В этой лемме полагаем  $B = E - \tau A$ . Заметим, что если  $A = A^*$ , то оператор  $B = E - \tau A$  также самосопряженный:

$$\begin{aligned} (Bx, y) &= ((E - \tau A)x, y) = (x, y) - \tau(Ax, y) = \\ &= (x, y) - \tau(x, Ay) = (x, (E - \tau A)y) = (x, By). \end{aligned}$$

Пусть  $\lambda_1, \lambda_2, \dots, \lambda_n$  — собственные числа оператора  $A$ , расположенные в порядке неубывания:

$$0 < \lambda_{\min} = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n = \lambda_{\max}, \quad (20)$$

$e_1, e_2, \dots, e_n$  — соответствующие собственные векторы, образующие ортонормированный базис:  $Ae_j = \lambda_j e_j$ . Тогда те же векторы являются собственными для оператора  $B$ , причем соответствующие собственные числа суть

$$\begin{aligned} \nu_j &= \nu_j(\tau) = 1 - \tau \lambda_j, \quad j = 1, 2, \dots, n, \\ Be_j &= (E - \tau A)e_j = e_j - \tau \lambda_j e_j = (1 - \tau \lambda_j)e_j = \nu_j e_j, \\ &\quad j = 1, 2, \dots, n. \end{aligned} \quad (21)$$

Очевидно, что в силу (20) собственные числа  $\nu_j$  расположены в порядке невозрастания (рис. 11):

$$\nu_1 \geq \nu_2 \geq \dots \geq \nu_n. \quad (22)$$

Легко видеть, что наибольшим среди чисел  $|\nu_j|$  ( $j = 1, 2, \dots, n$ ) является либо  $|\nu_1| = |1 - \tau \lambda_{\min}|$ , либо  $|\nu_n| = |1 - \tau \lambda_{\max}|$ , так что условие

$$q = \max_j |\nu_j| < 1 \quad (23)$$

леммы 1 совпадает с условием

$$q = \max \{|1 - \tau \lambda_{\min}|, |1 - \tau \lambda_{\max}|\} < 1. \quad (24)$$

Условие  $q < 1$  в случае  $\tau > 0$  выполняется в том и только том случае, если точка  $\nu_n$  на рис. 11 лежит правее точки  $-1$  (т. е. если

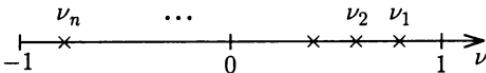


Рис. 11

$1 - \tau \lambda_{\max} > -1$ ). Это означает, что наряду с условием  $\tau > 0$  выполняется и второе неравенство (16). Если  $\tau \geq 2/\lambda_{\max}$ , то  $1 - \tau \lambda_{\max} \leq -1$  и  $q \geq 1$ . Итак, (16) равносильно требованию (9) леммы 1 в случае  $B = E - \tau A$ . Таким образом, доказаны утверждения 1°–3°.

Докажем теперь утверждение 4°. При очень маленьких  $\tau$  ( $\tau > 0$ ) все точки  $\nu_1, \nu_2, \dots, \nu_n$  расположены на рис. 11 левее точки 1, но вблизи нее, так что  $\max |\nu_j| = \nu_1 = 1 - \tau \lambda_{\min}$ . При увеличении  $\tau$  все точки смещаются влево, а число  $q = \max |\nu_j| = 1 - \tau \lambda_1$  уменьшается. Так будет происходить до тех пор, пока не наступит равенство  $|\nu_n| = \nu_1$ . При дальнейшем увеличении  $\tau$  число  $\nu_1$  будет продолжать уменьшаться. Однако теперь окажется, что  $|\nu_n| > \nu_1$ , так что число  $q = \max |\nu_j| = |\nu_n|$  с ростом  $\tau$  будет расти. Наименьшее значение  $q = q_{\text{опт}}$  будет при том значении  $\tau$ , при котором  $|\nu_n| = \nu_1$ , или  $-(1 - \tau \lambda_n) = 1 - \tau \lambda_1$ , т. е. при  $\tau = \tau_{\text{опт}} = 2/(\lambda_{\min} + \lambda_{\max})$ . Очевидно, что при этом

$$\begin{aligned} q = q_{\text{опт}} &= q(\tau_{\text{опт}}) = |\nu_n| = \nu_1 = 1 - \tau_{\text{опт}} \lambda_{\min} = \\ &= \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} = \frac{\lambda_{\max}/\lambda_{\min} - 1}{\lambda_{\max}/\lambda_{\min} + 1} = \frac{\mu(A) - 1}{\mu(A) + 1}. \quad \square \end{aligned}$$

Очевидно, что чем ближе  $\mu(A)$  к единице, тем ближе к нулю  $q_{\text{опт}}$  и тем быстрее в силу (17) убывает погрешность. С ростом числа обусловленности  $\mu(A)$  число  $q_{\text{опт}}$  увеличивается (оставаясь меньше единицы) и сходимость замедляется.

При оптимальном выборе  $\tau = \tau_{\text{опт}}$  справедлива оценка

$$\|\varepsilon^{(p)}\| \leq \left(\frac{1-\xi}{1+\xi}\right)^p \|\varepsilon^{(0)}\|, \quad \xi = \frac{\lambda_{\min}}{\lambda_{\max}}.$$

В силу леммы 1 существуют такие  $A, \varepsilon^{(0)}$ , что эта оценка достигается. Поэтому для того чтобы гарантировать оценку

$$\|\varepsilon^{(p)}\| \leq \varepsilon \|\varepsilon^{(0)}\|, \quad p = 0, 1, \dots, \quad (25)$$

при произвольном  $\varepsilon > 0$ , число  $p$  необходимо и достаточно выбрать из условия

$$\left(\frac{1-\xi}{1+\xi}\right)^p \leq \varepsilon, \quad p \geq -\frac{\ln \varepsilon}{\ln(1+\xi) - \ln(1-\xi)}.$$

Укажем более обозримую оценку для  $p$ . Заметим, что

$$\begin{aligned} \ln(1+\xi) - \ln(1-\xi) &= 2\xi \sum_{k=0}^{\infty} \frac{\xi^{2k}}{2k+1}, \\ 1 &\leq \sum_k \frac{\xi^{2k}}{2k+1} \leq \frac{1}{1-\xi^2}. \end{aligned}$$

Поэтому для гарантированной оценки (25) достаточно, чтобы значение  $p$  удовлетворяло оценке

$$p \geq -\frac{1}{2} \ln \varepsilon \cdot \mu, \quad \text{где } \mu = \frac{1}{\xi}, \quad (26)$$

и необходимо, чтобы

$$p \geq -\frac{1}{2} \ln \varepsilon \cdot (1 - \xi^2) \mu. \quad (27)$$

**Замечание.** Во многих случаях (например, при приближенной замене некоторых эллиптических краевых задач разностными) оператор  $A: \mathbb{R}^n \rightarrow \mathbb{R}^n$  возникающей линейной системы оказывается положительно определенным и самосопряженным ( $A = A^* > 0$ ) в смысле некоторого естественного скалярного умножения. Однако обычно не удается точно указать его наибольшее и наименьшее собственные значения. Удаётся указать лишь оценки границ спектра, т. е. такие числа  $a, b$ , чтобы выполнялись неравенства

$$0 < a \leq \lambda_{\min} \leq \lambda_{\max} \leq b.$$

В этом случае также можно воспользоваться методом простых итераций (15).

Зная вместо  $\lambda_{\min}, \lambda_{\max}$  лишь границы  $a, b$  спектра, можно воспользоваться значением  $\tau' = 2/(a+b)$ . При этом вместо

$$q_{\text{опт}} = \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}}$$

в гарантированной оценке (17) будет фигурировать число

$$q' = \max(|1 - \tau' \lambda_{\min}|, |1 - \tau' \lambda_{\max}|),$$

которое, вообще говоря, больше, чем  $q_{\text{опт}}$ . Как показано выше, при любом выборе  $\tau$  оценка (17) при некотором  $x^{(0)}$  превращается в точное равенство, причем  $q = q(\tau)$  определяется формулой (18). Поэтому при  $\tau = \tau' \neq \tau_{\text{опт}}$  получим неулучшаемую оценку (17), в которой будет  $q = q(\tau') > q_{\text{опт}}$ . Сходимость окажется тем медленнее, чем грубее известны границы  $a, b$  спектра.

**Пример 2.** Применим метод простых итераций к вычислению решения разностного аналога задачи Дирихле для уравнения Пуассона  $-\Delta^{(h)} u^{(h)} = f^{(h)}$ , сформулированной в § 7 из гл. 4.

Формула (15) примет вид

$$u^{(p+1)} = (E + \tau \Delta^{(h)}) u^{(p)} + \tau f, \quad p = 0, 1, \dots$$

Собственные числа оператора  $-\Delta^{(h)}$  найдены в (17) из § 7 гл. 4. Там же введено естественное скалярное умножение в пространстве  $U^{(h)}$ , где действует оператор  $-\Delta^{(h)}$ , установлена его самосопряженность и найдено число обусловленности  $\mu(-\Delta^{(h)}) = O(h^{-2})$ .

Поэтому в силу (25) при  $\tau_{\text{опт}}$  (укажите это  $\tau_{\text{опт}}$ ) количество итераций для уменьшения погрешности в  $e$  раз составит  $p \approx \frac{1}{2} \mu(-\Delta^{(h)}) = O(h^{-2})$ . Каждая итерация требует  $O(h^{-2})$  арифметических операций; их общее число  $O(h^{-4})$ .

В § 7 гл. 4 мы выписали точное решение задачи (20) в виде конечного ряда Фурье. Однако в случае непрямоугольной области или в случае, если вместо разностного аналога уравнения Пуассона рассматривался бы разностный аналог уравнения

$$\frac{\partial}{\partial x} \left( a \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( b \frac{\partial u}{\partial y} \right) = f \quad (28)$$

с переменными коэффициентами  $a = a(x, y) > 0$ ,  $b = b(x, y) > 0$ , мы не знали бы собственных функций и собственных чисел задачи и не могли бы воспользоваться конечными рядами Фурье. В то же время алгоритм простых итераций можно было бы построить и в этом случае вполне аналогично тому, как это сделано выше. Нужно лишь, чтобы разностное уравнение было самосопряженным и чтобы были известны не слишком грубые граници  $a, b$  спектра оператора задачи.

В дальнейшем мы укажем для уравнений  $Ax = f$  ( $A = A^* > 0$ ) с плохо обусловленным оператором  $A$  гораздо более эффективные итерационные алгоритмы, чем метод простой итерации.

**4. Переход от  $Ax = f$  ( $A = A^* > 0$ ) к лучше обусловленной системе с помощью энергетически эквивалентного оператора.** Мы видели, что скорость сходимости метода простых итераций (как и других, которые будут изложены в § 2) тем выше, чем число обусловленности системы  $\mu(A)$  меньше, т. е. чем это число ближе к единице.

В случае плохо обусловленной системы  $Ax = f$  иногда удается перейти к равносильной системе с оператором, который имеет меньшее число обусловленности, а затем решать эту систему методом итераций. Изложим этот прием.

Пусть  $B = B^* > 0$  — пока произвольный оператор. Умножим обе части уравнения  $Ax = f$  на  $B^{-1}$ . Получим равносильное уравнение

$$Cx = g, \quad C = B^{-1}A, \quad g = B^{-1}f. \quad (29)$$

Оператор  $C$  уже не является, вообще говоря, самосопряженным.

Введем новое скалярное умножение  $[x, y]_B = (Bx, y)$ . Оператор  $C$  оказывается самосопряженным в смысле нового скалярного умножения, т. е.  $[Cx, y]_B = [x, Cy]_B$ , а также положительно определенным, т. е.  $[Cx, x]_B > 0$ , если  $x \neq 0$ . Проверим это:

$$\begin{aligned} [Cx, y]_B &= (BCx, y) = (BB^{-1}Ax, y) = (Ax, y) = (x, Ay) = \\ &= (B^{-1}Bx, Ay) = (Bx, B^{-1}Ay) = (Bx, Cy) = [x, Cy]_B, \\ [Cx, x]_B &= (BCx, x) = (Ax, x) > 0, \quad x \neq 0. \end{aligned}$$

В нашем переходе от уравнения  $Ax = f$  к уравнению (29) выбор оператора пока произведен. В частности, если положить  $B = A$ , то оператор  $C = B^{-1}A$  окажется единичным, и решение  $x$  будет получено по формулам (29). Однако применение оператора  $A^{-1}$  равносильно точному решению уравнения  $Ax = f$ , которого мы как раз и хотим избежать за счет итерационного процесса, требующего умения вычислять при заданном  $z$  вектор  $Az$ , но не  $A^{-1}f$ . Поэтому имеет смысл выбирать оператор  $B$  лишь среди тех, для которых вычисление  $B^{-1}z$  по заданному  $z$  существенно проще, чем вычисление  $A^{-1}z$ . Если при

в этом удается выбрать  $B$  так, чтобы он был «похож» на оператор  $A$ , то можно надеяться, что оператор  $B^{-1}A$  будет «похож» на единичный, а его собственные числа  $\lambda_{\min}$ ,  $\lambda_{\max}$  и число обусловленности  $\mu(C)$  будут «ближе» к единице.

**Теорема 4.** Пусть  $B = B^* > 0$  и при заданных числах  $\gamma_1 > 0$ ,  $\gamma_2 > 0$  и всех  $x \in R^n$  справедливы неравенства

$$\gamma_1(Bx, x) \leq (Ax, x) \leq \gamma_2(Bx, x). \quad (30)$$

Тогда собственные числа  $\lambda_{\min}(C)$ ,  $\lambda_{\max}(C)$  и число обусловленности  $\mu_B(C)$  оператора  $C = B^{-1}A$  удовлетворяют неравенствам

$$\gamma_1 \leq \lambda_{\min}(C) \leq \lambda_{\max}(C) \leq \gamma_2, \quad \mu_B(C) \leq \gamma_2/\gamma_1. \quad (31)$$

**Доказательство.** Из курса линейной алгебры известно (и легко видеть непосредственно), что

$$\begin{aligned}\lambda_{\min}(C) &= \min_x \frac{[Cx, x]_B}{[x, x]_B} = \min_x \frac{(Ax, x)}{(Bx, x)}, \\ \lambda_{\max}(C) &= \max_x \frac{[Cx, x]_B}{[x, x]_B} = \max_x \frac{(Ax, x)}{(Bx, x)}.\end{aligned}$$

Отсюда в силу (30) следуют неравенства (31).  $\square$

Операторы  $A$ ,  $B$ , удовлетворяющие неравенствам (30), принято называть *эквивалентными по спектру*, или *энергетически эквивалентными* с константами эквивалентности  $\gamma_1$ ,  $\gamma_2$ .

Переход от  $Ax = f$  к  $Cx = g$  имеет смысл и существенно улучшает число обусловленности, если

$$\mu_B(C) \leq \frac{\gamma_2}{\gamma_1} \ll \mu(A).$$

При таком переходе увеличивается по сравнению с (17) скорость убывания погрешности  $\varepsilon^{(p)} = x - x^{(p)}$  в норме  $\|\cdot\|_B$ :

$$\begin{aligned}\|\varepsilon^{(p)}\|_B &= \|x - x^{(p)}\|_B \leq q_B^p \|x - x^{(0)}\|_B, \\ q_B &= \frac{\mu_B(C) - 1}{\mu_B(C) + 1}.\end{aligned}$$

Тогда при том же значении  $p$  будет  $q_B^p \ll q^p$ , так как

$$q_B \ll q, \quad q = \frac{\mu(A) - 1}{\mu(A) + 1}.$$

Типичная ситуация, в которой использование эквивалентных по спектру операторов дает большой эффект, впервые выделена и изучена Е.Г. Дьяконовым в начале 60-х годов и состоит в следующем. Пусть система алгебраических уравнений

$$A_n x = f, \quad f \in R^n, \quad x \in R^n,$$

где  $A_n: R^n \rightarrow R^n$ , возникла при дискретизации краевой задачи для эллиптического дифференциального уравнения. Размерность  $n$  пространства  $R^n$  тем выше, чем более точное приближение  $A_n$  исходной задачи мы используем при дискретизации. Таким образом, мы имеем дело с последовательностью пространств  $R^n$  ( $n \rightarrow \infty$ ), которые будем считать евклидовыми со скалярным умножением  $(x, y)^{(n)}$ .

Пусть  $\{A_n: R^n \rightarrow R^n\}$  — последовательность операторов, причем  $A_n = A_n^* > 0$ , а  $\{A_n x = f \ (x, f \in R^n)\}$  — последовательность уравнений. Пусть числа обусловленности  $\mu(A_n)$  растут с ростом  $n$  так, что  $\mu(A_n) \sim n^s$  ( $s > 0, s = \text{const}$ ). Тогда при отыскании решения с точностью  $\varepsilon > 0$  итерациями в соответствии с (27) потребуется  $O(n^s \ln \varepsilon)$  итераций.

Пусть оператор  $B_n: R^n \rightarrow R^n$  ( $B_n = B_n^* > 0$ ) энергетически эквивалентен оператору  $A_n$  с константами эквивалентности  $\gamma_1 > 0$ ,  $\gamma_2 > 0$ , которые не зависят от  $n$ . Тогда  $\mu_{B_n}(C_n) \leq \gamma_2/\gamma_1 = \text{const}$  ( $C_n = B_n^{-1}A_n$ ).

Перейдем к последовательности уравнений

$$C_n x = g_n, \quad C_n = B_n^{-1}A_n, \quad g_n = B_n^{-1}f^{(n)}.$$

Количество  $p$  итераций для уменьшения погрешности в  $\varepsilon$  раз:

$$\|x - x^{(p)}\|_{B_n} \leq \varepsilon \|x - x^{(0)}\|_{B_n}, \quad (32)$$

в силу ограниченности  $\mu_{B_n}(C_n) \leq \gamma_2/\gamma_1$  не возрастает с ростом  $n$  и имеет порядок  $O(\ln \varepsilon)$ .

Пусть нормы

$$\|x\| = [(x, x)^{(n)}]^{1/2}, \quad \|x\|_{B_n} = \{[x, x]_{B_n}\}^{1/2}$$

связаны неравенствами

$$n^{-l} \|x\|_{B_n} \leq \|x\| \leq n^l \|x\|_{B_n}, \quad l = \text{const} \geq 0.$$

Тогда для достижения оценки

$$\|x - x^{(p)}\| \leq \varepsilon \|x - x^{(0)}\|$$

в итерационном процессе для уравнения  $C_n x = g^{(n)}$  достаточно, чтобы выполнялось неравенство

$$\|x - x^{(p)}\|_{B_n} \leq \frac{\varepsilon}{n^l} \|x - x^{(0)}\|_{B_n}. \quad (33)$$

Это неравенство получается из (32) заменой  $\varepsilon$  на  $\varepsilon n^{-l}$ , так что потребуется лишь  $O\left(\ln \frac{\varepsilon}{n^l}\right) = O(l \ln n - \ln \varepsilon)$  итераций, в то время как без использования перехода от  $A_n x = f$  к  $C_n x = g_n$  требовалось, напомним,  $O(n^l \ln \varepsilon)$  итераций.

**5. Масштабирование как средство улучшения числа обусловленности.** Пусть требуется решить систему

$$Cx = \varphi, \quad (34)$$

заданную в каноническом виде:  $C = \{c_{ij}\}$ ,  $\det C \neq 0$ . Можно перейти к системе

$$Ax = f, \quad A = C^*C, \quad f = C^*\varphi,$$

причем  $A = A^* > 0$ , а затем воспользоваться методом итераций. Скорость их сходимости зависит от числа обусловленности  $\mu(A)$ . Легко показать, что в евклидовой норме, соответствующей скалярному умножению в пространстве  $R^n$ ,  $\mu(A) = \mu^2(C)$ . Прием масштабирования для уменьшения числа обусловленности заключается в следующем. Каждое из скалярных уравнений, составляющих (34), умножается на такой множитель, чтобы наибольший коэффициент нового уравнения оказался равным единице. От этой новой, масштабированной системы переходим к равносильной системе с симметрической и положительно определенной матрицей.

Прием масштабирования уменьшает число обусловленности используемой системы. В силу равенства  $\mu(A) = \mu^2(C)$  уменьшается также число  $\mu(A)$ , определяющее скорость сходимости итераций.

### Задачи

1. Пусть известно, что собственные значения оператора  $A: R^{100} \rightarrow R^{100}$  суть

$$\lambda_k = k^2, \quad k = 1, 2, \dots, 100.$$

Для приближенного вычисления решения системы  $Ax = f$  произвольно задается начальное приближение  $x^{(0)} \in R^{100}$ , а последующие приближения вычисляются по формулам

$$x^{(p+1)} = (E - \tau_p A)x^{(p)} + \tau_p f, \quad p = 0, 1, \dots, \quad (35)$$

где  $\tau_p$  — некоторые положительные числа.

Укажите такой набор итерационных параметров  $\tau_0, \tau_1, \dots, \tau_{99}$ , чтобы приближение  $x^{(100)}$  совпадало с точным решением  $x$  системы  $Ax = f$ .

2. Пусть в предыдущей задаче итерационные параметры  $\tau_p$  выбраны по формуле

$$\tau_p = \frac{1}{(p+1)^2}, \quad p = 0, 1, \dots, 99.$$

а) Проверить, что тогда  $x^{(100)} = x$ .

б) При реализации алгоритма (35) на компьютере с допустимым десятичным порядком чисел, скажем, до 10, вычисления натолкнутся на препятствие, которое состоит в том, что появятся числа слишком большого порядка. Объяснить механизм этого явления.

3. Пусть в задаче 1 итерационные параметры  $\tau$  выбраны по формуле

$$\tau_p = \frac{1}{(101-p)^2}, \quad p = 1, 2, \dots, 100.$$

а) Проверить, что тогда  $x^{(100)} = x$ .

б) При реализации этого алгоритма на компьютере с заданным количеством разрядов, скажем, 10, возникает большая погрешность. Объяснить механизм явления, препятствующего расчету.

## § 2. Метод Чебышёва и метод сопряженных градиентов

Для системы вида

$$Ax = f, \quad A = A^* > 0, \quad A: R^n \rightarrow R^n, \quad (1)$$

укажем два итерационных алгоритма для вычисления решения  $x = \bar{x}$ , в которых заданная точность достигается при меньшем объеме вычислений, чем в методе простых итераций, описанном в § 1. Обсудим также условия, при которых один из этих двух методов предпочтительнее другого. Обоснования указанных ниже алгоритмов см., напр., в [20].

В силу  $A = A^* > 0$  собственные числа  $\lambda_j$  ( $j = 1, 2, \dots, n$ ) оператора положительны. Будем считать, что они занумерованы в порядке возрастания и что известны такие числа  $a > 0, b > 0$ , что выполняется

$$a \leq \lambda_1 \leq \dots \leq \lambda_n \leq b. \quad (2)$$

Числа  $a, b$  называются *границами спектра* оператора  $A$ . Если  $a = \lambda_1, b = \lambda_n$ , то  $a, b$  называют *точными границами спектра*.

**1. Метод Чебышёва.** Зададим произвольно нулевое приближение  $x^0$  и будем вычислять последующие приближения по формулам

$$\begin{aligned} x^{(1)} &= (E - \tau A)x^{(0)} + \tau f, \\ x^{(p+1)} &= \alpha_{p+1}(E - \tau A)x^{(p)} + (1 - \alpha_{p+1})x^{(p-1)} + \tau \alpha_{p+1}f, \\ p &= 1, 2, \dots, \end{aligned} \quad (3)$$

где  $\tau$  и  $\alpha_k$  заданы формулами

$$\tau = \frac{2}{a+b}, \quad \alpha_1 = 2, \quad \alpha_{p+1} = \frac{4}{4 - \rho^2 \alpha_p}, \quad p = 1, 2, \dots, \quad \rho = \frac{b+a}{b-a}. \quad (4)$$

Можно показать, что погрешность  $\varepsilon^{(p)} = x - x^{(p)}$ ,  $\|\varepsilon^{(p)}\|^2 = (\varepsilon^{(p)}, \varepsilon^{(p)})$ , приближения  $x^{(p)}$  удовлетворяет оценке

$$\|\varepsilon^{(p)}\| \leq \frac{2p^q}{1 + q^{2p}} \|\varepsilon^{(0)}\|, \quad q = \frac{1 - \sqrt{\xi}}{1 + \sqrt{\xi}}, \quad \xi = \frac{a}{b}, \quad p = 0, 1, \dots. \quad (5)$$

Для уменьшения первоначальной погрешности в  $\varepsilon^{-1}$  раз, т. е. для достижения оценки  $\|\varepsilon^{(p)}\| \leq \varepsilon \|\varepsilon^{(0)}\|$ , достаточно, чтобы  $p$  было выбрано из условия

$$p \geq -\frac{1}{2} \left( \ln \frac{\varepsilon}{2} \right) \sqrt{\frac{b}{a}}.$$

Это число шагов  $p$  примерно в  $\nu = \frac{1}{2} \sqrt{\frac{b}{a}}$  раз меньше, чем число шагов, которого требует для достижения той же оценки  $\|\varepsilon^{(p)}\| \leq \varepsilon \|\varepsilon^{(0)}\|$  метод простой итерации. Экономия тем больше, чем больше число  $b/a = \xi^{-1}$ . В случае, если  $a, b$  — точные границы спектра, имеем  $\frac{b}{a} = \frac{\lambda_n}{\lambda_1} = \mu(A)$  и  $\nu = \frac{1}{2} \sqrt{\mu(A)}$ , т. е. экономия тем больше, чем больше число обусловленности системы (1). Отметим, что итерационный алгоритм (3) вычислительно устойчив.

Конструкция метода (3), доказательство оценки (5), свойства вычислительной устойчивости и других свойств алгоритма (3) основаны на использовании многочленов Чебышёва и некоторых других связанных с ними многочленов. Поэтому алгоритм (3) называют *трехслойным чебышёвским алгоритмом*.

Существует и используется на практике также так называемый *двухслойный чебышёвский итерационный алгоритм* для уравнения (1). Для вычисления  $x^{(p+1)}$  этот алгоритм использует и требует хранения в памяти машины только  $x^{(p)}$ , т. е., в отличие от (3),  $x^{(p-1)}$  хранить не нужно. Однако этот алгоритм требует дополнительной заботы об устойчивости счета, и мы его не излагаем.

**2. Метод сопряженных градиентов.** Формулы для вычисления приближений  $x^{(p+1)}$  методом сопряженных градиентов см. в § 6 гл. 4. В этих формулах не используются какие-либо границы  $a, b$  спектра, что является существенным преимуществом перед методом Чебышёва.

Гарантированная оценка скорости убывания погрешности зависит от числа обусловленности  $\mu(A)$  и не медленнее, чем в методе Чебышёва (3), построенного в случае известных  $a, b$ , т. е. в случае известных точных границ спектра.

В отличие от метода Чебышёва при использовании метода сопряженных градиентов может возникнуть вычислительная неустойчивость, которая проявляется тем сильнее, чем больше число обусловленности  $\mu(A)$ . Наиболее благоприятна для применения метода сопряженных градиентов ситуация, когда известно, что число обусловленности  $\mu(A)$  невелико, но границы спектра неизвестны, а порядок  $n$  системы много больше того числа итераций  $p$ , при котором погрешность  $\varepsilon^{(p)}$  удовлетворяет поставленному перед вычислителем требованию точности.

Скорость сходимости метода Чебышёва и метода сопряженных градиентов можно увеличить, если преобразовать систему (1) к системе вида

$$Cx = \varphi, \quad C = B^{-1}A, \quad B = B^* > 0,$$

уменьшив при этом число обусловленности. В этом случае полностью сохраняются построения и рассуждения из п. 4 § 1.

# ГЛАВА 6

## ПЕРЕОПРЕДЕЛЕННЫЕ СЛАУ. МЕТОД НАИМЕНЬШИХ КВАДРАТОВ

### § 1. Примеры задач, приводящих к переопределенным СЛАУ

**1. Задача обработки экспериментальных данных — эмпирические формулы.** Пусть известно, что величина  $y$  является некоторой функцией от аргумента  $t$ , причем в результате измерений получена таблица значений  $y_k = y(t_k)$  ( $k = 1, 2, 3, 4$ ) этой функции (рис. 12).

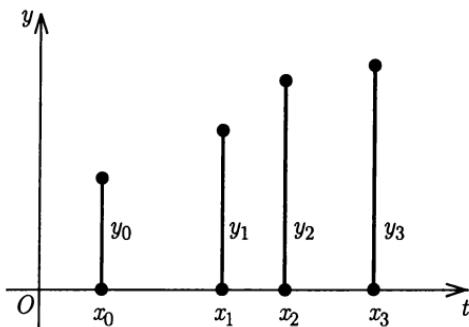


Рис. 12

Полученные измерения позволяют приближенно считать, что зависимость  $y = y(t)$  является линейной, т.е. имеет вид

$$y = x_1 t + x_2, \quad (1)$$

где  $x_1, x_2$  — некоторые числа. Числа  $x_1, x_2$  в эмпирической формуле (1) было бы желательно подобрать так, чтобы при  $t = t_k$  ( $k = 1, 2, 3, 4$ ) формула (1) давала те значения, которые получены измерениями, т.е. чтобы выполнялись уравнения

$$\begin{aligned} x_1 t_1 + x_2 &= y_1, \\ x_1 t_2 + x_2 &= y_2, \\ x_1 t_3 + x_2 &= y_3, \\ x_1 t_4 + x_2 &= y_4. \end{aligned} \quad (2)$$

Получилась система четырех линейных уравнений относительно двух неизвестных  $x_1, x_2$ . Эта переопределенная система не имеет (классического) решения, так как не существует прямой, проходящей через все четыре экспериментальные точки (см. рис. 12).

Введем обобщенное решение  $(x_1, x_2)$  системы (2), приняв за это решение ту пару чисел  $x_1, x_2$ , для которой невязки

$$x_1 t_k + x_2 - y_k, \quad k = 1, 2, 3, 4,$$

будут в некотором смысле как можно меньше. Можно, например, ввести функцию

$$\Phi(x_1, x_2) = \sum_{k=1}^4 (x_1 t_k + x_2 - y_k)^2, \quad (3)$$

равную сумме квадратов невязок, и принять за обобщенное решение  $(x_1, x_2)$  системы (2) ту пару чисел  $x_1, x_2$ , для которой  $\Phi(x_1, x_2)$  принимает наименьшее значение. Для вычисления этих чисел получим систему

$$\frac{\partial \Phi}{\partial x_1} = 0, \quad \frac{\partial \Phi}{\partial x_2} = 0 \quad (4)$$

двух уравнений, имеющую (обычное, классическое) решение  $(x_1, x_2)$ .

Выбор функции  $\Phi(x_1, x_2)$ , от которой зависит, какую именно пару чисел  $x_1, x_2$  принять за обобщенное решение системы (2), т. е. за коэффициенты эмпирической формулы (1), содержит произвол. Если бы, например, мы сочли разумным придать каждому  $k$ -му измерению свой вес  $b_k$  ( $k = 1, 2, 3, 4$ ), то вместо (3) было бы естественно принять функцию

$$\Phi(x_1, x_2) = \sum_{k=1}^4 b_k (x_1 t_k + x_2 - y_k)^2. \quad (5)$$

Соответственно изменилась бы и та пара чисел  $x_1, x_2$  (обобщенное решение системы (2)), которая минимизирует функцию  $\Phi(x_1, x_2)$ . Система (4) и ее решение  $(x_1, x_2)$  зависят от весов  $b_1, b_2, b_3, b_4$ .

Определение обобщенного решения системы (2) с помощью минимизации функции  $\Phi(x_1, x_2)$  от квадратов невязок есть пример метода наименьших квадратов для определения обобщенного решения переопределенной системы (2).

Можно было бы ввести меру невязки  $\Phi(x_1, x_2)$ , используя вместо (5) формулу

$$\Phi(x_1, x_2) = \sum_{k=1}^4 b_k |x_1 t_k + x_2 - y_k|, \quad (6)$$

т. е. использовать не квадраты, а модули невязок.

Отыскание минимума функции (6) — задача линейного программирования. Для ее решения нельзя воспользоваться уравнениями вида (4), так как функция (6), в отличие от функции (5), недифференцируема. Преимущество метода наименьших квадратов в том, что вычисление обобщенного решения, понимаемого в смысле метода наименьших квадратов, существенно проще.

**2. Задача уточнения результатов неточных измерений за счет числа этих измерений.** Пусть заранее известно, что некоторая величина  $y = y(t)$  зависит от своего аргумента  $t$  линейно, т. е. известно, что эта зависимость имеет вид (1). Задача состоит в том, чтобы по результатам измерения величины  $y(t)$  при нескольких значениях  $t$  вычислить

коэффициенты  $x_1, x_2$ . Допустим, что измерения произведены при  $t = t_k$  ( $k = 1, 2, 3, 4$ ) и их результаты приведены на рис. 12.

Для определения  $x_1, x_2$  снова получится та же переопределенная система (2). Эта система окажется несовместной, так как измерения производились с некоторой неизбежной погрешностью. Обобщенное решение снова можно определить как пару чисел  $x_1, x_2$ , придающую наименьшее значение функции вида (5), и воспользоваться для этого системой вида (4).

Подчеркнем, что переход от двух к четырем неточным измерениям и соответствующей переопределенной системе уравнений (2), вообще говоря, позволяет уменьшить влияние погрешностей, допускаемых при каждом измерении, на результат.

В терминах теории вероятностей существует точная постановка задачи о повышении надежности результатов, полученных при измерениях со случайными погрешностями, за счет числа этих измерений и их обработки с помощью метода наименьших квадратов (см., напр., [9]).

Переопределенные системы линейных уравнений могут возникать не только при построении эмпирических формул или при уточнении результатов неточных измерений за счет их числа (см., напр., гл. 13, § 2).

## § 2. Переопределенные СЛАУ и обобщенные решения в общем случае

1. Каноническая запись переопределенной системы линейных алгебраических уравнений имеет следующий вид:

$$\begin{aligned} a_{11}x_1 + \dots + a_{1s}x_s &= f_1, \\ \dots \dots \dots & \\ a_{n1}x_1 + \dots + a_{ns}x_s &= f_n, \\ n > s. \end{aligned} \tag{1}$$

Введем пространства  $R^s, R^n$ , состоящие из элементов вида  $x = \begin{bmatrix} x_1 \\ \dots \\ x_s \end{bmatrix}, f = \begin{bmatrix} f_1 \\ \dots \\ f_n \end{bmatrix}$  и имеющие размерности  $s, n$  ( $n > s$ ) соответственно. Обозначим через  $A$  прямоугольную матрицу коэффициентов системы (1):

$$A = \begin{bmatrix} a_{11} & \dots & a_{1s} \\ \dots \dots \dots & & \\ a_{n1} & \dots & a_{ns} \end{bmatrix}. \tag{2}$$

Тогда систему (1) можно коротко записать в виде

$$Ax = f, \quad x \in R^s, \quad f \in R^n. \tag{3}$$

Введем в  $R^n$  «основное» скалярное умножение, положив

$$(f, g)^{(n)} = \sum_{k=1}^n f_k g_k. \quad (4)$$

Наряду с (4) можно ввести скалярное произведение в  $R^n$  множеством других способов. Именно, произвольной симметрической и положительно определенной матрице  $B = B^* > 0$ , т. е.  $(Bf, f) > 0$  ( $f \neq 0$ ), соответствует скалярное умножение

$$[f, g]_B = (Bf, g), \quad f, g \in R^n; \quad (5)$$

обратно, любое скалярное умножение в пространстве  $R^n$  можно задать формулой вида (5), подобрав подходящую матрицу  $B = B^* > 0$ .

Система (1), вообще говоря, не имеет классического решения, т. е. не существует такого набора чисел  $x_1, x_2, \dots, x_s$ , который обращает каждое из  $n$  уравнений (1) в тождество.

**Определение.** Зафиксируем  $B: R^n \rightarrow R^n$ ,  $B = B^* > 0$ . Введем функцию от  $x \in R^n$ , положив

$$\Phi(x) = [Ax - f, Ax - f]_B. \quad (6)$$

Примем за *обобщенное решение*  $x_B$  системы (1) тот вектор  $x_B \in R^s$ , который придает наименьшее значение квадратичной функции (6).

**Замечание.** Выбор матрицы  $B = B^* > 0$  находится в руках исследователя. Она имеет смысл «весовой» матрицы и выбирается из тех или иных соображений о том, какую цену придать невязке системы (1) при заданных  $x_1, x_2, \dots, x_s$ .

**Теорема 1.** Пусть столбцы матрицы  $A$  линейно независимы, т. е. ранг  $A$  равен  $s$ . Тогда существует одно и только одно обобщенное решение  $x_B$  системы (1). Обобщенное решение системы (1) является классическим решением системы уравнений

$$A^*BAx = A^*Bf, \quad (7)$$

которая содержит  $s$  скалярных уравнений относительно  $s$  неизвестных  $x_1, x_2, \dots, x_s$ .

**Доказательство.** Введем обозначение  $a_k$  ( $a_k \in R^n$ ) для столбца с номером  $k$  ( $k = 1, 2, \dots$ ) матрицы  $A$ , так что

$$a_k = \begin{bmatrix} a_{1k} \\ \dots \\ a_{nk} \end{bmatrix}, \quad k = 1, 2, \dots, s.$$

Из формулы умножения матриц видим, что матрица  $C = A^*BA$  системы (7) есть квадратная ( $s \times s$ )-матрица. Элемент  $c_{ij}$  этой матрицы, стоящий на пересечении  $i$ -й строки и  $j$ -го столбца, есть число

$$c_{ij} = (a_i, Ba_j)^{(n)} = (Ba_i, a_j)^{(n)} = [a_i, a_j]_B. \quad (8)$$

Из (8) видно, что  $c_{ij} = c_{ji}$ , т. е.  $C = C^*$ .

Покажем, что матрица  $C$  невырождена и, более того, положительно определена:

$$(C\xi, \xi)^{(s)} > 0, \quad \xi \in R^s, \quad \xi \neq 0. \quad (9)$$

Отметим известную формулу

$$(f, A\xi)^{(n)} \equiv (A^* f, \xi)^{(s)}, \quad f \in R^n, \quad \xi \in R^s, \quad (10)$$

для проверки которой достаточно записать левую и правую части формулы (10) в развернутом виде.

Пусть  $\xi \in R^s$ ,  $\xi = \begin{bmatrix} \xi_1 \\ \dots \\ \xi_s \end{bmatrix} \neq \begin{bmatrix} 0 \\ \dots \\ 0 \end{bmatrix}$ . Тогда

$$A\xi = \xi_1 a_1 + \xi_2 a_2 + \dots + \xi_s a_s$$

есть линейная комбинация системы линейно независимых векторов  $a_k$ , не все коэффициенты которой равны нулю. Поэтому  $A\xi \neq 0 \in R^n$ . Но тогда согласно (10) скалярный квадрат вектора  $A\xi$  положителен:

$$0 < [A\xi, A\xi]_B = (BA\xi, A\xi)^{(n)} = (A^* BA\xi, \xi)^{(s)} = (C\xi, \xi)^{(s)},$$

и (9) доказано. В силу невырожденности матрицы  $C$  система (7) имеет одно и только одно решение  $x_B \in R^s$ .

Покажем, что  $x_B$  является единственным обобщенным решением системы (1). Это значит, что для любого  $x = x_B + \delta$  ( $\delta \in R^s$ ) справедливо строгое неравенство

$$\Phi(x_B + \delta) > \Phi(x_B), \quad \delta \neq 0. \quad (11)$$

Для доказательства этого неравенства предварительно заметим, что  $[Ax_B - f, A\delta]_B = (B(Ax_B - f), A\delta)^{(n)} = (A^* BAx_B - A^* Bf, \delta)^{(s)} = 0$ .

(12)

Это равенство справедливо, поскольку  $A^* BAx_B = A^* Bf$ .

Докажем теперь (11), опираясь на (9), (12):

$$\begin{aligned} \Phi(x_B + \delta) &= [A(x_B + \delta) - f, A(x_B + \delta) - f]_B = \\ &= [Ax_B - f, Ax_B - f]_B - 2[Ax_B - f, A\delta]_B + [A\delta, A\delta]_B = \\ &= \Phi(x_B) + (C\delta, \delta)^{(s)} > \Phi(x_B). \quad \square \end{aligned}$$

**2. Замечания о вычислении обобщенного решения.** Матрица системы (7) симметрическая и положительно определенная. Этим можно воспользоваться при вычислении решения методом итераций. В случае если  $a_1, a_2, \dots, a_s \in R^s$  — ортонормированная система векторов, т. е. если

$$[a_i, a_j]_B = \delta_{ij}, \quad i, j = 1, 2, \dots, s, \quad (13)$$

то матрица  $C$  оказывается единичной, так как было показано, что  $c_{ij} = [a_i, a_j]_B$ . В этом случае  $x_B = A^* Bf$ . Если условие (13) выполняется лишь приближенно, то матрица  $C$  оказывается близка

к единичной и потому хорошо обусловлена, так что решение (7) можно легко вычислить итерациями.

**3. Геометрический смысл метода наименьших квадратов.** Систему (1) можно записать так:

$$x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \dots + x_s \mathbf{a}_s = f, \quad (14)$$

где  $\mathbf{a}_k = \begin{bmatrix} a_{1k} \\ \dots \\ a_{nk} \end{bmatrix} \in R^n$ ,  $f = \begin{bmatrix} f_1 \\ \dots \\ f_n \end{bmatrix} \in R^n$ . Требуется найти коэффициенты

$x_1, x_2, \dots, x_s$  линейной комбинации  $x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \dots + x_s \mathbf{a}_s$  так, чтобы эта линейная комбинация наименее уклонялась от  $f$ :

$$\left\| f - \sum_{k=1}^s x_k \mathbf{a}_k \right\| \Rightarrow \min. \quad (15)$$

Обозначим через  $R^{(s)}(\mathbf{a}) \subset R^n$  подпространство размерности  $s$  пространства  $R^n$ , состоящее из всевозможных линейных комбинаций векторов  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_s$ . Покажем, что если  $x_1, x_2, \dots, x_s$  — обобщенное решение системы (1), то линейная комбинация  $\sum x_k \mathbf{a}_k$  есть ортогональная в смысле скалярного умножения  $[ , ]_B$  проекция вектора  $f \in R^n$  на подпространство  $R^s(\mathbf{a})$ .

В самом деле, любой вектор из  $R^s(\mathbf{a})$  имеет вид  $A\delta = \delta_1 \mathbf{a}_1 + \dots + \delta_s \mathbf{a}_s \in R^s(\mathbf{a})$  ( $\delta \in R^s$ ). Наименее уклоняется от  $f$  элемент  $\sum x_k \mathbf{a}_k$  подпространства  $R^s(\mathbf{a})$ , имеющий вид  $Ax_B$ , где  $x_B$  — решение системы (7). Очевидно, что в силу (12) элемент  $f - Ax_B$  ортогонален любому элементу  $A\delta \in R^s(\mathbf{a})$ .

Если в пространстве  $R^s(\mathbf{a})$  вместо  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_s$  выбрать какой-нибудь другой базис  $\mathbf{a}'_1, \mathbf{a}'_2, \dots, \mathbf{a}'_s$ , то система (7) заменится системой

$$C'x' = f \quad (16)$$

с матрицей  $C'[c'_{ij}]$ ,  $c'_{ij} = [A'\mathbf{a}'_i, \mathbf{a}'_j]_B$ ,  $i, j = 1, 2, \dots, s$ , где  $A'$  — матрица, столбцы которой суть  $\mathbf{a}'_k$ .

Вместо решения  $x_B$  системы (7) получим новое решение  $x'_B$  системы (16), но проекция  $f$  на  $R^s(\mathbf{a})$ , очевидно, останется прежней, так что справедливо равенство

$$\sum x_k \mathbf{a}_k = \sum x'_k \mathbf{a}'_k.$$

Если нас интересует проекция заданного  $f$  на заданное подпространство  $R^s \subset R^n$ , то естественно стремиться к выбору базиса  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_s$  этого подпространства, по возможности мало отличающегося от ортонормированного. Искомая проекция от выбора базиса в  $R^s(\mathbf{a})$  не зависит, а уравнение (7) в случае такого выбора базиса будет иметь хорошо обусловленную, близкую к единичной матрице (см. задачу 6 в конце параграфа).

#### 4. Переопределенные системы, заданные в операторной форме.

Наряду с канонической формой (1) переопределенная система может быть задана в операторной форме. Именно, пусть имеются два линейных пространства  $R^s$ ,  $R^n$  размерностей  $s$ ,  $n$  ( $s < n$ ) соответственно. Пусть заданы линейный оператор  $A: R^s \rightarrow R^n$  и вектор  $f \in R^n$ . Требуется найти такой элемент  $x \in R^s$ , чтобы выполнялись условия

$$Ax = f, \quad x \in R^s, \quad f \in R^n. \quad (17)$$

Эта запись по содержанию отличается от записи (3) тем, что здесь не предполагается, что  $x$ ,  $f$  заданы своими координатами и что  $A$  задан матрицей.

Система (17), вообще говоря, не имеет классического решения, так как вектор  $Ax \in R^n$  при любом  $x$  принадлежит  $s$ -мерному подпространству  $R(A) \subset R^n$ , в которое переходит  $R^s$  при преобразовании  $A: R^s \rightarrow R^n$ .

Для введения обобщенного решения надо определить в  $R^n$  скалярное умножение  $(f, g)^{(n)}$ . Все остальные скалярные умножения имеют вид

$$[f, g]_B = (Bf, g)^{(n)},$$

где  $B = B^* > 0$ . Введем функцию

$$\Phi_B(x) = [Ax - f, Ax - f]_B.$$

Примем за обобщенное решение  $x_B$  то значение  $x \in R^s$ , для которого эта функция принимает наименьшее значение:

$$\Phi_B(x) \Rightarrow \min.$$

**Пример.** Пусть в квадрате  $(0 \leq x \leq 1, 0 \leq y \leq 1)$  введены две сетки — крупная с шагом  $H = 0,1$  и с точками

$$(X_m, Y_n) = (mH, nH), \quad m, n = 0, 1, \dots, 10,$$

а также мелкая с шагом  $h = 10^{-2}$  и с точками

$$(x_m, y_n) = (mh, nh), \quad m, n = 0, 1, \dots, 10^2.$$

Введем линейные пространства  $U^{(H)}$ ,  $F^{(h)}$ . К первому отнесем все функции  $u_{mn}$ , определенные на крупной сетке, а ко второму — все функции  $f_{mn}$ , определенные на мелкой сетке. Зададим оператор  $A: U^{(H)} \rightarrow F^{(h)}$ , который по функции  $u_{mn}$  строит функцию  $f_{mn}$ , доопределяя  $u_{mn}$  в точках мелкой сетки с помощью линейной интерполяции по каждому из аргументов  $x$ ,  $y$ . Пусть путем приближенных измерений некоторого физического скалярного поля в точках мелкой сетки получена некоторая функция  $f \in F^{(h)}$ ,  $f = \{f_{mn}\}$ .

Поставим задачу — отыскать такую функцию  $u \in U^{(H)}$ , чтобы она была обобщенным решением следующей переопределенной системы:

$$Au = f, \quad u \in U^{(H)}, \quad f \in F^{(h)}. \quad (18)$$

Обобщенное решение будем понимать в смысле следующего скалярного умножения в  $F^{(h)}$ :

$$(f, g) = \sum_{m,n=0}^{100} f_{mn} g_{mn}, \quad f, g \in F^{(h)}.$$

Содержательный смысл этой задачи — произвести сглаживание экспериментальных данных, уменьшив влияние случайных погрешностей, допущенных при измерениях. В результате этого сглаживания вместо экспериментальной функции  $f \in F^{(h)}$  можно будет использовать функцию  $\tilde{f} = Au$ .

Мы не будем приводить здесь алгоритм отыскания обобщенного решения  $u \in U^{(H)}$  задачи (18), отнеся это к числу задач для самостоятельного решения. Подчеркнем только, что задача (18) поставлена не в канонической форме (1), а в операторной.

Можно многими способами перейти от операторной формы (18) к канонической форме (1) задания переопределенной системы (как это сделать?), однако это нецелесообразно.

### Задачи

**1.** Найти обобщенное решение (в смысле метода наименьших квадратов) переопределенной системы

$$x + y = 1, \quad x - y = 2, \quad 2x + y = 2,4.$$

**2.** Произведено некоторое число  $m$  приближенных измерений длины  $l$ . Получились следующие результаты:  $l = l_1, l = l_2, \dots, l = l_m$ , где  $l_j$  — некоторые числа.

Найти решение (в смысле метода наименьших квадратов) этой системы  $m$  уравнений относительно одного неизвестного  $l$ .

**3.** Электрическое сопротивление проволоки  $R$  линейно зависит от температуры  $t$ , так что  $R = a + bt$ .

Определить  $a, b$  по результатам следующих неточных измерений:

$t$	19,1	25,0	30,1	36,0
$R$	76,30	77,80	79,75	80,80

Найти сопротивление проволоки при  $t = 21, 28$ .

**4.** Скорость  $v$  корабля связана с мощностью  $p$  двигателя эмпирической формулой  $p = a + bv^3$ .

Определить  $a, b$  по данным измерений, приведенным в таблице:

$v$	6	8	10	12
$r$	420	805	1370	2370

**5.** Измерения углов  $\alpha, \beta, \gamma$  треугольника на плоскости привели к значениям  $\alpha = 52^\circ 5'$ ,  $\beta = 50^\circ 1'$ ,  $\gamma = 78^\circ 6'$ . Сумма результатов измерений есть  $180^\circ 12'$  и дает невязку  $12'$ , происходящую от погрешностей измерения.

Ликвидировать невязку, следя предписаниям метода наименьших квадратов.

**6.** Некоторая функция задана таблицей своих значений  $y_k = 1 + x_k + x_k^2 + \sin x_k$  на множестве точек

$$x_k = -1 + 10^{-2}k, \quad k = 0, 1, \dots, 200.$$

Требуется с помощью компьютера построить многочлен  $Q_m(x)$  степени не выше заданного  $m \leq 50$ , для которого сумма

$$\sum_{k=0}^{200} |Q_m(x_k) - y_k|^2$$

принимает наименьшее значение.

Рассмотреть и сравнить следующие три варианта решения и вычислить значения  $Q_m$  в точках

$$x_{k+1/2} = \frac{x_k + x_{k+1}}{2}, \quad k = 0, 1, \dots, 199.$$

**Вариант 1.** Ищем многочлен  $Q_m(x)$  в виде  $Q_m(x) = \sum_{s=0}^m c_s x^s$  с неопределенными коэффициентами  $c_s$ .

**Вариант 2.** Ищем многочлен  $Q_m(x)$  в виде

$$Q_m(x) = \sum_{s=0}^m c_s P_s(x),$$

где  $P_s(x)$  — многочлен Лежандра степени  $s$ . О многочленах Лежандра известно, что

$$\begin{aligned} P_0(x) &= 1, \quad P_1(x) = x, \\ P_{l+1}(x) &= \frac{2l}{l+1} x P_l(x) - \frac{l}{l+1} P_{l-1}(x), \quad l = 1, 2, \dots. \end{aligned}$$

Они ортогональны на отрезке  $-1 \leq x \leq 1$ , т. е.

$$\int_{-1}^1 P_k(x) P_l(x) dx = 0, \quad k \neq l.$$

Далее,

$$\int_{-1}^1 P_k^2(x) dx = \frac{2}{2k+1}.$$

**Вариант 3.** Ищем многочлен  $Q_m(x)$  в виде

$$Q_m(x) = \sum_{s=1}^m c_s \tilde{P}_s(x),$$

где

$$\tilde{P}_s(x) = \sqrt{\frac{2s+1}{2}} P_s(x).$$

Построение алгоритма отыскания коэффициентов  $c_s$  ( $s = 0, 1, \dots, m$ ) во всех трех вариантах осуществить по следующему плану.

1°. Вычислить матрицу  $A^*BA$  и матрицу  $A^*B$ , с помощью которых записывается система

$$A^*BAc = A^*By, \quad c = \begin{bmatrix} c_0 \\ \dots \\ c_m \end{bmatrix}, \quad y = \begin{bmatrix} y_0 \\ \dots \\ y_{200} \end{bmatrix}$$

метода наименьших квадратов.

2°. Написать программу решения полученной системы методом сопряженных градиентов, принимая за нулевое приближение  $\bar{c}^{(0)} = \begin{bmatrix} 0 \\ \dots \\ 0 \end{bmatrix}$  и заканчивая итерации, если  $c^{(p)}$  отличается от  $c^{(p-1)}$  «достаточно мало», а именно,

$$\max_s |c_s^{(p)} - c_s^{(p-1)}| \leq \varepsilon, \quad \varepsilon = 10^{-3}.$$

3°. Объяснить, почему система  $A^*BAc = A^*By$ , которая возникает в варианте 1, обусловлена плохо, а в варианте 3 ее число обусловленности не сильно превосходит единицу.

7\*. Составить программу для численного решения переопределенной системы, заданной в операторной форме в примере, который приведен в конце параграфа.

**Указание.** Ввести скалярное умножение в пространстве  $U^{(H)}$ ,  $F^{(h)}$ , положив

$$(f, g)^{(h)} = h^2 \sum_{m,n=0}^{10} f_{mn} g_{mn}, \quad f, g \in F^{(h)}.$$

Реализовать на компьютере программу решения итерациями системы линейных уравнений  $A^*Au = A^*f$ , заданной в операторной форме.

## ГЛАВА 7

### ЧИСЛЕННОЕ РЕШЕНИЕ НЕЛИНЕЙНЫХ СКАЛЯРНЫХ УРАВНЕНИЙ И СИСТЕМ УРАВНЕНИЙ

Здесь мы изложим основные способы отыскания решений нелинейных скалярных уравнений

$$F(x) = 0, \tag{1}$$

где  $F(x)$  — некоторая заданная функция, например,  $F(x) = \sin x - 0,5x$  или  $F(x) = e^{-x} + \cos x$ .

Будем рассматривать также системы скалярных уравнений

$$\mathbf{F}(\mathbf{x}) = \mathbf{0}, \tag{2}$$

где  $\mathbf{F}(\mathbf{x})$  — некоторая заданная вектор-функция от векторного аргумента  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ ,  $x_1, x_2, \dots, x_n$  — скалярные координаты вектора  $\mathbf{x}$ ; например, в случае

$$\mathbf{F}(\mathbf{x}) = \begin{cases} F_1(x_1, x_2) \\ F_2(x_1, x_2) \end{cases} = \begin{cases} x_1^2 + x_2^2 - 25 \\ x_2 - x_1^2 \end{cases}$$

система (2) в развернутом виде запишется так:

$$\begin{aligned} x_1^2 + x_2^2 - 25 &= 0, \\ x_2 - x_1^2 &= 0. \end{aligned} \tag{3}$$

При решении скалярных уравнений  $F(x) = 0$ , а также систем  $\mathbf{F}(\mathbf{x}) = \mathbf{0}$  таких уравнений возникают две задачи: задача указания областей, в каждой из которых расположено ровно одно решение (задача отделения корней), и задача отыскания корней с заданной точностью (задача уточнения корней).

Для отделения корней нет общих приемов. Используются графики, участки монотонности функций, на которых она меняет знак, и другие частные приемы. Имеется лишь один большой класс функций — многочлены с вещественными коэффициентами, для которого задача отделения вещественных корней решена в общем виде и полностью. Такое решение дает теорема Штурма, которую мы не приводим.

Для уточнения корней используются различные варианты последовательных приближений (итераций).

## § 1. Метод простых итераций

Рассмотрим сначала случай одного уравнения  $F(x) = 0$ , а затем случай системы  $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ .

**1. Случай одного скалярного уравнения.** Пусть известно, что интересующий нас корень  $x$  уравнения

$$F(x) = 0 \tag{1}$$

лежит в интервале  $a < x < b$ . Каким-нибудь способом приводим уравнение (1) к равносильному ему в интервале  $U = \{a < x < b\}$  уравнению вида  $x = f(x)$ . Можно положить

$$f(x) = x - \alpha(x)F(x), \tag{2}$$

где  $\alpha(x)$  — произвольная, не обращающаяся в нуль в точках  $U$  функция. Для отыскания решения  $\tilde{x}$ , принадлежащего интервалу  $a < x < b$ , зададим  $x_0$ , а затем будем последовательно вычислять  $x_1, x_2, \dots$  по формуле

$$x_{p+1} = f(x_p), \quad p = 0, 1, \dots \tag{3}$$

При этом предполагается, что каждое  $x_p$  принадлежит области определения функции  $f(x)$ , так что последовательное вычисление чисел  $x_1, x_2, \dots, x_p, \dots$  возможно.

**Теорема 1.** Пусть  $f(x)$  непрерывна в  $U$ , и пусть при заданном  $x_0$  последовательность (3) сходится к некоторому  $\tilde{x} \in U$ . Тогда  $\tilde{x}$  — корень уравнения (2), т. е. выполняется равенство  $\tilde{x} = f(\tilde{x})$ .

**Доказательство.** Очевидно, что  $\lim_{n \rightarrow \infty} x_{n+1} = \lim_{n \rightarrow \infty} f(x_n)$ . Но в силу определения непрерывности  $\lim_{n \rightarrow \infty} f(x_n) = f\left(\lim_{n \rightarrow \infty} x_n\right)$ . Поэтому  $\lim_{n \rightarrow \infty} x_{n+1} = f\left(\lim_{n \rightarrow \infty} x_n\right)$ , или  $\tilde{x} = f(\tilde{x})$ .  $\square$

О факте сходимости можно судить практически, наблюдая за последовательностью  $x_n$  в процессе вычисления ее членов на компьютере, «поручив» установление этого факта самому компьютеру. Если сходимость установлена, то за корень  $\tilde{x}$  можно приближенно принять член  $x^{(p)}$  последовательности (3) с достаточно большим номером  $p$ .

Однако последовательность (3) не всегда сходится. На рис. 13 изображены графики  $f(x)$ , для одного из которых при указанном на рисунке нулевом приближении  $x_0$  последовательность (3) сходится, а для другого — нет.

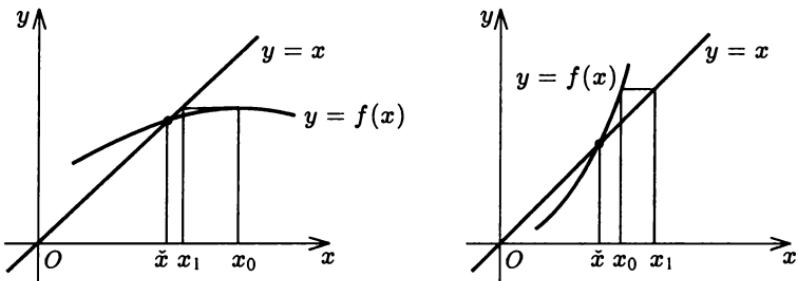


Рис. 13

Сформулируем условия, при которых гарантирована сходимость (3). Введем для этого следующее определение. Отображение (функция)  $f(x)$  называется *сжимающим* в области  $U$  с коэффициентом сжатия  $q$  ( $0 \leq q < 1$ ), если для любых двух  $x', x''$  из  $U$  выполнено неравенство

$$|f(x'') - f(x')| \leq q |x'' - x'|. \quad (4)$$

**Теорема 2.** Пусть уравнение  $x = f(x)$  имеет решение  $\tilde{x}$ , принадлежащее области  $U$ , и пусть отображение  $f(x)$  является сжимающим в области  $U$  с некоторым коэффициентом сжатия  $q$  ( $0 \leq q < 1$ ). Тогда решение  $\tilde{x}$  является единственным решением в области  $U$  и существует столь малое  $R > 0$ , что при выборе  $x_0$  из условия  $|\tilde{x} - x_0| < R$  все члены последовательности (3) будут определены, причем выполнются неравенства

$$|\tilde{x} - x_n| \leq q^n |\tilde{x} - x_0|, \quad n = 0, 1, \dots, \quad (5)$$

т. е. имеет место сходимость  $x_n$  к  $\tilde{x}$  со скоростью геометрической прогрессии со знаменателем  $q$ .

**Доказательство.** Сначала докажем единственность решения. Допустим, что наряду с  $\check{x}$  существует решение  $\check{x}'$ :

$$\check{x}' = f(\check{x}'), \quad \check{x}, \check{x}' \in U.$$

Тогда в силу того, что  $f(x)$  — сжимающее отображение,  $|f(\check{x}) - f(\check{x}')| \leq |q(\check{x} - \check{x}')|$ . Но в силу  $f(\check{x}) = \check{x}$  и  $f(\check{x}') = \check{x}'$  отсюда получаем

$$|\check{x}' - \check{x}| < |\check{x}' - \check{x}|.$$

Противоречие доказывает единственность.

Далее возьмем  $R > 0$  столь малым, чтобы совокупность  $U(\check{x}, R)$  точек, удовлетворяющих неравенству  $|x - \check{x}| < R$ , целиком принадлежала  $U$ . Пусть  $x_0 \in U(\check{x}, R)$ . Тогда

$$|\check{x} - x_1| = |f(\check{x}) - f(x_0)| \leq q|\check{x} - x_0| \leq qR < R,$$

так что  $x_1$  тоже принадлежит окрестности  $U(\check{x}, R)$ . Аналогично  $x_2, x_3, \dots, x_n, \dots$  лежат в этой окрестности.

При  $n = 0$  оценка (5) тривиальна. Допустим, она уже доказана для  $n = 0, 1, \dots, k$ . Докажем ее для  $n = k + 1$ :

$$|\check{x} - x_{k+1}| = |f(\check{x}) - f(x_k)| \leq q|\check{x} - x_k| \leq qq^k|\check{x} - x_0| = qq^{k+1}|\check{x} - x_0|.$$

Таким образом, неравенство (5) доказано для всех  $n$  по индукции.  $\square$

На первый взгляд доказанная теорема ничего не дает для отыскания  $\check{x}$ , поскольку мы не можем указать фактически то  $R > 0$  и ту окрестность  $U(\check{x}, R)$ , в которой надо выбирать  $x_0$ . Однако сам факт существования такой окрестности позволяет организовать автоматическое вычисление  $\check{x}$  на компьютере. Именно, зададим произвольно  $x_0 \in U$ . Если последовательность (3) удается вычислить и она сходится к некоторому  $\check{x}$ , то в силу теоремы 1 это и есть искомый корень. Если сходимости нет, то программа расчета должна взять другое значение  $x'_0$  и повторить расчет; если сходимости опять нет, то выбираем за начальное приближение третью точку  $x''_0$  и т. д. Если в процессе выбора этих начальных приближений  $x_0, x'_0, x''_0, \dots$  они все «гуше» располагаются в  $U$ , то на некотором шаге мы неизбежно попадем в достаточно малую окрестность  $U(\check{x}, R)$  корня, и итерационный процесс сойдется.

Укажем достаточный признак того, что отображение является сжимающим.

**Теорема 3.** Пусть  $f(x)$  имеет производную во всех точках области  $U$  (т. е. интервала  $a < x < b$ ), и пусть существует  $q$  ( $0 \leq q < 1$ ,  $q = \text{const}$ ), такое, что  $|f'(x)| \leq q$  всюду в  $U$ . Тогда отображение  $f(x)$  сжимающее в  $U$  с коэффициентом сжатия  $q$ .

**Доказательство.** По теореме Лагранжа о конечных приращениях для любых  $x', x''$  из  $U$

$$f(x'') - f(x') = (x'' - x')f'(\xi),$$

где  $\xi$  — некоторая точка, лежащая между  $x'$  и  $x''$ . Но тогда

$$|f(x'') - f(x')| = |x'' - x'||f'(\xi)| \leq q|x'' - x'|. \square$$

**2. Случай системы уравнений.** Пусть задана система уравнений

$$\mathbf{F}(\mathbf{x}) = \mathbf{0}, \quad (6)$$

причем вектор-функция  $\mathbf{F}(\mathbf{x})$  определена в некоторой области  $U$  пространства  $R^n$ ,  $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$ , и имеет решение  $\tilde{\mathbf{x}}$  в этой области. Приведем уравнение (6) к равносильному уравнению вида

$$\mathbf{x} = \mathbf{f}(\mathbf{x}). \quad (7)$$

Это, очевидно, можно сделать многими способами, например, положив  $\mathbf{f}(\mathbf{x}) = \mathbf{x} - \alpha(\mathbf{x})\mathbf{F}(\mathbf{x})$ , где  $\alpha(\mathbf{x})$  — произвольная матричная функция, определитель которой отличен от нуля всюду в  $U$ .

Зададим  $x_0 \in U$  и построим последовательность

$$\mathbf{x}_{n+1} = \mathbf{f}(\mathbf{x}_n), \quad n = 0, 1, \dots. \quad (8)$$

Будем считать пространство  $R^n$  нормированным. Тогда остаются справедливыми все определения и утверждения теоремы 1, относящиеся к случаю скалярного уравнения. Нужно только всюду вместо абсолютных величин (чисел) использовать нормы (векторов).

Признак того, что отображение  $\mathbf{f}: R^n \rightarrow R^n$  является сжимающим, сформулированный в теореме 2 для скалярной функции, т. е. в случае  $n = 1$ , в общем случае ( $n \geq 1$ ) формулируется несколько сложнее:

Область  $U$  пространства  $R^n$  называется *выпуклой*, если наряду с любыми двумя точками  $\mathbf{a} = (a_1, a_2, \dots, a_n)$ ,  $\mathbf{b} = (b_1, b_2, \dots, b_n)$ , принадлежащими области  $U$ , ей принадлежат и все точки отрезка  $[\mathbf{a}, \mathbf{b}]$ , содержащего эти точки, т. е. все точки вида

$$\mathbf{x} = \mathbf{a} + t(\mathbf{b} - \mathbf{a}), \quad (9)$$

где  $t$  изменяется от 0 до 1.

**Теорема 4.** Пусть область  $U \subset R^n$  выпуклая, и пусть компоненты  $f_i(x_1, x_2, \dots, x_n)$  вектор-функции

$$\mathbf{f}(\mathbf{x}) = \begin{cases} f_1(x_1, x_2, \dots, x_n), \\ \dots \dots \dots \\ f_n(x_1, x_2, \dots, x_n) \end{cases}$$

имеют равномерно непрерывные производные первого порядка в  $U$ . Рассмотрим матрицу

$$\frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \dots & \dots & \dots \\ \frac{\partial f_n}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_n} \end{bmatrix}.$$

Если норма этой матрицы при всех  $\mathbf{x} \in U$  не превосходит некоторого числа  $q$  ( $0 \leq q < 1$ ), то отображение  $\mathbf{f}(\mathbf{x})$  является сжимающим в области  $U$ :

$$\|\mathbf{f}(\mathbf{x}'') - \mathbf{f}(\mathbf{x}')\| \leq q \|\mathbf{x}'' - \mathbf{x}'\|,$$

где  $\mathbf{x}', \mathbf{x}''$  — произвольные точки из  $U$ .

Доказательства этой теоремы из математического анализа мы не приводим.

Для вычисления решения  $\tilde{x}$  уравнения  $\mathbf{F}(\mathbf{x}) = \mathbf{0}$  методом простых итераций (т. е. с помощью равносильного уравнения  $\mathbf{x} = \mathbf{f}(\mathbf{x})$  и последовательности  $\mathbf{x}_{n+1} = \mathbf{f}(\mathbf{x}_n)$ ,  $n = 0, 1, \dots$ ;  $\mathbf{x}_0$  задано) надо стараться выбрать это равносильное уравнение так, чтобы отображение  $f$  в области  $U$ , содержащей решение  $\tilde{x}$ , было сжимающим с коэффициентом сжатия  $q$ , как можно более близким к нулю.

## § 2. Метод линеаризации Ньютона

Сначала рассмотрим случай одного уравнения  $F(x) = 0$ , а потом — системы  $\mathbf{F}(\mathbf{x}) = \mathbf{0}$ .

**1. Метод линеаризации для одного уравнения.** Зададим  $x_0$ . Пусть некоторое приближение  $x_n$  к корню  $\tilde{x}$  уравнения  $F(x) = 0$  уже найдено. Воспользуемся приближенной формулой

$$F(x) \approx F(x_n) + F'(x_n)(x - x_n),$$

точность которой возрастает при приближении  $x$  к  $x_n$ . Вместо исходного уравнения  $F(x) = 0$  воспользуемся линейным уравнением

$$F(x_n) + F'(x_n)(x - x_n) = 0.$$

Решение этого уравнения примем за приближение  $x_{n+1}$ :

$$x_{n+1} = x_n - [F'(x_n)]^{-1} F(x_n), \quad (1)$$

$$n = 0, 1, \dots$$

Метод линеаризации Ньютона допускает простую геометрическую интерпретацию (рис. 14). Криволинейный график функции  $F(x)$  заменяется касательной к нему в точке  $(x_n, F(x_n))$ . За приближение  $x_{n+1}$  принимается точка пересечения этой касательной с осью абсцисс.

Формулу (1) можно интерпретировать следующим образом. Считая, что  $F'(x) \neq 0$ , переходим от уравнения  $F(x) = 0$  к равносильному уравнению

$$x = f(x),$$

где

$$f(x) = x - [F'(x)]^{-1} F(x),$$

и пользуемся простыми итерациями. Очевидно, что в точке  $\tilde{x}$ , где  $F(\tilde{x}) = 0$ ,

$$\frac{df}{dx} \Big|_{x=\tilde{x}} = 0.$$

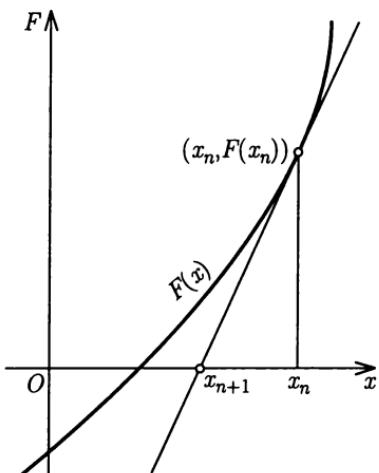


Рис. 14

Будем считать, что  $f''(x)$  непрерывна. Тогда для каждого  $0 < q < 1$  найдется (малая) окрестность точки  $\tilde{x}$ , в которой  $|f'(x)| \leq q$ .

В силу теоремы 3 из § 1 отображение  $f(x)$  в этой окрестности будет сжимающим с этим коэффициентом сжатия  $q$ , т. е. тем сильнее сжимающим, чем меньше окрестность корня  $\tilde{x}$ .

В соответствии с этим последовательность (1) (т. е. приближения, полученные методом линеаризации) будет сходиться к  $\tilde{x}$ , если точку  $x_0$  выбрать достаточно близко к  $\tilde{x}$ , причем скорость убывания погрешности  $|\tilde{x} - x_n|$  будет быстрее, чем скорость убывания членов геометрической прогрессии  $q^n$ , каково бы ни было фиксированное  $q$  ( $0 < q < 1$ ).

## 2. Метод линеаризации для системы.

Рассмотрим систему

$$\mathbf{F}(\mathbf{x}) = \mathbf{0}. \quad (2)$$

Пусть в некоторой области  $U$  есть решение  $\tilde{\mathbf{x}}$  системы (2). Будем считать, что компоненты  $F_j(x_1, x_2, \dots, x_n)$  ( $j = 1, 2, \dots, n$ ) функций, составляющих вектор-функцию  $\mathbf{F}(\mathbf{x})$ , являются достаточно гладкими функциями.

Пусть приближение  $\mathbf{x}_n$  к решению уже найдено. Воспользуемся приближенной формулой

$$\mathbf{F}(\mathbf{x}) \approx \mathbf{F}(\mathbf{x}_n) + \frac{d\mathbf{F}(\mathbf{x})}{d\mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}_n} (\mathbf{x} - \mathbf{x}_n)$$

и вместо уравнения (2) рассмотрим линейную систему уравнений

$$\mathbf{F}(\mathbf{x}_n) + \frac{d\mathbf{F}(\mathbf{x})}{d\mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}_n} (\mathbf{x} - \mathbf{x}_n) = \mathbf{0}. \quad (3)$$

Решение этой системы примем за  $\mathbf{x}_{n+1}$ .

Очевидно, что  $\mathbf{x}_{n+1}$  выражается формулой

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \left[ \frac{d\mathbf{F}(\mathbf{x}_n)}{d\mathbf{x}} \right]^{-1} \mathbf{F}(\mathbf{x}_n), \quad (4)$$

однако фактически решение линейной системы (3) можно искать одним из изложенных выше способов решения систем линейных уравнений (методы Гаусса и другие).

Для реализуемости метода Ньютона надо предполагать, что матрица  $d\mathbf{F}(\mathbf{x})/d\mathbf{x}$  в точке  $\mathbf{x} = \tilde{\mathbf{x}}$  невырождена, так что  $[d\mathbf{F}(\mathbf{x})/d\mathbf{x}]^{-1}$  существует. В силу непрерывности она невырождена и в малой окрестности точки  $\tilde{\mathbf{x}}$ .

Как и в случае одного уравнения, можно считать, что мы перешли от (2) к равносильному уравнению

$$\mathbf{x} = \mathbf{f}(\mathbf{x}),$$

где

$$\mathbf{f}(\mathbf{x}) = \mathbf{x} - \left[ \frac{d\mathbf{F}(\mathbf{x})}{d\mathbf{x}} \right]^{-1} \mathbf{F}(\mathbf{x}), \quad (5)$$

а затем воспользовались методом простых итераций. Можно проверить, что матрица  $d\mathbf{f}(\mathbf{x})/d\mathbf{x}$  в точке  $\tilde{\mathbf{x}}$  есть нулевая матрица. Поэтому, как и в случае одного уравнения, отображение  $\mathbf{f}(\mathbf{x})$ , задаваемое формулой (5), является сжимающим в окрестности точки  $\tilde{\mathbf{x}}$  с коэффициентом сжатия  $q$ , который тем ближе к нулю, чем меньше эта окрестность.

Доказательство аналогично случаю одного уравнения, но вместо теоремы 3 из § 1 опирается на теорему 4 из § 1 и более громоздко. В силу теоремы 4 из § 1 при выборе  $x_0$  из достаточно малой окрестности решения  $\tilde{x}$  процесс последовательных приближений (4) сходится, причем скорость убывания погрешности быстрее геометрической прогрессии  $q^p$ , как бы ни было мало фиксированное  $q$  ( $0 < q < 1$ ).

### Задачи

**1\***. Пусть уравнение  $F(x) = 0$  имеет корень  $\tilde{x}$  на отрезке  $[a, b]$ , причем на этом отрезке функция имеет ограниченную вторую производную, а  $F'(x) > 0$ .

Доказать, что погрешности  $|\tilde{x} - x_n|$  приближений  $x_n$ , полученные методом Ньютона, удовлетворяют оценкам

$$|\tilde{x} - x_{n+1}| \leq \text{const} |\tilde{x} - x_n|^2, \quad n = 0, 1, \dots,$$

если только нулевое приближение  $x_0$  выбрано достаточно близко к  $\tilde{x}$ .

**2.** Построить алгоритм для вычисления  $\sqrt{5}$  с заданным числом верных десятичных знаков, рассматривая  $\sqrt{5}$  как решение уравнения  $x^2 - 5 = 0$ . Воспользоваться методом Ньютона. Показать, что метод Ньютона сходится при произвольном выборе  $x_0 > 0$ .

**3.** Пусть графики некоторых заданных дифференцируемых функций  $y = \varphi(x)$ ,  $y = \psi(x)$  построены на плоскости  $Oxy$  и пересекаются в некоторой точке  $x = \tilde{x}$ . Для отыскания корня  $\tilde{x}$  уравнения  $F(x) \equiv \varphi(x) - \psi(x) = 0$  предполагается воспользоваться методом Ньютона. Пусть приближение  $x_n$  уже найдено и отмечено на оси абсцисс.

а) Построить следующее приближение  $x_{n+1}$  метода Ньютона, указав точку  $x_{n+1}$  на оси абсцисс.

б) Пусть графики  $y = \varphi(x)$ ,  $y = \psi(x)$  направлены один выпуклостью вверх, а другой выпуклостью вниз. Показать, что тогда метод Ньютона сходится для любого  $x_0 > \tilde{x}$ .

**4.** Построить алгоритмы для вычисления вещественных решений следующих систем скалярных уравнений и вычислить решения, получив результаты с пятью верными знаками:

а)  $\sin x - y = 1,30$ ,  $\cos y - x = -0,84$ ;

б)  $x^2 + 4y^2 = 1$ ,  $x^4 + y^4 = 0,5$ .

**5.** Для численного решения краевой задачи

$$\frac{d^2y}{dx^2} - y^3 = x^2, \quad 0 \leq x \leq 1, \quad y(0) = 1, \quad y(1) = 3,$$

на отрезке  $0 \leq x \leq 1$  введена сетка  $x_k = k/N$  ( $k = 0, 1, \dots, N$ ), а для неизвестных  $y_k = y(x_k)$  составлена разностная схема

$$\frac{y_{k+1} - 2y_k + y_{k-1}}{h^2} - y_k^3 = (kh)^2, \quad k = 1, 2, \dots, N-1,$$

$$y_0 = 1; \quad y_N = 3, \quad h = 1/N.$$

Для решения этой системы ( $N+1$ ) нелинейных уравнений относительно того же числа неизвестных  $y_0, y_1, \dots, y_N$  воспользоваться методом Ньютона.

а) Примем за начальное приближение  $y^{(0)} = \{y_0^{(0)}, y_1^{(0)}, \dots, y_N^{(0)}\}$  функцию  $y_k^{(0)} = 1 + 2kh$  ( $k = 0, 1, \dots, N$ ), определенную на сетке  $x_k$  и удовлетворяющую поставленным краевым условиям  $y_0 = 1, y_N = 3$ . Построить линейную систему метода Ньютона, позволяющую по уже найденному приближению  $y^{(n)} = \{y_k^{(n)}\}$  найти поправки  $\varepsilon^{(n)} = \{\varepsilon_k^{(n)}\}$  для получения следующего приближения:

$$y^{(n+1)} = \{y_k^{(n)} + \varepsilon_k^{(n)}\}.$$

б) Показать, что линейную систему уравнений для поправки  $\varepsilon^{(n)}$

$$\frac{\varepsilon_{k+1}^{(n)} - 2\varepsilon_k^{(n)} + \varepsilon_{k-1}^{(n)}}{h^2} - 3(y_k^{(n)})^2 \varepsilon_k^{(n)} = (kh)^2 + (y_k^{(n)})^3, \quad k = 1, 2, \dots, N-1,$$

$$\varepsilon_0^{(n)} = \varepsilon_N^{(n)} = 0,$$

можно решить методом прогонки.

в) Составить программу для вычисления решения исходной краевой задачи методом Ньютона на компьютере и найти таблицу значений решения, подобрав с помощью экспериментальных расчетов число  $N$  и число итераций так, чтобы таблица  $y^{(n)} = \{y_k^{(n)}\}$  имела четыре верных десятичных знака.

### ЧАСТЬ III

## МЕТОД КОНЕЧНЫХ РАЗНОСТЕЙ ДЛЯ ЧИСЛЕННОГО РЕШЕНИЯ ДИФФЕРЕНЦИАЛЬНЫХ УРАВНЕНИЙ

---

Самыми распространенными методами вычисления решений задач для обыкновенных дифференциальных уравнений и уравнений с частными производными являются метод конечных разностей (МКР) и его модификации. Во всех вариантах этого метода в области определения искомых функций вводится сетка и решение ищется на сетке. Для значений искомой сеточной функции строится система скалярных уравнений, решение которой и служит приближенной таблицей значений решения исходной задачи.

Простейший способ построения этой системы скалярных уравнений — разностной схемы — состоит в приближенной замене производных, входящих в дифференциальное уравнение и в краевые условия, разностными отношениями, чем и объясняется название метода.

### ГЛАВА 8

## ЧИСЛЕННОЕ РЕШЕНИЕ ЗАДАЧ ДЛЯ ОБЫКНОВЕННЫХ ДИФФЕРЕНЦИАЛЬНЫХ УРАВНЕНИЙ

Пусть на некотором отрезке  $a \leq x \leq b$  задано обыкновенное дифференциальное уравнение (или система уравнений) и требуется найти решение этого уравнения, удовлетворяющее дополнительным начальным или краевым условиям.

### § 1. Примеры разностных схем. Сходимость

Будем записывать краевые задачи символическим равенством

$$Lu = f. \quad (1)$$

Приведем примеры таких задач.

Пример задачи с начальными условиями для обыкновенного дифференциального уравнения первого порядка:

$$\frac{du}{dx} + \frac{x}{1+u^2} = \cos x, \quad 0 \leq x \leq 1, \quad u(0) = 3. \quad (2)$$

Пример задачи с краевыми условиями, поставленными в начале  $x = 0$  и в конце  $x = 1$  отрезка:

$$\frac{d^2u}{dx^2} = (1+x^2)u + \sqrt{x+1}, \quad 0 \leq x \leq 1, \quad (3)$$

$$u(0) = 2, \quad u(1) = 1. \quad (4)$$

Пример задачи с начальными условиями для системы двух уравнений первого порядка:

$$\begin{aligned} \frac{dv}{dx} + xvw &= x^2 - 3x + 1, \quad 0 \leq x \leq 1, \\ \frac{dw}{dx} + \frac{1}{1+x^2}(v+w) &= \cos^2 x, \quad 0 \leq x \leq 1, \\ v(0) &= 1, \quad w(0) = -3. \end{aligned} \quad (5)$$

Решением этой задачи является вектор-функция  $u = (v, w)$ , компоненты которой удовлетворяют равенствам (5).

Во всех примерах мы рассматриваем задачу на отрезке  $\overline{D} = \{0 \leq x \leq 1\}$ , а не на каком-либо другом, только для определенности.

**1. Примеры разностных схем.** Будем предполагать, что решение  $u(x)$  задачи (1) на отрезке  $0 \leq x \leq 1$  существует. Для вычисления этого решения с помощью метода конечных разностей (или метода сеток) надо прежде всего выбрать на отрезке  $\overline{D}$  конечное число точек, совокупность которых будем называть *сеткой* и обозначать через  $D_h$ , а затем считать искомым не решение  $u(x)$  задачи (1), а таблицу  $[u]_h$  значений этого решения в точках сетки  $D_h$ . Предполагается, что сетка  $D_h$  зависит от параметра  $h$ , который может принимать сколь угодно малые положительные значения. При стремлении шага сетки  $h$  к нулю сетка должна становиться все гуще. Например, можно положить  $h = 1/N$ , где  $N$  — какое-нибудь натуральное число, и принять за сетку  $D_h$  совокупность точек  $x_0 = 0, x_1 = h, x_2 = 2h, \dots, x_N = 1$ . В этом случае искомая сеточная функция  $[u]_h$  в точке  $x_n$  сетки принимает значение, которое будем обозначать  $u_n$ .

Для приближенного вычисления таблицы значений решения  $[u]_h$  в случае задачи (2) можно воспользоваться, например, системой уравнений

$$\frac{u_{n+1} - u_n}{h} + \frac{x_n}{1+u_n^2} = \cos x_n, \quad n = 0, 1, \dots, N-1, \quad u_0 = 3, \quad (6)$$

полученной в результате замены производной  $du/dx$  в точках сетки разностным отношением по приближенной формуле

$$\frac{du}{dx} \approx \frac{u(x+h) - u(x)}{h}.$$

Решение  $u^{(h)} = (u_0^{(h)}, u_1^{(h)}, \dots, u_N^{(h)})$  системы (6) определено на той же сетке  $D_h$ , что и искомая сеточная функция  $[u]_h$ . Его значения  $u_0^{(h)}, u_1^{(h)}, \dots, u_N^{(h)}$  в точках  $x_0, x_1, \dots, x_N$  последовательно вычисляются из системы (6) при  $n = 0, 1, \dots, N - 1$ .

(Для краткости в уравнениях (6) опущен значок  $h$  при  $u_n^{(h)}$ ; как правило, мы будем так же поступать в аналогичных случаях и в дальнейшем.)

В случае задачи (4) для отыскания сеточной функции  $u^{(h)}$ , приближенно совпадающей с искомой таблицей решения  $[u]_h$ , можно воспользоваться разностной схемой

$$\frac{u_{n+1} - 2u_n + u_{n-1}}{h^2} - (1 + x_n^2) u_n = \sqrt{x_n + 1}, \quad (7)$$

$$u_0 = 2, \quad u_N = 1, \quad n = 1, 2, \dots, N - 1.$$

Эта схема возникает в результате замены в точках сетки производной  $d^2u/dx^2$ , входящей в дифференциальное уравнение, по приближенной формуле

$$\frac{d^2u}{dx^2} \approx \frac{u(x+h) - 2u(x) + u(x-h)}{h^2}. \quad (8)$$

Для вычисления решения  $u^{(h)}$  этой задачи можно воспользоваться алгоритмом исключения (прогонки), описанным в § 4 гл. 4.

Выпишем еще разностную схему, пригодную для вычисления решения задачи (5):

$$\begin{aligned} \frac{v_{n+1} - v_n}{h} + x_n v_n w_n &= x_n^3 - 3x_n + 1, \\ \frac{w_{n+1} - w_n}{h} + \frac{1}{1+x_n^2} (v_n + w_n) &= \cos^2 x_n, \end{aligned} \quad n = 0, 1, \dots, N - 1, \quad (9)$$

$$v_0 = 1, \quad w_0 = -3.$$

Здесь  $u_0^{(h)} = \begin{bmatrix} v_0 \\ w_0 \end{bmatrix} = \begin{bmatrix} 1 \\ -3 \end{bmatrix}$  задано. При  $n = 0$  из уравнений (9) можно найти  $u_1^{(h)} = \begin{bmatrix} v_1 \\ w_1 \end{bmatrix}$ . Вообще, зная  $u_k = \begin{bmatrix} v_k \\ w_k \end{bmatrix}$  ( $k = 0, 1, \dots, n$ ), можно вычислить  $u_{n+1} = \begin{bmatrix} v_{n+1} \\ w_{n+1} \end{bmatrix}$ .

В рассмотренных примерах сетка  $D_h$  состоит из удаленных друг от друга на расстояние  $h$  точек. Ясно, что можно было бы расположить  $(N + 1)$  точек сетки  $D_h$  ( $h = 1/N$ ) на отрезке  $[0, 1]$  не равномерно, а так, чтобы выполнялись условия

$$x_0 = 0, \quad x_1 = x_0 + h_0, \quad x_2 = x_1 + h_1, \dots, \quad x_{n+1} = x_n + h_n, \dots, \quad x_N = 1,$$

где  $h_n > 0$  ( $n = 0, 1, \dots, N - 1$ ) — вообще говоря, не равные между собой числа, однако такие, что  $\max_n h_n \rightarrow 0$  при  $n \rightarrow \infty$ . Выбором

расположения узлов сетки  $D_h$  можно добиться того, чтобы искомая таблица  $[u]_h$  решения  $u(x)$  была подробнее при данном фиксированном  $N$  на тех участках, где  $u(x)$  изменяется быстрее. Такие участки иногда бывают заранее известны из физического смысла задачи или из предварительных грубых расчетов. Информация о скорости изменения  $u(x)$  выявляется также в ходе последовательного вычисления  $u_1, u_2, \dots, u_n$  и может быть учтена при выборе следующего узла сетки  $x_{n+1}$ .

Мы ограничимся приведенными примерами для иллюстрации понятий сетки и искомой сеточной функции (или вектор-функции) — таблицы значений решения  $[u]_h$ .

При измельчении сетки, т. е. при  $h \rightarrow 0$ , сеточная функция  $u^{(h)}$  является все более подробной таблицей искомого решения  $u(x)$  и дает о нем все более полное представление. Пользуясь интерполяцией, можно было бы с возрастающей при  $h \rightarrow 0$  точностью восстановить решение  $u$  всюду в области  $D$ . Ясно, что точность, с которой это можно сделать при заданных фиксированном числе и расположении узлов сетки  $D_h$ , зависит от дополнительно известных сведений о решении (типа оценок для его производных), а также от расположения узлов сетки.

Ограничимся этими беглыми замечаниями о восстановлении функции  $u$  по ее таблице  $u^{(h)}$ . Подробное рассмотрение вопросов восстановления функции по ее таблице составляет предмет теории интерполяции. Мы будем заниматься только задачей вычисления таблицы  $u^{(h)}$ . Поэтому условимся считать, что задача (1) решена точно, если найдена сеточная функция  $[u]_h$ . Однако в общем случае нам не удастся вычислить ее точно. Вместо сеточной функции  $[u]_h$ , будем искать другую сеточную функцию  $u^{(h)}$ , которая сходится к  $[u]_h$  при измельчении сетки. Для этой цели можно использовать разностные уравнения.

**2. Сходящиеся разностные схемы.** Способами построения и исследования сходящихся разностных схем мы будем заниматься на протяжении всей главы. Однако прежде надо придать точный смысл самому требованию сходимости  $u^{(h)} \rightarrow [u]_h$ , которое мы будем предъявлять к разностным схемам. Для этого рассмотрим линейное нормированное пространство функций  $U_h$ , определенных на сетке  $D_h$ . Норма  $\|u^{(h)}\|_{U_h}$  сеточной функции  $u^{(h)} \in U_h$  есть неотрицательное число, которое принимается за меру отклонения функции  $u^{(h)}$  от тождественного нуля.

Норма может быть определена различными способами. Можно, например, принять за норму функции точную верхнюю грань модуля ее значений в точках сетки, положив

$$\|u^{(h)}\|_{U_h} = \max_n |u_n|. \quad (10)$$

Если  $u^{(h)} = \begin{bmatrix} v_k \\ w_k \end{bmatrix}$  ( $k = 0, 1, \dots, N$ ), как в (9), то за норму, аналогичную (10), можно принять верхнюю грань модулей обеих функций  $v_k$ ,  $w_k$  на соответствующей им сетке.

Если  $u^{(h)}$  состоит из функций, определенных на сетке  $x_n = nh$  ( $n = 0, 1, \dots, N$ ), то часто используют норму, определенную равенством

$$\|u^{(h)}\|_{U_h} = \left( h \sum_{n=0}^N u_n^2 \right)^{1/2}.$$

Эта норма аналогична норме

$$\|u(x)\| = \left( \int_0^1 u^2 dx \right)^{1/2}$$

для функций  $u(x)$  с интегрируемым на отрезке  $0 \leq x \leq 1$  квадратом.

Всюду, где не оговорено противное, мы будем пользоваться нормой (10).

После того как введено нормированное пространство  $U_h$ , приобретает смысл понятие отклонения одной функции от другой. Если  $a^{(h)}$ ,  $b^{(h)}$  — произвольные сеточные функции из  $U_h$ , то мерой их отклонения друг от друга считается норма их разности, т. е. число

$$\|b^{(h)} - a^{(h)}\|_{U_h}.$$

Теперь можно перейти к строгому определению сходящейся разностной схемы.

Пусть для приближенного вычисления решения дифференциальной краевой задачи (1), т. е. для приближенного вычисления сеточной функции  $[u]_h$  на основе использования равенства (1), составлена некоторая система уравнений, которую будем символически записывать аналогично уравнению (1) в форме равенства

$$L_h u^{(h)} = f^{(h)}. \quad (11)$$

Примерами могут служить разностные схемы (6), (7), (9) для дифференциальных краевых задач (2), (4), (5) соответственно. Для записи схемы (6) в форме (11) можно положить

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - u_n}{h} + \frac{nh}{1 + u_n^2}, & n = 0, 1, \dots, N - 1, \\ u_0, & \end{cases}$$

$$f^{(h)} = \begin{cases} \cos nh, & n = 0, 1, \dots, N - 1, \\ 3. & \end{cases}$$

Схема (7) запишется в форме (11), если принять

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - 2u_n + u_{n-1}}{h^2} - [1 + (nh)^2]u_n, & n = 1, 2, \dots, N - 1, \\ u_0, \\ u_N, & \end{cases}$$

$$f^{(h)} = \begin{cases} \sqrt{1 + nh}, & n = 1, 2, \dots, N - 1, \\ 2, \\ 1. & \end{cases}$$

Запишем еще в виде (11) схему (9), приняв

$$L_h u^{(h)} = L_h \begin{bmatrix} v^{(h)} \\ w^{(h)} \end{bmatrix} =$$

$$= \begin{cases} \frac{v_{n+1} - v_n}{h} + nh v_n w_n, & n = 0, 1, \dots, N-1, \\ \frac{w_{n+1} - w_n}{h} + \frac{1}{1 + (nh)^2} (v_n + w_n), & n = 0, 1, \dots, N-1, \\ v_0, \\ w_0, \end{cases}$$

$$f^{(h)} = \begin{cases} (nh)^2 - 3nh + 1, & n = 0, 1, \dots, N-1, \\ \cos^2 nh, & n = 0, 1, \dots, N-1, \\ 1, \\ -3. \end{cases}$$

Система (11), как видим, зависит от  $h$  и должна быть выписана для всех тех  $h$ , для которых рассматриваются сетка  $D_h$  и сеточная функция  $[u]_h$ . Таким образом, разностная краевая задача (11) — это не одна система, а семейство систем, зависящее от параметра  $h$ . Будем предполагать, что при каждом рассматриваемом достаточно малом  $h$  существует решение  $u^{(h)}$  задачи (11), принадлежащее пространству  $U_h$ .

Будем говорить, что решение  $u^{(h)}$  разностной краевой задачи (11) при измельчении сетки *сходится* к решению дифференциальной краевой задачи (1), если

$$\|[u]_h - u^{(h)}\|_{U_h} \rightarrow 0, \quad h \rightarrow 0. \quad (12)$$

Если сверх того выполняется неравенство

$$\|[u]_h - u^{(h)}\|_{U_h} \leq c h^k, \quad (13)$$

где  $c > 0$ ,  $k > 0$  — некоторые постоянные, не зависящие от  $h$ , то будем говорить, что имеет место *сходимость порядка  $h^k$* , или что *разностная схема имеет  $k$ -й порядок точности*.

Сходимость является фундаментальным требованием, которое предъявляется к разностной схеме (11) для численного решения дифференциальной краевой задачи (1). Если она имеет место, то с помощью разностной схемы (11) можно вычислить решение  $[u]_h$  с любой наперед заданной точностью, выбирая для этого  $h$  достаточно малым. Мы точно сформулировали понятие сходимости и подошли к центральному вопросу о том, как построить сходящуюся разностную схему (11) для вычисления решения дифференциальной краевой задачи (1). Приведенные выше примеры схем дополняют рассмотренные в § 1 и дают представление о простейшем способе построения таких схем: следует выбрать сетку и заменить производные разностными отношениями. Однако для одной и той же дифференциальной краевой задачи можно получить различные разностные схемы (11), по-разному

выбирая сетку  $D_h$  и по-разному заменяя производные приближающими их разностными отношениями.

**3. Проверка сходимости разностной схемы.** Не будем пока заниматься построением разностных схем и поставим задачу несколько иначе. Пусть разностная схема  $L_h u^{(h)} = f^{(h)}$ , позволяющая надеяться, что сходимость

$$\|[u]_h - u^{(h)}\|_{U_h} \rightarrow 0, \quad h \rightarrow 0,$$

имеет место, на основании тех или иных соображений уже построена. Как проверить, является ли она в самом деле сходящейся?

Предположим, что разностная задача (11) имеет единственное решение  $u^{(h)}$ . Если бы при подстановке в левую часть (11) вместо  $u^{(h)}$  сеточной функции  $[u]_h$  равенство (11) оказалось в точности выполненным, то в силу единственности решения имело бы место равенство  $u^{(h)} = [u]_h$ , идеальное с точки зрения сходимости. Это означало бы, что решение  $u^{(h)}$  разностной задачи  $L_h u^{(h)} = f^{(h)}$  совпадает с исключенной сеточной функцией  $[u]_h$ , которую мы условились считать точным решением.

Однако, как правило, систему (11) не удается выбрать так, чтобы  $[u]_h$  в точности ей удовлетворяла. При подстановке  $[u]_h$  в уравнение (11) возникает некоторая невязка:

$$L_h [u]_h = f^{(h)} + \delta f^{(h)}. \quad (14)$$

Если эта невязка  $\delta f^{(h)}$  стремится к нулю при  $h \rightarrow 0$ , так что  $[u]_h$  удовлетворяет уравнению (11) все точнее, то будем говорить, что разностная схема  $L_h u^{(h)} = f^{(h)}$  аппроксимирует дифференциальную краевую задачу  $Lu = f$  на решении  $u$  последней.

В случае аппроксимации можно считать, что уравнение (14), которому удовлетворяет  $[u]_h$ , получается из уравнения (11) путем прибавления некоторой малой (при малом  $h$ ) добавки  $\delta f^{(h)}$  к правой части  $f^{(h)}$ . Следовательно, если решение  $u^{(h)}$  задачи (11) устойчиво относительно возмущения правой части  $f^{(h)}$ , т. е. мало изменяется при малом изменении правой части, то решение  $u^{(h)}$  задачи (11) и решение  $[u]_h$  задачи (14) различаются мало, так что из аппроксимации  $\delta f^{(h)} \rightarrow 0$  следует сходимость

$$u^{(h)} \rightarrow [u]_h, \quad h \rightarrow 0.$$

Намеченный нами путь проверки сходимости (12) состоит в том, чтобы разбить этот трудный вопрос на два более простых: сначала проверить, имеет ли место аппроксимация задачи (1) задачей (11), а затем выяснить, устойчива ли задача (11). В этом содержится и указание на способы построения сходящихся разностных схем для численного решения задачи (1): надо строить аппроксимирующую ее разностную схему; из многих возможных способов аппроксимации надо выбирать такие, при которых разностные схемы оказываются устойчивыми.

Изложенный общий план исследования сходимости, естественно, предполагает, что введены математически строгие понятия аппроксимации и устойчивости, позволяющие доказать, что из аппроксимации и устойчивости следует сходимость. Намеченные выше определения аппроксимации и устойчивости не являются строгими. Для определения аппроксимации надо еще уточнить, что такое невязка  $\delta f^{(h)}$  в общем случае и что такое ее величина, а для определения устойчивости — придать точный смысл словам «малому возмущению правой части соответствует малое возмущение решения разностной задачи  $L_h u^{(h)} = f^{(h)}$ ». Строгим определениям понятий аппроксимации и устойчивости мы посвятим отдельные параграфы.

## § 2. Апроксимация дифференциальной краевой задачи разностной схемой

**1. Невязка  $\delta f^{(h)}$ .** Придадим точный смысл понятию аппроксимации дифференциальной краевой задачи (1) из § 1:

$$Lu = f \quad (1)$$

на решении  $u$  разностной схемой (11) из § 1:

$$L_h u^{(h)} = f^{(h)}. \quad (2)$$

Для этого надо уточнить, что такое невязка  $\delta f^{(h)}$  в равенстве

$$L_h [u]_h = f^{(h)} + \delta f^{(h)}, \quad (3)$$

возникающая при подстановке сеточной функции  $[u]_h$  — таблицы исходного решения  $u(x)$  — в уравнение (2), а также, что такое ее величина. Стремление величины невязки  $\delta f^{(h)}$  к нулю при  $h \rightarrow 0$  мы и примем затем за определение аппроксимации.

Начнем с рассмотрения примера разностной схемы для численного решения дифференциальной краевой задачи

$$\begin{aligned} \frac{d^2u}{dx^2} + a(x) \frac{du}{dx} + b(x)u &= \cos x, \quad 0 \leq x \leq 1, \\ u(0) &= 1, \quad u'(0) = 2. \end{aligned} \quad (4)$$

За сетку  $D_h$  по-прежнему примем совокупность точек  $x_n = nh$  ( $n = 0, 1, \dots, N$ ,  $h = 1/N$ ). В качестве разностной схемы для приближенного вычисления  $[u]_h$  воспользуемся совокупностью равенств:

$$\begin{aligned} \frac{u_{n+1} - 2u_n + u_{n-1}}{h^2} + a(x_n) \frac{u_{n+1} - u_{n-1}}{2h} + b(x_n)u_n &= \cos x_n, \\ n &= 1, 2, \dots, N-1, \\ u_0 &= 1, \quad \frac{u_1 - u_0}{h} = 2, \end{aligned} \quad (5)$$

возникающей при замене производных в (4) по приближенным формулам:

$$\begin{aligned}\frac{d^2 u(x)}{dx^2} &\approx \frac{u(x+h) - 2u(x) + u(x-h)}{h^2}, \\ \frac{du(x)}{dx} &\approx \frac{u(x+h) - u(x-h)}{2h}, \\ \frac{du(0)}{dx} &\approx \frac{u(h) - u(0)}{h}.\end{aligned}\quad (6)$$

Разностная схема (5) записывается в форме (2), если обозначить

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - 2u_n + u_{n-1}}{h^2} + a(nh) \frac{u_{n+1} - u_{n-1}}{2h} + b(nh) u_n, & n = 1, 2, \dots, N-1, \\ u_0, \\ \frac{u_1 - u_0}{h}, \end{cases}$$

$$f^{(h)} = \begin{cases} \cos nh, & n = 1, 2, \dots, N-1, \\ 1, \\ 2. \end{cases} \quad (7)$$

Для вычисления и оценки величины невязки  $\delta f^{(h)}$ , возникающей при подстановке  $[u]_h$  в уравнение (2), уточним формулы (6). По формуле Тейлора имеем:

$$\begin{aligned}u(x+h) &= u(x) + hu'(x) + \frac{h^2}{2} u''(x) + \frac{h^3}{6} u'''(\xi_1), \\ u(x-h) &= u(x) - hu'(x) + \frac{h^2}{2} u''(x) - \frac{h^3}{6} u'''(\xi_2), \\ u(x+h) &= u(x) + hu'(x) + \frac{h^2}{2} u''(x) + \frac{h^3}{6} u'''(x) + \frac{h^4}{24} u^{(4)}(\xi_3), \\ u(x-h) &= u(x) - hu'(x) + \frac{h^2}{2} u''(x) - \frac{h^3}{6} u'''(x) + \frac{h^4}{24} u^{(4)}(\xi_4), \\ u(x+h) &= u(x) + hu'(x) + \frac{h^2}{2} u''(\xi_5).\end{aligned}$$

Здесь  $\xi_1, \xi_2, \xi_3, \xi_4, \xi_5$  — некоторые промежуточные точки отрезка  $[x-h, x+h]$ . Отсюда:

$$\begin{aligned}\frac{u(x+h) - u(x-h)}{2h} &= u'(x) + \frac{h^2}{12} [u'''(\xi_1) + u'''(\xi_2)], \\ \frac{u(x+h) - 2u(x) + u(x-h)}{h^2} &= u''(x) + \frac{h^2}{24} [u^{(4)}(\xi_3) + u^{(4)}(\xi_4)], \\ \frac{u(x+h) - u(x)}{h} &= u'(x) + \frac{h}{2} u''(\xi_5).\end{aligned}\quad (8)$$

**2. Вычисление невязки.** Будем считать, что решение  $u(x)$  задачи (4) имеет ограниченные производные до четвертого порядка. В силу формул (8) можно написать

$$\begin{aligned} \frac{u(x+h) - 2u(x) + u(x-h)}{h^2} + a(x) \frac{u(x+h) - u(x-h)}{2h} + b(x) u(x) = \\ = \frac{d^2 u(x)}{dx^2} + a(x) \frac{du(x)}{dx} + b(x) u(x) + \\ + h^2 \left[ \frac{u^{(4)}(\xi_3) + u^{(4)}(\xi_4)}{24} + a(x) \frac{u^{(3)}(\xi_1) + u^{(3)}(\xi_2)}{12} \right], \\ \xi_j \in [x-h, x+h]. \end{aligned}$$

Поэтому выражение

$$L_h[u]_h = \begin{cases} \frac{u(x_n+h) - 2u(x_n) + u(x_n-h)}{h^2} + \\ + a(x_n) \frac{u(x_n+h) - u(x_n-h)}{2h} + b(x_n) u(x_n), & n = 1, 2, \dots, N-1, \\ u(0), \\ \frac{u(h) - u(0)}{h} \end{cases}$$

можно переписать так:

$$L_h[u]_h = \begin{cases} \cos(nh) + h^2 \left[ \frac{u^{(4)}(\xi_3) + u^{(4)}(\xi_4)}{24} + a(nh) \frac{u^{(3)}(\xi_1) + u^{(3)}(\xi_2)}{12} \right], & n = 1, 2, \dots, N-1, \\ 1 + 0, \\ 2 + h \frac{u''(\xi_5)}{2}, \end{cases}$$

или

$$L_h[u]_h = f^{(h)} + \delta f^{(h)},$$

где

$$\delta f^{(h)}|_n = \begin{cases} h^2 \left[ \frac{u^{(4)}(\xi_3) + u^{(4)}(\xi_4)}{24} + a(nh) \frac{u^{(3)}(\xi_1) + u^{(3)}(\xi_2)}{12} \right], & n = 1, 2, \dots, N-1, \\ 0, \\ h \frac{u''(\xi_5)}{2}. \end{cases} \quad (9)$$

Удобно считать, что  $f^{(h)}$ ,  $\delta f^{(h)}$  принадлежат линейному нормированному пространству  $F_h$ , которое состоит из элементов вида

$$g^{(h)} = \begin{cases} \varphi_n, & n = 1, 2, \dots, N-1, \\ \psi_0, \\ \psi_1, \end{cases} \quad (10)$$

где  $\varphi_1, \varphi_2, \dots, \varphi_{N-1}$ , а также  $\psi_0, \psi_1$  — произвольная упорядоченная система чисел; можно считать, что  $g^{(h)}$  — это совокупность сеточной функции  $\varphi_n$  ( $n = 1, 2, \dots, N - 1$ ) и упорядоченной пары чисел  $\psi_0, \psi_1$ . Сложение двух элементов пространства  $F_h$  и умножение элементов  $g^{(h)}$  на числа производятся покомпонентно. Ясно, что в рассматриваемом примере  $F_h$  есть  $(N + 1)$ -мерное линейное пространство. Норма в  $F_h$  может быть введена многими способами. Если ввести в  $F_h$  норму равенством

$$\|g^{(h)}\|_{F_h} = \max \{|\psi_0|, |\psi_1|, \max_n |\varphi_n|\},$$

т.е. принять за норму максимум абсолютных величин всех компонент вектора  $g^{(h)}$ , то в силу (9) получим

$$\|\delta f^{(h)}\|_{F_h} \leq c h, \quad (11)$$

где  $c$  — некоторая постоянная, зависящая от  $u(x)$ , но не зависящая от  $h$ . Из этого неравенства следует стремление невязки  $\delta f^{(h)}$  к нулю при  $h \rightarrow 0$ .

В уравнении  $L_h u^{(h)} = f^{(h)}$ , подробно записанном равенствами (5), которое мы рассмотрели в качестве примера, на  $L_h$  можно смотреть как на оператор. Этот оператор каждой сеточной функции  $v^{(h)} = \{v_n\}$  ( $n = 0, 1, \dots, N$ ) из линейного пространства  $U_h$  функций, определенных на сетке  $D_h$ , ставит в соответствие некоторый элемент  $g^{(h)}$  вида (10) из линейного пространства  $F_h$  по формуле

$$L_h v^{(h)} = \begin{cases} \frac{v_{n+1} - 2v_n + v_{n-1}}{h^2} + a(x_n) \frac{v_{n+1} - v_{n-1}}{2h} + b(x_n)v_n, \\ v_0, \\ \frac{v_1 - v_0}{h}. \end{cases}$$

Условимся и в общем случае разностной краевой задачи (2) считать, что правые части тех скалярных уравнений, которые в совокупности записаны символическим равенством

$$L_h u^{(h)} = f^{(h)},$$

являются компонентами вектора  $f^{(h)}$  из некоторого линейного нормированного пространства  $F_h$ . Тогда  $L_h$  можно рассматривать как оператор, ставящий в соответствие каждой сеточной функции  $u^{(h)}$  из  $U_h$  некоторый элемент  $f^{(h)}$  из  $F_h$ . В таком случае осмысленно выражение  $L_h[u]_h$ , возникающее в результате применения оператора  $L_h$  к сеточной функции  $[u]_h$  из  $U_h$  и являющееся элементом пространства  $F_h$ . Невязка  $\delta f^{(h)} = L_h[u]_h - f^{(h)}$  принадлежит пространству  $F_h$  как разность двух элементов этого пространства. Под величиной невязки следует понимать  $\|\delta f^{(h)}\|_{F_h}$ .

### 3. Аппроксимация порядка $h^k$ .

**Определение.** Будем говорить, что разностная схема  $L_h u^{(h)} = f^{(h)}$  аппроксимирует задачу  $Lu = f$  на решении  $u$ , если  $\|\delta f^{(h)}\|_{F_h} \rightarrow 0$  при  $h \rightarrow 0$ . Если сверх того имеет место неравенство

$$\|\delta f^{(h)}\|_{F_h} \leq c h^k,$$

где  $c > 0$ ,  $k > 0$  — некоторые постоянные, то будем говорить, что имеет место *аппроксимация порядка  $h^k$* , или *порядка  $k$  относительно величины  $h$* .

То обстоятельство, что  $u$  является решением задачи (1), дает информацию о функции  $u$ , которую можно использовать для построения системы (2), а также для проверки факта аппроксимации. Поэтому в определении аппроксимации мы и упоминаем задачу (1). Однако подчеркнем, что приведенное определение аппроксимации задачи  $Lu = f$  на решении  $u$  разностной схемой  $L_h u^{(h)} = f^{(h)}$  само по себе не опирается на равенство  $Lu = f$  для функции  $u$ . Можно было бы говорить просто о том, что схема  $L_h u^{(h)} = f^{(h)}$  соответствует с порядком аппроксимации  $h^k$  функции  $u$ , не вникая в природу этой функции. В частности, если функция  $u$  является одновременно решением двух совсем различных задач вида (1), то одна и та же разностная схема  $L_h u^{(h)} = f^{(h)}$  одновременно аппроксимирует или не аппроксимирует каждую из этих задач на их общем решении  $u(x)$ .

### 4. Примеры.

Пример 1. Разностная схема (5) в силу оценки (11) аппрокси- мирует задачу (4) с первым порядком относительно  $h$ . Разностную схему (5) легко усовершенствовать так, чтобы аппроксимация стала порядка  $h^2$ . Из (9) видно, что все компоненты вектора  $\delta f^{(h)}$ , кроме последней, стремятся к нулю, как  $h^2$  (предпоследняя даже в точности равна нулю). Только последняя компонента вектора  $\delta f^{(h)}$ , т. е. невязка от подстановки  $u(x)$  в последнее уравнение  $(u_1 - u_0)/h = 2$  системы (5), стремится к нулю медленнее, а именно, как первая степень  $h$ . Это препятствие легко устранить.

По формуле Тейлора

$$\frac{u(h) - u(0)}{h} = u'(0) + \frac{h}{2!} u''(0) + \frac{h^2}{3!} u'''(\xi) = 2 + \frac{h}{2} u''(0) + \frac{h^2}{6} u'''(\xi),$$

$$0 \leq \xi \leq h.$$

Но из дифференциального уравнения (4) находим

$$u''(0) = -a(0)u'(0) - b(0)u(0) + \cos 0 = -2a(0) - b(0) + 1.$$

Поэтому, заменив последнее равенство (5) равенством

$$\frac{u_1 - u_0}{h} = 2 - \frac{h}{2} [2a(0) + b(0) - 1], \quad (12)$$

получим для  $f^{(h)}$  вместо (7) выражение

$$f^{(h)} = \begin{cases} \cos(nh), & n = 1, 2, \dots, N - 1, \\ 1, \\ 2 - \frac{h}{2} [2a(0) + b(0) - 1]. \end{cases}$$

Тогда вместо (9) будем иметь

$$\delta f^{(h)} = \begin{cases} \frac{h^2}{12} \left[ \frac{u^{(4)}(\xi_3) + u^{(4)}(\xi_4)}{2} + a(nh)(u'''(\xi_1) + u'''(\xi_2)) \right], \\ 0, \\ \frac{h^2}{6} u'''(\xi), \end{cases} \quad n = 1, 2, \dots, N - 1,$$

откуда  $\|\delta f^{(h)}\|_{F_h} \leq c h^2$ , где  $c$  — некоторая постоянная, не зависящая от  $h$ . Порядок аппроксимации станет вторым относительно  $h$ .

Подчеркнем, что для построения разностного граничного условия (12) мы использовали не только граничные условия задачи (4), но и само дифференциальное уравнение. Можно считать, что мы использовали граничное условие

$$u''(x) + a(x)u'(x) + b(x)u(x)|_{x=0} = \cos x|_{x=0},$$

которое является следствием дифференциального уравнения.

Пример 2. Выясним, каков порядок аппроксимации, который имеет разностная схема

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - u_{n-1}}{2h} + Au_n = 1 + x_n^2, & n = 1, 2, \dots, N - 1, \\ u_0 = b, \\ u_1 = b \end{cases} \quad (13)$$

на решении  $u$  задачи

$$\frac{du}{dx} + Au = 1 + x^2, \quad u(0) = b. \quad (14)$$

Здесь  $A$  — некоторое число,

$$f^{(h)} = \begin{cases} 1 + x_n^2, & n = 1, 2, \dots, N - 1, \\ b, \\ b. \end{cases}$$

Далее,

$$L_h[u]_h = \begin{cases} \frac{u(x_{n+1}) - u(x_{n-1})}{2h} + Au(x_n), & n = 1, 2, \dots, N - 1, \\ u(0), \\ u(h), \end{cases}$$

или

$$L_h[u]_h = \begin{cases} \frac{du(x_n)}{dx} + Au(x_n) + \frac{h^2}{12} [u'''(\xi_n^{(1)}) + u'''(\xi_n^{(2)})], \\ u(0), \quad n = 1, 2, \dots, N - 1, \quad \xi_n^{(i)} \in [x_{n-1}, x_{n+1}] \\ u(0) + h \frac{du(\xi_0)}{dx}, \quad 0 \leq \xi_0 \leq h. \end{cases}$$

Поскольку для решения  $u(x)$  выполняется равенство

$$\frac{du(x_n)}{dx} + Au(x_n) = 1 + x_n^2,$$

то невязка  $\delta f^{(h)}$  имеет вид

$$\delta f^{(h)} = \begin{cases} \frac{h^2}{12} [u'''(\xi_n^{(1)}) + u'''(\xi_n^{(2)})], \quad n = 1, 2, \dots, N - 1, \\ 0, \quad \xi_n^{(i)} \in [x_{n-1}, x_{n+1}], \\ hu'(\xi_0), \quad 0 \leq \xi_0 \leq h. \end{cases}$$

Аппроксимация задачи (14) схемой (13) имеет первый относительно  $h$  порядок. Бросается в глаза, что компоненты невязки, как и в примере 1, имеют различные порядки относительно  $h$ . Разностное уравнение

$$\frac{u_{n+1} - u_{n-1}}{2h} + Au_n = 1 + x_n^2, \quad n = 1, 2, \dots, N - 1, \quad (15)$$

при подстановке  $[u]_h$  удовлетворяется с невязкой  $\frac{h^2}{6} u''(\xi_n)$  порядка  $h^2$ .

Первое граничное условие  $u_0 = b$  при подстановке  $[u]_h$  выполняется точно, а второе условие  $u_1 = b$  — с невязкой  $hu'(\xi_0)$  порядка первой степени  $h$ .

**5. Замена производных разностными отношениями.** В рассмотренных примерах для получения разностных схем мы заменили производные в дифференциальном уравнении разностными отношениями. Этот прием универсален и позволяет построить для любой дифференциальной краевой задачи, имеющей достаточно гладкое решение  $u(x)$ , разностную схему с любым наперед заданным порядком аппроксимации.

**6. Другие способы построения разностных схем.** Замена производных разностными отношениями — не единственный, а часто и не лучший способ построения разностных схем. Некоторым другим способам, приводящим к наиболее употребительным разностным схемам, будет посвящен § 4. Здесь ограничимся примером.

Простейшая разностная схема

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - u_n}{h} - G(x_n, u_n) = 0, \quad n = 0, 1, \dots, N - 1, \\ u_0 = a, \quad h = 1/N, \end{cases}$$

называемая *схемой Эйлера*, аппроксимирует задачу

$$\frac{du}{dx} - G(x, u) = 0, \quad 0 \leq x \leq 1, \quad u(0) = a \quad (16)$$

(где  $G(x, u)$  — заданная функция) с первым порядком относительно  $h$ . При известном  $u_n$  значение  $u_{n+1}$  вычисляется по формуле  $u_{n+1} = u_n + hG(x_n, u_n)$ .

Схема

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - u_n}{h} - \frac{1}{2} [G(x_n, u_n) + G(x_{n+1}, \tilde{u})] = 0, \\ u_0 = a, \end{cases}$$

где  $\tilde{u} = u_n + hG(x_n, u_n)$ , называется *схемой Эйлера с пересчетом*. Она же является одной из схем Рунге–Кутты второго порядка аппроксимации, о которых будет подробно рассказано в § 4. Если  $u_n$  уже вычислено, то по схеме Эйлера вычисляем значение

$$\tilde{u} = u_n + hG(x_n, u_n),$$

а потом осуществляем уточнение найденного  $\tilde{u}$ , полагая

$$u_{n+1} = u_n + \frac{h}{2} [G(x_n, u_n) + G(x_{n+1}, \tilde{u})].$$

На практике вычисление решения задачи Коши, т. е. задачи вида (16), для обыкновенных дифференциальных уравнений без особенностей обычно производится по одной–двум довольно универсальным, хорошо проверенным схемам, для которых на современных компьютерах имеются стандартные программы.

Если приходится с очень большой точностью решать задачи специального вида, то применяются многочисленные специальные схемы, приспособленные именно для этих задач, но уступающие универсальным схемам при решении задач другого вида.

### Задачи

- Проверить, что схема Эйлера с пересчетом аппроксимирует задачу (16) на гладком решении  $u(x)$  со вторым относительно  $h$  порядком.
- Усовершенствовать второе начальное условие  $u_1 = b$  в схеме (13) так, чтобы схема стала аппроксимировать с порядком  $h^2$ .

### § 3. Определение устойчивости разностной схемы.

#### Сходимость как следствие аппроксимации и устойчивости

Выше мы построили ряд разностных схем, аппроксимирующих некоторые краевые задачи для обыкновенных дифференциальных уравнений. Можно показать, что эти схемы являются сходящимися, причем порядок сходимости совпадает с порядком аппроксимации.

Однако можно сконструировать примеры аппроксимирующих, но не сходящихся разностных схем. Легко проверить, что разностная схема

$$L_h u^{(h)} = \begin{cases} 4 \frac{u_{n+1} - u_{n-1}}{2h} - 3 \frac{u_{n+1} - u_n}{h} + Au_n = 0, & n = 1, 2, \dots, N-1, \\ u_0 = b, \\ u_1 = be^{-Ah} \end{cases}$$

аппроксимирует задачу Коши

$$\begin{aligned} \frac{du}{dx} + Au &= 0, \quad 0 \leq x \leq 1, \\ u(0) &= b, \quad A = \text{const}, \end{aligned}$$

на решении  $u$  с первым порядком относительно  $h$ . Однако решение  $u^{(h)}$ , доставляемое этой разностной схемой, не стремится к  $[u]_h$  и даже не остается ограниченным при  $h \rightarrow 0$ .

Действительно, общее решение этого разностного уравнения имеет вид

$$u_n = c_1 q_1^n + c_2 q_2^n,$$

где  $q_1, q_2$  — корни квадратного уравнения  $q^2 - (3 + Ah)q + 2 = 0$ . Выбирая произвольные постоянные так, чтобы получить частное решение, удовлетворяющее заданным начальным условиям, получаем формулу

$$u_n = u_0 \left( \frac{q_2}{q_2 - q_1} q_1^n - \frac{q_1}{q_2 - q_1} q_2^n \right) + u_1 \left( -\frac{1}{q_2 - q_1} q_1^n + \frac{1}{q_2 - q_1} q_2^n \right).$$

Анализ этой формулы показывает, что  $\max_{0 \leq nh \leq 1} |u_n| \rightarrow \infty$  при  $h \rightarrow 0$ .

Таким образом, аппроксимации, вообще говоря, недостаточно для сходимости. Нужна еще устойчивость.

**1. Определение устойчивости.** Пусть для приближенного вычисления решения  $u$  дифференциальной краевой задачи

$$Lu = f \tag{1}$$

составлена разностная схема

$$L_h u^{(h)} = f^{(h)}, \tag{2}$$

которая аппроксимирует задачу (1) на решении  $u$  с некоторым порядком  $h^k$ . Это значит, что невязка  $\delta f^{(h)}$ :

$$L_h [u]_h = f^{(h)} + \delta f^{(h)}, \tag{3}$$

возникающая при подстановке таблицы  $[u]_h$  решения  $u$  в уравнение (2), удовлетворяет оценке вида

$$\|\delta f^{(h)}\|_{F_h} \leq c_1 h^k, \tag{4}$$

где  $c_1$  — некоторая постоянная, не зависящая от  $h$ .

**Определение 1.** Будем называть разностную схему (2) *устойчивой*, если существуют числа  $h_0 > 0$ ,  $\delta > 0$ , такие, что при любом  $h < h_0$  и любом  $\varepsilon^{(h)} \in F_h$ ,  $\|\varepsilon^{(h)}\|_{F_h} < \delta$ , разностная задача

$$L_h z^{(h)} = f^{(h)} + \varepsilon^{(h)}, \quad (5)$$

полученная из задачи (2) добавлением к правой части возмущения  $\varepsilon^{(h)}$ , имеет одно и только одно решение  $z^{(h)}$ , причем отклонение  $z^{(h)} - u^{(h)}$  этого решения от решения  $u^{(h)}$  невозмущенной задачи (2) есть сеточная функция, удовлетворяющая оценке

$$\|z^{(h)} - u^{(h)}\|_{U_h} \leq c \|\varepsilon^{(h)}\|_{F_h}, \quad (6)$$

где  $c$  — некоторая постоянная, не зависящая от  $h$ .

В частности, неравенство (6) означает, что малое возмущение  $\varepsilon^{(h)}$  правой части разностной схемы (2) вызывает равномерно относительно  $h$  малое отклонение  $z^{(h)} - u^{(h)}$  решения.

Пусть оператор  $L_h$ , отображающий  $U_h$  в  $F_h$ , линейный. Тогда приведенное выше определение устойчивости равносильно следующему определению.

**Определение 2.** Будем называть разностную схему (2) с линейным оператором  $L_h$  *устойчивой*, если при любом  $f^{(h)} \in F_h$  уравнение  $L_h u^{(h)} = f^{(h)}$  имеет единственное решение  $u^{(h)} \in U_h$ , причем

$$\|u^{(h)}\|_{U_h} \leq c \|f^{(h)}\|_{F_h}, \quad (7)$$

где  $c$  — некоторая постоянная, не зависящая от  $h$ .

Докажем равносильность обоих определений устойчивости в случае линейного оператора.

Сначала установим, что из устойчивости разностной схемы (2) в смысле определения 2 следует устойчивость в смысле определения 1. Пусть линейная задача (2) при всех рассматриваемых  $h > 0$  и произвольном  $f^{(h)} \in F_h$  имеет единственное решение, причем выполняется оценка (7). Вычитая из равенства (5) равенство (2), получаем

$$L_h(z^{(h)} - u^{(h)}) = \varepsilon^{(h)},$$

откуда в силу (7) следует оценка (6) при произвольном  $\varepsilon^{(h)} \in F_h$ , а значит, и устойчивость в смысле определения 1.

Покажем теперь, что устойчивость в смысле определения 1 влечет устойчивость в смысле определения 2. В силу определения 1 при некоторых  $h_0 > 0$ ,  $\delta > 0$  и при произвольных  $h < h_0$ ,  $\varepsilon^{(h)} \in F_h$ ,  $\|\varepsilon^{(h)}\|_{F_h} \leq \delta$ , существуют и единственны решения уравнений

$$L_h z^{(h)} = f^{(h)} + \varepsilon^{(h)}, \quad L_h u^{(h)} = f^{(h)}.$$

Положим  $w^{(h)} = z^{(h)} - u^{(h)}$  и вычтем эти равенства почленно. Получим, что существует и единственное решение уравнения

$$L_h w^{(h)} = \varepsilon^{(h)},$$

причем в силу (6)

$$\|w^{(h)}\|_{U_h} \leq c \|\varepsilon^{(h)}\|_{F_h}.$$

Очевидно, что, изменив обозначения решения и правой части уравнения  $L_h w^{(h)} = \varepsilon^{(h)}$ , последний результат можно сформулировать так: при произвольных  $h < h_0$ ,  $f^{(h)} \in F_h$ ,  $\|f^{(h)}\|_{F_h} < \delta$ , задача (2) имеет единственное решение  $u^{(h)}$ , причем это решение удовлетворяет оценке (7). Однако в таком случае задача (2) имеет единственное решение  $u^{(h)}$  и выполняется оценка (7) не только для всех  $f^{(h)}$ , удовлетворяющих оценке  $\|f^{(h)}\|_{F_h} < \delta$ , но и вообще для всех  $f^{(h)} \in F_h$ , т. е. имеет место устойчивость в смысле определения 2.

В самом деле, пусть  $\|f^{(h)}\|_{F_h} < \delta$ . Докажем однозначную разрешимость и оценку (7) в этом случае. Положим

$$u^{(h)} = \frac{2\|f^{(h)}\|_{F_h}}{\delta} \tilde{u}^{(h)}, \quad f^{(h)} = \frac{2\|f^{(h)}\|_{F_h}}{\delta} \tilde{f}^{(h)}. \quad (8)$$

Для  $\tilde{u}^{(h)}$  получим уравнение

$$L_h \tilde{u}^{(h)} = \tilde{f}^{(h)},$$

причем

$$\|\tilde{f}^{(h)}\|_{F_h} = \frac{\delta}{2\|f^{(h)}\|_{F_h}} \|f^{(h)}\|_{F_h} = \frac{\delta}{2} < \delta.$$

Поэтому уравнение  $L_h \tilde{u}^{(h)} = \tilde{f}^{(h)}$  однозначно разрешимо, причем

$$\|\tilde{u}^{(h)}\|_{U_h} \leq c \|\tilde{f}^{(h)}\|_{F_h}.$$

В силу формул (8) отсюда следуют однозначная разрешимость задачи (2) и справедливость оценки (7) при произвольном рассматриваемом  $f^{(h)} \in F_h$ .

**2. Зависимость между аппроксимацией, устойчивостью и сходимостью.** Докажем теперь, что из аппроксимации и устойчивости следует сходимость.

Теорема 1. Пусть разностная схема  $L_h u^{(h)} = f^{(h)}$  аппроксимирует задачу  $Lu = f$  на решении  $u$  с порядком  $h^k$  и устойчива. Тогда решение  $u^{(h)}$  разностной задачи  $L_h u^{(h)} = f^{(h)}$  сходится к  $[u]_h$ , причем имеет место оценка

$$\|[u]_h - u^{(h)}\|_{U_h} \leq c c_1 h^k,$$

где  $c, c_1$  — числа, входящие в оценки (4), (6).

Доказательство. Положим  $\varepsilon^{(h)} = \delta f^{(h)}$ ,  $[u]_h = z^{(h)}$ . Тогда оценка (6) примет вид

$$\|[u]_h - u^{(h)}\|_{U_h} \leq c \|\delta f^{(h)}\|_{F_h}.$$

Учитывая (4), сразу получаем доказываемое неравенство.  $\square$

Пример. Докажем устойчивость разностной схемы Эйлера

$$\begin{aligned} \frac{u_{n+1} - u_n}{h} - G(x_n, u_n) &= \varphi(x_n), \quad n = 0, 1, \dots, N-1, \\ u_0 &= \psi, \end{aligned} \quad (9)$$

$x_n = nh$ ,  $h = 1/N$ , для численного решения дифференциальной краевой задачи

$$\begin{aligned} \frac{du}{dx} - G(x, u) &= \varphi(x), \quad 0 \leq x \leq 1, \\ u|_{x=0} &= \psi. \end{aligned} \quad (10)$$

Будем предполагать функцию  $G(x, u)$  двух аргументов и функцию  $\varphi(x)$  такими, что существует решение  $u(x)$ , имеющее ограниченную вторую производную. Кроме того, будем считать, что  $G(x, u)$  имеет ограниченную производную по  $u$ :

$$\left| \frac{\partial G}{\partial u} \right| < M. \quad (11)$$

Положим  $\|f^{(h)}\|_{F_h} = \max\{|\psi|, \max_n |\varphi(x_n)|\}$ .

Читателю рекомендуется проверить, что разностная схема (9) аппроксимирует задачу (10) на решении  $u(x)$  с первым относительно  $h$  порядком. (Разностное уравнение соответствует задаче с первым порядком, а граничное условие  $u_0 = \psi$  соответствует точно.)

Определим норму

$$\|u^{(h)}\|_{U_h} = \max_n |u_n|$$

и проверим устойчивость разностной схемы (9). Запишем ее в форме (2), положив

$$\begin{aligned} L_h u^{(h)} &= \begin{cases} \frac{u_{n+1} - u_n}{h} - G(x_n, u_n), & n = 0, 1, \dots, N-1, \\ u_0, & \end{cases} \\ f^{(h)} &= \begin{cases} \varphi(x_n), & n = 0, 1, \dots, N-1, \\ \psi. & \end{cases} \end{aligned}$$

Задача

$$L_h z^{(h)} = f^{(h)} + \varepsilon^{(h)}$$

в подробной записи имеет вид

$$\begin{aligned} \frac{z_{n+1} - z_n}{h} - G(x_n, z_n) &= \varphi(x_n) + \varepsilon_n, \quad n = 0, 1, \dots, N-1, \\ z_0 &= \psi + \varepsilon, \end{aligned} \quad (12)$$

где

$$\varepsilon^{(h)} = \begin{cases} \varepsilon_n, & n = 0, 1, \dots, N-1, \\ \varepsilon. & \end{cases}$$

Вычтем из уравнений (12) соответствующие уравнения (9) почленно. Обозначим  $z_n - u_n = w_n$  и учтем, что

$$G(x_n, z_n) - G(x_n, u_n) = \frac{\partial G(x_n, \xi_n)}{\partial u} w_n \equiv M_n^{(h)} w_n,$$

где  $\xi_n$  — некоторое число, заключенное между числами  $z_n$  и  $u_n$ . Получим следующую систему уравнений для определения  $w^{(h)} = \{w_0, w_1, \dots, w_N\}$ :

$$\frac{w_{n+1} - w_n}{h} - M_n^{(h)} w_n = \varepsilon_n, \quad n = 0, 1, \dots, N - 1, \quad (13)$$

$$w_0 = \varepsilon.$$

Учитывая, что  $|M_n^{(h)}| < M$  в силу (11) и что  $(n + 1)h \leq 1$ , получаем

$$\begin{aligned} |w_{n+1}| &= |(1 + hM_n^{(h)}) w_n + h\varepsilon_n| \leq (1 + Mh)|w_n| + h|\varepsilon_n| \leq \\ &\leq (1 + Mh)^2 |w_{n-1}| + h(1 + Mh)|\varepsilon_{n-1}| + h|\varepsilon_n| \leq \\ &\leq (1 + Mh)^2 |w_{n-1}| + 2h(1 + Mh)\|\varepsilon^{(h)}\|_{F_h} \leq \\ &\leq (1 + Mh)^3 |w_{n-2}| + 3h(1 + Mh)^2 \|\varepsilon^{(h)}\|_{F_h} \leq \dots \\ \dots &\leq (1 + Mh)^{n+1} |w_0| + (n + 1)h(1 + Mh)^n \|\varepsilon^{(h)}\|_{F_h} \leq \\ &\leq 2(1 + Mh)^N \|\varepsilon^{(h)}\|_{F_h} \leq 2e^M \|\varepsilon^{(h)}\|_{F_h}. \end{aligned}$$

Из доказанного неравенства

$$|w_{n+1}| \leq 2e^M \|\varepsilon^{(h)}\|_{F_h}$$

следует оценка вида (6):

$$\|w^{(h)}\|_{U_h} \leq 2e^M \|\varepsilon^{(h)}\|_{F_h},$$

означающая устойчивость с постоянной  $c = 2e^M$ . В силу теоремы 1 разностная схема (9) является сходящейся с первым относительно  $h$  порядком.

Исследуем теперь сходимость разностной схемы (7) из § 1:

$$\frac{u_{n+1} - 2u_n + u_{n-1}}{h^2} - (1 + x_n^2) u_n = \sqrt{1 + x_n}, \quad n = 1, 2, \dots, N - 1, \quad (14)$$

$$u_0 = 2, \quad u_N = 1$$

для дифференциальной краевой задачи (3), (4) из § 1. Апроксимация со вторым относительно  $h$  порядком задачи (3), (4) из § 1 задачей (14) благодаря формуле

$$\frac{u(x + h) - 2u(x) + u(x - h)}{h^2} = u''(x) + \frac{h^2}{12} u^{(4)}(\xi)$$

здесь очевидна.

Проверим устойчивость. Рассматриваемая задача линейна относительно  $u$ . Поэтому проверка устойчивости состоит в том, чтобы установить существование единственного решения задачи

$$\frac{u_{n+1} - 2u_n + u_{n-1}}{h^2} - (1 + x_n^2)u_n = g_n, \quad n = 1, 2, \dots, N-1, \quad (15)$$

$$u_0 = \alpha, \quad u_N = \beta$$

при любых  $\{g_n\}$ ,  $\alpha$ ,  $\beta$ , и в том, чтобы получить оценку

$$\max |u_n| \leq c \max \{|\alpha|, |\beta|, \max |g_n|\}. \quad (16)$$

Задачу вида (15) мы рассматривали в связи с прогонкой. Для задачи вида

$$a_n u_{n-1} + b_n u_n + c_n u_{n+1} = g_n,$$

$$u_0 = \alpha, \quad u_N = \beta$$

в предположении

$$|b_n| > |a_n| + |c_n| + \delta, \quad \delta > 0,$$

были доказаны ее однозначная разрешимость и оценка

$$|u_n| \leq \max \left\{ |\alpha|, |\beta|, \frac{1}{\delta} \max |g_m| \right\}. \quad (17)$$

В случае задачи (15)

$$a_n = \frac{1}{h^2}, \quad c_n = \frac{1}{h^2}, \quad b_n = \frac{2}{h^2} + 1 + x_n^2 \geq |a_n| + |b_n| + 1.$$

Поэтому оценка (17) влечет при  $\delta = 1$  оценку (16) с постоянной  $c = 1$ . Устойчивость доказана.

В заключение подчеркнем, что схема доказательства сходимости решения задачи  $L_h u^{(h)} = f^{(h)}$  к решению задачи  $Lu = f$  путем проверки аппроксимации и устойчивости носит общий характер. Под  $Lu = f$  можно понимать любое функциональное уравнение, а не только краевую задачу для обыкновенного дифференциального уравнения. Само по себе неважно, решением какой задачи является функция  $u$ . Уравнение  $Lu = f$  используется только для конструирования разностного уравнения  $L_h u^{(h)} = f^{(h)}$ . Поясним эту мысль.

**3. Сходящаяся разностная схема для интегрального уравнения.** Построим и исследуем разностную схему для вычисления решения интегрального уравнения

$$Lu \equiv u(x) - \int_0^1 K(x, y)u(y) dy = f(x), \quad 0 \leq x \leq 1.$$

Будем предполагать, что  $|K(x, y)| < \rho < 1$ .

Зададим  $N$ , положим  $h = 1/N$  и будем искать таблицу  $[u]_h$  значений решения на сетке  $x_n = nh$  ( $n = 0, 1, \dots, N$ ). Для получения разностной схемы в равенстве

$$u(x_n) - \int_0^1 K(x_n, y)u(y) dy = f(x_n), \quad n = 0, 1, \dots, N,$$

приближенно заменим интеграл суммой, пользуясь квадратурной формулой трапеций.

Напомним эту формулу. Для произвольной, дважды дифференцируемой на отрезке  $0 \leq y \leq 1$  функции  $\varphi(y)$  справедливо приближенное равенство

$$\int_0^1 \varphi(y) dy \approx h \left( \frac{\varphi_0}{2} + \varphi_1 + \dots + \varphi_{N-1} + \frac{\varphi_N}{2} \right), \quad h = \frac{1}{N},$$

причем погрешность есть величина  $O(h^2)$ . После указанной замены интеграла получим

$$u_n - h \left( \frac{K(x_n, 0)}{2} u_0 + K(x_n, h) u_1 + \dots + K(x_n, (N-1)h) u_{N-1} + \frac{K(x_n, 1)}{2} u_N \right) = f_n, \quad n = 0, 1, \dots, N. \quad (18)$$

Выписанная система равенств записывается в форме  $L_h u^{(h)} = f^{(h)}$ , если положить

$$L_h u^{(h)} = \begin{cases} g_0, \\ g_1, \\ \dots \\ g_N, \end{cases} \quad f^{(h)} = \begin{cases} f(0), \\ f(h), \\ \dots \\ f(Nh), \end{cases}$$

где

$$g_n = u_n - h \left[ \frac{K(x_n, 0)}{2} u_0 + K(x_n, h) u_1 + \dots + \frac{K(x_n, 1)}{2} u_N \right], \quad n = 0, 1, \dots, N.$$

Построенная разностная схема  $L_h u^{(h)} = f^{(h)}$  аппроксимирует задачу  $Lu = f$  на решении  $u$  со вторым порядком относительно шага  $h$ , поскольку квадратурная формула трапеций имеет второй порядок точности.

Проверим устойчивость. Пусть  $u^{(h)} = (u_0, u_1, \dots, u_N)$  — какое-нибудь решение системы (18), и пусть  $u_s$  — одна из тех компонент решения, которые по модулю не меньше каждой из остальных:

$$|u_s| \geq |u_n|, \quad n = 0, 1, \dots, N.$$

Из уравнения с номером  $n = s$  системы (18) следует неравенство

$$|f(x_s)| = \left| u_s - h \left( \frac{K(x_s, 0)}{2} u_0 + K(x_s, h) u_1 + \dots + \frac{K(x_s, 1)}{2} u_N \right) \right| \geq |u_s| - h \left( \frac{\rho}{2} + \rho + \dots + \rho + \frac{\rho}{2} \right) |u_s| = (1 - Nh\rho) |u_s| = (1 - \rho) |u_s|.$$

Поэтому

$$\|u^{(h)}\|_{U_h} = \max_n |u_n| = |u_s| \leq \frac{1}{1-\rho} |f(x_s)| \leq \frac{1}{1-\rho} \|f^{(h)}\|_{F_h}. \quad (19)$$

В частности, при  $f(x_n) \equiv 0$  отсюда следует, что система (17) не имеет нетривиальных решений, а следовательно, однозначно разрешима при любой правой части. Неравенство (19) означает устойчивость (6) с постоянной  $c = 1/(1 - \rho)$ . Решение  $u^{(h)}$  задачи  $L_h u^{(h)} = f^{(h)}$  в силу теоремы о сходимости удовлетворяет неравенству

$$\|[u]_h - u^{(h)}\|_{U_h} = \max_n |u(nh) - u_n| \leq ch^2,$$

где  $c$  — некоторая постоянная.

#### § 4. Схемы Рунге–Кутты

Изложим здесь некоторые употребительные разностные схемы решения задачи Коши для дифференциального уравнения первого порядка

$$\frac{du}{dx} - G(x, u) = 0, \quad 0 \leq x \leq 1, \quad u(0) = a. \quad (1)$$

В конце параграфа эти схемы будут перенесены на системы дифференциальных уравнений первого порядка, к которым сводится общий случай уравнений и систем любого порядка.

Выберем на отрезке  $0 \leq x \leq 1$  сетку точек

$$0 = x_0 < x_1 < \dots < x_{N-1} < x_N = 1, \quad x_n = nh, \quad h = \frac{1}{N},$$

и будем составлять разностные схемы для приближенного отыскания таблицы  $[u]_h$  значений решения  $u(x)$  на выбранной сетке.

С простейшей употребительной схемой мы уже встречались. Это — схема Эйлера

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - u_n}{h} - G(x_n, u_n) = 0, & n = 0, 1, \dots, N-1, \\ u_0 = a, \end{cases} \quad (2)$$

обладающая первым порядком аппроксимации (и точности).

Вычисления по этой схеме имеют простой геометрический смысл. Если  $u_n$  уже вычислено, то вычисление

$$u_{n+1} = u_n + hG(x_n, u_n)$$

равносильно сдвигу из точки  $(x_n, u_n)$  в точку  $(x_{n+1}, u_{n+1})$  на плоскости  $Oxy$  по касательной к интегральной кривой  $u = u(x)$  дифференциального уравнения  $du/dx - G(x, u) = 0$ , проходящей через точку  $(x_n, u_n)$ .

Среди схем более высокого порядка аппроксимации наиболее употребительны различные варианты схем Рунге–Кутты.

**1. Схемы Рунге–Кутты.** Пусть значение  $u_n$  приближенного решения в точке  $x_n$  уже найдено и требуется вычислить  $u_{n+1}$  в точке  $x_{n+1} = x_n + h$ . Задаем целое  $l$  и выписываем выражения:

$$\begin{aligned} k_1 &= G(x_n, u_n), \\ k_2 &= G(x_n + \alpha h, u_n + \alpha h k_1), \\ k_3 &= G(x_n + \beta h, u_n + \beta h k_2), \\ &\dots \\ k_l &= G(x_n + \gamma h, u_n + \gamma h k_{l-1}). \end{aligned}$$

Затем полагаем

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - u_n}{h} - (p_1 k_1 + p_2 k_2 + \dots + p_l k_l) = 0, \\ u_0 = a. \end{cases} \quad n = 0, 1, \dots, N - 1,$$

Коэффициенты  $\alpha, \beta, \dots, \gamma, p_1, p_2, \dots, p_l$  подбираем так, чтобы получить при заданном  $l$  аппроксимацию возможно более высокого порядка. Зная  $u_n$ , можно вычислить  $k_1, k_2, \dots, k_l$ , а затем  $u_{n+1} = u_n + h(p_1 k_1 + p_2 k_2 + \dots + p_l k_l)$ .

Простейшей схемой Рунге–Кутты является схема Эйлера ( $l = 1$ ).

Схема Рунге–Кутты

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - u_n}{h} - \frac{1}{6} (k_1 + 2k_2 + 2k_3 + k_4) = 0, \\ u_0 = a, \end{cases} \quad n = 0, 1, \dots, N - 1, \quad (3)$$

где

$$\begin{aligned} k_1 &= G(x_n, u_n), & k_2 &= G\left(x_n + \frac{h}{2}, u_n + \frac{k_1 h}{2}\right), \\ k_3 &= G\left(x_n + \frac{h}{2}, u_n + \frac{k_2 h}{2}\right), & k_4 &= G(x_n + h, u_n + k_3 h), \end{aligned}$$

имеет четвертый порядок аппроксимации.

Схема Рунге–Кутты

$$L_h u^{(h)} = \begin{cases} \frac{u_{n+1} - u_n}{h} - \left[ \frac{2\alpha - 1}{2\alpha} k_1 + \frac{1}{2\alpha} k_2 \right] = 0, \\ u_0 = a, \end{cases} \quad n = 0, 1, \dots, N - 1, \quad (4)$$

где  $k_1 = G(x_n, u_n)$ ,  $k_2 = G(x_n + \alpha h, u_n + \alpha h k_1)$ , при любом фиксированном  $\alpha \neq 0$  имеет второй порядок аппроксимации.

Мы докажем только утверждение о схеме (4). Доказательство утверждения о схеме (3) аналогично, но более громоздко.

Решение  $u(x)$  уравнения  $u' = G(x, u)$  удовлетворяет тождествам

$$\begin{aligned}\frac{du}{dx} &\equiv G(x, u(x)), \\ \frac{d^2u}{dx^2} &\equiv \frac{d}{dx} G(x, u) = \frac{\partial G}{\partial x} + \frac{\partial G}{\partial u} G.\end{aligned}$$

Поэтому из формулы Тейлора

$$\frac{u(x_n + h) - u(x_n)}{h} = u'(x_n) + \frac{h}{2} u''(x_n) + O(h^2)$$

для решения  $u(x)$  следует равенство

$$\frac{u(x_{n+1}) - u(x_n)}{h} - \left[ G + \frac{h}{2} \left( \frac{\partial G}{\partial x} + \frac{\partial G}{\partial u} G \right) \right]_{\substack{x=x_n \\ u=u(x_n)}} = O(h^2). \quad (5)$$

С другой стороны, разлагая по  $h$  функцию двух переменных по формуле Тейлора и удерживая члены первой степени, получаем

$$\begin{aligned}\frac{2\alpha - 1}{2\alpha} k_1 + \frac{1}{2\alpha} k_2 &= \\ &= \frac{2\alpha - 1}{2\alpha} G + \frac{1}{2\alpha} \left[ G + \frac{\partial G}{\partial x} \alpha h + \frac{\partial G}{\partial u} \alpha h G + O(h^2) \right]_{\substack{x=x_n \\ u=u(x_n)}} = \\ &= G + \frac{h}{2} \left( \frac{\partial G}{\partial x} + \frac{\partial G}{\partial u} G \right) \Big|_{\substack{x=x_n \\ u=u(x_n)}} + O(h^2).\end{aligned} \quad (6)$$

Поэтому при подстановке в левую часть равенства (4) вместо  $u_n$ ,  $u_{n+1}$  соответственно значений  $u(x_n)$ ,  $u(x_{n+1})$  решения  $u(x)$  получится выражение, совпадающее с левой частью равенства (5) с точностью  $O(h^2)$ . Следовательно, это выражение имеет второй порядок аппроксимации относительно  $h$ . Поскольку значение  $u_0 = a$  задано точно, этим завершается доказательство того, что схема (4) имеет второй порядок аппроксимации.

**2. Обобщение на системы уравнений.** Все описанные схемы численного решения задачи Коши для дифференциального уравнения первого порядка (1) автоматически переносятся на системы уравнений первого порядка. Для этого в записи (1) надо понимать под  $u(x) = \mathbf{u}(x)$  и  $G(x, u) = \mathbf{G}(x, \mathbf{u})$  вектор-функции, а под  $a = \mathbf{a}$  — заданный вектор. Тогда схемы Рунге–Кутты (3), (4) сохранят смысл и останутся применимыми.

Например, система уравнений

$$\begin{aligned}\frac{dv}{dx} - (x + v^2 + \sin w) &= 0, \\ \frac{dw}{dx} + xvw &= 0, \\ v(0) = a_1, \quad w(0) &= a_2\end{aligned}$$

запишется в форме

$$\frac{d\mathbf{u}}{dx} - \mathbf{G}(x, \mathbf{u}) = 0,$$

$$\mathbf{u}(0) = \mathbf{a},$$

если положить

$$\mathbf{u}(x) = \begin{bmatrix} v(x) \\ w(x) \end{bmatrix}, \quad \mathbf{G}(x, \mathbf{u}) = \begin{bmatrix} x + v^2 + \sin w \\ -xvw \end{bmatrix}, \quad \mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}.$$

Формула для  $\mathbf{u}_{n+1}$  в схеме Эйлера:

$$\mathbf{u}_{n+1} = \mathbf{u}_n + h\mathbf{G}(x_n, \mathbf{u}_n),$$

подробно запишется так:

$$v_{n+1} = v_n + h(x_n + v_n^2 + \sin w_n),$$

$$w_{n+1} = w_n + h(-x_n v_n w_n).$$

Все рассуждения о порядке аппроксимации в случае одного уравнения сохраняются и в случае системы уравнений. При этом в формуле (6) под производной вектора  $\mathbf{G} = (G_1, G_2, \dots, G_k)$  по вектору  $\mathbf{u} = (u_1, u_2, \dots, u_k)$  надо понимать матрицу

$$\frac{\partial \mathbf{G}}{\partial \mathbf{u}} = \begin{bmatrix} \frac{\partial G_1}{\partial u_1} & \dots & \frac{\partial G_1}{\partial u_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial G_k}{\partial u_1} & \dots & \frac{\partial G_k}{\partial u_k} \end{bmatrix}.$$

Произвольная система дифференциальных уравнений, разрешенных относительно старших производных, сводится к системе уравнений первого порядка  $d\mathbf{u}/dx = \mathbf{G}(x, \mathbf{u})$  путем замены искомых функций. Как это делается, ясно из следующего примера.

Задача

$$\frac{d^2v}{dx^2} + \sin(xv' + v^2 + w) = 0,$$

$$\frac{dw}{dx} + \sqrt{x^2 + v^2 + (v')^2 + w^2} = 0,$$

$$v(0) = a, \quad v'(0) = b, \quad w(0) = c$$

приводится к требуемому виду, если положить

$$u_1 = v(x), \quad u_2 = \frac{dv}{dx}, \quad u_3 = w(x).$$

Получим

$$\frac{du_1}{dx} - u_2 = 0,$$

$$\frac{du_2}{dx} + \sin(xu_2 + u_1^2 + u_3) = 0,$$

$$\frac{du_3}{dx} + \sqrt{x^2 + u_1^2 + u_2^2 + u_3^2} = 0,$$

$$u_1(0) = a, \quad u_2(0) = b, \quad u_3(0) = c.$$

**Замечание.** Разработаны разностные схемы типа схем Рунге–Кутты, применимые непосредственно для уравнения второго порядка и не требующие предварительного сведения этих уравнений к системам первого порядка.

### Задача\*

Дифференциальное уравнение

$$\frac{dy}{dx} = \frac{16x - 5y + xy}{x + y - xy}$$

имеет особую точку (седло) при  $x = y = 0$ .

Численно построить сепаратрисы, объединяя аналитические средства для выхода из особой точки с расчетами с помощью разностных схем.

## § 5. Методы решения краевых задач

Примером краевой задачи является задача

$$\begin{aligned} y'' &= f(x, y, y'), \quad 0 \leq x \leq 1, \\ y(0) &= Y_0, \quad y(1) = Y_1 \end{aligned} \tag{1}$$

с граничными условиями на концах отрезка  $0 \leq x \leq 1$ , на котором надо найти решение  $y(x)$ . На этом примере схематически изложим некоторые способы численного решения краевых задач.

**1. Метод стрельбы.** В § 4 указаны удобные способы численного решения задачи Коши, т. е. задачи вида

$$\begin{aligned} y'' &= f(x, y, y'), \quad 0 \leq x < 1, \\ y(0) &= Y_0, \quad \left. \frac{dy}{dx} \right|_{x=0} = \operatorname{tg} \alpha, \end{aligned} \tag{2}$$

где  $Y_0$  — ордината точки  $(0, Y_0)$ , из которой выходит интегральная кривая, а  $\alpha$  — угол наклона интегральной кривой к оси  $Ox$  при выходе из точки  $(0, Y_0)$  (рис. 15). При фиксированном  $Y_0$  решение задачи (2) имеет вид  $y = y(x, \alpha)$ . При  $x = 1$  решение  $y(x, \alpha)$  зависит только от  $\alpha$ :

$$y(x, \alpha)|_{x=1} = y(1, \alpha).$$

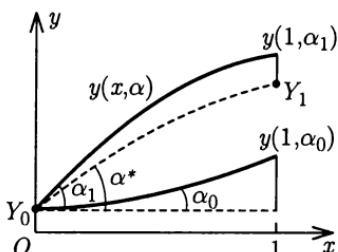


Рис. 15

ки  $(0, Y_0)$  под углом  $\alpha$  к оси абсцисс, попадет в точку  $(1, Y_1)$ :

$$y(1, \alpha) = Y_1. \tag{3}$$

Решение задачи (2) при этом  $\alpha = \alpha^*$  совпадает с искомым решением задачи (1). Дело сводится, таким образом, к решению уравнения (3). Уравнение (3) есть уравнение вида  $F(\alpha) = 0$ , где  $F(\alpha) = y(1, \alpha) - Y_1$ . Оно отличается от привычных уравнений лишь тем, что функция  $F(\alpha)$  задана не аналитическим выражением, а с помощью алгоритма решения задачи (2). Сведение решения краевой задачи (1) к решению задачи Коши (2) и составляет сущность метода стрельбы.

Для решения уравнения (3) можно использовать метод деления отрезка пополам, метод хорд, метод касательных (метод Ньютона) и т. п. Например, при использовании метода деления отрезка пополам мы задаем  $\alpha_0, \alpha_1$  так, чтобы разности  $y(1, \alpha_0) - Y_1, y(1, \alpha_1) - Y_1$  имели разные знаки. Затем полагаем

$$\alpha_2 = \frac{\alpha_0 + \alpha_1}{2}.$$

Вычисляем  $y(1, \alpha_2)$ . Вычисляем затем  $\alpha_3$  по одной из формул

$$\alpha_3 = \frac{\alpha_1 + \alpha_2}{2}, \quad \alpha_3 = \frac{\alpha_0 + \alpha_2}{2}$$

в зависимости от того, имеют ли разности  $y(1, \alpha_2) - Y_1, y(1, \alpha_1) - Y_1$  соответственно разные или одинаковые знаки. Затем вычисляем  $y(1, \alpha_3)$ . Процесс продолжается до тех пор, пока не будет достигнута требуемая точность  $|y(1, \alpha_n) - Y_1| < \varepsilon$ .

В случае использования метода хорд задаем  $\alpha_0, \alpha_1$ , а затем вычисляем последующие  $\alpha_k$  по рекуррентной формуле

$$\alpha_{n+1} = \alpha_n - \frac{F(\alpha_n)}{F'(\alpha_n) - F'(\alpha_{n-1})} (\alpha_n - \alpha_{n-1}), \quad n = 1, 2, \dots$$

Метод стрельбы, сводящий решение краевой задачи (1) к вычислению решений задачи Коши (2), хорошо работает в том случае, если решение  $y(x, \alpha)$  «не слишком сильно» зависит от  $\alpha$ . В противном случае он становится вычислительно неустойчивым, даже если решение задачи (1) зависит от входных данных «умеренно».

Поясним взятые в кавычки слова на примере следующей линейной краевой задачи:

$$\begin{aligned} y'' - a^2 y = 0, \quad 0 \leq x \leq 1, \\ y(0) = Y_0, \quad y(1) = Y_1 \end{aligned} \tag{1'}$$

при постоянном  $a^2$ . Выпишем решение этой задачи:

$$y(x) = \frac{e^{-ax} - e^{-a(2-x)}}{1 - e^{-2a}} Y_0 + \frac{e^{-\alpha(1-x)} - e^{-\alpha(1+x)}}{1 - e^{-2a}} Y_1.$$

Коэффициенты при  $Y_0, Y_1$  с ростом  $a$  остаются ограниченными на отрезке  $0 \leq x \leq 1$  функциями; при всех  $a > 0$  они не превосходят единицы. Поэтому небольшие погрешности при задании  $Y_0, Y_1$  ведут к столь же небольшим погрешностям в решении.

Рассмотрим теперь задачу Коши

$$\begin{aligned} y'' - a^2 y = 0, \quad 0 \leq x \leq 1, \\ y(0) = Y_0, \quad y'(0) = \operatorname{tg} \alpha. \end{aligned} \quad (2')$$

Ее решение имеет вид

$$y(x) = \frac{aY_0 + \operatorname{tg} \alpha}{2a} e^{ax} + \frac{aY_0 - \operatorname{tg} \alpha}{2a} e^{-ax}.$$

Если при задании  $\operatorname{tg} \alpha$  допущена погрешность  $\varepsilon$ , то значение решения при  $x = 1$  получит отклонение

$$\Delta y(1) = \frac{\varepsilon}{2a} e^a - \frac{\varepsilon}{2a} e^{-a}. \quad (4)$$

При больших  $a$  вычитаемое в этом равенстве пренебрежимо мало, но коэффициент при  $\varepsilon$  в первом слагаемом становится большим. Поэтому метод стрельбы при решении задачи (1'), будучи формально приемлемой процедурой, при больших  $a$  становится практически непригодным.

**2. Метод прогонки.** Для решения краевой задачи

$$\begin{aligned} y'' - p(x)y = f(x), \quad 0 \leq x \leq 1, \\ y(0) = Y_0, \quad y(1) = Y_1 \end{aligned}$$

при  $p(x) \gg 1$  можно воспользоваться разностной схемой

$$\begin{aligned} \frac{y_{n+1} - 2y_n + y_{n-1}}{h^2} - p(x_n)y_n = f_n, \quad n = 1, 2, \dots, N-1, \\ y(0) = Y_0, \quad y(1) = Y_1 \end{aligned}$$

и решать разностную задачу прогонкой. Условия применимости прогонки при  $p(x) > 0$ , как легко проверит читатель, выполнены.

**3. Метод Ньютона.** Метод стрельбы при решении хорошо поставленной краевой задачи может оказаться, как мы видели, неприменимым из-за вычислительной неустойчивости. А метод прогонки даже формально можно применять только для решения линейных задач.

Метод Ньютона сводит решение нелинейной задачи к серии линейных задач и состоит в следующем. Пусть известна некоторая функция  $y_0(x)$ , удовлетворяющая граничным условиям (1) и грубо приближенно равная искомому решению  $y(x)$ . Положим

$$y(x) = y_0(x) + v(x), \quad (5)$$

где  $v(x)$  — поправка к нулевому приближению  $y_0(x)$ . Подставим (5) в уравнение (1) и линеаризуем задачу, используя равенства

$$y''(x) = y_0''(x) + v''(x),$$

$$f(x, y_0 + v, y'_0 + v'_0) =$$

$$= f(x, y_0, y'_0) + \frac{\partial f(x, y_0, y'_0)}{\partial y} v + \frac{\partial f(x, y_0, y'_0)}{\partial y'} v' + O(v^2 + |v'|^2).$$

Отбрасывая остаточный член  $O(v^2 + |v'|^2)$ , получаем линейную задачу для поправки  $\bar{v}$ :

$$\bar{v}'' = p\bar{v}' + q\bar{v} + \varphi(x), \quad \bar{v}(0) = \bar{v}(1) = 0, \quad (6)$$

где

$$p = p(x) = \frac{\partial f(x, y_0, y'_0)}{\partial y'}, \quad q = q(x) = \frac{\partial f(x, y_0, y'_0)}{\partial y},$$

$$\varphi(x) = f(x, y_0, y'_0) - y''_0.$$

Решая линейную задачу (6) аналитически или каким-либо численным методом, найдем приближенно поправку  $v$  и примем

$$y_1 = y_0(x) + \bar{v}(x)$$

за следующее приближение.

Описанная процедура может применяться к нелинейной разностной краевой задаче, возникшей при аппроксимации задачи (1) (см. задачу 5 в § 2 гл. 7).

### Задача

Численно найти наименьшее, а также третье по величине положительные числа  $\lambda_1, \lambda_3$ , при которых задача

$$y'' + (x - \lambda)y = 0, \quad 0 \leq \lambda \leq 1,$$

$$y(0) = y(1) = 0$$

имеет отличные от тождественного нуля решения  $y_1(x), y_3(x) \neq 0$  соответственно.

## ГЛАВА 9

### РАЗНОСТНЫЕ СХЕМЫ

### ДЛЯ УРАВНЕНИЙ С ЧАСТНЫМИ ПРОИЗВОДНЫМИ

Выше в связи с разностными схемами для обыкновенных дифференциальных уравнений мы определили понятия сходимости, аппроксимации и устойчивости. Мы доказали теорему о том, что если разностная краевая задача аппроксимирует дифференциальную задачу и устойчива, то при измельчении сетки решение разностной задачи сходится к решению дифференциальной задачи. В этой теореме содержится указание на способы построения сходящихся разностных схем для численного решения дифференциальных краевых задач: надо строить аппроксимирующие разностные схемы и выбирать среди них устойчивые.

Определения сходимости, аппроксимации и устойчивости и теорема о связи между этими понятиями носят общий характер. Эти понятия

имеют смысл для любых функциональных уравнений. Мы иллюстрировали их примерами разностных схем для обыкновенных дифференциальных уравнений и для интегрального уравнения. Здесь мы проиллюстрируем некоторые основные способы построения разностных схем и проверки их устойчивости примерами разностных схем для уравнений с частными производными. При этом обнаружится много важных и существенно новых по сравнению со случаем обыкновенных дифференциальных уравнений обстоятельств. Главные из них — разнообразие сеток и способов аппроксимации, неустойчивость большинства взятых наудачу аппроксимирующих схем, сложность исследований устойчивости и трудности вычисления решений разностных краевых задач, требующие специальных усилий для их преодоления.

## § 1. Основные определения и их иллюстрация

**1. Определение сходимости.** Пусть требуется приближенно вычислить решение  $u$  дифференциальной краевой задачи

$$Lu = f, \quad (1)$$

поставленной в некоторой области  $D$  с границей  $\Gamma$ . Для этого следует выбрать дискретное множество точек  $D_h$  (сетку), принадлежащее  $D \cup \Gamma$ , ввести линейное нормированное пространство  $U_h$  функций, определенных на сетке  $D_h$ , установить соответствие между решением  $u$  и функцией  $[u]_h \in U_h$ , которую будем считать искомой таблицей решения  $u$ . Для приближенного отыскания таблицы  $[u]_h$ , которую мы условились считать точным решением задачи (1), надо на основе этой задачи составить такую систему разностных уравнений

$$L_h u^{(h)} = f^{(h)} \quad (2)$$

относительно функции  $u^{(h)}$  из  $U_h$ , чтобы имела место сходимость

$$\|[u]_h - u^{(h)}\|_{U_h} \rightarrow 0, \quad h \rightarrow 0. \quad (3)$$

Если для решения разностной краевой задачи (2) выполняется неравенство

$$\|[u]_h - u^{(h)}\|_{U_h} \leq ch^k, \quad c = \text{const},$$

то говорят, что *сходимость* (в смысле выбранной нами нормы) *имеет порядок  $k$  относительно  $h$* .

Задачу построения сходящейся разностной схемы (2) разбивают на две — построение разностной схемы (2), аппроксимирующей задачу (1) на решении  $u$  и последней, и проверку устойчивости схемы (2).

**2. Определение аппроксимации.** Чтобы понятие аппроксимации имело смысл, надо ввести норму в пространстве  $F_h$ , которому принадлежит правая часть  $f^{(h)}$  уравнения (2). По определению разностная задача (2) *аппроксимирует* задачу (1) на решении  $u$ , если в равенстве

$$L_h [u]_h = f^{(h)} + \delta f^{(h)}$$

невязка  $\delta f^{(h)}$ , возникающая при подстановке  $[u]_h$  в разностную краевую задачу (2), стремится к нулю при  $h \rightarrow 0$ :

$$\|\delta f^{(h)}\|_{F_h} \rightarrow 0.$$

Если

$$\|\delta f^{(h)}\|_{F_h} \leq ch^k,$$

где  $c$  не зависит от  $h$ , то аппроксимация имеет порядок  $k$  относительно  $h$ .

Построим, например, для задачи Коши

$$\begin{aligned} \frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} &= \varphi(x, t), \quad -\infty < x < \infty, \quad 0 \leq t \leq T, \\ u(x, 0) &= \psi(x), \quad -\infty < x < \infty, \end{aligned} \tag{4}$$

одну из аппроксимирующих ее разностных схем. Задача (4) записывается в форме (1), если положить

$$\begin{aligned} Lu &\equiv \begin{cases} \frac{\partial u}{\partial t} - \frac{\partial u}{\partial x}, & -\infty < x < \infty, \quad 0 \leq t \leq T, \\ u(x, 0), & -\infty < x < \infty, \end{cases} \\ f &= \begin{cases} \varphi(x, t), & -\infty < x < \infty, \quad 0 \leq t \leq T, \\ \psi(x), & -\infty < x < \infty. \end{cases} \end{aligned}$$

В качестве сетки  $D_h$  используем совокупность точек пересечения прямых

$$x = mh, \quad t = p\tau, \quad m = 0, \pm 1, \dots, \quad p = 0, 1, \dots, [T/\tau],$$

где  $h > 0$ ,  $\tau > 0$  — некоторые числа, а  $[T/\tau]$  — целая часть дроби  $T/\tau$ . Будем считать, что шаг  $\tau$  связан с шагом  $h$  зависимостью  $\tau = rh$ , где  $r = \text{const}$ , так что сетка  $D_h$  зависит только от одного параметра  $h$ . Искомой сеточной функцией является таблица  $[u]_h = \{u(mh, p\tau)\}$  значений решения  $u(x, t)$  задачи (4) в точках сетки  $D_h$ .

Перейдем к построению аппроксимирующей задачи (4) разностной схемы (2). Значения сеточной функции  $u^{(h)}$  в точке  $(x_m, t_p) = (mh, p\tau)$  сетки  $D_h$  будем обозначать  $u_m^p$ . Схему (2) получим, приблизив производные  $\partial u / \partial t$ ,  $\partial u / \partial x$  разностными отношениями

$$\begin{aligned} \left. \frac{\partial u}{\partial t} \right|_{x,t} &\approx \frac{u(x, t + \tau) - u(x, t)}{\tau}, \\ \left. \frac{\partial u}{\partial x} \right|_{x,t} &\approx \frac{u(x + h, t) - u(x, t)}{h}. \end{aligned}$$

Эта схема имеет вид

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} &= \varphi_m^p, \\ u_m^0 &= \psi_m, \quad \varphi_m^p \equiv \varphi(mh, p\tau), \quad \psi_m = \psi(mh), \\ p &= 0, 1, \dots, [T/\tau] - 1, \quad m = 0, \pm 1, \dots. \end{aligned} \tag{5}$$

Оператор  $L_h$  и правая часть  $f^{(h)}$  для схемы (5) задаются соответственно равенствами

$$L_h u^{(h)} = \begin{cases} \frac{u_m^{p+1} - u_m^p}{h} - \frac{u_{m+1}^p - u_m^p}{h}, & p = 0, 1, \dots, [\frac{T}{\tau}] - 1, \\ u_m^0, & m = 0, \pm 1, \dots; \end{cases}$$

$$f^{(h)} = \begin{cases} \varphi_m^p, & m = 0, \pm 1, \dots, p = 0, 1, \dots, [\frac{T}{\tau}] - 1, \\ \psi_m, & m = 0, \pm 1, \dots. \end{cases}$$

Таким образом,  $f^{(h)}$  — это пара сеточных функций  $\varphi_m^p$ ,  $\psi_m$ , одна из которых задана на двумерной сетке  $(x_m, t_r) = (mh, r\tau)$ ,  $m = 0, \pm 1, \dots$ ,  $p = 0, 1, \dots, [T/\tau]$ , а другая — на одномерной сетке

$$(x_m, 0) = (mh, 0), \quad m = 0, \pm 1, \dots.$$

Разностное уравнение (5) можно разрешить относительно  $u_m^{p+1}$ , получив

$$u_m^{p+1} = (1 - r)u_m^p + ru_{m+1}^p + \tau\varphi_m^p. \quad (6)$$

Итак, зная значения  $u_m^p$  ( $m = 0, \pm 1, \dots$ ) решения  $u^{(h)}$  в точках сетки при  $t = r\tau$ , можно вычислить значения  $u_m^{p+1}$  в точках сетки при  $t = (p+1)\tau$ . Поскольку значения  $u_m^0$  при  $t = 0$  заданы равенствами  $u_m^0 = \psi_m$ , мы можем шаг за шагом вычислить значения решения  $u_m^p$  в точках сетки на прямых  $t = \tau, 2\tau, \dots$ , т. е. всюду на  $D_h$ .

Перейдем к выяснению порядка аппроксимации, которым обладает схема (5). За  $F_h$  можно принять линейное пространство всех пар ограниченных функций  $g^{(h)} = \begin{bmatrix} \varphi_m^p \\ \psi_m \end{bmatrix}$ , положив

$$\|g^{(h)}\|_{F_h} = \max_{m,p} |\varphi_m^p| + \max_m |\psi_m|^*. \quad (7)$$

Вообще говоря, норма, в которой рассматривается аппроксимация, может быть выбрана многими способами, и выбор этот небезразличен. Будем иметь в виду всюду в этом параграфе именно норму (7).

Предположим, что решение  $u(x, t)$  задачи (4) имеет ограниченные вторые производные. Тогда по формуле Тейлора

$$\frac{u(x_m + h, t_p) - u(x_m, t_p)}{h} = \frac{\partial u(x_m, t_p)}{\partial x} + \frac{h}{2} \cdot \frac{\partial^2 u(x_m + \xi, t_p)}{\partial x^2},$$

$$\frac{u(x_m, t_p + \tau) - u(x_m, t_p)}{\tau} = \frac{\partial u(x_m, t_p)}{\partial t} + \frac{\tau}{2} \cdot \frac{\partial^2 u(x_m, t_p + \eta)}{\partial t^2},$$

где  $\xi, \eta$  — некоторые числа, зависящие от  $m, p, h$  и удовлетворяющие неравенствам  $0 \leq \xi \leq h, 0 \leq \eta \leq \tau$ .

<sup>\*)</sup> Если  $\max |\varphi_m^p|$  или  $\max |\psi_m|$  не достигается, то имеется в виду точная верхняя грань  $\sup |\varphi_m^p|$  или  $\sup |\psi_m|$  соответственно.

С помощью этих формул выражение

$$L_h[u]_h = \begin{cases} \frac{u(x_m, t_p + \tau) - u(x_m, t_p)}{\tau} - \frac{u(x_m + h, t_p) - u(x_m, t_p)}{h}, \\ u(x_m, 0) \end{cases}$$

можно переписать в виде

$$L_h[u]_h = \begin{cases} \left( \frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} \right) \Big|_{x_m, t_p} + \frac{\tau}{2} \cdot \frac{\partial^2 u(x_m, t_p + \eta)}{\partial t^2} - \frac{h}{2} \cdot \frac{\partial^2 u(x_m + \xi, t_p)}{\partial x^2}, \\ u(x_m, 0) + 0, \end{cases}$$

или

$$L_h[u]_h = f^{(h)} + \delta f^{(h)},$$

где

$$\delta f^{(h)} = \begin{cases} \frac{\tau}{2} \cdot \frac{\partial^2 u(x_m, t_p + \eta)}{\partial t^2} - \frac{h}{2} \cdot \frac{\partial^2 u(x_m + \xi, t_p)}{\partial x^2}, \\ 0. \end{cases}$$

Следовательно,

$$\|\delta f^{(h)}\|_{F_h} \leq \frac{1}{2} \left( r \sup \left| \frac{\partial^2 u}{\partial t^2} \right| + \sup \left| \frac{\partial^2 u}{\partial x^2} \right| \right) h.$$

Таким образом, рассматриваемая разностная схема (5) имеет первый порядок аппроксимации относительно  $h$  на решении  $u(x, t)$ , обладающем ограниченными вторыми производными.

**3. Определение устойчивости.** Разностная краевая задача (2) по определению *устойчива*, если существуют числа  $\delta > 0$ ,  $h_0 > 0$ , такие, что при любом  $h < h_0$  и любом  $\varepsilon^{(h)}$  из  $F_h$ , удовлетворяющем неравенству  $\|\varepsilon^{(h)}\|_{F_h} < \delta$ , разностная краевая задача

$$L_h z^{(h)} = f^{(h)} + \varepsilon^{(h)}$$

имеет одно и только одно решение, причем выполняется условие

$$\|z^{(h)} - u^{(h)}\|_{U_h} \leq c \|\varepsilon^{(h)}\|_{F_h},$$

где  $c$  — некоторая постоянная, не зависящая от  $h$ .

В § 3 гл. 8, где введено понятие устойчивости, показано, что в случае линейного оператора  $L_h$  сформулированное определение равносильно следующему.

Разностная краевая задача (2) *устойчива*, если существует  $h_0 > 0$ , такое, что при  $0 < h < h_0$  и любом  $f^{(h)} \in F_h$  она однозначно разрешима, причем

$$\|u^{(h)}\|_{U_h} \leq c \|f^{(h)}\|_{F_h}, \quad (8)$$

где  $c$  — некоторая постоянная, не зависящая от  $h$ ,  $f^{(h)}$ .

Свойство устойчивости можно трактовать как равномерную относительно  $h$  чувствительность решения разностной краевой задачи (2) к возмущениям  $\varepsilon^{(h)}$  правой части.

Подчеркнем, что в силу приведенного определения устойчивость есть некоторое внутреннее свойство разностной краевой задачи. Оно формулируется независимо от какой-либо связи с дифференциальной краевой задачей, в частности, независимо от аппроксимации и сходимости. Напомним, что имеется связь между устойчивостью, аппроксимацией и сходимостью. Эту связь устанавливает следующая

**Теорема.** *Если разностная краевая задача (2) аппроксимирует на решении и дифференциальную краевую задачу (1) и устойчива, то имеет место сходимость (3). При этом порядок относительно  $h$  скорости сходимости совпадает с порядком аппроксимации.*

Несколько более сильная формулировка и доказательство этой важной теоремы приведены в § 3 гл. 8 (теорема 1).

Покажем, что при  $r = \tau/h \leq 1$  разностная схема (5) устойчива. При этом норму  $\|u^{(h)}\|_{U_h}$  определим равенством

$$\|u^{(h)}\|_{U_h} = \max_p \sup_m |u_m^p|. \quad (9)$$

Норму  $\|\cdot\|_{F_h}$  будем понимать как

$$\|f^{(h)}\|_{F_h} = \max_p \sup_m |\varphi_m^p| + \sup_m |\psi_m|.$$

Разностную задачу

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} = \varphi_m^p, \quad u_m^0 = \psi_m, \quad m = 0, \pm 1, \dots, \\ p = 0, 1, \dots, [T/\tau] - 1, \quad (5')$$

которая отличается от задачи (5) только тем, что  $\varphi_m^p$ ,  $\psi_m$  — произвольные правые части, вообще говоря, не совпадающие с  $\varphi(mh, pt)$ ,  $\psi(mh)$  соответственно, перепишем в форме

$$u_m^{p+1} = (1 - r)u_m^p + ru_{m+1}^p + \tau\varphi_m^p, \quad u_m^0 = \psi_m. \quad (6')$$

Поскольку  $r \leq 1$ , то  $1 - r \geq 0$ . В этом случае справедлива оценка

$$|(1 - r)u_m^p + ru_{m+1}^p| \leq |(1 - r) + r| \max \{|u_m^p|, |u_{m+1}^p|\} = \\ = \max \{|u_m^p|, |u_{m+1}^p|\} \leq \sup_m |u_m^p|.$$

Используя эту оценку, выводим из (6') неравенство

$$|u_m^{p+1}| \leq \sup_m |u_m^p| + \tau \sup_m |\varphi_m^p| \leq \sup_m |u_m^p| + \tau \sup_{m,p} |\varphi_m^p|. \quad (6'')$$

Отметим, что в случае  $\varphi_m^p \equiv 0$  из неравенства (6'') следует, что  $\sup_m |u_m^p|$  не возрастает с ростом  $p$ . Отмеченное свойство разностной схемы принято называть *принципом максимума*. Для краткости будем иногда пользоваться этим названием для общего случая неравенства (6'').

Правая часть этого неравенства не зависит от  $m$ , так что в левой части вместо  $|u_m^{p+1}|$  можно написать  $\sup_m |u_m^{p+1}|$ , получив неравенство

$$\sup_m |u_m^{p+1}| \leq \sup_m |u_m^p| + \tau \sup_{p,m} |\varphi_m^p|.$$

Аналогично получаем неравенства

$$\sup_m |u_m^p| \leq \sup_m |u_m^{p-1}| + \tau \sup_{p,m} |\varphi_m^p|,$$

.....

$$\sup_m |u_m^1| \leq \sup_m |u_m^0| + \tau \sup_{p,m} |\varphi_m^p|.$$

После почленного сложения всех этих неравенств и приведения подобных членов получаем

$$\sup_m |u_m^{p+1}| \leq \sup_m |u_m^0| + (p+1)\tau \sup_{p,m} |\varphi_m^p|.$$

Отсюда непосредственно следует неравенство

$$\sup_m |u_m^{p+1}| \leq \sup_m |\psi_m| + T \sup_{p,m} |\varphi_m^p| \leq (1+T) \|f^{(h)}\|_{F_h}.$$

Доказанное неравенство имеет место для всех  $p < T/\tau$ , так что оно останется справедливым, если вместо  $\sup_m |u_m^{p+1}|$  написать

$$\max_p \sup_m |u_m^p| = \|u^{(h)}\|_{U_h}:$$

$$\|u^{(h)}\|_{U_h} \leq (1+T) \|f^{(h)}\|_{F_h}. \quad (10)$$

Неравенство (10) означает устойчивость линейной задачи (5), поскольку существование и единственность решения задачи (6') при произвольных ограниченных  $\varphi_m^p$ ,  $\psi_m$  имеют место. Роль постоянной  $c$  из неравенства (8) играет здесь число  $1+T$ .

Не следует думать, что одна только аппроксимация дифференциальной краевой задачи (1) разностной краевой задачей (2) обеспечивает сходимость (3). Мы убедились в этом в § 3 гл. 1 с помощью специально сконструированного примера аппроксимирующей, но расходящейся разностной схемы.

В случае уравнений с частными производными взятая наудачу аппроксимирующая разностная схема обычно непригодна, а выбор устойчивой (и, следовательно, сходящейся) разностной схемы — постоянная забота вычислителя.

Напомним, например, что доказательство устойчивости разностной схемы (5) мы провели в предположении, что  $r \leq 1$ . В случае же  $r > 1$  разностная задача (5) по-прежнему аппроксимирует задачу (4), но наше доказательство устойчивости не проходит. Покажем, что в этом случае нет сходимости решения  $u^{(h)}$  разностной задачи (5) к решению  $[u]_h$  дифференциальной задачи (4), а значит, не может быть и устойчивости, так как аппроксимация и устойчивость влечет бы за собой сходимость.

**4. Условие Куранта, Фридрихса и Леви.** Покажем, что в случае  $\tau/h > 1$  разностная схема (5) не может сходиться для произвольной функции  $\psi(x)$ . Пусть для определенности  $\varphi(x, t) \equiv 0$ , так что

$$\varphi(mh, rt) \equiv 0; \text{ пусть, далее, } T = 1.$$

Шаг  $h$  будем выбирать так, чтобы точка  $(0, 1)$  на плоскости  $Oxt$  принадлежала сетке, т. е. чтобы число

$$N = \frac{1}{\tau} = \frac{1}{rh}$$

было целым (рис. 16). В силу разностного уравнения имеем

$$u_m^{p+1} = (1 - r)u_m^p + ru_{m+1}^p.$$

Значение  $u_0^{p+1}$  решения  $u^{(h)}$  в точке  $(0, 1)$  сетки выражается через значения  $u_0^p, u_1^p$  решения в точках

$(0, 1 - \tau), (h, 1 - \tau)$  сетки. Значения  $u_0^p, u_1^p$  выражаются через значения  $u_0^{p-1}, u_1^{p-1}, u_2^{p-1}$  решения в трех точках сетки:  $(0, 1 - 2\tau), (h, 1 - 2\tau), (2h, 1 - 2\tau)$ . Значения решения  $u_0^{p-1}, u_1^{p-1}, u_2^{p-1}$  в свою очередь выражаются через значения решения в четырех точках:  $(0, 1 - 3\tau), (h, 1 - 3\tau), (2h, 1 - 3\tau), (3h, 1 - 3\tau)$ , и т. д. В конечном счете значение  $u_0^{p+1}$  выражается через значения  $u_m^0 = \psi_m$  решения в точках сетки  $(0, 0), (h, 0), (2h, 0), \dots, (Nh, 0)$ . Все эти точки лежат на отрезке

$$0 \leq x \leq \frac{h}{\tau} = \frac{1}{r}$$

прямой  $t = 0$  (см. рис. 16), где задано начальное условие

$$u(x, 0) = \psi(x)$$

для дифференциального уравнения.

Таким образом, решение разностного уравнения в точке  $(0, 1)$  сетки не зависит от значений функции  $\psi(x)$  в точках  $x$ , лежащих вне отрезка  $0 \leq x \leq 1/r$ . Но решением задачи

$$\begin{aligned} \frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} &= 0, \quad -\infty < x < \infty, \quad t > 0, \\ u(x, 0) &= \psi(x), \quad -\infty < x < \infty, \end{aligned} \tag{11}$$

как легко проверить, является функция

$$u(x, t) \equiv \psi(x + t).$$

Она постоянна на каждой прямой  $x + t = \text{const}$  (характеристике уравнения (11)) и, в частности, на прямой  $x + t = 1$ , которая проходит через точки  $(0, 1), (1, 0)$  (см. рис. 16) и в точке  $(0, 1)$  принимает значение  $\psi(1)$ . Отсюда видно, что в случае  $r > 1$  сходимости, вообще говоря, быть не может.

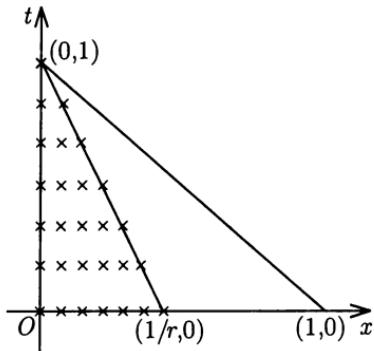


Рис. 16

Действительно, в этом случае отрезок оси абсцисс  $0 \leq x \leq 1/r < 1$  не содержит точку  $(1, 0)$ . Если бы для какой-нибудь функции  $\psi(x)$  сходимость имела место (на самом деле это практически невозможно), то, не меняя значения  $\psi(x)$  на отрезке  $0 \leq x \leq 1/r$  и не меняя, таким образом, значения решения разностного уравнения в точке  $(0, 1)$ , мы могли бы нарушить сходимость, изменив  $\psi(x)$  в точке  $x = 1$  и ее окрестности, что отразилось бы на значении  $u(0, 1)$  решения дифференциального уравнения.

Изменение  $\psi(x)$  в точке  $x = 1$  и ее окрестности можно внести так, чтобы не нарушить существования вторых производных функции  $\psi(x)$  и решения  $u(x, t) = \psi(x + t)$ , так что аппроксимация на решении  $u(x, t)$  будет иметь место. В этих условиях из устойчивости схемы (5) вытекала бы сходимость. Но поскольку при  $r > 1$  нет сходимости, то нет и устойчивости.

Соображение, которое мы использовали для доказательства непригодности схемы (5) в случае  $r = \tau/h > 1$ , носит общий характер. Оно состоит в следующем.

Допустим, что в постановке исходной задачи участвует некоторая функция  $\psi$  (см., например, задачу (4)). Выберем произвольную точку  $P$ , принадлежащую области определения решения  $u$ . Пусть значение решения  $u(P)$  зависит от значений функции  $\psi$  в точках некоторого множества  $G_\psi = G_\psi(P)$ , принадлежащего области определения функции  $\psi$ . Это значит, что, изменяя значения  $\psi$  в малой окрестности любой точки  $Q$  из области  $G_\psi(P)$ , можно вызвать изменение значения  $u(P)$ . Множество  $G_\psi(P)$  будем называть *областью влияния* значений функции  $\psi$  на значение  $u(P)$  решения  $u$ .

Допустим, что для вычисления решения  $u$  используется некоторая разностная схема  $L_h u^{(h)} = f^{(h)}$ , причем значение решения  $u^{(h)}$  в ближайшей к  $P$  точке сетки  $P^{(h)}$  полностью определяется значениями функции  $\psi$  на некотором множестве  $G_\psi^{(h)}(P)$ .

Условие Куранта, Фридрихса и Леви формулируется следующим образом: *для того чтобы имела место сходимость  $u^{(h)} \rightarrow u$  при  $h \rightarrow 0$ , разностная схема должна бытьстроена так, чтобы в произвольной окрестности любой точки области  $G_\psi(P)$  при достаточно малом  $h$  имелась точка множества  $G_\psi^{(h)}(P)$ .*

Объясним, почему в случае невыполнения сформулированного условия Куранта, Фридрихса и Леви сходимости ожидать не приходится. Пусть это условие не выполнено, так что в некоторой фиксированной окрестности некоторой точки  $Q$  из области  $G_\psi(P)$  при всех достаточно малых  $h$  нет точек из множества  $G_\psi^{(h)}(P)$ . Если сходимость  $u^{(h)} \rightarrow u$  при данной функции  $\psi$  имеет место, то изменим  $\psi$  в указанной окрестности точки  $Q$  так, чтобы изменилось значение  $u$ , оставляя вне этой окрестности функцию  $\psi$  неизменной. Сходимость  $u^{(h)} \rightarrow u$  для измененной функции  $\psi$  уже не может иметь места: значение  $u(P)$  изменилось, в то время как значения  $u^{(h)}$  в ближайшей к  $P$  точке

сетки  $P^{(h)}$  остались при малых  $h$  неизменными, поскольку функция  $\psi$  в точках множества  $G_\psi^{(h)}(P)$  осталась неизменной.

Условию Куранта, Фридрихса и Леви нетрудно придать форму теоремы, а проведенные рассуждения превратить в ее доказательство, но мы не будем этого делать.

Подчеркнем еще раз, что в случае, если разностная схема  $L_h u^{(h)} = f^{(h)}$  аппроксимирует исходную задачу на решении  $u$ , то условие Куранта, Фридрихса и Леви является не только необходимым условием сходимости, но также и необходимым условием устойчивости схемы  $L_h u^{(h)} = f^{(h)}$ : в случае невыполнения этого условия устойчивости быть не может, так как аппроксимация и устойчивость влекли бы за собой сходимость в силу теоремы, сформулированной в п. 3.

**5. Механизм неустойчивости.** Проведенное нами доказательство неустойчивости схемы (5) основано на использовании условия Куранта, Фридрихса и Леви, необходимого для сходимости и устойчивости. Оно носит, таким образом, косвенный характер. Интересно проследить непосредственно, как оказывается неустойчивость при  $r > 1$  разностной схемы (5) на чувствительности решения  $u^{(h)}$  к погрешностям в задании  $f^{(h)}$ . Ведь именно равномерная относительно  $h$  чувствительность решения к погрешностям при задании  $f^{(h)}$  и определена выше как устойчивость.

Допустим, что при всех  $h$  выполняются тождества  $\varphi_m^p \equiv 0$ ,  $\psi_m \equiv 0$ , так что

$$f^{(h)} = \begin{Bmatrix} \varphi_m^p \\ \psi_m \end{Bmatrix} = 0,$$

и решение  $u^{(h)} = \{u_m^p\}$  задачи (5) есть тождественный нуль:  $u_m^p \equiv 0$ . Допустим, далее, что при задании начальных данных допущена погрешность и вместо  $\psi_m = 0$  задано  $\tilde{\psi}_m = (-1)^m \varepsilon$  ( $\varepsilon = \text{const}$ ), так что вместо  $f^{(h)} = 0$  задано

$$f^{(h)} = \begin{Bmatrix} 0 \\ \tilde{\psi}_m \end{Bmatrix}, \quad \|f^{(h)}\|_{F_h} = \varepsilon.$$

Будем обозначать получающееся при этом решение через  $\tilde{u}^{(h)}$ . В силу уравнений

$$\tilde{u}_m^{p+1} = (1 - r)\tilde{u}_m^p + r\tilde{u}_{m+1}^p, \quad \tilde{u}_m^0 = (-1)^m \varepsilon$$

для  $\tilde{u}_m^1$  получим  $\tilde{u}_m^1 = (1 - r)\tilde{u}_m^0 + r\tilde{u}_{m+1}^0 = (1 - 2r)\tilde{u}_m^0$ .

Мы видим, что допущенная при  $p = 0$  погрешность умножилась на число  $1 - 2r$ . При переходе к  $\tilde{u}_m^2$  получим

$$\tilde{u}_m^2 = (1 - 2r)\tilde{u}_m^1 = (1 - 2r)^2 \tilde{u}_m^0.$$

Вообще,

$$\tilde{u}_m^p = (1 - 2r)^p \tilde{u}_m^0.$$

При  $r > 1$  будем иметь  $1 - 2r < -1$ , так что исходная погрешность

$$\tilde{u}_m^0 = (-1)^m \varepsilon$$

при каждом переходе от одного слоя  $t = pt$  сетки к следующему умножается на отрицательное число, превосходящее единицу по модулю. Для последнего слоя:  $p = [T/\tau]$ , и

$$|\tilde{u}_m^p| = |1 - 2r|^{[T/\tau]} |\tilde{u}_m^0|.$$

Отсюда

$$\|\tilde{u}^{(h)}\|_{U_h} = \sup_{m,r} |\tilde{u}_m^p| = |1 - 2r|^{[T/(rh)]} \sup_m |\tilde{\psi}_m| = |1 - 2r|^{[T/(rh)]} \|f^{(h)}\|_{F_h}.$$

При фиксированном  $T$  первоначально допущенная в начальных данных погрешность  $(-1)^m \varepsilon$  увеличивается в очень быстро возрастающее при  $h \rightarrow 0$  число раз, равное  $|1 - 2r|^{[T/(rh)]}$ .

**6. Об экономичности разностных схем.** Остановимся теперь кратко на критике принятого нами способа оценки качества аппроксимации сравнением величины нормы невязки  $\|\delta f^{(h)}\|_{F_h}$  с той или иной степенью  $h$ . Как мы знаем, для устойчивых схем порядок аппроксимации совпадает с порядком погрешности  $\|[u]_h - u^{(h)}\|_{U_h}$  в решении. Качество схемы естественно оценивать по количеству вычислительной работы, необходимой для получения заданной точности. Количество же работы, как правило, пропорционально числу точек  $N$  использованной разностной сетки. Для обыкновенных дифференциальных уравнений  $N$  обратно пропорционально шагу  $h$ . Поэтому, когда мы говорим, что  $\varepsilon$  — норма погрешности, где  $\varepsilon \sim h^k$ , мы тем самым утверждаем, что  $\varepsilon = N^{-k}$ , т. е. что, например, уменьшение погрешности вдвое требует увеличения работы в  $\sqrt[4]{2}$  раз. Таким образом, в случае обыкновенных разностных уравнений порядок аппроксимации относительно  $h$  характеризует объем работы.

Для уравнений с частными производными дело обстоит уже не так. В рассмотренном нами примере задачи с двумя переменными  $x, t$  сетка задается двумя шагами  $h, \tau$ . Число  $N$  точек сетки, помещающихся в ограниченной области на плоскости  $Oxt$ , имеет порядок  $1/(\tau h)$ . Это число также может применяться для оценки количества работы, затрачиваемой при решении разностных уравнений. Пусть  $\tau = rh$ . В этом случае  $N \sim h^{-2}$ , и утверждение, что  $\varepsilon \sim h^k$ , эквивалентно утверждению  $\varepsilon \sim N^{-k/2}$ . Если  $\tau = rh^2$ , то  $N \sim h^{-3}$ , и утверждение  $\varepsilon \sim h^k$  эквивалентно тому, что  $\varepsilon \sim N^{-1/3}$ .

Мы видим, что в случае уравнений с частными производными порядок погрешности естественнее было бы измерять не в степенях  $h$ , а в степенях  $N^{-1}$ . Мы все же остановимся на описанном выше способе оценки аппроксимации степенями, так как это удобнее при проведении выкладок. Читатель, однако, должен при оценке качества разностных схем иметь в виду отмеченное обстоятельство.

Надо еще заметить, что утверждение о пропорциональности вычислительной работы числу  $N$  точек сетки тоже не всегда является верным. Можно привести примеры разностных схем, для которых вычисление решения требует произвести  $O(N^{1+q})$  арифметических операций, где  $q$  равно  $1/2$ ,  $1$  или даже  $2$ . С этим приходится встречаться при решении разностных краевых задач, аппроксимирующих эллиптические уравнения, или при решении задач в случае трех или более независимых переменных (например,  $u(t, x, y)$ ).

При реальных расчетах на компьютере для сравнительной оценки используемых алгоритмов за меру качества схемы обычно естественно принять машинное время. Машинное время не обязательно пропорционально числу арифметических действий. Оказывают влияние (иногда превалирующее) также затраты времени на пересылку информации из одного блока машинной памяти в другой. Может играть роль и время, расходуемое на логические операции.

**7. Библиографическая справка.** Понятие устойчивости разностных схем относительно ошибок округления при задании начальных данных впервые описано Дж. фон Нейманом и Р.Д. Рихтмайером в 1950 г. (см.: Механика. — 1951. — Вып. 1) в работе, посвященной расчету газодинамических скачков. Первая система определений устойчивости и аппроксимации, при которой сходимость является следствием аппроксимации и устойчивости, была предложена В.С. Рябеньким (О применении метода конечных разностей к решению задачи Коши // ДАН СССР. — 1952. — Т. 86, № 6) в случае разностных аналогов задачи Коши для систем уравнений с частными производными с коэффициентами, зависящими только от  $t$ .

Принятая в нашей книге система основных определений и теорема о том, что из аппроксимации и устойчивости следует сходимость, близки к предложенным А.Ф. Филипповым (ДАН СССР. — 1955. — Т. 100, № 6; см. также [18]). Отличие состоит главным образом в том, что мы используем более универсальное определение аппроксимации.

Существуют другие естественные системы определений основных понятий, при которых аппроксимация и устойчивость обеспечивают сходимость. Среди них наиболее известна система определений П.Д. Лакса. В теории Лакса рассматриваются разностные схемы для нестационарных задач, причем предполагается, что эти разностные схемы действуют не в пространстве сеточных функций, а в том же функциональном пространстве, что и дифференциальное уравнение. При этом предположении доказывается, что для аппроксимирующей разностной схемы устойчивость и сходимость имеют место одновременно.

В последующие годы А.А. Самарский предложил и развил в соавторстве с А.В. Гулиным теорию устойчивости, применимую к широкому классу разностных схем [19].

Следует сказать, что в работе 1928 г. Р. Куранта, К. Фридрихса и Г. Леви (О разностных уравнениях математической физики // УМН. — 1940. — Т. 8. — С. 125–160) и во многих других работах, где метод

конечных разностей используется для доказательства существования решений дифференциальных уравнений, устанавливаются неравенства, которые в современной терминологии можно истолковать как устойчивость в тех или иных нормах. Однако понятие устойчивости возникло в связи с использованием разностных схем для приближенного вычисления решений в предположении, что эти решения существуют. Поэтому устойчивость изучается обычно в более слабых нормах, чем это нужно для доказательства существования.

Отметим, что впервые метод конечных разностей для доказательства существования решений уравнений с частными производными был использован в 1924 г. Л.А. Люстерником, который рассматривал уравнение Лапласа (см.: УМН. — 1940. — Т. 8).

### Задачи

1. Для задачи Коши (4) исследовать следующую разностную схему:

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_m^p - u_{m-1}^p}{h} = \varphi_m^p, \quad m = 0, \pm 1, \dots,$$

$$u_m^0 = \psi_m, \quad p = 0, 1, \dots, \left[ \frac{T}{\tau} \right] - 1,$$

где  $\tau = rh$ ,  $r = \text{const}$ .

а) Выписать оператор  $L_h$  и правую часть  $f^{(h)}$ , возникающие при записи этой схемы в виде  $L_h u^{(h)} = f^{(h)}$ .

б) Изобразить взаимное расположение трех точек сетки (шаблон), значения  $u^{(h)}$  в которых связывает разностное уравнение при фиксированных  $m, p$ .

в) Показать, что разностная схема аппроксимирует дифференциальную задачу с первым относительно  $h$  порядком на решении  $u(x, t)$ , имеющем ограниченные вторые производные.

г) Выяснить, устойчива ли исследуемая разностная схема при каком-либо выборе  $r$ .

2. Для задачи Коши

$$u_t + u_x = \varphi(x, t), \quad -\infty < x < \infty,$$

$$u(x, 0) = \psi(x), \quad 0 < t < T,$$

исследовать по предложенному в задаче 1 плану каждую из следующих разностных схем:

$$\frac{u_m^{p+1} - u_m^p}{\tau} + \frac{u_m^p - u_{m-1}^p}{h} = \varphi_m^p, \quad u_m^0 = \psi_m;$$

$$\frac{u_m^{p+1} - u_m^p}{\tau} + \frac{u_{m+1}^p - u_m^p}{h} = \varphi_m^p, \quad u_m^0 = \psi_m.$$

3. Для задачи Коши для уравнения теплопроводности

$$\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = \varphi(x, t), \quad -\infty < x < \infty, \quad 0 \leq t \leq T, \quad (12)$$

$$u(x, 0) = \psi(x), \quad -\infty < x < \infty,$$

рассмотреть разностную схему

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} &= \varphi(mh, p\tau), \\ m = 0, \pm 1, \dots, \quad p = 0, 1, \dots, \left[ \frac{T}{\tau} \right] - 1, \\ u_m^0 &= \psi(mh), \quad m = 0, \pm 1, \dots \end{aligned} \tag{13}$$

Нормы в  $U_h$ ,  $F_h$  понимаем в смысле (9), (7) соответственно.

а) Проверить, что при  $\tau/h^2 = r = \text{const}$  разностная схема (13) аппроксимирует задачу (12) с порядком  $O(h^2)$ .

б)\* Показать, что в случае  $r \leq 1/2$ ,  $\varphi(x, t) \equiv 0$  справедлив следующий принцип максимума:

$$\sup_m |u_m^{p+1}| \leq \sup_m |u_m^p|.$$

в)\* Опираясь на принцип максимума, доказать, что в случае  $r \leq 1/2$  разностная схема (13) устойчива.

**4.** Решение задачи Коши (12) для уравнения теплопроводности имеет вид

$$u(x, t) = \frac{1}{2\sqrt{\pi}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{t}} \psi(s) \exp \left\{ -\frac{(x - \xi)^2}{4t} \right\} d\xi.$$

Может ли оказаться устойчивой и сходящейся «явная» разностная схема (12), аппроксимирующая эту задачу, если  $\tau = h$ ,  $h \rightarrow 0$ ?

**Указание.** Сопоставить области влияния для дифференциальной и разностной задач и воспользоваться условием Куранта, Фридрихса и Леви (см. п. 4).

**5.** Система уравнений акустики

$$\begin{aligned} \frac{\partial v}{\partial t} &= \frac{\partial w}{\partial x}, \quad \frac{\partial w}{\partial t} = \frac{\partial v}{\partial x}, \quad -\infty < x < \infty, \quad 0 \leq t \leq T, \\ v(x, 0) &= \varphi(x), \quad w(x, 0) = \psi(x), \end{aligned}$$

имеет решение вида

$$\begin{aligned} v(x, t) &= \frac{\varphi(x - t) - \psi(x - t) + \varphi(x + t) + \psi(x + t)}{2}, \\ w(x, t) &= \frac{-\varphi(x - t) + \psi(x - t) + \varphi(x + t) + \psi(x + t)}{2}. \end{aligned}$$

Может ли оказаться сходящейся разностная схема вида

$$\begin{aligned} \frac{v_m^{p+1} - v_m^p}{\tau} + \frac{w_{m+1}^p - w_m^p}{h} &= 0, \quad m = 0, \pm 1, \dots, p \geq 0, \\ \frac{w_m^{p+1} - w_m^p}{\tau} + \frac{v_{m+1}^p - v_m^p}{h} &= 0, \quad v_m^0 = \varphi(mh), \quad w_m^0 = \psi(mh)? \end{aligned}$$

Сопоставить области влияния начальных данных для разностной и дифференциальной задач.

**6\*. Задача Коши**

$$\frac{\partial u}{\partial t} - \frac{\partial u}{\partial x}, \quad t > 0, \quad -\infty < x < \infty,$$

$$u(x, 0) = e^{i\alpha x}, \quad 0 < \alpha < 2\pi, \quad \alpha = \text{const},$$

имеет решение

$$u(x, t) = e^{i\alpha t} e^{i\alpha x}.$$

Соответствующая разностная схема

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} = 0, \quad p = 0, 1, \dots, \quad m = 0, \pm 1, \dots,$$

$$u_m^0 = e^{i\alpha(mh)}$$

имеет решение

$$u_m^p = (1 - r + re^{i\alpha h})^p e^{i\alpha(mh)}, \quad \text{где } p = t/\tau, \quad m = x/h,$$

которое стремится к решению дифференциальной задачи при  $h \rightarrow 0$ , каково бы ни было фиксированное  $r = \tau/h$ . Между тем при  $r > 1$  разностная схема не удовлетворяет условию Куранта, Фридрихса и Леви, необходимому для сходимости. Объясните кажущийся парадокс.

## § 2. Некоторые приемы построения аппроксимирующих разностных схем

**1. Замена производных разностными отношениями.** Простейший прием построения разностных краевых задач, аппроксимирующих дифференциальные, состоит в замене производных соответствующими разностными отношениями. Приведем несколько примеров разностных схем, полученных таким способом. В этих примерах будут использованы приближенные формулы:

$$\begin{aligned} \frac{df(z)}{dz} &\approx \frac{f(z + \Delta z) - f(z)}{\Delta z}, \quad \frac{df(z)}{dz} \approx \frac{f(z) - f(z - \Delta z)}{\Delta z}, \\ \frac{df(z)}{dz} &\approx \frac{f(z + \Delta z) - f(z - \Delta z)}{2\Delta z}, \\ \frac{d^2 f(z)}{dz^2} &\approx \frac{f(z + \Delta z) - 2f(z) + f(z - \Delta z)}{(\Delta z)^2}. \end{aligned} \tag{1}$$

Предполагая функцию  $f(z)$  имеющей достаточное число ограниченных производных, можно выписать выражение для остаточных членов этих формул. По формуле Тейлора

$$\begin{aligned} f(z + \Delta z) &= f(z) + \Delta z f'(z) + \frac{(\Delta z)^2}{2!} f''(z) + \\ &\quad + \frac{(\Delta z)^3}{3!} f'''(z) + \frac{(\Delta z)^4}{4!} f^{(4)}(z) + o[(\Delta z)^4], \end{aligned} \tag{2}$$

$$\begin{aligned} f(z - \Delta z) &= f(z) - \Delta z f'(z) + \frac{(\Delta z)^2}{2!} f''(z) - \\ &\quad - \frac{(\Delta z)^3}{3!} f'''(z) + \frac{(\Delta z)^4}{4!} f^{(4)}(z) + o[(\Delta z)^4]. \end{aligned}$$

Используя разложения (2), можно получить выражения для остаточных членов приближенных формул (1). Именно, справедливы равенства

$$\begin{aligned}\frac{f(z + \Delta z) - f(z)}{\Delta z} &= f'(z) + \left[ \frac{\Delta z}{2} f''(z) + o(\Delta z) \right], \\ \frac{f(z) - f(z - \Delta z)}{\Delta z} &= f'(z) + \left[ -\frac{\Delta z}{2} f''(z) + o(\Delta z) \right], \\ \frac{f(z + \Delta z) - 2f(z) + f(z - \Delta z)}{(\Delta z)^2} &= f''(z) + \left[ \frac{(\Delta z)^2}{12} f^{(4)}(z) + o((\Delta z)^2) \right].\end{aligned}\quad (3)$$

Остаточные члены приближенных формул (1) входят в соответствующие равенства (3) в виде выражений в квадратных скобках.

Очевидно, что формулы (1) и выражения остаточных членов, выписанные в формулах (3), можно использовать и при замене частных производных разностными отношениями. Например,

$$\frac{\partial u(x, t)}{\partial t} \approx \frac{u(x, t + \Delta t) - u(x, t)}{\Delta t},$$

причем

$$\frac{u(x, t + \Delta t) - u(x, t)}{\Delta t} = \frac{\partial u(x, t)}{\partial t} + \left[ \frac{\Delta t}{2} \frac{\partial^2 u(x, t)}{\partial t^2} + o(\Delta t) \right].$$

Точно так же справедлива формула

$$\frac{\partial u(x, t)}{\partial x} \approx \frac{u(x + \Delta x, t) - u(x, t)}{\Delta x},$$

и при этом

$$\frac{u(x + \Delta x, t) - u(x, t)}{\Delta x} = \frac{\partial u(x, t)}{\partial x} + \left[ \frac{\Delta x}{2} \frac{\partial^2 u(x, t)}{\partial x^2} + o(\Delta x) \right],$$

и т. д.

**Пример 1.** Вернемся к задаче Коши (4) из § 5:

$$\begin{aligned}\frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} &= \varphi(x, t), \quad -\infty < x < \infty, \quad 0 \leq t \leq T, \\ u(x, 0) &= \psi(x), \quad -\infty < x < \infty.\end{aligned}\quad (4)$$

Для аппроксимации этой задачи Коши построим две разностные схемы. В этих схемах используем сетку  $D_h$ , образованную точками пересечения прямых  $x = mh$ ,  $t = p\tau$ , попавшими в полосу  $0 \leq t \leq T$ . Значения  $\tau$ ,  $h$  будем считать связанными соотношением  $\tau = rh$ , где  $r$  — некоторая положительная постоянная. Простейшая из этих схем имеет вид

$$L_h^{(1)} u^{(h)} = \begin{cases} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} = \varphi(mh, p\tau), \\ u_m^0 = \psi(mh) \end{cases} \quad (5)$$

и получается при замене производных  $\partial u / \partial t$ ,  $\partial u / \partial x$  по приближенным формулам

$$\frac{\partial u(x, t)}{\partial t} \approx \frac{u(x, t + \tau) - u(x, t)}{\tau},$$

$$\frac{\partial u(x, t)}{\partial x} = \frac{u(x + h, t) - u(x, t)}{h}.$$

Невязка  $\delta f^{(h)}$ , возникающая при подстановке решения  $[u]_h$  дифференциальной задачи в левую часть разностной задачи:

$$L_h[u]_h = f^{(h)} + \delta f^{(h)},$$

выражается формулой

$$\delta f^{(h)} = \begin{cases} \left( \frac{\tau}{2} u_{tt} - \frac{h}{2} u_{xx} \right)_m^p + O(\tau + h), \\ 0. \end{cases}$$

За норму элемента  $f^{(h)}$  пространства  $F_h$  примем в этом параграфе максимум всех компонент элемента  $f^{(h)} \in F_h$ . Тогда очевидно, что

$$\| \delta f^{(h)} \|_{F_h} = O(\tau + h) = O(rh + h) = O(h),$$

и порядок аппроксимации получается первый.

Вторая схема получается при использовании другой формулы для замены  $\partial u(x, t) / \partial x$ :

$$\frac{\partial u(x, t)}{\partial x} \approx \frac{u(x, t) - u(x - h, t)}{h}.$$

Она имеет вид

$$L_h^{(2)} u^{(h)} = \begin{cases} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_m^p - u_{m-1}^p}{h} = \varphi(mh, p\tau), \\ u_m^0 = \psi_m. \end{cases}$$

Порядок аппроксимации снова получается первый.

Вторая схема, казалось бы, совсем несущественно отличается от первой. В действительности, однако, вторая схема непригодна для счета: для нее при любом  $r = \tau/h$  не выполнено условие Куранта, Фридрихса и Леви.

Чтобы различить операторы  $L_h$  этих двух схем, мы снабдили их номерами и написали  $L_h^{(1)} u^{(h)} = f^{(h)}$ ,  $L_h^{(2)} u^{(h)} = f^{(h)}$ . Для облегчения запоминания разностной схемы ее обычно принято сопоставлять с картинкой, на которой изображено взаимное расположение точек сетки (шаблон), значения в которых связывает разностное уравнение при некоторых фиксированных значениях  $m$ ,  $p$ . Для двух рассмотренных схем эти картинки изображены на рис. 17.

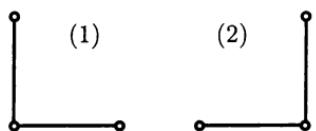


Рис. 17

Пример 2. Приведем две разностные схемы, аппроксимирующие задачу Коши для уравнения теплопроводности:

$$\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = \varphi(x, t), \quad -\infty < x < \infty, \quad 0 < t < T,$$

$$u(x, 0) = \psi(x), \quad -\infty < x < \infty.$$

Простейшая из них:

$$L_h^{(1)} u^{(h)} = \begin{cases} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} = \varphi(mh, p\tau), \\ u_m^0 = \psi(mh), \end{cases}$$

$$f^{(h)} = \begin{cases} \varphi(mh, p\tau), \\ \psi(mh) \end{cases}$$

получается при замене производных  $u_t$ ,  $u_{xx}$  разностными отношениями по формулам

$$u_t(x, t) \approx \frac{u(x, t + \tau) - u(x, t)}{\tau},$$

$$u_{xx}(x, t) \approx \frac{u(x + h, t) - 2u(x, t) + u(x - h, t)}{h^2}.$$

Если для замены  $u_{xx}(x, t)$  использовать другую формулу:

$$u_{xx}(x, t) \approx \frac{u(x + h, t + \tau) - 2u(x, t + \tau) + u(x - h, t + \tau)}{h^2},$$

придем к другой схеме для того же уравнения:

$$L_h^{(2)} u^{(h)} = \begin{cases} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^{p+1} - 2u_m^{p+1} + u_{m-1}^{p+1}}{h^2} = \varphi(mh, p\tau), \\ u_m^0 = \psi(mh). \end{cases}$$

Шаблоны, соответствующие этим разностным схемам, изображены на рис. 18.



Рис. 18

Эти схемы существенно различаются. Вычисление решения по первой из них не представляет труда и проводится по явной формуле

$$u_m^{p+1} = (1 - 2r)u_m^p + r(u_{m+1}^p + u_{m-1}^p) + \tau\varphi(mh, p\tau),$$

где  $r = \tau/h^2$ . Эта формула получена из разностного уравнения в результате решения его относительно  $u_m^{p+1}$ . Зная значения решения  $u_m^p$  ( $m = 0, \pm 1, \dots$ ) на слое  $t = t_p = p\tau$  сетки, мы можем вычислять его значения  $u_m^{p+1}$  на следующем слое:  $t = t_{p+1} = (p + 1)\tau$ .

Вторая схема лишена этого удобного свойства. В этом случае разностное уравнение, выписанное при фиксированных  $m, p$ , нельзя разрешить относительно  $u_m^{p+1}$ , выразив это значение через известные значения  $u_{m-1}^p, u_m^p, u_{m+1}^p$  с предыдущего слоя. Дело в том, что в это уравнение входит не только неизвестное значение  $u_m^{p+1}$ , но также и неизвестные  $u_{m+1}^{p+1}, u_{m-1}^{p+1}$ . Поэтому для определения  $u_m^{p+1}$  ( $m = 0, \pm 1, \dots$ ) придется решать разностное уравнение относительно сеточной функции  $u_m^{p+1}$  аргумента  $m$ . Тем не менее в дальнейшем будет показано, что схема  $L_h^{(2)}u^{(h)} = f^{(h)}$ , как правило, удобнее схемы  $L_h^{(1)}u^{(h)} = f^{(h)}$ .

При  $\tau = rh^2$  ( $r = \text{const}$ ) обе схемы имеют второй порядок аппроксимации относительно  $h$ . Вычислим невязку  $\delta f^{(h)}$  и оценим порядок аппроксимации для второй из этих схем. Пользуясь формулами (3), можно написать

$$L_h^{(2)}[u]_h = \begin{cases} (u_t - u_{xx}) \Big|_{\substack{x=mh \\ t=(p+1)\tau}} - \frac{\tau}{2} u_{tt}(x, t_{p+1}) - \\ \quad - \frac{h^2}{12} \cdot \frac{\partial^4 u(x, t_{p+1})}{\partial x^4} + O(\tau + h^2), \\ u(mh, 0). \end{cases}$$

Отсюда с учетом  $\tau = rh^2$  можно написать

$$\begin{aligned} L_h^{(2)}[u]_h &= \begin{cases} \varphi(mh, (p+1)\tau) + O(h^2), \\ \psi(mh) + 0; \end{cases} \\ \delta f^{(h)} &= \begin{cases} \varphi(x_m, t_p) - \varphi(x_m, t_{p+1}) + O(h^2), \\ 0. \end{cases} \end{aligned}$$

Но

$$\begin{aligned} \varphi(x_m, t_{p+1}) &= \varphi(x_m, t_p) + [\varphi(x_m, t_{p+1}) - \varphi(x_m, t_p)] = \\ &= \varphi(x_m, t_p) + O(\tau) = \varphi(x_m, t_p) + O(h^2). \end{aligned}$$

Поэтому

$$\|\delta f^{(h)}\|_{F_h} = O(h^2).$$

**Пример 3.** Рассмотрим простейшую разностную схему, аппроксимирующую задачу Дирихле для уравнения Пуассона в квадрате  $D$  ( $0 < x < 1, 0 < y < 1$ ) с границей  $\Gamma$  (см. рис. 8 в гл. 4, § 1):

$$\begin{aligned} \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} &= \varphi(x, y), \quad (x, y) \in D, \\ u|_{\Gamma} &= \psi(x, y), \quad (x, y) \in \Gamma. \end{aligned}$$

Построим сетку  $D_h$ , отнеся к ней те точки  $(x_m, y_n) = (mh, nh)$ , которые попали внутрь квадрата или на его границу. Шаг  $h$  выберем так, чтобы число  $1/h$  было целым.

Разностную схему  $L_h u^{(h)} = f^{(h)}$  зададим равенствами

Здесь

$$f^{(h)} = \begin{cases} \varphi(mh, nh), & (mh, nh) \in D, \\ \psi(mh, nh), & (mh, nh) \in \Gamma. \end{cases}$$

Невязка  $\delta f^{(h)}$  имеет в силу формул (3) вид

$$\delta f^{(h)} = \begin{cases} \frac{h^2}{12} \left( \frac{\partial^4 u}{\partial x^4} + \frac{\partial^4 u}{\partial y^4} \right) + O(h^2), \\ 0, \end{cases}$$

так что аппроксимация имеет второй порядок. Пятиточечный шаблон, отвечающий использованному разностному уравнению, изображен на рис. 19.

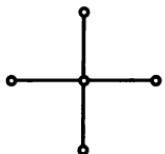


Рис. 19

Разностные схемы, построенные выше, получались путем замены каждой производной в дифференциальном уравнении тем или иным разностным отношением.

**2. Метод неопределенных коэффициентов.** Более общий способ построения разностных схем состоит в том, что приближается не каждая производная в отдельности, а сразу весь дифференциальный оператор.

Разъясним этот способ на примерах разностных схем для задачи Коши (4). Сначала рассмотрим схему первого порядка аппроксимации (5). Она связывает значения искомой функции в трех точках, изображенных на рис. 17 слева. Разностное уравнение

$$\Lambda_h u^{(h)} = \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} = \varphi(mh, p\tau),$$

используемое в этой схеме, имеет вид

$$\Lambda_h u^{(h)} \equiv a^1 u_m^{p+1} + a_0 u_m^p + a_1 u_{m+1}^p = \varphi(mh, p\tau).$$

Забудем на время, что нам уже известна разностная схема (5), для которой

$$a^1 = \frac{1}{\tau}, \quad a_0 = \frac{1}{h} - \frac{1}{\tau}, \quad a_1 = -\frac{1}{h},$$

и, считая эти коэффициенты неопределенными, постараемся подобрать их так, чтобы имело место равенство

$$\Lambda_h[u]_h \Big|_{\substack{x=mh \\ t=p\tau}} = \left( \frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} \right) \Big|_{\substack{x=mh \\ t=p\tau}} + O(h),$$

или

$$\Lambda_h[u]_h \Big|_{\substack{x=mh \\ t=p\tau}} = \Lambda u \Big|_{\substack{x=mh \\ t=p\tau}} + O(h), \quad (6)$$

где

$$\Lambda u = \frac{\partial u}{\partial t} - \frac{\partial u}{\partial x}. \quad (7)$$

Воспользуемся формулой Тейлора:

$$u(mh, (p+1)\tau) = u(mh, p\tau) + \tau \frac{\partial u(mh, p\tau)}{\partial t} + O(\tau^2),$$

$$u((m+1)h, p\tau) = u(mh, p\tau) + h \frac{\partial u(mh, p\tau)}{\partial x} + O(h^2).$$

Подставив эти выражения в правую часть равенства

$$\Lambda_h[u]_h \Big|_{\substack{x=mh \\ t=p\tau}} \equiv a^1 u(mh, (p+1)\tau) + a_0 u(mh, p\tau) + a_1 u((m+1)h, p\tau),$$

получим

$$\Lambda_h[u]_h \Big|_{\substack{x=mh \\ t=p\tau}} = (a^1 + a_0 + a_1) u(mh, p\tau) +$$

$$+ a^1 \tau \frac{\partial u(mh, p\tau)}{\partial t} + a_1 h \frac{\partial u(mh, p\tau)}{\partial x} + O(a^1 \tau^2) + O(a_1 h^2). \quad (8)$$

Поскольку нашей целью является такой подбор коэффициентов  $a^1$ ,  $a_0$ ,  $a_1$ , чтобы выполнялось условие аппроксимации (6), то естественно так предварительно сгруппировать слагаемые в правой части равенства (8), чтобы выделился член (7). Тогда остальные слагаемые образуют остаточный член аппроксимации, который должен быть мал. Чтобы выделить член  $\Lambda u$ , можно заменить в правой части равенства (8) производные  $\partial u / \partial t$  или  $\partial u / \partial x$  соответственно по одной из формул

$$\frac{\partial u}{\partial t} \equiv \Lambda u + \frac{\partial u}{\partial x} \quad \text{или} \quad \frac{\partial u}{\partial x} = \frac{\partial u}{\partial t} - \Lambda u.$$

Для определенности воспользуемся первой из них.

Кроме того, подчиним шаги  $\tau$ ,  $h$  связи  $\tau = rh$ , где  $r$  — какая-нибудь постоянная. После этого равенство (8) примет следующий вид:

$$\Lambda_h[u]_h \Big|_{(mh, p\tau)} = a^1 rh \Lambda u \Big|_{(mh, p\tau)} + (a^1 + a_0 + a_1) u(mh, p\tau) +$$

$$+ (a^1 r + a_1) h \frac{\partial u(mh, p\tau)}{\partial x} + O(a^1 r^2 h^2) + O(a_1, h^2). \quad (9)$$

Среди всех гладких функций  $u(x, t)$  можно указать такие, для которых  $u$ ,  $\partial u / \partial x$ ,  $\partial u / \partial t$  в любой заранее заданной фиксированной точке принимают любые независимые друг от друга значения. Следовательно, и значения

$$u, \quad \frac{\partial u}{\partial x}, \quad \Lambda u \equiv \frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} = \varphi(x, t)$$

также можно считать независимыми друг от друга. В силу этого из равенства (9) следует, что для выполнения при любой правой части  $\varphi(x, t)$  задачи (4) условия аппроксимации

$$\Lambda_h[u]_h|_{(mh, p\tau)} = (\Lambda u)_{(mh, p\tau)} + O(h)$$

необходимо, чтобы выполнялись равенства

$$a^1 rh = 1 + O_1(h), \quad a^1 + a_0 + a_1 = 0 + O_2(h), \quad a^1 r + a_1 = 0 + O_3(h),$$

где  $O_1(h)$ ,  $O_2(h)$ ,  $O_3(h)$  — какие-нибудь произвольные величины порядка  $O(h)$ . Положим  $O_1(h) = O_2(h) = O_3(h) = 0$ . Получающаяся при этом система

$$a^1 rh = 1, \quad a^1 + a_0 + a_1 = 0, \quad a^1 r + a_1 = 0$$

имеет единственное решение

$$a^1 = \frac{1}{rh} = \frac{1}{\tau}, \quad a_0 = \frac{r - 1}{rh} = \frac{1}{h} - \frac{1}{\tau}, \quad a_1 = -\frac{1}{h},$$

которое приводит к уже известной схеме (5).

Теперь мы дополнительно узнали, что среди разностных схем вида

$$L_h u^{(h)} = \begin{cases} a^1 u_m^{p+1} + a_0 u_m^p + a_1 u_{m+1}^p = \varphi(mh, p\tau), \\ u_m^0 = \psi(mh) \end{cases}$$

она является единственной, аппроксимирующей рассматриваемую задачу Коши. Говоря о единственности, мы пренебрегаем тем произволом, который вносит свобода выбора функций  $O_1(h)$ ,  $O_2(h)$ ,  $O_3(h)$ . Всюду в дальнейших примерах также будем пренебречь подобного рода очевидным произволом и даже не всегда будем вводить произвольные величины, аналогичные функциям  $O_1(h)$ ,  $O_2(h)$ ,  $O_3(h)$ , с самого начала полагая их равными нулю.

Посмотрим теперь, как можно строить для задачи (4) разностные схемы

$$L_h u^{(h)} \equiv \begin{cases} a^1 u_m^{p+1} + a_0 u_m^p + a_{-1} u_{m-1}^p + a_1 u_{m+1}^p = \varphi(mh, p\tau), \\ u_m^0 = \psi(mh) \end{cases} \quad (10)$$

более общего вида, связывающие значения искомой функции в четырех точках. Шаги сетки снова связем равенством  $\tau = rh$  ( $r = \text{const}$ ) и введем обозначение  $\Lambda_h$ , положив

$$\Lambda_h u^{(h)} \equiv a^1 u_m^{p+1} + a_0 u_m^p + a_{-1} u_{m-1}^p + a_1 u_{m+1}^p. \quad (11)$$

Для всякой достаточной гладкой функции  $u(x, t)$  с помощью формулы Тейлора можно записать

$$\begin{aligned} \Lambda_h[u]_h &= (a^1 + a_0 + a_{-1} + a_1) u(mh, p\tau) + a^1 rh \frac{\partial u(mh, p\tau)}{\partial t} + \\ &+ (a_1 - a_{-1}) h \frac{\partial u(mh, p\tau)}{\partial x} + \frac{1}{2} a^1 r^2 h^2 \frac{\partial^2 u(mh, p\tau)}{\partial t^2} + \\ &+ \frac{1}{2} (a_1 + a_{-1}) h^2 \frac{\partial^2 u(mh, p\tau)}{\partial x^2} + O(a^1 r^3 h^3 + a_1 h^3 + a_{-1} h^3). \end{aligned}$$

Выделим в правой части этого равенства член  $\Lambda u \equiv \partial u / \partial t - \partial u / \partial x$ , воспользовавшись для этого тождеством  $u_t = u_x + \Lambda u$ . Тогда

$$\begin{aligned} \Lambda_h[u]_h &= a^1 r h \Lambda u|_{(mh, p\tau)} + (a^1 + a_0 + a_1 + a_{-1}) u(mh, p\tau) + \\ &+ (a^1 r + a_1 - a_{-1}) h \frac{\partial u(mh, p\tau)}{\partial x} + \frac{1}{2} a^1 r^2 h^2 \frac{\partial^2 u(mh, p\tau)}{\partial t^2} + \\ &+ \frac{1}{2} (a_1 + a_{-1}) h^2 \frac{\partial^2 u(mh, p\tau)}{\partial x^2} + O(a^1 r^3 h^3 + a_1 h^3 + a_{-1} h^3). \end{aligned} \quad (12)$$

Если предположить, что величина  $O(a^1 r^3 h^3 + a_1 h^3 + a_{-1} h^3)$  достаточно мала (это предположение подтверждается в дальнейшем), то для выполнения условия аппроксимации

$$\Lambda_h[u]_h|_{(mh, p\tau)} = \Lambda u|_{(mh, p\tau)} + O(h)$$

необходимо, чтобы четыре числа  $a^1, a_0, a_1, a_{-1}$  удовлетворяли следующим трем равенствам:

$$\begin{aligned} a^1 r h &= 1, \\ a^1 + a_0 + a_1 + a_{-1} &= 0, \\ a^1 r + a_1 - a_{-1} &= 0. \end{aligned} \quad (13)$$

Система (13) имеет много решений — семейство решений, зависящее от одного параметра. Одно из этих решений

$$a^1 = \frac{1}{rh}, \quad a_0 = \frac{r-1}{h}, \quad a_{-1} = 0, \quad a_1 = -\frac{1}{h}$$

дает уже рассмотренную схему (5). Решению

$$a^1 = \frac{1}{rh}, \quad a_0 = -\frac{1}{rh}, \quad a_{-1} = \frac{1}{2h}, \quad a_1 = -\frac{1}{2h}$$

соответствует схема

$$L_h u^{(h)} \equiv \begin{cases} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_{m-1}^p}{2h} = \varphi(mh, p\tau), \\ u_m^0 = \psi(mh). \end{cases} \quad (14)$$

Выбрав какое-либо решение системы (13), надо его подставить в остаточный член и убедиться, что он мал. Для двух сейчас приведенных решений подстановка чисел  $a^1, a_0, a_1, a_{-1}$  дает остаточные члены

$$\frac{a^1 r^2 h^2}{2} \cdot \frac{\partial^2 u}{\partial t^2} + \frac{a_1 + a_{-1}}{2} h^2 \frac{\partial^2 u}{\partial x^2} + O(a^1 r^3 h^3 + a_1 h^3 + a_{-1} h^3)$$

порядка  $O(h)$ .

В самом деле, среди гладких функций  $u(x, t)$  имеются многочлены второй степени, для которых  $\partial^2 u / \partial t^2, \partial^2 u / \partial x^2$  принимают в любой фиксированной точке любые независимые наперед заданные значения. При этом член  $O(a^1 r^3 h^3, a_1 h^3, a_{-1} h^3)$ , в который входят третий производные многочлена  $u(x, t)$ , обращается в нуль. Поэтому для того чтобы при любом выборе гладкой функции  $u(x, t)$  остаточный член был

порядка  $h$ , необходимо и достаточно, чтобы хотя бы один из коэффициентов при  $\partial^2 u / \partial t^2$ ,  $\partial^2 u / \partial x^2$  был порядка  $h$ . Поскольку из первого уравнения (13) имеем  $a^1 = 1/(rh)$ , то коэффициент при  $\partial^2 u / \partial t^2$  есть  $rh/2$ , и порядок остаточного члена всегда не выше первого.

Мы установили, что в общем случае нельзя построить разностную схему вида (10), которая аппроксимирует задачу (4) с порядком  $h^2$ . Для увеличения порядка аппроксимации пришлось бы увеличить число точек разностной сетки, используемых в конструируемой схеме.

Однако при дополнительном предположении, что  $\varphi(x, t) \equiv 0$ , существует одна и только одна схема вида (10), аппроксимирующая задачу (4) со вторым порядком относительно  $h$ . Построим ее, одновременно убеждаясь в ее единственности. Заметим, что из тождества  $\partial u / \partial t - \partial u / \partial x = 0$  вытекает тождество  $\partial^2 u / \partial t^2 = \partial^2 u / \partial x^2$ . Действительно,

$$\frac{\partial^2 u}{\partial t^2} = \frac{\partial}{\partial t} \left( \frac{\partial u}{\partial t} \right) = \frac{\partial}{\partial t} \left( \frac{\partial u}{\partial x} \right) = \frac{\partial}{\partial x} \left( \frac{\partial u}{\partial t} \right) = \frac{\partial}{\partial x} \left( \frac{\partial u}{\partial x} \right).$$

Поэтому при выполнении условий (13), необходимых для аппроксимации первого порядка, первая строка правой части равенства (14) примет вид

$$\frac{1}{2} (a^1 r^2 + a_1 + a_{-1}) h^2 u_{xx} \Big|_{(mh, pr)} + O(a^1 r^3 h^3, a_1 h^3, a_{-1} h^3). \quad (15)$$

Для получения схемы второго порядка с учетом  $a^1 = 1/(rh)$  и  $\varphi(x, t) \equiv 0$  необходимо, чтобы

$$a^1 r^2 + a_1 + a_{-1} = 0. \quad (16)$$

Равенства (13), (16) образуют систему, которая имеет единственное решение

$$a^1 = \frac{1}{rh}, \quad a_0 = -\frac{1}{rh} + \frac{r}{h}, \quad a_{-1} = \frac{1-r}{2h}, \quad a_1 = -\frac{1+r}{2h}. \quad (17)$$

Второе слагаемое в (15) есть  $O(h^3)$ , так что схема (10) с коэффициентами (17) есть единственная схема второго порядка аппроксимации для задачи (4) при  $\varphi(x, t) \equiv 0$ . Соответствующая разностная схема имеет вид

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_{m-1}^p}{2h} - \frac{r}{2h} (u_{m+1}^p - 2u_m^p + u_{m-1}^p) = 0. \quad (18)$$

Методом неопределенных коэффициентов можно убедиться в единственности этих формул: с точностью до несущественного произвола есть только один набор коэффициентов  $a_{-1}$ ,  $a_0$ ,  $a_1$ , при котором для любой достаточно гладкой функции  $u(x, y)$  имеет место формула

$$\begin{aligned} \frac{\partial^2 u(x_m, y_n)}{\partial x^2} &= a_{-1} u(x_{m-1}, y_n) + a_0 u(x_m, y_n) + \\ &\quad + a_1 u(x_{m+1}, y_n) + O(\max(\Delta x_{m-1}, \Delta x_m)) \end{aligned}$$

с остаточным членом первого порядка малости относительно  $\max(\Delta x_{m-1}, \Delta x_m)$ . Формулы такого вида с остаточным членом второго порядка малости при  $\Delta x_m \neq \Delta x_{m-1}$  не существует.

Для более точной замены производной разностным отношением здесь необходимо привлечь более трех точек сетки.

Во всех рассмотренных до сих пор в этой главе примерах разностных схем  $L_h u^{(h)} = f^{(h)}$  оператор  $L_h$ , который отображает пространство  $U_h$  в пространство  $F_h$ , задается явными формулами. Но часто оказываются полезными разностные схемы, в которых оператор  $L_h$  описывается более сложно. В дальнейшем мы еще встретимся с задачами, где такие схемы возникают естественным образом.

Изложенные приемы построения разностных схем остаются применимыми и в случаях задач с переменными коэффициентами, нелинейных задач, сеток с переменным шагом. Например, если для замены производных, входящих в дифференциальное уравнение  $u_{xx} + u_{yy} = \varphi(x, y)$ , разностными отношениями используется неравномерная прямоугольная сетка, то для построения разностной схемы можно применить формулы

$$\frac{\partial^2 u}{\partial x^2} \Big|_{x_m, y_n} \approx \frac{\frac{u(x_{m+1}, y_n) - u(x_m, y_n)}{\Delta x_m} - \frac{u(x_m, y_n) - u(x_{m-1}, y_n)}{\Delta x_{m-1}}}{\frac{\Delta x_m + \Delta x_{m-1}}{2}},$$

$$\frac{\partial^2 u}{\partial y^2} \Big|_{x_m, y_n} \approx \frac{\frac{u(x_m, y_{n+1}) - u(x_m, y_n)}{\Delta y_n} - \frac{u(x_m, y_n) - u(x_m, y_{n-1})}{\Delta y_{n-1}}}{\frac{\Delta y_n + \Delta y_{n-1}}{2}}.$$

**3. Идея метода А.С. Холодова.** В этом методе даются некоторая точная постановка и способы решения следующей задачи. Рассматривается класс схем, имеющих заданный шаблон. Пусть поставлены два требования, которым должна удовлетворять искомая схема, причем одно из них считается основным, а второе — дополнительным. Эти требования могут быть несовместимыми.

Задача состоит в том, чтобы среди схем из заданного класса, удовлетворяющих основному требованию, найти ту (или те), для которой дополнительное требование нарушается «меньше», чем для всех других.

Для изложения идеи воспользуемся классом схем вида

$$u_m^{p+1} = \sum_{-j_{лев} \leq j \leq j_{прав}} b_j u_{m+j}^p, \quad (19)$$

где  $j_{лев} \geq 0$ ,  $j_{прав} \geq 0$  — целые числа. Схема вида (10) в случае  $\varphi(x, t) \equiv 0$  приводится к виду (19), если положить  $j_{лев} = 1$ ,  $j_{прав} = 1$ :

$$u_m^{p+1} = b_{-1} u_{m-1}^p + b_0 u_m^p + b_1 u_{m+1}^p,$$

где

$$b_{-1} = -\frac{a_{-1}}{a^1}, \quad b_0 = -\frac{a_0}{a^1}, \quad b_1 = -\frac{a_1}{a^1}.$$

Отождествим с каждой схемой вида (19) точку  $\mathbf{b} = (b_{-j_{\text{лев}}}, b_{-j_{\text{лев}}+1}, \dots, b_{j_{\text{прав}}})$  с координатами  $\{b_j\}$  в  $k$ -мерном линейном пространстве,  $k = j_{\text{лев}} + j_{\text{прав}} + 1$ .

Обозначим основное требование  $A_{\text{осн}}$ , а дополнительное —  $A_{\text{доп}}$ .

Через  $M_0$  обозначим совокупность точек в пространстве разностных схем, для которых выполнено условие  $A_{\text{осн}}$ , а через  $M_D$  — совокупность точек, удовлетворяющих условию  $A_{\text{доп}}$ .

А.С. Холодов вводит в пространстве разностных схем ту или иную норму, после чего за меру невыполнения условия  $A_{\text{доп}}$  для заданной схемы  $\mathbf{b} \in M_0$  принимает расстояние  $\text{dist}(\mathbf{b}, M_D) = \inf_{x \in M_D} \|\mathbf{b} - x\|$  от точки  $\mathbf{b} \in M_0$  до множества  $M_D$ . Сформулированная выше задача приобретает тогда точную постановку: среди разностных схем, удовлетворяющих условию  $A_{\text{осн}}$ , т. е. среди точек  $\mathbf{b}$ , принадлежащих множеству  $M_0$ , найти ту точку, которая наименее удалена от множества  $M_D$ . Эта задача решается средствами линейного программирования, оптимизации и другими.

Выбор нормы остается в руках исследователя и осуществляется с учетом содержательного смысла дополнительного условия и простоты отыскания нужной разностной схемы. Подробное изложение метода и приложения см. в [12].

**4. Схемы с пересчетом, или схемы «предиктор-корректор».** При построении разностных схем, аппроксимирующих нестационарные задачи, может быть использована та же идея, которая лежит в основе конструкции схем Рунге–Кутты для обыкновенных дифференциальных уравнений,— идея пересчета. Пересчет позволяет повысить порядок аппроксимации, получаемый по исходной схеме, не использующей пересчета. Кроме того, в случае квазилинейных дифференциальных уравнений пересчет дает дополнительную возможность получения так называемых дивергентных схем, о которых будет идти речь в гл. 10.

Напомним идею пересчета на примере простейшей из схем Рунге–Кутты численного решения задачи Коши для обыкновенного дифференциального уравнения

$$\frac{dy}{dt} = f(t, y), \quad y(0) = \psi, \quad 0 < t < T. \quad (20)$$

Если значение  $y_p$  в точке  $t_p = p\tau$  уже вычислено, то для вычисления  $y_{p+1}$  находим предварительно вспомогательную величину  $\tilde{y}_{p+1/2}$ , пользуясь простейшей схемой Эйлера (схема «предиктор»)

$$\frac{\tilde{y}_{p+1/2} - y_p}{\tau/2} = f(t_p, y_p), \quad (21)$$

а затем осуществляем пересчет по схеме «корректор»

$$\frac{y_{p+1} - y_p}{\tau} = f(t_{p+1/2}, \tilde{y}_{p+1/2}). \quad (22)$$

Вспомогательная величина  $\tilde{y}_{p+1/2}$ , найденная по схеме первого порядка точности, позволяет приближенно найти угловой коэффициент

интегральной кривой в середине отрезка  $[t_p, t_{p+1}]$  и получить  $y_{p+1}$  по формуле (22) с большей точностью, чем это было бы сделано по схеме Эйлера (21).

Согласно п. 2 в § 4 гл. 8, все соображения остаются в силе, если  $y$ ,  $y_p$ ,  $\tilde{y}_{p+1/2}$  будут конечномерными векторами, а  $f$  — вектор-функцией. Но можно пойти и дальше, а именно, считать  $y$ ,  $y_p$ ,  $\tilde{y}_{p+1/2}$  элементами функционального пространства, а  $f$  — операторами в этом пространстве.

Например, задачу Коши

$$\begin{aligned} \frac{\partial u}{\partial t} + A \frac{\partial u}{\partial x} &= 0, \quad -\infty < x < \infty, \quad 0 < t < T, \\ u(x, 0) &= \psi(x), \quad -\infty < x < \infty, \quad A = \text{const}, \end{aligned} \quad (23)$$

можно считать задачей вида (20), если положить  $y(t) = u(x, t)$ , так что при каждом  $t$  под  $y$  надо понимать функцию аргумента  $x$ , а под операцией  $f$  — оператор  $-A\partial/\partial x$ . Приведем пример разностной схемы с пересчетом для задачи (23).

Пример. Пусть сеточная функция  $u^p = \{u_m^p\}$  ( $m = 0, \pm 1, \dots$ ) при данном  $p$  уже вычислена. Найдем вспомогательную сеточную функцию  $\tilde{u}^{p+1/2} = \{u_m^{p+1/2}\}$  ( $m = 0, \pm 1, \dots$ ), отнесенную к моменту времени  $t_{p+1/2} = (p + 1/2)\tau$  и к точкам  $x_{m+1/2} = (m + 1/2)h$ , воспользовавшись следующей схемой первого порядка точности:

$$\frac{\tilde{u}_{m+1/2}^{p+1/2} - (u_{m+1}^p + u_m^p)/2}{\tau/2} + A \frac{u_{m+1}^p - u_m^p}{h} = 0, \quad m = 0, \pm 1, \dots. \quad (24)$$

Затем осуществим коррекцию и найдем  $u^{p+1}$  с помощью схемы

$$\frac{u_m^{p+1} - u_m^p}{\tau} + A \frac{\tilde{u}_{m+1/2}^{p+1/2} - \tilde{u}_{m-1/2}^{p+1/2}}{h} = 0, \quad m = 0, \pm 1, \dots. \quad (25)$$

Исключая  $\tilde{u}^{p+1/2}$  из уравнений (24), (25), получаем схему

$$\frac{u_m^{p+1} - u_m^p}{\tau} + A \frac{u_{m+1}^p - u_{m-1}^p}{2h} - A^2 \frac{\tau}{2} \cdot \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} = 0, \quad (26)$$

$$u_m^0 = \psi(mh), \quad m = 0, \pm 1, \dots; \quad p = 0, 1, \dots, [T/\tau] - 1.$$

Эта схема при  $A = -1$  совпадает со схемой (18). Общий случай  $A \neq 1$  несущественно отличается от разобранного. Схема (26), а значит, и схема с пересчетом (24), (25) имеют второй порядок аппроксимации по  $h$  ( $t/h = \text{const}$ ).

### Задачи

- Для задачи Коши

$$\begin{aligned} \frac{\partial u}{\partial t} - \left( \frac{\partial u}{\partial x} + \frac{\partial u}{\partial y} \right) &= \varphi(x, y, t), \quad -\infty < x, y < \infty, \quad 0 \leq t < T, \\ u(x, y, 0) &= \varphi(x, y), \quad -\infty < x, y < \infty, \end{aligned}$$

воспользоватьсяся сеткой  $x_m = mh$ ,  $y_n = nh$ ,  $t_p = p\tau$  и построить какую-либо аппроксимирующую ее разностную схему.

**2.** Для задачи о теплопроводности

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2}, \quad -\infty < x, y < \infty, \quad 0 \leq t \leq T, \quad (27)$$

$$u(x, 0) = \psi(x), \quad -\infty < x < \infty,$$

рассмотреть разностную схему

$$\frac{u_m^{p+1} - u_m^p}{\tau} = \sigma \frac{u_{m+1}^{p+1} - 2u_m^{p+1} + u_{m-1}^{p+1}}{h^2} + (1 - \sigma) \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2},$$

$$u_m^0 = \psi(mh),$$

где  $\sigma$  — параметр.

а) Показать, что при любом  $\sigma$  имеет место аппроксимация на гладком решении  $u(x, t)$  с порядком  $O(\tau + h^2)$ .

б) Подобрать  $\sigma$  так, чтобы аппроксимация была  $O(\tau^2 + h^2)$ .

в) Связав шаги сетки соотношением  $\tau/h^2 = r = \text{const}$ , подобрать затем  $\sigma$  так, чтобы получить аппроксимацию порядка  $h^4$ .

г) При  $\sigma = 0$  подобрать число  $r = \tau/h^2$  так, чтобы аппроксимация имела порядок  $h^4$ .

д)\* Можно ли за счет выбора  $\sigma$  при фиксированном  $\tau/h^2 = r$  добиться того, чтобы аппроксимация на любом гладком решении была порядка выше четвертого?

**3.** Для задачи о теплопроводности

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left[ a(x, t) \frac{\partial u}{\partial x} \right], \quad -\infty < x < \infty, \quad 0 < t < T,$$

$$u(x, 0) = \psi(x), \quad -\infty < x < \infty,$$

пользуясь сеткой  $x_m = mh$ ,  $t_p = p\tau$ , построить аппроксимирующую ее разностную схему.

**4.** Для нелинейной задачи о теплопроводности

$$\frac{\partial u}{\partial t} = \frac{\partial}{\partial x} \left[ a(u) \frac{\partial u}{\partial x} \right], \quad -\infty < x < \infty, \quad 0 \leq t \leq T,$$

$$u(x, 0) = \psi(x),$$

пользуясь сеткой  $x_m = mh$ ,  $t_p = p\tau$ , построить аппроксимирующую ее явную разностную схему. Выписать формулы для вычисления  $u^{(h)}$  по этой схеме.

5. Доказать, что схема с пересчетом для задачи (27), в которой значения решения  $\tilde{u}_m^{p+1/2}$  на промежуточном слое определяются по неявной схеме порядка аппроксимации  $O(\tau + h^2)$

$$\frac{\tilde{u}_m^{p+1/2} - u_m^p}{\tau/2} - \frac{\tilde{u}_{m+1}^{p+1/2} - 2\tilde{u}_m^{p+1/2} + \tilde{u}_{m-1}^{p+1/2}}{h^2} = 0, \quad m = 0, \pm 1, \dots,$$

а решение  $\{u_m^{p+1}\}$  — по схеме

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{\tilde{u}_{m+1}^{p+1/2} - 2\tilde{u}_m^{p+1/2} + \tilde{u}_{m-1}^{p+1/2}}{h^2} = 0,$$

$$u_m^0 = \psi(mh),$$

имеет порядок аппроксимации  $O(\tau^2 + h^2)$  на гладком решении  $u$ .

**6.** Изменим разностную схему (10), отвечающую коэффициентам (17), заменив в правой части  $\varphi(mh, p\tau)$  на

$$\varphi(mh, p\tau) + \frac{rh}{2}(\varphi_t + \varphi_x) \Big|_{\substack{x=mh \\ t=p\tau}}.$$

Показать, что при произвольной правой части  $\varphi(x, t)$ , имеющей ограниченные производные достаточно высокого порядка, полученная схема аппроксимирует задачу (4) с порядком  $h^2$ .

**7.** Записать схему (10) при  $\varphi(x, t) = 0$  в виде

$$u_m^{p+1} = b_{-1}u_{m-1}^p + b_0u_m^p + b_1u_{m+1}^p. \quad (10')$$

Исходную схему (10) будем называть *монотонной*, если  $b_j \geq 0$  ( $j = -1, 0, 1$ ).

Будем считать, что условие  $A_0 = A_{\text{осн}}$  состоит в том, чтобы схема (10) имела порядок аппроксимации не ниже первого и была монотонна.

Описать в трехмерном пространстве  $R$  множество  $M_0$  точек  $(b_{-1}, b_0, b_1)$ , которому соответствуют схемы, удовлетворяющие условию  $A_0$ .

Ответ. В случае  $0 < r \leq 1$ , где  $r = \tau/h$ , множество  $M_0$  есть отрезок с концами  $\left(\frac{1-r}{2}, 0, \frac{1+r}{2}\right)$ ,  $(0, 1-r, r)$ ; в частности, в случае  $r = 1$  оно состоит из одной точки  $(0, 0, 1)$ . В случае  $r > 1$  множество  $M_0$  пусто.

**8.** Для задачи Коши для уравнения колебаний струны

$$\begin{aligned} \frac{\partial^2 u}{\partial t^2} - \frac{\partial^2 u}{\partial x^2} &= \varphi(x, t), \quad -\infty < x < \infty, \quad 0 \leq t \leq T, \\ u(x, 0) &= \psi^{(0)}(x), \quad \frac{\partial u(x, 0)}{\partial t} = \psi^{(1)}(x), \quad -\infty < x < \infty, \end{aligned}$$

исследовать аппроксимацию, которой обладает на достаточно гладком решении  $u(x, t)$  разностная схема

$$L_h u^{(h)} = \begin{cases} \frac{u_m^{p+1} - 2u_m^p + u_m^{p-1}}{\tau^2} - \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} = \varphi(mh, p\tau), \\ u_m^0 = \psi^{(0)}(mh), \\ \frac{u_m^1 - u_m^0}{\tau} = \psi^{(1)}(mh). \end{cases}$$

За норму  $\|f^{(h)}\|_{F_h}$  принять максимум модулей всех компонент элемента

$$f^{(h)} = \begin{cases} \varphi_m^p, \\ \psi_m^{(0)}, \\ \psi_m^{(1)}. \end{cases}$$

Показать, что аппроксимация имеет первый порядок относительно  $h$ , если  $\tau = rh$  ( $r = \text{const}$ ).

Как следует задать значения  $\psi_m^{(1)}$ , используя заданные функции  $\varphi(x, y, t)$ ,  $\psi^{(0)}(x)$ ,  $\psi^{(1)}(x)$ , чтобы порядок аппроксимации оказался вторым?

9. Для задачи о распространении тепла на отрезке:

$$\frac{\partial u}{\partial t} - \frac{\partial^2 u}{\partial x^2} = \varphi(x, t), \quad 0 \leq x \leq 1, \quad 0 \leq t \leq T,$$

$$u(x, 0) = \psi^{(0)}(x), \quad 0 \leq x \leq 1,$$

$$u(0, t) = \psi^{(1)}(t), \quad 0 \leq t \leq 1,$$

$$u(1, t) = \psi^{(2)}(t), \quad 0 \leq t \leq 1,$$

рассмотреть разностную схему вида

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} &= \varphi(mh, p\tau), \\ m = 1, 2, \dots, M-1; \quad h = M^{-1}, \quad p = 0, 1, \dots, \left[\frac{T}{\tau}\right] - 1, \\ u_m^0 &= \psi^{(0)}(mh), \quad m = 0, 1, \dots, M, \\ u_0^p &= \psi^{(1)}(p\tau), \quad p = 0, 1, \dots, \left[\frac{T}{\tau}\right], \\ u_M^p &= \psi^{(2)}(p\tau), \quad p = 0, 1, \dots, \left[\frac{T}{\tau}\right]. \end{aligned} \quad (11)$$

За норму  $\|\cdot\|_{F_h}$  принять максимум абсолютных величин правых частей уравнений, составляющих в совокупности рассматриваемую разностную схему.

Шаги  $\tau, h$  считать связанными равенством  $\tau = rh^2$ .

Показать, что схема (11) обладает порядком аппроксимации  $h^2$ .

### § 3. Спектральный признак устойчивости разностной задачи Коши

Изложим широко применяемый способ Неймана исследования разностных задач с начальными данными. В этом параграфе ограничимся случаем разностной задачи Коши с постоянными коэффициентами, а в § 4 частично распространим результаты на случай переменных коэффициентов и на случай смешанных задач.

**1. Устойчивость по начальным данным.** Простейшим примером разностной задачи Коши может служить неоднократно встречавшаяся выше задача

$$L_h u^{(h)} \equiv \begin{cases} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} = \varphi_m^p, & p = 0, 1, \dots, \left[\frac{T}{\tau}\right] - 1, \\ u_m^0 = \psi_m, & m = 0, \pm 1, \dots. \end{cases} \quad (1)$$

Положив

$$f^{(h)} = \begin{cases} \varphi_m^p, & p = 0, 1, \dots, [T/\tau] - 1, \\ \psi_m, & m = 0, \pm 1, \dots, \end{cases}$$

запишем задачу (1) в форме

$$L_h u^{(h)} = f^{(h)}. \quad (2)$$

Определим нормы  $\|u^{(h)}\|_{U_h}$  и  $\|f^{(h)}\|_{F_h}$  соответственно равенствами

$$\|u^{(h)}\|_{U_h} = \sup_{m,p} |u_m^p|, \quad \|f^{(h)}\|_{F_h} = \sup_{m,p} |\varphi_m^p| + \sup_m |\psi_m|.$$

Тогда условие устойчивости задачи (2)

$$\|u^{(h)}\|_{U_h} \leq c \|f^{(h)}\|_{F_h} \quad (3)$$

примет вид

$$\max_m |u_m^p| \leq c \left[ \max_m |\psi_m| + \max_{m,k} |\varphi_m^k| \right], \quad p = 0, 1, \dots, [T/\tau], \quad (4)$$

где  $c$  не зависит от  $h$  (и от  $\tau = rh$ ,  $r = \text{const}$ ). Условие (4) должно выполняться при произвольных  $\{\psi_m\}$ ,  $\{\varphi_m^p\}$ . В частности, для устойчивости необходимо, чтобы оно выполнялось при произвольных  $\{\psi_m\}$  и  $\varphi_m^p \equiv 0$ , т. е. чтобы решение задачи

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} &= 0, \quad p = 0, 1, \dots, \left[ \frac{T}{\tau} \right] - 1, \\ u_m^0 &= \psi_m, \quad m = 0, \pm 1, \dots, \end{aligned} \quad (5)$$

удовлетворяло условию

$$\max_m |u_m^p| \leq c \max_m |\psi_m|, \quad p = 0, 1, \dots, [T/\tau], \quad (6)$$

при произвольной ограниченной функции  $\psi_m$ .

Свойство (6), необходимое для устойчивости (4) задачи (1), называют *устойчивостью задачи (1) относительно возмущения начальных данных*. Оно означает, что возмущение  $\psi_m$ , внесенное в начальные данные задачи (1), вызовет возмущение решения задачи (1), которое в силу (6) не более, чем в  $c$  раз, где  $c$  не зависит от  $h$ , превосходит возмущение начальных данных.

**2. Необходимое спектральное условие устойчивости.** Для устойчивости задачи Коши (1) по начальным данным необходимо, чтобы условие (6) выполнялось, в частности, если  $\psi_m$  есть какая-нибудь гармоника:

$$u_m^0 = \psi_m = e^{i\alpha m}, \quad m = 0, \pm 1, \dots, \quad (7)$$

где  $\alpha$  — вещественный параметр. Но решение задачи (5) при начальном условии (7) имеет вид

$$u_m^p = \lambda^p e^{i\alpha m}, \quad (8)$$

где  $\lambda = \lambda(\alpha)$  определяется путем подстановки выражения (8) в однородное разностное уравнение задачи (5):

$$\lambda(\alpha) = 1 - r + re^{i\alpha}, \quad r = \frac{\tau}{h} = \text{const}. \quad (9)$$

Для решения (8) справедливо равенство

$$\max_m |u_m^p| = |\lambda(\alpha)|^p \max_m |e^{i\alpha m}| = |\lambda(\alpha)|^p \max_m |u_m^0| = |\lambda(\alpha)|^p \max_m |\psi_m|.$$

Поэтому для выполнения условия (6) необходимо, чтобы при всех вещественных  $\alpha$  выполнялось неравенство

$$|\lambda(\alpha)|^p \leq c, \quad p = 0, 1, \dots, [T/\tau],$$

или

$$|\lambda(\alpha)| \leq 1 + c_1 \tau, \quad (10)$$

где  $c_1$  — некоторая постоянная, не зависящая от  $\alpha, \tau$ . Это и есть необходимое спектральное условие Неймана применительно к рассматриваемому примеру. Спектральным оно называется по следующей причине.

Существование решения вида (8) показывает, что гармоника  $\{e^{i\alpha m}\}$  является собственной функцией оператора перехода

$$u_m^{p+1} = (1 - r)u_m^p + ru_{m+1}^p, \quad m = 0, \pm 1, \dots,$$

который в силу разностного уравнения (5) ставит в соответствие сеточной функции  $\{u_m^p\}$  ( $m = 0, \pm 1, \dots$ ), определенной на слое  $t = t_p$ , сеточную функцию  $\{u_m^{p+1}\}$ , определенную на слое  $t = t_{p+1}$ . Число  $\lambda(\alpha) = 1 - r + re^{i\alpha}$  является соответствующим этой гармонике  $\{e^{i\alpha m}\}$  собственным числом оператора перехода. Линия, которую пробегает точка  $\lambda = \lambda(\alpha)$  на комплексной плоскости, когда  $\alpha$  пробегает вещественную ось, вся состоит из собственных значений и является спектром оператора перехода.

Таким образом, необходимое условие устойчивости (10) можно сформулировать так: спектр оператора перехода, соответствующего разностному уравнению задачи (5), должен лежать в круге радиуса  $1 + c_1 \tau$  на комплексной плоскости. В нашем примере спектр (9) не зависит от  $\tau$ . Поэтому условие (10) равносильно требованию, чтобы спектр  $\lambda(\alpha)$  лежал в единичном круге:

$$|\lambda(\alpha)| \leq 1. \quad (11)$$

Воспользуемся сформулированным признаком для анализа устойчивости задачи (1). Спектр (9) представляет собой окружность с центром в точке  $1 - r$  и радиусом  $r$  на комплексной плоскости. В случае  $r < 1$  эта окружность лежит внутри единичного круга (касаясь единичной окружности в точке  $\lambda = 1$ ), при  $r = 1$  совпадает с единичной окружностью, а при  $r > 1$  лежит вне единичного круга (касаясь его в точке  $\lambda = 1$ ) (рис. 20). Соответственно необходимое условие устойчивости (11) выполняется при  $r \leq 1$  и не выполняется при  $r > 1$ . В § 1, п. 3 мы исследовали рассматриваемую разностную схему и показали, что при  $r \leq 1$  она устойчива, а при  $r > 1$  неустойчива. Таким образом, необходимое условие устойчивости Неймана оказалось в данном случае достаточно чувствительным, чтобы в точности отделить случай устойчивости от случая неустойчивости.

В общем случае задачи Коши для разностных уравнений и систем разностных уравнений необходимый спектральный признак устойчиво-

сти Неймана состоит в том, что спектр  $\lambda = \lambda(\alpha, h)$  разностной задачи при всех достаточно малых  $h$  должен лежать в круге

$$|\lambda| \leq 1 + \varepsilon \quad (12)$$

на комплексной плоскости, как бы мало ни было заранее выбранное положительное число  $\varepsilon$ .

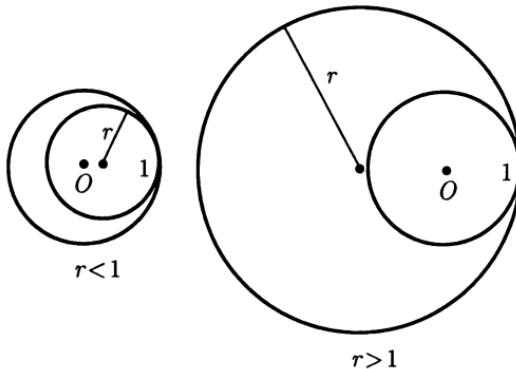


Рис. 20

Заметим, что если для рассматриваемой разностной задачи спектр окажется не зависящим от  $h$  (и от  $\tau$ ), то условие (12) равносильно требованию, чтобы спектр  $\lambda = \lambda(\alpha, h) = \lambda(\alpha)$  лежал в единичном круге (11).

Под *спектром* разностной задачи в условии (12) понимается совокупность всех  $\lambda = \lambda(\alpha, h)$ , при которых соответствующее однородное разностное уравнение (или система уравнений) имеет решение вида

$$u_m^p = \lambda^p (u^0 e^{i\alpha m}), \quad (13)$$

где  $u^0$  — число (единица), если речь идет о скалярном разностном уравнении, и числовая вектор, если речь идет о векторном разностном уравнении, т. е. о системе скалярных разностных уравнений.

Если необходимое условие Неймана (12) не выполняется, то ни при каком разумном выборе норм нельзя ожидать устойчивости, а в случае его выполнения можно надеяться, что при некотором разумном выборе норм устойчивость имеет место.

**3. Примеры.** Рассмотрим ряд интересных разностных задач Коши и применим для анализа их устойчивости спектральный признак Неймана. Начнем с разностных схем, аппроксимирующих дифференциальную задачу Коши

$$\frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} = \varphi(x, t), \quad -\infty < x < \infty, \quad 0 < t < T, \quad (14)$$

$$u(x, 0) = \psi(x), \quad -\infty < x < \infty.$$

Пример 1. Рассмотрим разностную схему

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_m^p - u_{m-1}^p}{h} = \varphi(mh, p\tau), \quad m = 0, \pm 1, \dots,$$

$$u_m^0 = \psi(mh), \quad p = 0, 1, \dots, [T/\tau] - 1.$$

Подставляя выражение вида (8) в соответствующее однородное разностное уравнение, после простых преобразований получаем

$$\lambda(\alpha) = 1 + r - re^{-i\alpha}.$$

В силу этой формулы спектр представляет собой окружность с центром в точке  $1 + r$  и радиусом  $r$  (рис. 21). Ни при каком  $r$  спектр не лежит в единичном круге. Условие устойчивости (12) всегда не выполняется. Этого можно было ожидать, так как при любом  $r = \tau/h$  не выполняется необходимое условие сходимости (и устойчивости) Куранта, Фридрихса и Леви.

Рис. 21

При  $r = \tau/h$  спектр лежит на единичной окружности, а не внутри ее. Поэтому условие устойчивости Куранта, Фридрихса и Леви не выполняется.

Пример 2. Рассмотрим разностную схему

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_{m-1}^p}{2h} - \frac{\tau}{2h^2} (u_{m-1}^p - 2u_m^p + u_{m+1}^p) = \varphi(mh, p\tau),$$

$$u_m^0 = \psi(mh), \quad (15)$$

аппроксимирующую задачу (14) при  $\varphi \equiv 0$  со вторым порядком относительно  $h$  (проверьте это). Для нее  $\lambda = \lambda(\alpha)$  определяется из уравнения

$$\frac{\lambda - 1}{\tau} - \frac{e^{i\alpha} - e^{-i\alpha}}{2h} - \frac{\tau}{2h^2} (e^{i\alpha} - 2 + e^{-i\alpha}) = 0.$$

Обозначим по-прежнему  $r = \tau/h$ . Заметив, что

$$\frac{e^{i\alpha} - e^{-i\alpha}}{2i} = \sin \alpha, \quad \frac{e^{i\alpha} - 2 + e^{-i\alpha}}{4} = -\left(\frac{e^{i\alpha/2} - e^{-i\alpha/2}}{2i}\right)^2 = -\sin^2 \frac{\alpha}{2},$$

получим

$$\lambda(\alpha) = 1 + ir \sin \alpha - 2r^2 \sin^2 \frac{\alpha}{2}, \quad (16)$$

$$|\lambda(\alpha)|^2 = \left(1 - 2r^2 \sin^2 \frac{\alpha}{2}\right)^2 + r^2 \sin^2 \alpha.$$

После простых преобразований найдем

$$1 - |\lambda|^2 = 4r^2 \sin^4 \frac{\alpha}{2} (1 - r^2). \quad (17)$$

Условие Неймана выполняется, если правая часть неотрицательна ( $r \leq 1$ ), и не выполняется при  $r > 1$ .

Пример 3. Рассмотрим разностную схему

$$L_h u^{(h)} = \begin{cases} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_{m-1}^p}{2h} = \varphi(mh, p\tau), \\ u_m^0 = \psi(mh) \end{cases} \quad (18)$$

для той же задачи Коши (14). Подставляя в уравнение (18) выражение (8), после сокращений получаем уравнение для  $\lambda$ :

$$\frac{\lambda - 1}{\tau} - \frac{e^{i\alpha} - e^{-i\alpha}}{2h} = 0,$$

или

$$\lambda(\alpha) = 1 + i\left(\frac{\tau}{h} \sin \alpha\right).$$

Спектр  $\lambda = \lambda(\alpha)$  есть вертикальный отрезок длины  $2\tau/h$ , проходящий через точку  $\lambda = 1$  (рис. 22).

Если  $\tau/h = r = \text{const}$ , то условие (10) не выполняется, спектр не лежит в нужном круге  $|\lambda| \leq 1 + ct$ . Отметим, что признак Куранта, Фридрихса и Леви позволяет утверждать (проверьте), что схема неустойчива, только при  $\tau/h = r > 1$ , а при  $\tau/h = r \leq 1$  суждений об устойчивости не дает и оказывается слабее признака Неймана.

Если при  $h \rightarrow 0$  шаг  $\tau$  изменяется как  $h^2$ , так что  $r = \tau/h^2$ , то самая далекая от точки  $\lambda = 0$  точка  $\lambda(\alpha)$  имеет модуль

$$|\lambda(\alpha)|_{\alpha=\pi/2} = \sqrt{1 + \left(\frac{\tau}{h}\right)^2} = \sqrt{1 + \tau r} < 1 + \frac{r}{2}\tau.$$

Условие  $|\lambda(\alpha)| \leq 1 + ct$  в этом случае выполняется при  $c = r/2$ .

Ясно, что требование  $\tau = rh^2$  является гораздо более жестким условием на убывание шага  $\tau$  при стремлении шага  $h$  к нулю, чем требование  $\tau = rh$ , которого было достаточно для выполнения признака Неймана для разностных схем (5), (15), аппроксимирующих ту же задачу Коши (14).

Рассмотрим теперь две разностные схемы, аппроксимирующие задачу Коши для уравнения теплопроводности

$$\begin{aligned} \frac{\partial u}{\partial t} - a^2 \frac{\partial^2 u}{\partial x^2} &= \varphi(x, t), \quad -\infty < x < \infty, \quad 0 \leq t \leq T, \\ u(x, 0) &= \psi(x), \quad -\infty < x < \infty. \end{aligned} \tag{19}$$

Пример 4. Явная разностная схема

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - a^2 \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} &= \varphi(mh, p\tau), \\ u_m^0 &= \psi(mh), \quad m = 0, \pm 1, \dots, \quad p = 0, 1, \dots, [T/\tau] - 1, \end{aligned}$$

позволяющая вычислить  $\{u_m^{p+1}\}$  в явном виде через  $\{u_m^p\}$ :

$$u_m^{p+1} = (1 - 2r)u_m^p + r(u_{m+1}^p + u_{m-1}^p) + \tau \varphi_m^p,$$

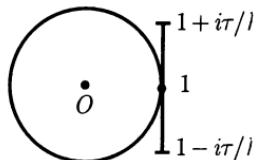


Рис. 22

при подстановке  $u^p = \lambda^p e^{i\alpha m}$  в соответствующее однородное разностное уравнение приводит к соотношению

$$\frac{\lambda - 1}{\tau} - a^2 \frac{e^{-i\alpha} - 2 + e^{i\alpha}}{h^2} = 0.$$

Заметив, что

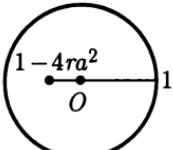
$$e^{i\alpha} - 2 + e^{-i\alpha} = -4 \sin^2 \frac{\alpha}{2},$$

получим

$$\lambda(\alpha) = 1 - 4ra^2 \sin^2 \frac{\alpha}{2}, \quad r = \frac{\tau}{h^2}.$$

При изменении  $\alpha$  число  $\lambda(\alpha)$  пробегает отрезок  $1 - 4ra^2 \leq \lambda \leq 1$  вещественной оси (рис. 23).

Для устойчивости необходимо, чтобы левый конец этого отрезка лежал в единичном круге  $1 - 4ra^2 \geq -1$ , или



$$r \leq \frac{1}{2a^2}. \quad (20)$$

В случае, если  $r > 1/(2a^2)$ , точка

$$\lambda(\alpha) = 1 - 4ra^2 \sin^2 \left( \frac{\alpha}{2} \right),$$

Рис. 23

отвечающая  $\alpha = \pi$ , лежит левее точки  $-1$ . Гармоника  $e^{i\pi m} = (-1)^m$  порождает решение

$$u_m^p = (1 - 4ra^2)^p (-1)^m,$$

не удовлетворяющее условию (6) ни при какой постоянной  $c$ .

Пример 5. Вторая схема — неявная:

$$L_h u^{(h)} = \begin{cases} \frac{u_m^{p+1} - u_m^p}{\tau} - a^2 \frac{u_{m+1}^{p+1} - 2u_m^{p+1} + u_{m-1}^{p+1}}{h^2} = \varphi(mh, p\tau), \\ u_m^0 = \psi(mh), \\ m = 0, \pm 1, \dots; \quad p = 0, 1, \dots, [T/\tau] - 1. \end{cases} \quad (21)$$

Здесь  $\{u_m^{p+1}\}$  не выражается через  $\{u_m^p\}$  явной формулой, так как в уравнение при данном  $m$  входит не одно, а три неизвестных:  $u_{m-1}^{p+1}$ ,  $u_m^{p+1}$ ,  $u_{m+1}^{p+1}$ . Поэтому схему (21) называют *неявной*.

Аналогично предыдущему примеру получаем выражение

$$\lambda(\alpha) = \frac{1}{1 + 4ra^2 \sin^2(\alpha/2)}, \quad r = \frac{\tau}{h^2}. \quad (22)$$

Спектр этой задачи есть отрезок

$$(1 + 4ra^2)^{-1} \leq \lambda \leq 1$$

вещественной оси, и условие  $|\lambda| < 1$  выполняется при любом  $r$ .

Спектральный признак Неймана применим для исследования разностной задачи Коши и в случае, если пространственных переменных две или более.

Пример 6. Для задачи

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}, \quad -\infty < x, y < \infty, \quad 0 < t < T,$$

$$u(x, y, 0) = \psi(x, y)$$

возьмем сетку  $(x_m, y_n, t_p) = (mh, nh, pt)$ . Заменяя производные разностными отношениями, построим разностную схему

$$L_h u^{(h)} = \begin{cases} \frac{u_{mn}^{p+1} - u_{mn}^p}{\tau} - \frac{u_{m+1,n}^p - 2u_{mn}^p + u_{m-1,n}^p}{h^2} + \\ \qquad \qquad \qquad + \frac{u_{m,n+1}^p - 2u_{mn}^p + u_{m,n-1}^p}{h^2} = 0, \\ u_{mn}^0 = \psi(mh, nh), \\ m, n = 0, \pm 1, \dots, \quad p = 0, 1, \dots, [T/\tau] - 1. \end{cases} \quad (23)$$

Задавая  $u_{mn}^0$  в виде  $e^{i(\alpha m + \beta n)}$ , т. е. в виде двумерной гармоники, зависящей от двух вещественных параметров  $\alpha, \beta$ , находим решение вида

$$u_{mn}^p = \lambda^p(\alpha, \beta) e^{i(\alpha m + \beta n)}.$$

Подставляя это выражение в разностное уравнение, после сокращений и тождественных преобразований находим

$$\lambda(\alpha, \beta) = 1 - 4r \left( \sin^2 \frac{\alpha}{2} + \sin^2 \frac{\beta}{2} \right).$$

При изменении вещественных  $\alpha, \beta$  точка  $\lambda = \lambda(\alpha, \beta)$  пробегает отрезок  $1 - 8r \leq \lambda \leq 1$  вещественной оси. Условие устойчивости выполняется, если  $1 - 8r \geq -1$ , или  $r \leq 1/4$ .

Приведем пример, иллюстрирующий применение признака Неймана для разностных уравнений, связывающих значения искомой функции не на двух, а на трех временных слоях.

Пример 7. Задачу Коши для волнового уравнения

$$\frac{\partial^2 u}{\partial t^2} - \frac{\partial^2 u}{\partial x^2} = 0, \quad -\infty < x < \infty, \quad 0 \leq t \leq T,$$

$$u(x, 0) = \psi^{(0)}(x), \quad \frac{\partial u(x, 0)}{\partial t} = \psi^{(1)}(x), \quad -\infty < x < \infty,$$

аппроксимируем разностной схемой

$$\frac{u_m^{p+1} - 2u_m^p + u_{m-1}^{p-1}}{\tau^2} - \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} = 0,$$

$$u_m^0 = \psi^{(0)}(mh), \quad \frac{u_m^1 - u_m^0}{\tau} = \psi^{(1)}(mh), \quad (24)$$

$$m = 0, \pm 1, \dots, \quad p = 0, 1, \dots, [T/\tau] - 1.$$

Подставляя в разностное уравнение решение вида (8), получаем после простых преобразований следующее квадратное уравнение для определения  $\lambda$ :

$$\lambda^2 - 2\left(1 - 2r^2 \sin^2 \frac{\alpha}{2}\right)\lambda + 1 = 0, \quad r = \frac{\tau}{h}.$$

Произведение корней этого уравнения равно единице. Если дискриминант

$$d(\alpha) = 4r^2 \sin^2 \frac{\alpha}{2} \left(r^2 \sin^2 \frac{\alpha}{2} - 1\right)$$

квадратного уравнения отрицателен, то корни  $\lambda_1(\alpha), \lambda_2(\alpha)$  комплексно сопряженные и равные единице по модулю. В случае  $r < 1$  дискриминант остается отрицательным при всех  $\alpha$ . На рис. 24, а изображен

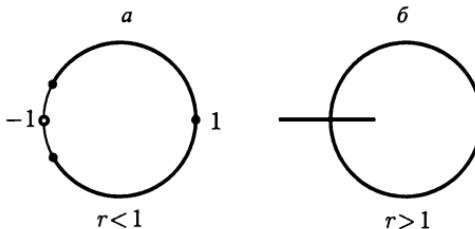


Рис. 24

спектр для этого случая. Он представляет собой часть единичной окружности. В случае  $r = 1$  спектр заполняет всю окружность. При  $r > 1$  по мере увеличения  $\alpha$  от 0 до  $\pi$  корни  $\lambda_1(\alpha), \lambda_2(\alpha)$  движутся из точки  $\lambda = 1$  по единичной окружности: один по часовой стрелке, а другой против часовой стрелки, пока не сольются в точке  $\lambda = -1$ , а затем один из корней пойдет по вещественной оси из точки  $\lambda = -1$  влево, а другой вправо, так как они вещественны и  $\lambda_1 \lambda_2 = 1$  (рис. 24, б). Условие устойчивости выполняется при  $r \leq 1$ .

Рассмотрим задачу Коши для следующей гиперболической системы дифференциальных уравнений, описывающей распространение звука:

$$\begin{aligned} \partial v / \partial t &= \partial w / \partial x, & \partial w / \partial t &= \partial v / \partial x, \\ -\infty < x < \infty, \quad 0 &\leq t \leq T, \\ v(x, 0) &= \psi^{(1)}(x), & w(x, 0) &= \psi^{(2)}(x). \end{aligned} \tag{25}$$

Положим

$$u(x, t) = \begin{bmatrix} v(x, t) \\ w(x, t) \end{bmatrix}, \quad \psi(x) = \begin{bmatrix} \psi^{(1)}(x) \\ \psi^{(2)}(x) \end{bmatrix}$$

и запишем (25) в векторной форме:

$$\begin{aligned} \frac{\partial u}{\partial t} - A \frac{\partial u}{\partial x} &= 0, & -\infty < x < \infty, \quad 0 \leq t \leq T, \\ u(x, 0) &= \psi(x), & -\infty < x < \infty, \end{aligned} \tag{25'}$$

где

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

Исследуем две разностные схемы, аппроксимирующие задачу (25').

Пример 8. Рассмотрим разностную схему

$$\frac{u_m^{p+1} - u_m^p}{\tau} - A \frac{u_{m+1}^p - u_m^p}{h} = 0, \quad u_m^0 = \psi(mh). \quad (26)$$

Ищем решение векторного однородного разностного уравнения в виде (13):

$$u_m^p = \lambda^p \begin{bmatrix} v^0 \\ w^0 \end{bmatrix} e^{i\alpha m}.$$

Подставляя это выражение в разностное уравнение (26), приходим к равенству

$$\frac{\lambda - 1}{\tau} u^0 - A \frac{e^{i\alpha} - 1}{h} u^0 = 0,$$

или

$$(\lambda - 1)u^0 - r(e^{i\alpha} - 1)Au^0, \quad r = \frac{\tau}{h}, \quad (27)$$

которое можно рассматривать как векторную запись системы линейных уравнений для определения компонент вектора  $u^0$ .

Запишем систему (27) в развернутой форме:

$$\begin{bmatrix} \lambda - 1 & -r(e^{i\alpha} - 1) \\ -r(e^{i\alpha} - 1) & \lambda - 1 \end{bmatrix} \begin{bmatrix} v^0 \\ w^0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}. \quad (28)$$

Эта система линейных уравнений имеет нетривиальное решение  $u^0 = \begin{bmatrix} v^0 \\ w^0 \end{bmatrix}$  лишь при тех  $\lambda = \lambda(\alpha)$ , при которых определитель системы (28) обращается в нуль:

$$(\lambda - 1)^2 = r^2(e^{i\alpha} - 1)^2.$$

Отсюда

$$\begin{aligned} \lambda_1(\alpha) &= 1 - r + re^{i\alpha}, \\ \lambda_2(\alpha) &= 1 + r - re^{i\alpha}. \end{aligned} \quad (29)$$

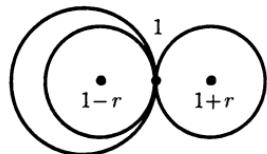


Рис. 25

Корни  $\lambda_1(\alpha)$ ,  $\lambda_2(\alpha)$  пробегают окружности радиуса  $r$  с центрами в точках  $1 - r$ ,  $1 + r$  соответственно (рис. 25). Условие устойчивости Неймана не выполняется ни при каком  $r$ .

Пример 9. Рассмотрим разностную схему

$$\frac{u_m^{p+1} - u_m^p}{\tau} - A \frac{u_{m+1}^p - u_{m-1}^p}{2h} - \frac{\tau}{2h^2} A^2(u_{m+1}^p - 2u_m^p + u_{m-1}^p) = 0,$$

$$u_m^0 = \psi(mh), \quad m = 0, \pm 1, \dots, \quad p = 0, 1, \dots, [T/\tau] - 1,$$

аппроксимирующую задачу (25') со вторым порядком и аналогичную схеме (15) для скалярного случая (14). Условие существования нетривиального решения вида (13) у векторного уравнения (25) состоит, как и в примере 8, в том, чтобы обращался в нуль определитель системы, получающейся для определения  $u^0 = \begin{bmatrix} v^0 \\ w^0 \end{bmatrix}$ .

Приравняв этот определитель нулю, получаем квадратное уравнение относительно  $\lambda = \lambda(\alpha)$ , из которого находим:

$$\begin{aligned}\lambda_1(\alpha) &= 1 + ir \sin \alpha - 2r^2 \sin^2 \frac{\alpha}{2}, \\ \lambda_2(\alpha) &= 1 - ir \sin \alpha - 2r^2 \sin^2 \frac{\alpha}{2}.\end{aligned}\quad (30)$$

Эти формулы аналогичны (16), и, как в (17), получим

$$1 - |\lambda_{1,2}(\alpha)|^2 = 4r^2 \sin^4 \frac{\alpha}{2} (1 - r^2).$$

Спектр, задаваемый формулами (30), лежит в единичном круге при  $r \leq 1$ .

### Задачи

1. При каких значениях параметра  $\sigma > 0$  разностная схема, аппроксимирующая задачу Коши для уравнения теплопроводности

$$\begin{aligned}\frac{u_m^{p+1} - u_m^p}{\tau} &= \sigma \frac{u_{m+1}^{p+1} - 2u_m^{p+1} + u_{m-1}^{p+1}}{h^2} + (1 - \sigma) \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2}, \\ u_m^0 &= \psi(mh), \quad m = 0, \pm 1, \dots, \quad p \geq 0,\end{aligned}$$

удовлетворяет спектральному признаку устойчивости Неймана при любом  $r = \tau/h^2 = \text{const}$ ?

2. Удовлетворяет ли спектральному признаку устойчивости следующая разностная схема, аппроксимирующая задачу Коши (19) для уравнения теплопроводности с порядком  $O(\tau^2 + h^2)$ :

$$\begin{aligned}\frac{u_m^{p+1} - u_m^{p-1}}{2\tau} - \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} &= \varphi(mh, p\tau), \\ u_m^0 &= \psi(mh), \quad m = 0, \pm 1, \dots, \\ u_m^1 &= \tilde{\psi}(mh), \quad p \geq 0,\end{aligned}$$

где

$$\begin{aligned}\tilde{\psi}(mh) &= u(x, 0) + \tau \left. \frac{\partial u(x, 0)}{\partial t} \right|_{x=mh} = u(x, 0) + \tau \frac{\partial^2 u(x, 0)}{\partial x^2} = \\ &= \psi(mh) + \tau \psi''(mh)?\end{aligned}$$

## § 4. Принцип замороженных коэффициентов

Здесь мы изложим прием, существенно расширяющий класс нестационарных разностных задач, для исследования которых можно пользоваться спектральным признаком устойчивости. Этот необходимый признак устойчивости, изложенный в § 3 для исследования разностной задачи Коши с постоянными коэффициентами, можно применять и в случае разностной задачи Коши с непрерывными, но не постоянными коэффициентами, а также для задач в ограниченных областях, когда граничные условия задаются не только при  $t = 0$ , но и на боковых

границах. Этим приемом можно пользоваться и для исследования нелинейных задач.

**1. Замораживание коэффициентов во внутренних точках.** Сформулируем принцип замороженных коэффициентов, воспользовавшись в качестве примера следующей разностной краевой задачей:

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - a(x_m, t_p) \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} &= 0, \\ u_m^0 = \psi(mh), \quad m &= 1, \dots, M-1, \\ l_1 u_0^{p+1} = 0, \quad l_2 u_M^{p+1} = 0, \quad Mh &= 1, \quad p \geq 0. \end{aligned} \quad (1)$$

В этой разностной краевой задаче  $l_1 u_0^{p+1} = 0$ ,  $l_2 u_M^{p+1} = 0$  — некоторые условия, задаваемые соответственно на левой и правой границах сеточного отрезка  $\{0, h, 2h, \dots, Mh\}$ .

Выберем произвольную внутреннюю точку  $(\tilde{x}, \tilde{t})$  области  $0 < x < 1$ ,  $t > 0$ , где рассматривается задача (1), и «заморозим» коэффициенты задачи (1) в этой точке. Возникающее разностное уравнение с постоянными коэффициентами

$$\frac{u_m^{p+1} - u_m^p}{\tau} - a(\tilde{x}, \tilde{t}) \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} = 0, \quad p \geq 0, \quad (2)$$

будем рассматривать теперь не при  $m = 1, 2, \dots, M-1$ , а при всех целочисленных  $m$ .

Сформулируем теперь следующий принцип замороженных коэффициентов. Для устойчивости задачи (1) необходимо, чтобы задача Коши для разностного уравнения с постоянными коэффициентами (2) удовлетворяла спектральному признаку устойчивости Неймана.

В обоснование принципа замороженных коэффициентов приведем следующее рассуждение на эвристическом уровне строгости. При измельчении сетки коэффициент  $a(x, t)$  в окрестности точки  $(\tilde{x}, \tilde{t})$  за любое фиксированное число шагов сетки длины  $h$  по пространству и длины  $\tau$  по времени в силу непрерывности функции  $a(x, t)$  меняется все меньше и все меньше отличается от значения  $a(\tilde{x}, \tilde{t})$ . Добавим к этому, что расстояние от точки  $(\tilde{x}, \tilde{t})$  до границ  $t = 0$ ,  $x = 0$  и  $x = 1$  отрезка, измеренное числом шагов сетки, стремится к бесконечности. Поэтому при мелкой сетке возмущения, наложенные на решение задачи (1) в окрестности точки  $(x = \tilde{x}, t = \tilde{t})$ , развиваются (за малое время) примерно так же, как для задачи (2).

Понятно, что это рассуждение носит общий характер. Оно не зависит от числа пространственных переменных, числа искомых функций и вида разностного уравнения или системы уравнений.

В § 3 мы рассмотрели задачу Коши для уравнения вида (2) и нашли, что для выполнения условия Неймана отношение  $\tau/h^2$  шагов сетки должно удовлетворять условию

$$r \leq \frac{1}{2a(\tilde{x}, \tilde{t})}. \quad (3)$$

Поскольку в силу принципа замороженных коэффициентов для устойчивости задачи (1) это условие должно выполняться при любых  $\tilde{x}$ ,  $\tilde{t}$ , отношение  $\tau/h^2 = r$  шагов сетки должно быть подчинено условию

$$r \leq \frac{1}{2 \max_{x,t} a(x, t)}.$$

Принцип замороженных коэффициентов позволяет ориентироваться на эвристическом уровне строгости и при исследовании устойчивости нелинейных задач. Поясним это на следующей нелинейной задаче:

$$\begin{aligned} u_t &= (1 + u^2)u_{xx} = 0, \quad 0 < x < 1, \quad 0 < t < T, \\ u(x, 0) &= \psi(x), \quad 0 < x < 1, \\ u(0, t) &= \psi^{(1)}(t), \quad u(1, t) = \psi^{(2)}(t). \end{aligned}$$

Используем разностную схему

$$L_h u^{(h)} = \begin{cases} \frac{u_m^{p+1} - u_m^p}{\tau_p} - [1 + (u_m^p)^2] \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} = 0, \\ p = 0, 1, \dots, [T/\tau] - 1, \quad 1 \leq m \leq M - 1, \\ u_m^0 = \psi(mh), \quad 0 \leq m \leq M, \\ u_0^p = \psi^{(1)}(t_p), \\ u_M^p = \psi^{(2)}(t_p), \quad p = 0, 1, \dots, [T/\tau]. \end{cases}$$

В ней допускается изменение шага  $\tau_p$  от слоя к слою. Эта схема позволяет последовательно, слой за слоем, вычислить  $u_m^1$  ( $m = 0, 1, \dots, M$ ), затем  $u_m^2$  ( $m = 0, 1, \dots$ ) и т. д.

Допустим, что мы уже добрались до слоя  $t = t_p$ ,  $p < [T/\tau]$ , вычислили  $u_m^p$  ( $m = 0, 1, \dots, M$ ) и хотим продолжать счет. Как выбрать следующий шаг  $\tau = \tau_p$ ? Можно считать, что нам предстоит найти решение линейного разностного уравнения

$$\frac{u_m^{p+1} - u_m^p}{\tau_p} - a(x_m, t_p) \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} = 0$$

с заданным переменным коэффициентом  $a(x_m, t_p) \equiv 1 + (u_m^p)^2$ . Действительно, естественно считать, что значения  $u_m^p$  близки к значениям  $u(mh, t_p)$  гладкого решения  $u(x, t)$  дифференциальной задачи. Коэффициент тогда близок к непрерывной функции  $a(x, t) = 1 + u^2(x, t)$ , которая на протяжении нескольких временных шагов мало изменяется.

Применение признака Неймана к уравнению с переменным коэффициентом  $a(x_m, t_p)$  дает ограничение (3) на соотношение шагов сетки, необходимое для устойчивости:

$$\frac{\tau_p}{h^2} = r_p \leq \frac{1}{2 \max_{x=mh} a(x, t_p)} = \frac{1}{2 \max_m (1 + |u_m^p|^2)}.$$

Отсюда следует рекомендация выбрать очередной шаг  $\tau_p$  из условия

$$\tau_p \leq \frac{h^2}{2 \max_m (1 + |u_m^p|^2)}.$$

Численный эксперимент на компьютере подтверждает правильность этих эвристических рассуждений.

Если необходимое условие устойчивости, полученное путем рассмотрения задачи Коши с замороженными в произвольной точке области коэффициентами, окажется нарушенным, то устойчивости нельзя ожидать ни при каком задании граничных условий. Подчеркнем, однако, что принцип замороженных коэффициентов не учитывает влияния граничных условий. В случае выполнения необходимого условия устойчивости, вытекающего из принципа замороженных коэффициентов, устойчивость может иметь место при одних граничных условиях и не иметь места при других. Теперь изложим признак Бабенко–Гельфандса, учитывающий влияние границ в случае задачи на отрезке.

**2. Признак Бабенко–Гельфандса.** При рассмотрении задачи (1) мы полагали, что возмущения, приданые ее решению в окрестности произвольной внутренней точки  $(\tilde{x}, \tilde{t})$ , при мелкой сетке развиваются примерно так же, как такие же возмущения, приданые решению задачи Коши (2) с замороженными в точке  $(\tilde{x}, \tilde{t})$  коэффициентами. В обоснование этого принципа мы принимали во внимание, что расстояния от внутренней точки  $(\tilde{x}, \tilde{t})$  до границ, измеренные числом шагов сетки, при измельчении сетки неограниченно возрастают. Но если точка  $(\tilde{x}, \tilde{t})$  лежит на боковой границе  $x = 0$  или  $x = 1$ , то это эвристическое рассуждение теряет убедительность. Пусть, например,  $\tilde{x} = 0$ . Тогда расстояние от точки  $(\tilde{x}, \tilde{t})$  до любой фиксированной точки  $(x, \tilde{t})$ , где  $x > 0$  (в частности, до точки  $(1, \tilde{t})$ , лежащей на правой границе), измеренное числом шагов сетки, при  $h \rightarrow 0$  по-прежнему неограниченно возрастает, но число шагов до точки  $(0, \tilde{t})$  на левой границе не меняется и остается равным нулю.

Поэтому возмущение решения задачи (1) вблизи левой границы  $x = 0$  за малое время должно развиваться подобно возмущению решения задачи

$$\frac{u_m^{p+1} - u_m^p}{\tau} - a(0, \tilde{t}) \frac{u_{m+1}^{p+1} - 2u_m^p + u_{m-1}^p}{h^2} = 0, \quad m = 1, 2, \dots, \quad (4)$$

$$l_1 u_0^{p+1} = 0.$$

Эта задача получилась из исходной задачи (1) при замораживании коэффициента  $a(x, t)$  в левом конце отрезка  $0 \leq x \leq 1$  и одновременном удалении правой границы в плюс бесконечность. Задачу (4) естественно рассматривать только на тех функциях  $u^p = \{u_0^p, u_1^p, \dots\}$ , для которых

$$u_m^p \rightarrow 0, \quad m \rightarrow +\infty.$$

Только в этом случае возмущение сосредоточено вблизи границы  $x = 0$ , и только относительно возмущений такого вида задача (1) и задача (4) вблизи левой границы  $x = 0$  сходны между собой.

Точно так же развитие возмущений решения задачи (1) вблизи правой границы  $x = 1$  должно быть похоже на развитие таких же возмущений для задачи

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - a(1, \tilde{t}) \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} &= 0, \\ m = \dots, -1, 0, 1, \dots, M-1, \quad p \geq 0, \\ l_2 u_M^{p+1} &= 0 \end{aligned} \quad (5)$$

с одной только правой границей. Эта задача возникла из исходной задачи (1) при замораживании коэффициента  $a(x, t)$  в правом конце  $x = 1$  и при удалении левой границы в минус бесконечность. Задачу (5) надо рассматривать на сеточных функциях  $u^p = \{ \dots, u_{-1}^p, u_0^p, u_1^p, \dots, u_M^p \}$ , удовлетворяющих условию  $u_m^p \rightarrow 0$  при  $m \rightarrow -\infty$ .

Задачи (2), (4), (5) проще исходной задачи (1) в том смысле, что при фиксированном  $r = \tau/h^2$  они не зависят от  $h$  и являются задачами с постоянными коэффициентами.

Таким образом, процедура исследования устойчивости, учитывающая влияние границ, применительно к задаче (1) состоит в следующем. Надо составить вспомогательные задачи (2), (4), (5). Для каждой из этих трех задач, не зависящих от  $h$ , надо найти все те числа  $\lambda$  (собственные числа оператора перехода от  $u^p$  к  $u^{p+1}$ ), при которых существуют решения вида

$$u_m^p = \lambda^p u_m^0.$$

При этом в случае задачи (2) функция  $u^0 = \{u_m^0\}$  ( $m = 0, \pm 1, \dots$ ) должна быть ограничена. В случае задачи (4) для  $u^0 = \{u_m^0\}$  ( $m \geq 0$ ) должно быть  $u_m^0 \rightarrow 0$  ( $m \rightarrow +\infty$ ), а в случае задачи (5) должно быть  $u^0 = \{u_m\}$  ( $m \leq M$ ),  $u_m^0 \rightarrow 0$  ( $m \rightarrow -\infty$ ).

Для устойчивости задачи (1) совокупность собственных чисел каждой из трех задач (2), (4), (5) должна лежать в единичном круге  $|\lambda| \leq 1$ . (Задача (2) рассматривается при любом фиксированном  $\tilde{x}$  ( $0 < \tilde{x} < 1$ ).)

Продолжим рассмотрение задачи (1). Будем считать в дальнейшем, что  $a(x, t) \equiv 1$ , и вычислим спектры для всех трех задач (2), (4), (5) при различных краевых условиях  $l_1 u_0^{p+1} = 0$ ,  $l_2 u_M^{p+1} = 0$ .

Подставляя решение  $u_m^p = \lambda^p u_m^0$  в разностное уравнение (2), получаем

$$(\lambda - 1)u_m - r(u_{m+1} - 2u_m + u_{m-1}) = 0, \quad r = \frac{\tau}{h^2},$$

или

$$u_{m+1} - \frac{-2r + 1 - \lambda}{r} u_m + u_{m-1} = 0. \quad (6)$$

Это — обыкновенное разностное уравнение второго порядка.

Чтобы написать общее решение уравнения (6), составим характеристическое уравнение

$$q^2 - \frac{-2r + 1 - \lambda}{r} q + 1 = 0. \quad (7)$$

Если  $q$  — корень этого уравнения, то сеточная функция

$$u_m = \lambda^p q^m$$

есть одно из решений уравнения

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} = 0.$$

Если  $|q| = 1$ , т. е.  $q = e^{i\alpha}$ , то ограниченная при  $m \rightarrow +\infty$  и при  $m \rightarrow -\infty$  сеточная функция

$$u_m^p = \lambda^p e^{i\alpha m},$$

как мы видели в § 3, является решением при

$$\lambda = 1 - 4r \sin^2 \left( \frac{\alpha}{2} \right), \quad 0 \leq \alpha < 2\pi.$$

Эти  $\lambda = \lambda(\alpha)$  заполняют отрезок  $1 - 4r \leq \lambda \leq 1$  на вещественной оси. Этот отрезок и есть спектр задачи (2). Собственных значений  $\lambda$ , не лежащих на этом отрезке, задача (2) не имеет, так как в случае отсутствия у характеристического уравнения (7) корня  $q$ , по модулю равного единице, уравнение (6) не имеет ограниченного при  $m \rightarrow \pm\infty$  решения.

Если  $\lambda$  не лежит на отрезке  $1 - 4r \leq \lambda \leq 1$ , то оба корня характеристического уравнения (7) отличны по модулю от единицы, но их произведение равно свободному члену квадратного уравнения (7), т. е. единице. Поэтому среди корней уравнения (7) один по модулю больше, а другой меньше единицы. Пусть для определенности  $|q_1(\lambda)| < 1$ ,  $|q_2(\lambda)| > 1$ . Тогда общее решение уравнения (6), убывающее по модулю при  $m \rightarrow +\infty$ , имеет вид

$$u_m = cq_1^m, \quad q_1 = q_1(\lambda) < 1,$$

а общее решение уравнения (6), стремящееся к нулю при  $m \rightarrow -\infty$ , имеет вид

$$u_m = cq_2^m, \quad q_2 = q_2(\lambda) > 1.$$

Для определения собственных значений задачи (4) надо подставить  $u_m = cq_1^m$  в левое граничное условие  $l_1 u = 0$  и найти те  $\lambda$ , при которых оно выполняется. Это и будут все собственные значения задачи (4). Если, например,  $l_1 u_0 \equiv u_0 = 0$ , то условие  $cq_1^0 = 0$  не выполняется ни при каком  $c \neq 0$ , так что собственных значений нет. Если  $l_1 u_0 = u_1 - u_0 = 0$ , то условие  $c(q_1 - q_1^0) = c(q_1 - 1) = 0$  в силу  $q_1 \neq 1$  приводит к  $c = 0$ , так что собственных значений опять нет. В случае  $l_1 u_0 = 2u_1 - u_0 = 0$  условие  $c(2q_1 - q_1^0) = 0$  выполняется при  $c \neq 0$ , если  $q_1 = 1/2 < 1$ .

Из уравнения (7) находим, что при  $q_1 = 1/2$  число  $\lambda$  есть

$$\lambda = 1 + r \left( q_1 - 2 + \frac{1}{q_1} \right) = 1 + r \frac{1 - 4 + 4}{2} = 1 + \frac{r}{2}.$$

Это и есть единственное собственное значение задачи (4). Оно лежит вне единичного круга, и устойчивости нет. Аналогично вычисляются собственные значения задачи (5). Они получаются из уравнения

$$l_2 u_M = 0$$

при

$$u_m = q_2^m, \quad q_2 = q_2(\lambda), \quad m = M, M-1, \dots$$

Рассмотрим в качестве еще одного примера разностную схему

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} &= 0, \quad p = 0, 1, \dots, \left[ \frac{T}{\tau} \right] - 1, \\ m &= 0, 1, \dots, M-1, \quad Mh = 1, \\ u_m^0 &= \psi(mh), \quad u_M^{p+1} = 0, \end{aligned} \tag{8}$$

которая аппроксимирует задачу

$$\begin{aligned} \frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} &= 0, \quad 0 \leq x \leq 1, \quad 0 < t < T, \\ u(x, 0) &= \psi(x), \quad u(1, t) = 0. \end{aligned}$$

Применим для исследования ее устойчивости признак Бабенко-Гельфанда. Сопоставим схеме (8) три задачи: задачу без боковых границ

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} = 0, \quad m = 0, \pm 1, \dots, \tag{9}$$

задачу с одной только левой границей

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} = 0, \quad m = 0, 1, \dots, \tag{10}$$

и задачу с одной только правой границей

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_m^p}{h} &= 0, \quad m = M-1, M-2, \dots, \\ u_M^{p+1} &= 0. \end{aligned} \tag{11}$$

В случае задачи (10) — с одной только левой боковой границей — граничного условия нет, так как его не было в исходной задаче (8).

Надо найти совокупность собственных чисел всех трех операторов перехода от  $u^p$  к  $u^{p+1}$ , соответствующих каждой из трех вспомогательных задач (9)–(11), и выяснить, при каких условиях все они лежат в круге  $|\lambda| \leq 1$ .

Решение вида

$$u_m^p = \lambda^p u_m^0$$

при подстановке в разностное уравнение

$$u_m^{p+1} = (1 - r)u_m^p + ru_{m+1}^p, \quad r = \frac{\tau}{h},$$

приводит к следующему обыкновенному разностному уравнению первого порядка для собственной функции:

$$(\lambda - 1 + r)u_m - ru_{m+1} = 0. \quad (12)$$

Соответствующее характеристическое уравнение

$$\lambda - 1 + r - rq = 0 \quad (13)$$

дает связь между  $\lambda$  и  $q$ . Общее решение уравнения (12) есть

$$u_m = cq^m = c \left( \frac{\lambda - 1 + r}{r} \right)^m, \quad m = 0, \pm 1, \dots$$

При  $|q| = 1$ ,  $q = e^{i\alpha}$  ( $0 \leq \alpha \leq 2\pi$ )

$$\lambda = 1 - r + re^{i\alpha}.$$

Точка  $\lambda = \lambda(\alpha)$  пробегает окружность с центром в точке  $1 - r$  и радиусом  $r$ . Это и есть множество собственных значений задачи (9) (рис. 26, а).

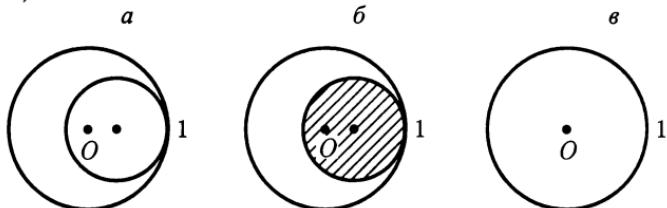


Рис. 26

Убывающее при  $m \rightarrow +\infty$  нетривиальное решение задачи (10) существует при любом  $q$  ( $|q| < 1$ ). Соответствующие  $\lambda = 1 - r + rq$  заполняют, очевидно, всю внутренность круга, ограниченного окружностью  $\lambda = 1 - r + re^{i\alpha}$  (рис. 26, б).

Наконец, решения задачи (11)  $u_m^p = \lambda^p u_m$ , убывающие при  $m \rightarrow -\infty$ , должны иметь вид

$$u_m^p = c\lambda^p q^m, \quad |q| > 1,$$

где  $\lambda$ ,  $q$  связаны равенством (13).

Из граничного условия  $u_m^p = 0$  следует, что нетривиальное решение существует только при  $\lambda = \lambda(q) = 0$ , т.е. при  $q = (r - 1)/r$ . Эта величина  $q$  по модулю больше единицы в случае выполнения одного из неравенств, а именно,  $(r - 1)/r > 1$  или  $(r - 1)/r < -1$ . Первое неравенство решений не имеет. Решение второго имеет вид  $r < 1/2$ . Итак, при  $r < 1/2$  задача (11) имеет собственное значение  $\lambda = 0$  (рис. 26, в).

На рис. 27 изображены объединения собственных значений всех трех задач соответственно для случаев  $r < 1/2$ ,  $1/2 < r < 1$ ,  $r > 1$ . Ясно, что объединение собственных значений всех трех задач лежит в круге  $|\lambda| \leq 1 + c\tau$ , где  $c$  не зависит от  $\tau$ , в том и только том случае, если  $r \leq 1$ .

Изложенный здесь признак устойчивости нестационарных разностных задач на отрезке, учитывающий влияние граничных условий, применим и в случае краевых задач на отрезке для систем разностных уравнений. В этом случае естественные на первый взгляд схемы, удовлетворяющие признаку Неймана, часто оказываются неустойчивыми из-за неудачной аппроксимации граничных условий, и важно уметь подбирать схемы, свободные от этого недостатка.

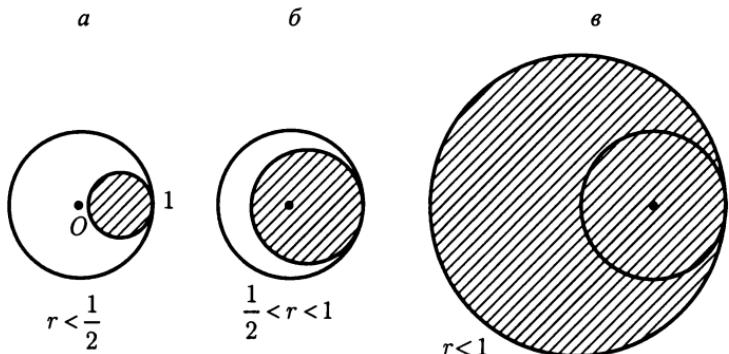


Рис. 27

В [5, 6] спектральный признак Бабенко–Гельфанда обсуждается с более общей точки зрения — на основе введенного там понятия *спектра семейства операторов*. В частности, С.К. Годуновым и В.С. Рябеньким строго показано, что его выполнение необходимо для устойчивости и что в случае его выполнения устойчивость не может грубо нарушаться.

### Задачи

1. Выяснить условия выполнения спектрального признака устойчивости для разностной схемы

$$\frac{u_m^{p+1} - u_m^p}{\tau} - \frac{u_{m+1}^p - u_{m-1}^p}{2h} - \frac{\tau}{2} \cdot \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} = 0, \quad (14)$$

$$p \geq 0, \quad m = 1, 2, \dots, M-1, \quad Mh = 1,$$

$$u_m^0 = \psi(mh), \quad m = 0, 1, \dots, M, \quad (15)$$

$$u_M^{p+1} = 0, \quad p = 0, 1, \dots, [T/\tau] - 1, \quad (16)$$

$$\frac{u_0^{p+1} - u_0^p}{\tau} - \frac{u_1^p - u_0^p}{h} = 0, \quad (17)$$

аппроксимирующей дифференциальную задачу

$$\frac{\partial u}{\partial t} - \frac{\partial u}{\partial x} = 0, \quad 0 \leq x \leq 1, \quad 0 \leq t \leq T, \quad (18)$$

$$u(x, 0) = \varphi(x), \quad u(1, t) = 0 \quad (19)$$

на гладком решении  $u(x, t)$  со вторым порядком относительно  $h$ .

Ответ.  $\tau/h \leq 1$ .

2. Для построения разностной схемы, аппроксимирующей следующую краевую задачу для гиперболической системы дифференциальных уравнений:

$$\begin{aligned}\frac{\partial v}{\partial t} &= \frac{\partial w}{\partial x}, \quad \frac{\partial w}{\partial t} = \frac{\partial v}{\partial x}, \quad 0 \leq x \leq 1, \quad 0 \leq t \leq T, \\ v(x, 0) &= \psi^{(1)}(x), \\ w(x, 0) &= \psi^{(2)}(x), \\ v(0, t) &= w(1, t) = 0,\end{aligned}$$

положим  $u(x, t) = \begin{bmatrix} v(x, t) \\ w(x, t) \end{bmatrix}$ ,  $\psi(x) = \begin{bmatrix} \psi^{(1)}(x) \\ \psi^{(2)}(x) \end{bmatrix}$  и запишем ее в матричной форме:

$$\begin{aligned}\frac{\partial u}{\partial t} - A \frac{\partial u}{\partial x} &= 0, \\ u(x, 0) &= \psi(x), \\ v(0, t) &= w(1, t) = 0,\end{aligned}$$

где  $A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$ . Выберем сетку  $(x_m, t_p) = (mh, p\tau)$ ,  $h = M^{-1}$ ,  $M$  — натуральное число. Зададим разностные уравнения:

$$\begin{aligned}\frac{u_m^{p+1} - u_m^p}{\tau} - A \frac{u_{m+1}^{p+1} - u_{m-1}^p}{2h} - \frac{A^2 \tau}{2} \cdot \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} &= 0, \\ m &= 1, 2, \dots, M-1, \quad p \geq 0, \\ u_m^0 &= \psi(mh), \quad v_0^{p+1} = w_M^{p+1} = 0.\end{aligned}$$

Для завершения построения схемы надо задать дополнительные граничные условия на левой и правой боковых границах. Заметив, что при любых  $\alpha, \beta$  из системы дифференциальных уравнений следуют равенства

$$\begin{aligned}\frac{\partial(v + \alpha w)}{\partial t} - \frac{\partial(w + \alpha v)}{\partial x} \Big|_{x=0} &= 0, \\ \frac{\partial(v + \beta w)}{\partial t} - \frac{\partial(w + \beta v)}{\partial x} \Big|_{x=1} &= 0,\end{aligned}$$

зададим дополнительные разностные краевые условия, положив

$$\begin{aligned}\frac{(v_0^{p+1} + \alpha w_0^{p+1}) - (v_0^p + \alpha w_0^p)}{\tau} - \frac{(w_1^p + \alpha v_1^p) - (w_0^p + \alpha v_0^p)}{h} &= 0, \\ \frac{(v_M^{p+1} + \beta w_M^{p+1}) - (v_M^p + \beta w_M^p)}{\tau} - \frac{(w_M^p + \beta v_M^p) - (w_{M-1}^p - \beta v_{M-1}^p)}{h} &= 0.\end{aligned}$$

При условии  $r = \tau/h \leq 1$  показать, что:

а) если  $\alpha = 1, \beta = -1$ , то спектральный признак устойчивости выполняется;

б) если  $\alpha = -1$ , то независимо от выбора числа  $\beta$  спектральный признак устойчивости не выполняется;

в) найти условия, которым должны удовлетворять  $\alpha, \beta$ , чтобы выполнялся спектральный признак устойчивости, учитывающий влияние граничных условий.

## § 5. Явные и неявные разностные схемы для уравнения теплопроводности

Рассмотрим следующую задачу для модельного уравнения теплопроводности:

$$\begin{aligned} \frac{\partial u}{\partial t} - a(x, t) \frac{\partial^2 u}{\partial x^2} &= 0, \quad a(x, t) > 0, \quad 0 \leq x \leq 1, \\ u(x, 0) &= \psi(x), \\ u|_{x=0} &= \varphi_{\text{лев}}(t), \quad u|_{x=1} = \varphi_{\text{прав}}(t). \end{aligned} \quad (1)$$

Для численного решения этой задачи рассмотрим и сравним две разностные схемы — явную и неявную. При этом выяснится, что в некоторых случаях неявная схема обладает преимуществами перед явной, несмотря на то, что сам алгоритм вычислений по явной схеме проще. Это преимущество связано с ее устойчивостью при любом соотношении шагов сетки по времени и по пространству.

**1. Явная разностная схема.** Воспользуемся следующим разностным аналогом (см. рис. 18):

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau_p} - a(x_m, t_p) \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2} &= 0, \\ u_m^0 &= \psi_m, \quad m = 0, 1, \dots, M, \quad Mh = 1, \\ u_0^{p+1} &= \psi_{\text{лев}}(t_{p+1}), \quad u_M^{p+1} = \varphi_{\text{прав}}(t_{p+1}), \\ t_0 &= 0, \quad t_p = \tau_0 + \tau_1 + \dots + \tau_{p-1}, \quad p = 1, 2, \dots. \end{aligned} \quad (2)$$

Пусть решение  $u_m^k$  ( $m = 0, 1, \dots, M$ ) уже найдено при  $k = 0, 1, \dots, p$ . Тогда вычисление  $u_m^{p+1}$  на временному слое  $t = t_{p+1} = t_p + \tau_p$  осуществляется в силу (2) по явной формуле

$$\begin{aligned} u_m^{p+1} &= u_m^p + \frac{\tau_p}{h^2} a(x_m, t_p) (u_{m+1}^p - 2u_m^p + u_{m-1}^p), \\ m &= 1, 2, \dots, M-1. \end{aligned} \quad (3)$$

Поэтому разностная схема (2) называется *явной разностной схемой*.

В силу принципа замороженных коэффициентов можно ожидать устойчивости лишь в том случае, если

$$\tau_p \leq \frac{h^2}{2 \max_m a(x_m, t_p)}. \quad (4)$$

Отсюда следует, что в случае больших значений  $a(x_m, t_p)$  в окрестности какой-нибудь точки  $(\tilde{x}, \tilde{t})$  для вычисления решения на слое  $t = t_{p+1}$  придется воспользоваться очень маленьким значением  $\tau$ , так что продвижение до заданного  $t = T$  может потребовать очень большого числа шагов по времени и оказаться практически слишком трудоемким.

Заметим, что это не связано с физическим смыслом задачи. В самом деле, большое значение коэффициента теплопроводности  $a(x_m, t_p)$  в окрестности точки  $(\tilde{x}, \tilde{t})$  физически означает просто, что эту окрестность можно «вырезать» из теплопроводного стержня, не изменив картину распространения тепла: можно считать, что этот кусочек стержня сделан из материала с нулевой теплоемкостью.

**2. Неявная разностная схема.** Вместо схемы (2) воспользуемся неявной разностной схемой

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} - a(x_m, t_p) \frac{u_{m+1}^{p+1} - 2u_m^{p+1} + u_{m-1}^{p+1}}{h^2} = 0, \\ m = 1, 2, \dots, M-1, \\ u_m^0 = \psi_m, \quad m = 0, 1, \dots, M, \\ u_0^{p+1} = \psi_{лев}(t_{p+1}), \quad u_M^{p+1} = \psi_{прав}(t_{p+1}), \quad p \geq 0. \end{aligned} \quad (5)$$

Пусть решение на слое  $t = t_p$  уже найдено, так что значения  $u_m^p$  ( $m = 0, 1, \dots, M$ ) уже вычислены. Для вычисления решения  $u_m^{p+1}$  ( $m = 0, 1, \dots, M$ ) надо решить систему (5) вида

$$\begin{aligned} u_0 = \varphi_{лев}(t_{p+1}), \\ \alpha_m u_{m-1} + \beta_m u_m + \gamma_m u_{m+1} = f_m, \quad m = 1, 2, \dots, M-1, \\ u_M = \psi_{прав}(t_{p+1}), \end{aligned}$$

где

$$\begin{aligned} \alpha_m = \gamma_m = -\frac{a(x_m, t_p)}{h^2} \tau, \\ \beta_m = \frac{2a(x_m, t_p)}{h^2} \tau + 1, \quad f_m = u_m^p. \end{aligned}$$

Очевидно, что

$$|\beta_m| > |\alpha_m| + |\gamma_m| + \delta, \quad \delta = 1 > 0,$$

так что выполнены условия применимости метода прогонки.

В § 3 (см. пример 5) мы рассмотрели неявную разностную схему для уравнения теплопроводности с постоянным коэффициентом и установили, что спектральный признак устойчивости выполняется при любом отношении  $\tau/h^2 = r$  шагов сетки по времени и пространству. В силу принципа замороженных коэффициентов спектральный признак не накладывает ограничений на шаги сетки и в случае переменного коэффициента теплопроводности. Это делает неявную схему пригодной и в тех случаях, когда  $a(x, t)$  в отдельных местах принимает очень большие значения.

В заключение отметим, что условие устойчивости (4) для явной схемы (2) и абсолютная устойчивость схемы (5) установлены нами на эвристическом уровне строгости, поскольку мы опирались на принцип замороженных коэффициентов, носящий эвристический характер. Однако наши выводы о свойствах явной и неявной разностных схем (2), (5) действительно справедливы и могут быть строго доказаны.

**Задача**

Пусть коэффициент теплопроводности в задаче (1) задается формулой  $a = 1 + u^2$ , так что задача (1) становится нелинейной.

а) Предложить явную и неявную разностные схемы для этой задачи.

б) Рассмотреть явную схему

$$\frac{u_m^{p+1} - u_m^p}{\tau_p} = [1 + (u_m^p)^2] \frac{u_{m+1}^p - 2u_m^p + u_{m-1}^p}{h^2},$$

$$m = 1, 2, \dots, M - 1,$$

$$u_0^p = \varphi_{лев}^p, \quad u_m^0 = \psi_m, \quad u_M^p = \varphi_{прав}^p, \quad p = 0, 1, \dots$$

Как выбрать  $\tau_p$  по значениям  $u_m^p$  решения на очередном слое  $p$ ?

в) Рассмотреть неявную разностную схему на основе следующего разностного уравнения:

$$\frac{u_m^{p+1} - u_m^p}{\tau} = [1 + (u_m^{p+1})^2] \frac{u_{m+1}^{p+1} - 2u_m^{p+1} + u_{m-1}^{p+1}}{h^2}.$$

Как следует усовершенствовать это уравнение, чтобы при переходе от  $u_m^p$  ( $m = 0, 1, \dots, M$ ) к  $u_m^{p+1}$  ( $m = 0, 1, \dots, M$ ) можно было воспользоваться методом прогонки?

**ГЛАВА 10****ПОНЯТИЕ О РАЗРЫВНЫХ РЕШЕНИЯХ  
И СПОСОБАХ ИХ ВЫЧИСЛЕНИЯ**

Такие дифференциальные уравнения математической физики, как уравнения газовой динамики, упругости, теплопроводности и др., являются следствиями физических законов сохранения величин (массы, импульса, энергии и т. п.) [22], записываемых в некоторой интегральной форме. Эти дифференциальные уравнения равносильны исходным интегральным соотношениям только в том случае, если искомые физические поля (т. е. искомые функции) дифференцируемы.

Во всех рассмотренных до сих пор примерах мы предполагали, что существуют достаточно гладкие решения дифференциальных краевых задач, а в основу построения разностных схем клади приближенную замену производных в дифференциальном уравнении разностными отношениями. Однако дифференцируемых функций недостаточно для описания многих важных процессов физики. Так, например, физические эксперименты показывают, что распределения давления, плотности и скорости в сверхзвуковом течении невязкого газа описываются функциями, имеющими скачки. Скачки (разрывы) могут возникать с течением времени при гладких начальных данных.

Интегральная форма записи законов сохранения имеет смысл и для разрывных решений, так как разрывные функции можно интегрировать. Однако теряется возможность перейти к равносильной дифференциальной постановке задачи, так как разрывные функции нельзя дифференцировать. Вместе с тем теряются удобства, которыми мы пользовались при исследовании свойств решений дифференциальных краевых задач и при построении разностных схем. Поэтому прежде чем переходить к построению алгоритмов вычисления разрывных решений задач, сформулированных в терминах интегральных законов сохранения, полезно так обобщить понятие решения дифференциальной краевой задачи, чтобы она сохраняла смысл и оставалась равносильной исходному интегральному закону сохранения даже в случае разрывных решений.

Мы изложим способ перехода от постановки задачи в терминах интегрального закона сохранения к равносильной дифференциальной краевой задаче, имеющей лишь обобщенное (разрывное) решение, а также способы вычисления обобщенного решения на примере следующей задачи.

В полосе  $0 \leq t \leq T$  найти функцию  $u(x, t)$ , для которой выполняется равенство вида

$$\int_{\Gamma} \frac{u^k}{k} dx - \frac{u^{k+1}}{k+1} dt \equiv 0, \quad (\text{I})$$

где  $\Gamma$  — произвольный замкнутый контур,  $k$  — некоторое фиксированное натуральное число, и которая удовлетворяет начальному условию

$$u(x, 0) = \psi(x), \quad -\infty < x < \infty. \quad (\text{II})$$

Левую часть уравнения (I) можно понимать как поток вектора

$$\Phi(x, t) = \begin{bmatrix} u^k/k \\ u^{k+1}/(k+1) \end{bmatrix}$$

через контур  $\Gamma$ . Требование, чтобы поток этого вектора через произвольный контур  $\Gamma$  был равен нулю, будем понимать как некоторый закон сохранения, записанный в интегральной форме.

Задача (I), (II) является простейшей среди имеющихся разрывные решения при гладких начальных данных. Поэтому она может служить модельной для понимания способов решения подобных задач в газовой динамике.

## § 1. Дифференциальная формулировка интегрального закона сохранения

**1. Дифференциальное уравнение в случае гладких решений.** Предположим сначала, что задача (I), (II) имеет решение  $u(x, t)$ , непрерывное вместе со своими производными всюду в полосе  $0 \leq t \leq T$ .

Оказывается, что в таком случае она равносильна следующей задаче Коши:

$$\begin{aligned} \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} &= 0, \quad 0 \leq t \leq T, \quad -\infty < x < \infty, \\ u(x, 0) &= \psi(x). \end{aligned} \quad (1)$$

Уравнение (1) называется *уравнением Хопфа*.

Докажем сформулированное утверждение о равносильности. Для этого напомним формулу Остроградского–Грина. Пусть  $D$  – произвольная область с границей  $\Gamma$  на плоскости  $Oxt$ , и пусть функции  $\Phi_1(x, t)$ ,  $\Phi_2(x, t)$  имеют в области  $D$  непрерывные вплоть до границы  $\Gamma$  частные производные. Тогда справедлива следующая *формула Остроградского–Грина*:

$$\iint_D \left( \frac{\partial \Phi_1}{\partial t} + \frac{\partial \Phi_2}{\partial x} \right) dx dt = \int_{\Gamma} \Phi_1 dx - \Phi_2 dt. \quad (2)$$

Это тождество означает, что интеграл от дивергенции  $\partial \Phi_1 / \partial t + \partial \Phi_2 / \partial x$  векторного поля  $\Phi = \begin{bmatrix} \Phi_1 \\ \Phi_2 \end{bmatrix}$  по области  $D$  равен потоку вектора  $\Phi$  через границу  $\Gamma$  этой области. Применив эту формулу к (I), можно написать

$$\int_{\Gamma} \frac{u^k}{k} dx - \frac{u^{k+1}}{k+1} dt = \iint_D \left[ \frac{\partial}{\partial t} \left( \frac{u^k}{k} \right) + \frac{\partial}{\partial x} \left( \frac{u^{k+1}}{k+1} \right) \right] dx dt. \quad (3)$$

Отсюда видно, что если гладкая функция  $u(x, t)$  удовлетворяет уравнению (1), то выполнено и равенство (I). В самом деле, если выполнено (1), то также

$$u^{k-1} \left( \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} \right) \equiv \frac{\partial}{\partial t} \left( \frac{u^k}{k} \right) + \frac{\partial}{\partial x} \left( \frac{u^{k+1}}{k+1} \right) = 0. \quad (4)$$

Поэтому правая часть (3) обращается в нуль.

Справедливо и обратное: если гладкая функция  $u(x, t)$  удовлетворяет интегральному закону сохранения (I), то в каждой точке  $(x_0, t_0)$  полосы  $0 \leq t \leq T$  выполняется уравнение (4), а вместе с тем и (1). Для доказательства предположим противное, и пусть для определенности в некоторой точке  $(x_0, t_0)$

$$\left. \frac{\partial}{\partial t} \left( \frac{u^k}{k} \right) + \frac{\partial}{\partial x} \left( \frac{u^{k+1}}{k+1} \right) \right|_{\substack{x=x_0 \\ t=t_0}} > 0.$$

Тогда в силу непрерывности можно найти столь малый круг  $D$  с границей  $\Gamma$  и с центром в точке  $(x_0, t_0)$ , в котором всюду

$$\frac{\partial}{\partial t} \left( \frac{u^k}{k} \right) + \frac{\partial}{\partial x} \left( \frac{u^{k+1}}{k+1} \right) > 0.$$

Получим

$$0 = \int_{\Gamma} \frac{u^k}{k} dx - \frac{u^{k+1}}{k+1} dt = \iint_D \left[ \frac{\partial}{\partial t} \left( \frac{u^k}{k} \right) + \frac{\partial}{\partial x} \left( \frac{u^{k+1}}{k+1} \right) \right] dx dt > 0.$$

Возникшее противоречие  $0 > 0$  доказывает сделанное утверждение: для гладких функций  $u(x, t)$  задача (I), (II) и задача Коши (1) равносильны.

**2. Механизм возникновения разрывов.** Предположим сначала, что задача (1) имеет гладкое решение  $u(x, t)$ . Рассмотрим на плоскости  $Oxt$  линии  $x = x(t)$  — графики решения дифференциального уравнения

$$\frac{dx}{dt} = u(x, t). \quad (5)$$

Эти линии называются *характеристиками* уравнения  $u_t + uu_x = 0$ .

Вдоль каждой характеристики  $x = x(t)$  решение  $u(x, t)$  можно считать функцией, зависящей только от  $t$ :

$$u(x, t) = u(x(t), t).$$

Тогда

$$\frac{du}{dt} = \frac{\partial u}{\partial t} + \frac{\partial u}{\partial x} \cdot \frac{dx}{dt} = \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = 0.$$

Поэтому вдоль характеристики решение постоянно:  $u(x, t) = C$ . Но в силу уравнения (5) из этого следует, что характеристика есть прямая линия

$$x = Ct + C_1.$$

Пусть  $x_0 \in (-\infty, \infty)$  — абсцисса точки на оси  $Ox$ , через которую проходит характеристика. Подставив  $t = 0$  в уравнение прямой, получаем  $C_1 = x_0$ , а начальное условие при  $x = x_0$  даст  $C = u(x_0, 0) = \psi(x_0)$ , так что уравнение характеристики имеет вид  $x = \psi(x_0)t + x_0$ , где  $\psi(x_0)$  — угловой коэффициент ее наклона к оси  $Ot$ . Заданием начальной функции  $u(x, 0) = \psi(x)$ , таким образом, наглядно определяются и картина характеристик, и значения решения  $u(x, t)$  в каждой точке полуплоскости  $t > 0$  (рис. 28).

Заметим сразу же, что в предположении существования гладкого решения  $u(x, t)$  характеристики не могут пересекаться, так как каждая характеристика приносила бы в точку пересечения свое значение решения, так что решение не было бы однозначной функцией. При монотонно возрастающей функции  $\psi(x)$  с ростом  $x_0$  угол  $\alpha$  увеличивается, характеристики не пересекаются (рис. 29). Но в случае убывания функции  $\psi(x)$  характеристики

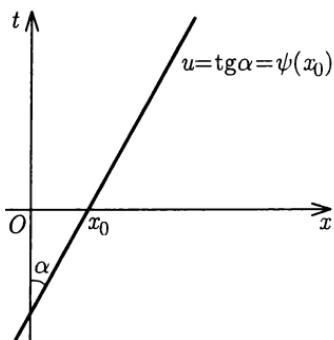


Рис. 28

сходятся и независимо от гладкости функции  $u(x, t)$  пересечения неизбежны. Гладкое решение задачи (1) перестает существовать с момента  $t = \tilde{t}$ , когда хотя бы две характеристики пересекутся (рис. 30); графики функции  $u(x, t) = u$  при  $t = 0$ ,  $t = \tilde{t}/2$ ,  $t = \tilde{t}$  изображены на рис. 31.

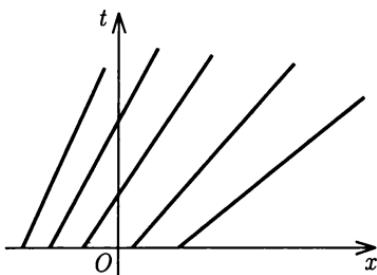


Рис. 29

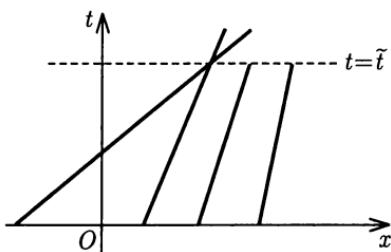


Рис. 30

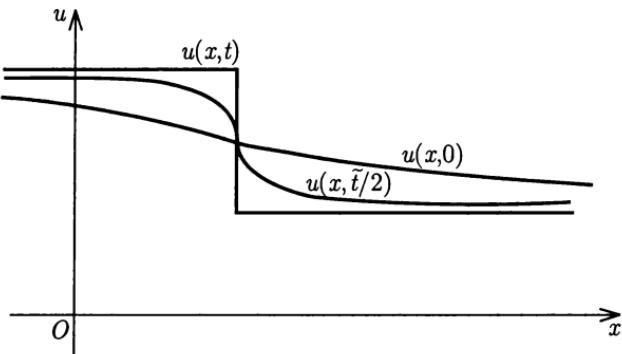


Рис. 31

**3. Условие на линии разрыва решения.** Пусть внутри области, где отыскивается решение, имеется линия  $x = x(t)$ , на которой решение  $u(x, t)$  задачи (I), (II) терпит разрыв первого рода. Пусть при приближении к этой линии слева или справа получаем на ней соответственно  $u_{\text{лев}}(x, t)$ ,  $u_{\text{прав}}(x, t)$ . Оказывается, что значения  $u_{\text{лев}}(x, t)$ ,  $u_{\text{прав}}(x, t)$  и скорость движения точки разрыва  $\dot{x} = dx/dt$  не могут быть произвольными: они связаны между собой некоторым соотношением.

Пусть  $L$  — линия разрыва (рис. 32). Интеграл

$$\int_{\Gamma} \frac{u^k}{k} dx - \frac{u^{k+1}}{k+1} dt$$

по контуру  $ABCDA$ , как и по любому другому контуру, обращается в нуль. Когда отрезки  $BC$  и  $DA$  стя-

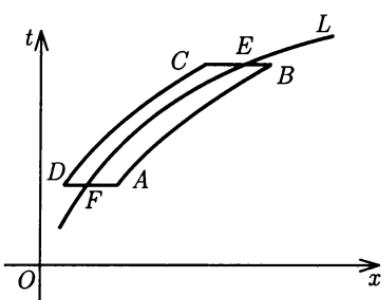


Рис. 32

гиваются к точкам  $E$  и  $F$  соответственно, интегралы по ним обращаются в нуль и получается равенство

$$\int_E^F \left[ \frac{u^k}{k} \right] dx - \left[ \frac{u^{k+1}}{k+1} \right] dt = 0,$$

или в силу произвольности контура  $ABCDA$ ,

$$\int_{L'} \left( \left[ \frac{u^k}{k} \right] \frac{dx}{dt} - \left[ \frac{u^{k+1}}{k+1} \right] \right) dt = 0,$$

где  $[z] = z_{\text{прав}} - z_{\text{лев}}$  — скачок величины  $z$  на линии разрыва, а  $L'$  — произвольный участок этой линии.

В силу произвольности участка  $L'$  подынтегральная функция в каждой точке должна обращаться в нуль:

$$\left[ \frac{u^k}{k} \right] \frac{dx}{dt} - \left[ \frac{u^{k+1}}{k+1} \right] = 0.$$

Тогда

$$\frac{dx}{dt} = \left[ \frac{u^{k+1}}{k+1} \right] / \left[ \frac{u^k}{k} \right]. \quad (6)$$

Отсюда при различных  $k$  на линии разрыва получаются различные условия. В частности, в случае  $k = 1$

$$\frac{dx}{dt} = \frac{u_{\text{лев}} + u_{\text{прав}}}{2}, \quad (7)$$

а в случае  $k = 2$

$$\frac{dx}{dt} = \frac{2}{3} \cdot \frac{u_{\text{лев}}^2 + u_{\text{лев}}u_{\text{прав}} + u_{\text{прав}}^2}{u_{\text{лев}} + u_{\text{прав}}}.$$

Отсюда видно, что условия на разрыве, которым должно удовлетворять разрывное решение задачи (I), (II), зависят от  $k$ .

**4. Обобщенное (разрывное) решение дифференциальной краевой задачи.** Введем понятие *обобщенного решения* задачи Коши (1), отождествив это решение с решением задачи (I), (II).

В случае решения, имеющего, как мы видели, всюду непрерывные производные, это решение не зависит от  $k$  и совпадает с обычным решением задачи Коши (1); решение в этом случае — это дифференцируемая функция, обращающая уравнение  $\partial u / \partial t + \partial u / \partial x = 0$  в тождество и удовлетворяющая начальному условию  $u(x, 0) = \psi(x)$ . В плодотворности рассмотрения наряду с задачей (I), (II) равносильной ей задачи Коши (1) мы убедились, занимаясь выяснением механизма возникновения разрывов. Этот механизм трудно было бы выяснить, оставаясь в рамках постановки задачи (I), (II) и не используя задачу Коши (1).

В случае разрывного решения приведенное нами определение обобщенного решения задачи (1) пока ничем не обогащает постановку

задачи (I), (II), являясь лишь ее переименованием. Поэтому придадим определению обобщенного решения задачи (1), приведенному выше, другую форму.

При этом мы будем рассматривать лишь те решения задачи (I), (II), которые в полосе  $0 \leq t \leq T$  имеют непрерывные первые производные всюду, кроме некоторых гладких кривых  $x = x(t)$ , на которых они могут претерпевать разрывы первого рода (скакки).

Будем называть  $u(x, t)$  обобщенным решением задачи Коши (1), соответствующим интегральному закону сохранения (I), если:

а) функция  $u(x, t)$  удовлетворяет дифференциальному уравнению (1) в каждой точке полосы  $0 \leq t \leq T$ , не лежащей на линиях разрыва;

б) на линиях разрыва выполняется условие (6);

в) для каждого  $x$ , при котором  $\psi(x)$  непрерывна, функция  $u(x, t)$  непрерывна в точке  $(x, 0)$  и удовлетворяет начальному условию  $u(x, 0) = \psi(x)$ .

Мы не будем доказывать равносильность двух определений обобщенного решения, приведенных выше, предоставляя это читателю в качестве упражнения.

Подчеркнем, что обобщенное решение задачи Коши (1) задается в случае разрывного решения не только равенствами (1), но и указанием того, какому интегральному закону сохранения это обобщенное решение соответствует. В нашей постановке задачи (I), (II) нет оснований предпочесть какое-либо значение  $k$ .

В задачах математической физики интегральный закон сохранения, аналогичный условию (I), описывает некоторое конкретное физическое явление и вполне определен.

В дальнейших рассмотрениях мы будем считать для определенности, что выполняются интегральный закон сохранения (I), отвечающий

значению  $k = 1$ , и вытекающее из него условие (7) на разрыве.

**5. Распад произвольного разрыва.** Пусть заданы разрывные начальные данные

$$u = \begin{cases} 2, & x < 0, \\ 1, & x > 0. \end{cases}$$

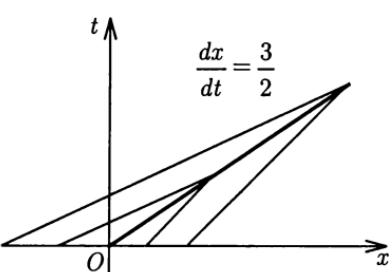


Рис. 33

Построенное по этим начальным данным решение изображено на рис. 33. Тангенс угла наклона линии

разрыва  $dx/dt = (2+1)/2 = 3/2$  является средним арифметическим из тангенсов углов наклона характеристик по обе стороны от нее.

Зададим теперь в начальных условиях другой разрыв:

$$u = \begin{cases} 1/2, & x < 0, \\ 3/2, & x > 0. \end{cases}$$

На рис. 34 показаны два возможных способа построения решения. При первом способе (рис. 34, а) мы получаем непрерывное решение,

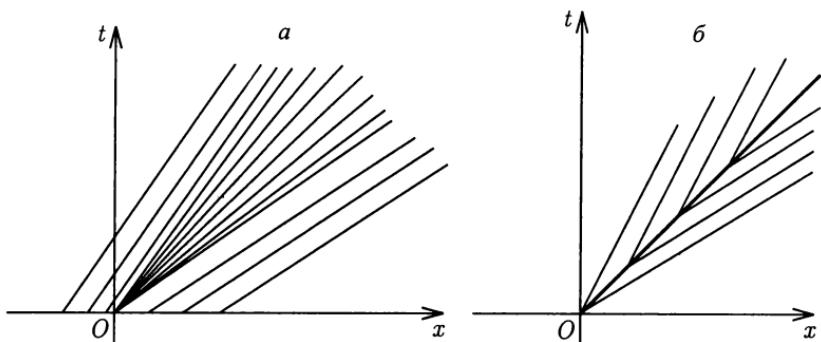


Рис. 34

а при втором (рис. 34, б) — разрывное решение при  $t > 0$ . Следует предпочтеть непрерывное решение. В пользу этого говорит следующее рассуждение. Если несколько изменить начальные данные, задав их формулой

$$u = \begin{cases} 1/2, & x \leq 0, \\ 3/2, & x \geq \varepsilon, \\ 1 + x/\varepsilon, & 0 \leq x \leq \varepsilon, \end{cases}$$

то решение  $u$  определится однозначно. Оно изображено на рис. 35. При стремлении  $\varepsilon$  к нулю это решение переходит в непрерывное решение, изображенное на рис. 34, а.

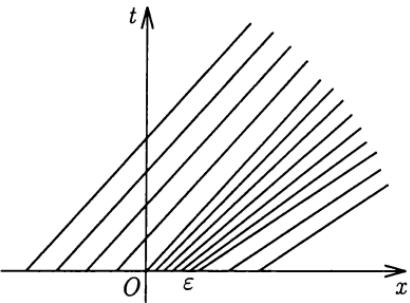


Рис. 35

Запрет решения, изображенного на рис. 34, б, по причине его неустойчивости относительно возмущения начальных данных аналогичен запрету ударных волн разрежения при математическом описании течения идеального газа.

### Задача\*

Рассмотрим вспомогательную задачу

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = \mu \frac{\partial^2 u}{\partial x^2}, \quad u(x, 0) = \psi(x), \quad -\infty < x < \infty, \quad (8)$$

где  $\mu > 0$  — некоторый параметр, который аналогичен вязкости в случае задач газовой динамики. Уравнение (8) параболического типа и имеет гладкое решение при любой гладкой функции  $\psi(x)$ .

Доказать, что при  $\mu \rightarrow 0$  в случае  $\psi(x) = 2$  при  $x < 0$  и  $\psi(x) = 1$  при  $x > 0$  решение стремится к обобщенному в смысле данного выше определения решению, отвечающему задаче (I), (II) при  $k = 1$ .

## § 2. Построение разностных схем

Перейдем теперь к вопросу о построении разностных схем для вычисления обобщенного решения задачи

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = 0, \quad u(x, 0) = \psi(x), \quad (1)$$

соответствующего интегральному закону сохранения (3) из § 1 при  $k = 1$ .

Будем предполагать для определенности, что  $\psi(x) > 0$ . Тогда  $u(x, t) > 0$ . Первое, что кажется естественным, это рассмотреть разностную схему

$$\begin{aligned} \frac{u_m^{p+1} - u_m^p}{\tau} + u_m^p \frac{u_m^p - u_{m-1}^p}{h} &= 0, \\ m = 0, \pm 1, \dots, \quad p = 0, 1, \dots, \\ u_m^0 &= \psi(mh). \end{aligned} \quad (2)$$

Замораживая коэффициент  $u_m^p$  в точке  $m = m_0$ , мы видим, что для возникающего уравнения с постоянными коэффициентами при переходе на слой  $t = (p+1)\tau$  выполняется принцип максимума (см. гл. 9, § 1), если шаг  $\tau = \tau_p$  выбран из условия

$$r_p = \frac{\tau_p}{h} \leq \frac{1}{\max_m |u_m^p|}.$$

Поэтому можно ожидать устойчивости. Если решение задачи (1) гладкое, то аппроксимация задачи (1) задачей (2) не вызывает сомнения. Действительно, в этом случае экспериментальные расчеты заранее известных гладких решений подтверждают сходимость.

Однако если задача (1) имеет разрывное решение, то ожидать сходимости к обобщенному решению в каком-либо разумном смысле нет оснований. Ведь в используемую разностную схему (2) не заложена информация о том, какой именно закон сохранения положен в основу определения обобщенного решения. Здесь нарушено условие Куранта, Фридрихса и Леви в том смысле, что, в отличие от решения разностной схемы, обобщенное решение зависит от числа  $k$ , определяющего интегральный закон сохранения (I), а решение разностной схемы не зависит.

Изложим два подхода к построению алгоритмов для вычисления обобщенного решения.

**1. Метод характеристик.** В этом методе развитие возникающих в процессе расчета, т. е. при увеличении времени  $t$ , разрывов ведется по особым формулам, основанным на условии (7) из § 1, которое должно выполняться на разрыве. В точках областей гладкости используется запись закона сохранения в дифференциальной форме, т. е. уравнение

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} = 0.$$

Принципиальная схема одного из вариантов метода характеристик применительно к нашему примеру состоит в следующем. Отметим на оси  $Ox$  точки  $x_m = mh$ . Будем считать для определенности, что начальное условие  $u(x, 0) = \psi(x)$  задается гладкой функцией  $\psi(x)$ . Из каждой точки  $(x_m, 0)$  проведем характеристику уравнения  $u_t + uu_x = 0$ .

Предположим, чтобы не осложнять изложение, что для заданной функции  $\psi(x)$  значение  $\tau$  выбрано столь малым, что на любом отрезке времени длины  $\tau$  каждая характеристика пересекается не более чем с одной из соседних характеристик. Проведем прямые  $t = t_p = p\tau$ . Рассмотрим точки пересечения характеристик, выходящих из точек  $(x_m, 0)$ , с прямой  $t = \tau$  и перенесем в эти точки значения решения  $u(x_m, 0) = \psi(x_m)$  по характеристикам.

Если на участке  $0 \leq t \leq \tau$  никакие две характеристики не пересекались, то делаем следующий шаг: продолжаем характеристики до пересечения с прямой  $t = 2\tau$  и переносим по характеристикам значения решения в точки пересечения. Если пересечения характеристик за время  $\tau \leq t \leq 2\tau$  опять не было, то делаем следующий шаг, и так до тех пор, пока на некотором участке  $t_p \leq t \leq t_{p+1}$  не пересекутся две характеристики (например, выходящие из точек  $(x_m, 0)$  и  $(x_{m+1}, 0)$ ) (рис. 36). Тогда середину отрезка  $Q_{m+1}^{p+1}Q_m^{p+1}$  будем считать точкой, из

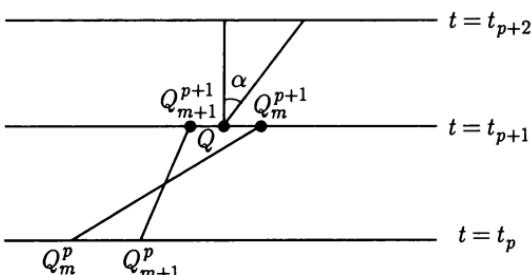


Рис. 36

которой выходит зарождающийся разрыв. Точки  $Q_{m+1}^{p+1}$  и  $Q_m^{p+1}$  заменим одной точкой  $Q$ , приписывая ей два значения решения:

$$u_{\text{лев}} = u(Q_m^{p+1}), \quad u_{\text{прав}} = u(Q_{m+1}^{p+1}).$$

Из точки  $Q$  проводим линию разрыва до пересечения с прямой  $t = t_{p+2}$ . Угловой коэффициент линии разрыва определяем из условия на разрыве

$$\operatorname{tg} \alpha = \frac{u_{\text{лев}} + u_{\text{прав}}}{2}.$$

Из точки пересечения линии разрыва с прямой  $t = t_{p+2}$  проводим в обратном направлении характеристики с угловыми коэффициентами  $u_{\text{лев}}$ ,  $u_{\text{прав}}$  с предыдущего слоя до пересечения с прямой  $t = t_{p+1}$ . В точках пересечения этих характеристик с прямой  $t = t_{p+1}$  с помощью

интерполяции по  $x$  находим значения  $u$  и принимаем их за левое и правое значения решения в точке разрыва, лежащей на прямой  $t = t_{p+2}$ . Это позволяет определить новый наклон линии разрыва как среднее арифметическое найденных значений слева и справа и продолжить разрыв еще на шаг  $\tau$  по времени.

Достоинство метода характеристик состоит в том, что он позволяет следить за разрывами и аккуратно их рассчитывать. Однако в процессе счета в общем случае возникают все новые разрывы; в частности, малосущественные разрывы могут пересекаться, так что с течением времени картина усложняется. Логика расчета усложняется, требования к памяти компьютера и его быстродействию возрастают. В этом заключается недостаток метода характеристик, в котором разрывы выделены и считаются нестандартным образом.

**2. Дивергентные разностные схемы. Схема Годунова.** Разностные схемы, не использующие искусственно введенную вязкость (см. задачу в конце § 1) и не использующие условия на разрыве, как было выяснено, должны опираться на интегральные законы сохранения.

Проведем на плоскости  $Oxt$  сетку прямых  $t = p\tau$ ,  $x = (m + 1/2)h$  ( $p = 0, 1, \dots$ ,  $m = 0, \pm 1, \dots$ ). Отметим на сторонах возникающих прямоугольников их середины (рис. 37); пусть  $D_h$  — сетка, образованная отмеченными точками (оси координат на рисунке не показаны).

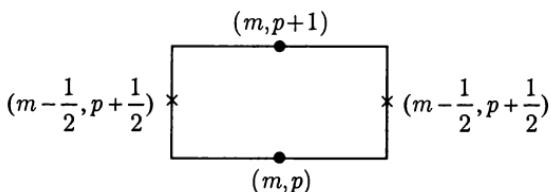


Рис. 37

Искомой функцией  $[u]_h$  будем считать сеточную функцию, определенную в точках сетки  $D_h$  путем усреднения решения  $u(x, t)$  по той стороне сеточного прямоугольника, которому принадлежит рассматриваемая точка сетки:

$$[u]_h \Big|_{\substack{x=x_m \\ t=t_p}} \equiv \tilde{u}_m^p = \frac{1}{h} \int_{x_{m-1/2}}^{x_{m+1/2}} u(x, t_p) dx,$$

$$[u]_h \Big|_{\substack{x=x_{m+1/2} \\ t=t_{p+1/2}}} \equiv \tilde{U}_{m+1/2}^{p+1/2} = \frac{1}{\tau} \int_{t_p}^{t_{p+1}} u(x_{m+1/2}, t) dt.$$

Приближенное решение  $u^{(h)}$  задачи определено на той же сетке  $D_h$ . Значения  $u^{(h)}$  в точках  $(x_m, t_p)$ , лежащих на горизонтальных сторонах прямоугольников, будем обозначать через  $u_m^p$ , а в точках  $(x_{m+1/2}, t_{p+1/2})$  вертикальных сторон — через  $U_{m+1/2}^{p+1/2}$ .

Величину  $u_m^p$  можно считать продолженной на всю сторону  $t = t_p$ ,  $x_{m-1/2} < x < x_{m+1/2}$  прямоугольника, которому принадлежит точка  $(x_m, t_p)$ . Аналогично будем считать, что  $U_{m+1/2}^{p+1/2}$  определена во всем вертикальном промежутке

$$x = x_{m+1/2}, \quad t_p < t < t_{p+1}.$$

Таким образом,  $u^{(h)}$  будет функцией, определенной на прямых  $x = (m + 1/2)h$ ,  $t = p\tau$ . Связь между величинами  $u_m^p$  и  $U_{m+1/2}^{p+1/2}$  ( $p = 0, 1, \dots$ ,  $m = 0, \pm 1, \dots$ ) установим, исходя из интегрального закона сохранения

$$\oint_{\Gamma} u \, dx - \frac{u^2}{2} \, dt = 0.$$

Рассмотрим в качестве контура  $\Gamma$  элементарный прямоугольник сетки и потребуем, чтобы

$$\oint_{\Gamma} u^{(h)} \, dx - \frac{(u^{(h)})^2}{2} \, dt = 0.$$

Перепишем последнее равенство в следующем эквивалентном виде:

$$-\oint_{\Gamma} u^{(h)} \, dx - \frac{(u^{(h)})^2}{2} \, dt = 0, \quad (3)$$

или, в развернутом виде,

$$h(u_m^{p+1} - u_m^p) + \frac{\tau}{2} [(U_{m+1/2}^{p+1/2})^2 - (U_{m-1/2}^{p+1/2})^2] = 0. \quad (4)$$

Если будет указано правило вычисления величин  $U_{m+1/2}^{p+1/2}$  ( $m = 0, \pm 1, \dots$ ) по уже известным величинам  $u_m^p$  ( $m = 0, \pm 1, \dots$ ), то формула (4) позволит вычислить величины  $u_m^{p+1}$  ( $m = 0, \pm 1, \dots$ ), т. е. продвинуться на один шаг по времени. Однако независимо от конкретного способа, который мы изберем для вычисления величины  $U_{m+1/2}^{p+1/2}$ , разностная схема вида (4) обладает свойством дивергентности, которое состоит в следующем.

Проведем в полуплоскости  $t \geq 0$  какой-либо замкнутый несамопересекающийся контур, целиком состоящий из сторон сеточных прямоугольников (рис. 38). Этот контур  $g_h$  ограничит некоторую область  $G_h$ , составленную из сеточных прямоугольников. Просуммируем почленно все уравнения (4), относя-

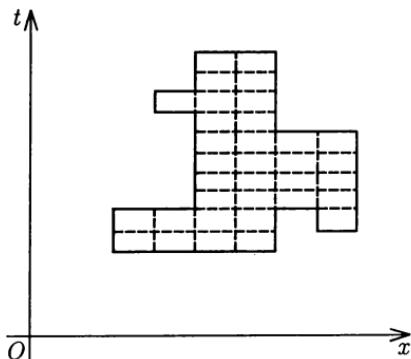


Рис. 38

щиеся к прямоугольникам, составляющим область  $G_h$ . Уравнения (4) и (3) отличаются только формой записи. Поэтому можно считать, что мы суммируем уравнения (3). Получим

$$-\oint_{g_h} u^{(h)} dx - \frac{(u^{(h)})^2}{2} dt = 0. \quad (5)$$

Интегралы по тем сторонам прямоугольников, которые не лежат на границе  $g_h$  области  $G_h$ , но входят в выражение (3), после суммирования уравнений (3) взаимно уничтожаются. Действительно, каждая

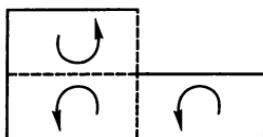


Рис. 39

из этих сторон принадлежит двум соседним прямоугольникам, так что интегрирование функции  $u^{(h)}$  по такой стороне встречается дважды и ведется в противоположных направлениях (рис. 39).

Разностные схемы, при суммировании которых по точкам сеточной области  $G_h$  остаются только алгебраические суммы значений неизвестных или функций от них вдоль границы области, называют *дивергентными*, или *консервативными*. Такие схемы аналогичны дифференциальным уравнениям дивергентного вида

$$\operatorname{div} \Phi = \frac{\partial \Phi_1}{\partial t} + \frac{\partial \Phi_2}{\partial x} = 0,$$

при почленном интегрировании которых по двумерной области  $D$  в левой части возникает контурный интеграл (2) из § 1. Разностная схема (2) недивергентна, схема (4) дивергентна.

Заметим следующее. Пусть сеточная функция  $u^{(h)}$ , удовлетворяющая уравнению (4), при  $h \rightarrow 0$  равномерно сходится к некоторой кусочно непрерывной функции  $u(x, t)$  во всякой замкнутой области, не содержащей линий разрыва, и пусть  $u^{(h)}$  равномерно по  $h$  ограничена. Тогда  $u(x, t)$  удовлетворяет интегральному закону сохранения

$$\oint_g u dx - \frac{u^2}{2} dt = 0,$$

где  $g$  — произвольный кусочно гладкий контур, т. е.  $u^{(h)}$  сходится к обобщенному решению. Это непосредственно следует из возможности приблизить контур  $g$  контуром  $g_h$ , из равенства (5) и предположения о сходимости \*).

\* ) Функция  $u = u(x, t)$  определена почти всюду, а функция  $u^{(h)}(x, t)$  — лишь на сетке прямых. Это формальное несоответствие можно преодолеть, считая, что при уменьшении  $h$  каждая новая сетка построена так, что прежняя является ее подсеткой, и говоря о сходимости в точках сетки, построенной для любого фиксированного  $h$  из числа допустимых.

Чтобы схема (4) приобрела смысл, надо указать способ вычисления величин  $U_{m+1/2}^{p+1/2}$  по величинам  $u_m^p$ . В схеме Годунова, которую мы возьмем для иллюстрации понятия дивергентных схем, для этого используется решение следующей задачи о распаде разрыва.

Пусть в начальный момент решение  $u(x, 0)$  задано условиями

$$u(x, 0) = \begin{cases} u_{\text{лев}}, & x < 0, \\ u_{\text{прав}}, & x > 0, \end{cases}$$

где  $u_{\text{лев}} = \text{const}$ ,  $u_{\text{прав}} = \text{const}$ . Тогда можно найти соответствующее обобщенное решение. Как это делается, мы видели в § 1 при разборе примера  $u_{\text{лев}} = 1/2$ ,  $u_{\text{прав}} = 3/2$  и примера  $u_{\text{лев}} = 2$ ,  $u_{\text{прав}} = 1$ . Нам важно знать значение  $U = u(0, t)$  решения  $u(x, t)$  при  $x = 0$ .

Читатель, построив картинки типа рис. 28, изображающие решение  $u(x, t)$ , легко проверит, что на прямой  $x = 0$  решение принимает значения  $u_{\text{лев}}$ ,  $u_{\text{прав}}$  или 0 в зависимости от заданных начальных данных, и для каждой конкретной пары чисел  $u_{\text{лев}}$ ,  $u_{\text{прав}}$  выяснит, какое именно. Например, при  $u_{\text{лев}} > 0$ ,  $u_{\text{прав}} > 0$  будет  $u(0, t) = u_{\text{лев}}$ , а при  $u_{\text{лев}} < 0$ ,  $u_{\text{прав}} < 0$  будет  $u(0, t) = u_{\text{прав}}$ .

Величину  $U_{m+1/2}^{p+1/2}(= U)$  в схеме (4) будем определять из задачи о распаде разрыва, возникающего на границе  $x = x_{m+1/2}$  каждого из двух участков, где заданы постоянные значения  $u_m^p (= u_{\text{лев}})$ ,  $u_{m+1}^p (= u_{\text{прав}})$ .

Если, например,  $u_m^p > 0$  ( $m = 0, \pm 1, \dots$ ), то  $U_{m+1/2}^{p+1/2} = u_{\text{лев}} = u_m^p$  ( $m = 0, \pm 1, \dots$ ), и схема (4) примет вид

$$\frac{u_m^{p+1} - u_m^p}{\tau} + \frac{1}{h} \left[ \frac{(u_m^p)^2}{2} - \frac{(u_{m-1}^p)^2}{2} \right] = 0,$$

$$u_m^0 = \frac{1}{h} \int_{x_{m-1/2}}^{x_{m+1/2}} \psi(x) dx,$$

или

$$\frac{u_m^{p+1} - u_m^p}{\tau} + \frac{u_{m-1}^p + u_m^p}{2} \cdot \frac{u_m^p - u_{m-1}^p}{h} = 0.$$

Легко видеть, что при

$$r \leq \frac{\tau}{h} \leq \frac{1}{\max_m |u_m^p|}$$

имеет место принцип максимума

$$\max_m |u_m^{p+1}| \leq \max_m |u_m^p| \leq \dots \leq \max_m |u_m^0| \leq \max_x |\psi(x)|.$$

Отсюда видно, что при  $\tau \leq \frac{1}{\max_x |\psi(x)|} h$  можно надеяться, что при некотором разумном выборе норм полученная разностная схема устойчива. Мы не будем фактически указывать эти нормы: экспериментальные расчеты подтверждают, что при измельчении сетки решение  $u^{(h)}$

задачи (4) с кусочно монотонными и кусочно гладкими начальными данными  $\psi(x)$  сходится к некоторой функции  $u(x, t)$ , имеющей конечное число разрывов, причем вне любой окрестности разрывов сходимость равномерная.

Схема (4) с вычислением  $U_{m+1/2}^{p+1/2}$  путем использования распада разрыва, т. е. схема Годунова, не является единственной дивергентной схемой для задачи (1); см. [4, 9, 12].

## ГЛАВА 11

### РАЗНОСТНЫЕ МЕТОДЫ ДЛЯ ЭЛЛИПТИЧЕСКИХ ЗАДАЧ

Простейшим эллиптическим уравнением является уравнение Пуассона

$$\Delta u = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = \varphi(x, y). \quad (1)$$

Пусть задана некоторая область  $D$  на плоскости, а на ее границе  $\Gamma = \partial D$  поставлено краевое условие вида

$$\left( au - b \frac{\partial u}{\partial n} \right) \Big|_{\Gamma} = \psi(s), \quad (2)$$

где  $\partial/\partial n$  — производная в направлении внутренней нормали,  $a \geq 0$ ,  $b \geq 0$  ( $a^2 + b^2 = 1$ ) — некоторые числа,  $s$  — длина дуги, отсчитываемая вдоль границы  $\Gamma$ . Функции  $f(x, y)$ ,  $\psi(s)$  считаются заданными. Требуется вычислить решение краевой задачи (1), (2). В случае  $a = 1$ ,  $b = 0$  возникающая задача называется *первой краевой задачей*, или *задачей Дирихле*. В случае  $a = 0$ ,  $b = 1$  — это *вторая краевая задача*, или *задача Неймана*. В случае  $a > 0$ ,  $b > 0$  — это *третья краевая задача*.

Перечисленные три краевые задачи являются основными для уравнения Пуассона, но в математической физике возникают также и другие задачи для уравнения Пуассона. Наряду с уравнением Пуассона могут возникать уравнения с переменными коэффициентами вида

$$\frac{\partial}{\partial x} \left( a \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( b \frac{\partial u}{\partial y} \right) = \varphi(x, y), \quad (3)$$

где  $a = a(x, y) > 0$ ,  $b = b(x, y) > 0$  — функции, для которых также корректно поставлены задачи с краевыми условиями вида (2).

Эллиптические уравнения и системы эллиптических уравнений возникают при математическом моделировании многих стационарных состояний. Так, например, уравнение Пуассона может описывать потенциал электрического поля, потенциал скоростей установившегося потока несжимаемой жидкости, установившуюся температуру в однородном

теплопроводном теле; система уравнений упругости Ламе описывает смещения в находящемся под действием стационарных сил теле, и т. п.

Численное решение краевых задач для эллиптических уравнений осуществляется во многих случаях с помощью разностных схем.

В простейших случаях, когда решение во всей рассматриваемой области меняется более или менее равномерно, а сама область не имеет узких «перешейков», можно пользоваться регулярными сетками, а для получения разностных схем заменять производные разностными отношениями.

В противном случае регулярная сетка становится практически непригодной, так как места быстрого изменения решения или места, где область  $D$  имеет узкие «горловины», диктуют слишком мелкую сетку. А в случае нерегулярной сетки построение разностной схемы путем замены производных разностными отношениями становится невозможным. Для построения разностных схем на нерегулярных сетках большое распространение получили так называемые вариационно-разностные и проекционно-разностные схемы. Среди прикладников соответствующие схемы называют обычно *методом конечных элементов*.

В этой главе мы остановимся на трех вопросах. Во-первых, проверим, что простейшая пятиточечная разностная схема для задачи Дирихле в прямоугольной области, с которой мы уже не раз встречались, обладает свойствами аппроксимации и устойчивости, а значит, является сходящейся при измельчении сетки. Во-вторых, дадим представление о методе конечных элементов. Наконец, объясним, как можно вычислять решения больших систем линейных уравнений, которыми являются при мелких сетках как простейшие разностные схемы, так и системы, возникающие в методе конечных элементов.

При этом в дополнение к изложенным в части II точным и итерационным методам решения систем линейных уравнений будет изложен многосеточный метод Федоренко.

## § 1. Апроксимация и устойчивость простейшей разностной схемы

Рассмотрим задачу Дирихле для уравнения Пуассона в квадратной области  $\overline{D} = (0 \leq x, y \leq 1)$  с границей  $\Gamma$ :

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = \varphi(x, y), \quad 0 \leq x, y \leq 1, \quad (1)$$

$$u|_{\Gamma} = \psi(s).$$

Совокупность точек  $(x_m, y_n) = (mh, nh)$  сетки (где  $h = 1/M$ ,  $M$  целое), попавших внутрь квадрата или на его границу, обозначим  $D_h$ .

Точки  $D_h$ , лежащие строго внутри квадрата  $D$ , будем называть *внутренними точками* сеточного квадрата  $D_h$ ; совокупность

внутренних точек обозначим  $D_h^0$ . Совокупность точек  $D_h$ , попавших на границу квадрата  $\bar{D}$ , будем обозначать  $\Gamma_h$ .

Рассмотрим разностную схему

$$L_h u^{(h)} = f^{(h)}; \quad (2)$$

здесь

$$L_h u^{(h)} = \begin{cases} \frac{u_{m+1,n} - 2u_{mn} + u_{m-1,n}}{h^2} + \frac{u_{m,n+1} - 2u_{mn} + u_{m,n-1}}{h^2} = \\ u_{mn} = \psi(s_{mn}), \quad (x_m, y_n) \in \Gamma_h, \end{cases} = \varphi(x_m, y_n), \quad (3)$$

где  $\psi(s_{mn})$  — значение функции  $\psi(s)$  в точке  $(x_m, y_n)$ , принадлежащей границе  $\Gamma_h$ .

**1. Аппроксимация.** Правая часть  $f^{(h)}$  разностной схемы (2) имеет вид

$$f^{(h)} = \begin{cases} \varphi(x_m, y_n), & (x_m, y_n) \in D_h^0, \\ \psi(s_{mn}), & (x_m, y_n) \in \Gamma_h. \end{cases} \quad (4)$$

В предположении, что решение  $u(x, y)$  задачи (1) имеет ограниченные четвертые производные, с помощью формулы Тейлора устанавливается равенство

$$\frac{u(x+h, y) - 2u(x, y) + u(x-h, y)}{h^2} + \frac{u(x, y+h) - 2u(x, y) + u(x, y-h)}{h^2} = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} + O(h^2). \quad (5)$$

Поэтому для решения  $u(x, y)$  задачи (1) имеем

$$L_h[u]_h = \begin{cases} \varphi(x_m, y_n) + O(h^2), & (x_m, y_n) \in D_h^0, \\ \psi(s_{mn}) + 0, & (x_m, y_n) \in \Gamma_h. \end{cases} \quad (6)$$

Таким образом, невязка  $\delta f^{(h)}$ , возникающая при подстановке  $[u]_h$  в левую часть разностной схемы (2), имеет вид

$$\delta f^{(h)} = \begin{cases} O(h^2), & (x_m, y_n) \in D_h^0, \\ 0, & (x_m, y_n) \in \Gamma_h. \end{cases} \quad (7)$$

В пространстве  $F_h$ , состоящем из элементов вида

$$f^{(h)} = \begin{cases} \varphi_{mn}, & (mh, nh) \in D_h^0, \\ \psi_{mn}, & (mh, nh) \in \Gamma_h, \end{cases}$$

введем норму

$$\|f^{(h)}\|_{F_h} = \max_{(mh, nh) \in D_h^0} |\varphi_{mn}| + \max_{(mh, nh) \in \Gamma_h} |\psi_{mn}|. \quad (8)$$

Тогда

$$\|\delta f^{(h)}\|_{F_h} = O(h^2).$$

Таким образом, разностная краевая задача (3) аппроксимирует задачу Дирихле (1) со вторым порядком относительно  $h$ .

**2. Устойчивость.** Определим норму в пространстве  $U_h$  функций, заданных на сетке  $D_h$ , положив

$$\|u^{(h)}\|_{U_h} = \max_{(mh,nh) \in D_h} |u_{mn}|. \quad (9)$$

Для доказательства устойчивости разностной схемы (3), к которому мы приступаем, в соответствии с определением устойчивости надо установить, что задача (2) однозначно разрешима при произвольной правой части  $f^{(h)}$  (это свойство не зависит от выбора норм) и что выполняется оценка вида

$$\|u^{(h)}\|_{U_h} \leq c \|f^{(h)}\|_{F_h}, \quad (10)$$

где  $c$  не зависит ни от  $h$ , ни от  $f^{(h)}$ .

**Лемма 1.** Пусть функция  $v^{(h)} = \{v_{mn}\}$  определена на сетке  $D_h$  и во внутренних точках  $(x_m, y_n) = (mh, nh) \in D_h^0$  удовлетворяет условию

$$\Lambda_h v^{(h)}|_{(mh,nh)} \geq 0, \quad (mh, nh) \in D_h^0, \quad (11)$$

где

$$\Lambda_h v^{(h)}|_{(mh,nh)} \equiv \frac{v_{m+1,n} - 2v_{mn} + v_{m-1,n}}{h^2} + \frac{v_{m,n+1} - 2v_{mn} + v_{m,n-1}}{h^2}.$$

Тогда наибольшее на сетке  $D_h$  значение  $v^{(h)}$  достигается хотя бы в одной точке  $\Gamma_h$ .

**Доказательство.** Допустим противное. Выберем среди точек сетки  $D_h$ , в которых  $v^{(h)}$  достигает своего наибольшего значения, какую-нибудь одну точку  $(x_m, y_n)$ , имеющую самую большую абсциссу. По нашему предположению  $(x_m, y_n)$  — внутренняя точка, причем  $v_{mn} > v_{m+1,n}$ . В точке  $(mh, nh)$

$$\begin{aligned} \Lambda_h v^{(h)}|_{(mh,nh)} &\equiv \\ &\equiv \frac{(v_{m+1,n} - v_{mn}) + (v_{m-1,n} - v_{mn}) + (v_{m,n+1} - v_{mn}) + (v_{m,n-1} - v_{mn})}{h^2} < 0, \end{aligned}$$

поскольку первая скобка в числителе отрицательна, а остальные скобки неположительны. Получаем противоречие с (11).  $\square$

**Лемма 2.** Пусть функция  $v^{(h)} = \{v_{mn}\}$  определена на сетке  $D_h$  и во внутренних точках удовлетворяет условию

$$\Lambda_h v^{(h)}|_{(mh,nh)} \leq 0, \quad (mh, nh) \in D_h^0. \quad (12)$$

Тогда наименьшее на сетке  $D_h$  значение  $v^{(h)}$  достигается хотя бы в одной точке границы.

Лемма 2 доказывается аналогично лемме 1.

**Теорема 1** (принцип максимума). *Каждое решение разностного уравнения*

$$\Lambda_h v^{(h)}|_{(mh, nh)} = 0, \quad (mh, nh) \in D_h^0, \quad (13)$$

*достигает своих наибольшего и наименьшего значений в некоторых точках  $\Gamma_h$ .*

Доказательство получается объединением утверждений лемм 1, 2.

Указанное свойство решений разностного уравнения (13) аналогично свойству решений  $v(x, y)$  уравнения Лапласа  $v_{xx} + v_{yy} = 0$  принимать наибольшее и наименьшее значения на границе области, где эти решения определены.

Из принципа максимума следует, что задача

$$L_h u^{(h)} = \begin{cases} \Lambda_h u^{(h)}|_{(mh, nh)} = 0, & (mh, nh) \in D_h^0, \\ u^{(h)}|_{(mh, nh)} = 0, & (mh, nh) \in \Gamma_h, \end{cases} \quad (14)$$

имеет только нулевое решение  $u^{(h)} = 0$ , поскольку наибольшее и наименьшее значения этого решения получаются в точках  $\Gamma_h$ , где  $u_{mn} = 0$ . Следовательно, разностная краевая задача (2) однозначно разрешима при произвольной правой части.

Переходим к доказательству оценки (10). В силу (5) для произвольного многочлена второй (и даже третьей) степени  $P(x, y)$  выполняется равенство

$$\Lambda_h P = \frac{\partial^2 P}{\partial x^2} + \frac{\partial^2 P}{\partial y^2}, \quad (15)$$

так как четвертые производные от  $P(x, y)$ , входящие в выражение остаточного члена формулы (5), обращаются в нуль.

Используя функции  $\varphi_{mn} = \varphi(x_m, y_n)$ ,  $\psi_{mn} = \psi(s_{mn})$  из правой части системы (3) и фиксируя  $R = \sqrt{2}$ , строим вспомогательную функцию

$$P^{(h)}(x, y) = \frac{1}{4} [R^2 - (x^2 + y^2)] \max_{(mh, nh) \in D_h^0} |\varphi_{mn}| + \max_{(mh, nh) \in \Gamma_h} |\psi_{mn}|.$$

Будем рассматривать только ее в точках сетки  $D_h$ , что отражено значком  $(h)$  в обозначении  $P^{(h)}(x, y)$ . В силу (15)

$$\Lambda_h P^{(h)} \Big|_{\substack{x= mh \\ y= nh}} = - \max_{(rh, sh) \in D_h^0} |\varphi_{rs}|, \quad (mh, nh) \in D_h^0.$$

Поэтому разность решения  $u^{(h)}$  задачи (3) и функции  $P^{(h)}$  удовлетворяет в точках  $D_h^0$  равенствам

$$\Lambda_h (u^{(h)} - P^{(h)}) = \Lambda_h u^{(h)} - \Lambda_h P^{(h)} = \varphi_{mn} + \max_{i,j} |\varphi_{ij}| \geq 0.$$

В силу леммы 1 разность  $u^{(h)} - P^{(h)}$  принимает свое наибольшее значение на границе  $\Gamma_h$ . Но на границе  $\Gamma_h$  эта разность

$$\begin{aligned} u^{(h)} \Big|_{\Gamma_h} - P^{(h)} \Big|_{\Gamma_h} &= \psi_{mn} - P^{(h)}(mh, nh) = \\ &= (\psi_{mn} - \max_{(rh, sh) \in \Gamma_h} |\psi_{rs}|) + \frac{1}{4} |x^2 + y^2 - R^2| \max_{(rh, sh) \in D_h^0} |\varphi_{rs}| \end{aligned}$$

неположительна, так как в квадрате  $D$  всюду  $x^2 + y^2 \leq R^2$  и обе скобки в правой части неположительны. Поскольку наибольшее значение  $u^{(h)} - P^{(h)}$  неположительно, то в силу леммы 1 всюду на  $D_h^0$

$$[u^{(h)} - P^{(h)}]_{mh, nh} \leq 0,$$

или

$$u^{(h)} \leq P^{(h)}.$$

Аналогично для функции  $u^{(h)} + P^{(h)}$  в точках  $D_h^0$

$$\Lambda_h(u^{(h)} + P^{(h)}) \leq 0,$$

а в точках  $\Gamma_h$  сумма  $u^{(h)} + P^{(h)}$  неотрицательна. В силу леммы 2 всюду на  $D_h^0$

$$u^{(h)} + P^{(h)} \geq 0,$$

или

$$-P^{(h)} \leq u^{(h)}.$$

Таким образом, всюду на  $D_h$

$$|u_{mn}| \leq |P^{(h)}(mh, nh)| \leq \frac{1}{4} R^2 \max_{(rh, sh) \in D_h^0} |\varphi_{rs}| + \max_{(rh, sh) \in \Gamma_h} |\psi_{rs}|.$$

Отсюда вытекает неравенство (10):

$$\begin{aligned} \|u^{(h)}\|_{U_h} &= \max |u_{mn}| \leq c \left( \max_{(rh, sh) \in D_h^0} |\varphi_{rs}| + \max_{(rh, sh) \in \Gamma_h} |\psi_{rs}| \right) = \\ &= c \|f^{(h)}\|_{F_h}, \end{aligned}$$

где

$$c = \max \left( 1, \frac{R^2}{4} \right) = 1;$$

этим завершается доказательство устойчивости.

В случае задачи Дирихле для эллиптического уравнения с переменными коэффициентами

$$\frac{\partial}{\partial x} \left( a \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( b \frac{\partial u}{\partial y} \right) = \varphi(x, y), \quad (x, y) \in D,$$

$$u|_{\Gamma} = \psi(s),$$

где  $a = a(x, y)$ ,  $b = b(x, y)$  — положительные в прямоугольнике  $D$  гладкие функции, разностную схему можно построить аналогично.

Используя во внутренних точках сетки замену выражений  $\frac{\partial}{\partial x} \left( a \frac{\partial u}{\partial x} \right)$ ,  $\frac{\partial}{\partial y} \left( b \frac{\partial u}{\partial y} \right)$  разностными отношениями по приближенным формулам

$$\frac{\partial}{\partial x} \left[ a(x, y) \frac{\partial u(x, y)}{\partial x} \right] \approx \tilde{\Lambda}_{xx} u(x, y) = \frac{1}{h} \left[ a \left( x + \frac{h}{2}, y \right) \times \right. \\ \left. \times \frac{u(x+h, y) - u(x, y)}{h} - a \left( x - \frac{h}{2}, y \right) \frac{u(x, y) - u(x-h, y)}{h} \right],$$

$$\frac{\partial}{\partial y} \left[ b(x, y) \frac{\partial u(x, y)}{\partial y} \right] \approx \tilde{\Lambda}_{yy} u(x, y) = \frac{1}{h} \left[ b \left( x, y + \frac{h}{2} \right) \times \right. \\ \left. \times \frac{u(x, y+h) - u(x, y)}{h} - b \left( x, y - \frac{h}{2} \right) \frac{u(x, y) - u(x, y-h)}{h} \right],$$

получаем разностную схему вида (2), где

$$L_h u^{(h)} = \begin{cases} \tilde{\Lambda}_{xx} u^{(h)} + \tilde{\Lambda}_{yy} u^{(h)} = \varphi(mh, nh), & (mh, nh) \in D_h^0, \\ u|_{\Gamma_h} = \psi(s_{mn}), & (mh, nh) \in \Gamma_h. \end{cases} \quad (16)$$

Пользуясь формулой Тейлора, можно убедиться в том, что имеет место второй порядок аппроксимации. Можно было бы доказать устойчивость построенной схемы, преодолевая некоторые дополнительные трудности по сравнению с рассмотренными нами при разборе примера.

На практике при решении конкретных задач обычно ограничиваются обоснованиями принципиального характера на модельных задачах типа приведенных выше. Конкретные суждения о погрешности получаются, как правило, не из теоретических оценок, а из сравнения результатов расчетов, выполненных на сетках с различными значениями шага  $h$ .

После того как разностная краевая задача, аппроксимирующая дифференциальную, построена, нужно еще указать не слишком трудоемкий способ ее решения. Ведь при малом  $h$  задача (16) есть система скалярных уравнений очень высокого порядка. В разобранном нами примере решение разностных уравнений — сложная и интересная задача, но мы отложим ее рассмотрение до § 3.

### Задачи

1. Доказать, что если во внутренних точках области  $D_h$  функция  $u^{(h)}$  удовлетворяет уравнению

$$\Lambda_h u^{(h)}|_{(mh, nh)} = 0, \quad (mh, nh) \in D_h^0,$$

то либо  $u^{(h)}$  принимает всюду на  $D_h$  одинаковые значения, либо наибольшее и наименьшее значения функции  $u^{(h)}$  не достигаются ни в одной внутренней точке сетки  $D_h$  (*усиленный принцип максимума*).

2. Если во всех внутренних точках области  $D_h$  выполнено условие  $\Lambda_h u^{(h)} \geq 0$ , причем хотя бы в одной точке неравенство строгое, то  $u^{(h)}$  не достигает своего наибольшего значения ни в одной внутренней точке.

3. Рассмотрим разностную схему  $L_h u^{(h)} = f^{(h)}$  вида

$$L_h u^{(h)} = \begin{cases} \Lambda_h u^{(h)}|_{(mh,nh)} = \varphi(mh, nh), & (mh, nh) \in D_h^0, \\ u_{mn} = \psi_1(s_{mn}), & m = M \text{ или } n = 0, \text{ или } n = M, \\ \frac{u_{1,n} - u_{0,n}}{h} = \psi_2(s_{mn}), & n = 1, 2, \dots, M - 1. \end{cases}$$

Эта разностная схема аппроксимирует задачу

$$\begin{aligned} \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} &= \varphi(x, y), \quad (x, y) \in D, \\ u(x, y) &= \psi_1(s), \quad x = 1 \text{ или } y = 0, \text{ или } y = 1, \\ \frac{\partial u}{\partial x} &= \psi_2(s), \quad x = 0. \end{aligned}$$

a) Доказать, что при любых  $\varphi(mh, nh) = 0$ ,  $\psi_1(s_{mn})$ ,  $\psi_2(s_{mn})$  задача  $L_h u^{(h)} = f^{(h)}$  имеет единственное решение.

б)\* Доказать, что если  $\varphi(mh, nh)$  неотрицательно, а  $\psi_1(s_{mn})$ ,  $\psi_2(s_{mn})$  неположительны, то  $u^{(h)}$  неположительно.

в)\* Доказать, что при любых  $\varphi$ ,  $\psi_1$ ,  $\psi_2$  имеет место оценка вида

$$\max_{(mh,nh) \in D_h} |u_{mn}| \leq c \left( \max_{(mh,nh) \in D_h^0} |\varphi_{mn}| + \max_{m,n} |\psi_1(s_{mn})| + \max_n |\psi_2(s_{0n})| \right),$$

где  $c$  — некоторая постоянная, не зависящая от  $h$ . Вычислить  $c$ .

## § 2. Понятие о методе конечных элементов

В этом параграфе мы дадим представление о методе конечных элементов для построения системы алгебраических уравнений, связывающих значения искомого решения дифференциальной краевой задачи на нерегулярной сетке. Возникающая благодаря методу конечных элементов свобода в выборе сеток позволяет использовать нерегулярные сетки и располагать их узлы гуще в тех частях области определения искомого решения, где это решение ведет себя более сложно или где нас интересуют более мелкие детали его поведения. Обычный для регулярных сеток способ построения разностных схем путем аппроксимации производных в точках сетки разностными отношениями в случае нерегулярной сетки становится затруднительным.

Мы дадим понятие о методе конечных элементов на характерном примере его применения к численному решению задачи Дирихле для уравнения Пуассона

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = f(x, y), \quad (x, y) \in D, \tag{1}$$

$$u|_{\Gamma} = 0. \tag{2}$$

Для простоты будем предполагать, что  $D$  — заданная выпуклая ограниченная область на плоскости  $Oxy$  с кусочно гладкой границей  $\Gamma$ ,

функция  $f(x, y)$  — заданная функция, причем решение задачи (1), (2) имеет ограниченные вторые производные всюду в замкнутой области  $\bar{D} = D \cup \Gamma$ .

Линейное пространство всех функций, дважды непрерывно дифференцируемых всюду в замкнутой области  $\bar{D}$ , обозначим  $W$ .

Для дискретизации и численного решения задачи (1), (2) построим сетку (рис. 40).

Для построения этой сетки сначала впишем в контур  $\Gamma$  некоторую несамопересекающуюся замкнутую ломаную, вершины которой обозначим через  $Q_m$ ,  $m = 1, 2, \dots, M$ ; эта ломаная ограничивает некоторый многоугольник  $Q$ , имеющий  $M$  вершин. Разобьем многоугольник  $Q$  на треугольники так, чтобы каждое звено замкнутой ломаной  $Q_1 Q_2 \dots Q_M Q_1$ .

оказалось стороной одного из треугольников разбиения, причем любые два треугольника этого разбиения могли бы пересекаться лишь по общей стороне либо общей вершине. Общее число тех вершин  $P_1, \dots, P_N$  треугольников, которые лежат внутри  $D$ , обозначим  $N$ .

Совокупность точек  $P_1, P_2, \dots, P_N$  и точек  $Q_1, Q_2, \dots, Q_M$  примем за сетку, в точках  $P_n$  которой надо определить численно значения искомого решения  $u(x, y)$  задачи (1), (2). Значения искомого решения в точках  $Q_m \in \Gamma$  в силу условия (2) заданы и равны нулю.

Заметим, что при заданных  $M$  и  $N$  в расположении узлов  $Q_m$  и  $P_n$  остается большой произвол, которым можно пользоваться для сгущения сетки в окрестности любой заданной точки области  $D$ . С возрастанием числа точек  $Q_m$  и  $P_n$  сетка может быть сделана сколь угодно мелкой всюду в  $D$ .

Один из двух основных подходов к построению вычислительных схем метода конечных элементов опирается на использование вариационной постановки исходной краевой задачи (Ритц, 1908). Второй подход не требует перехода к вариационной постановке краевой задачи и является одной из конкретизаций проекционного метода Галёркина (1916).

В п. 1–3 мы изложим вариационный подход, а в п. 4–5 проекционный подход к построению вычислительных схем метода конечных элементов.

**1. Вариационная постановка краевых задач.** Начнем с предварительных построений. В результате построения сетки область  $D$  разбивается на треугольники с вершинами  $Q_m, P_n$ ,  $m = 1, 2, \dots, M$ ,  $n = 1, 2, \dots, N$ . Если все три или хотя бы две вершины треугольника являются внутренними точками  $P_i$  и  $P_j$  сетки, то треугольник ограничен отрезками прямых; если же какие-нибудь две вершины этого

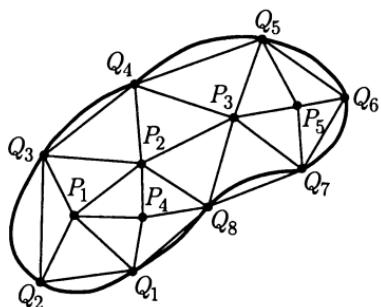


Рис. 40

треугольника  $Q_i$  и  $Q_j$  лежат на границе  $\Gamma$ , то одна из сторон возникающего криволинейного треугольника есть отрезок  $Q_i Q_j$  криволинейной границы  $\Gamma$ .

Перенумеруем произвольным образом треугольники  $D_k, k = 1, 2, \dots, K$ , на которые мы разбили исходную область  $D$  при построении сетки.

Введем теперь линейное пространство  $W^N$  функций  $w(x, y)$ , заданных и непрерывных всюду на  $\overline{D}$ , обращающихся в нуль на  $\Gamma$  и имеющих непрерывные производные первого порядка на замкнутом треугольнике  $\overline{D}_k, k = 1, 2, \dots, K$ , разбиения. При этом в точках границы  $\Gamma_k = \partial D_k, k = 1, 2, \dots, K$ , имеются в виду односторонние производные.

При переходе из одного треугольника разбиения в соседний с ним треугольник эти производные, возможно, имеют скачок.

Заметим, что запас функций  $w^N$ , образующих пространство  $W^N$ , зависит от выбора разбиения области  $\overline{D}$  на треугольники  $D_k$ , но всегда  $W^N \subset W$ .

**Лемма 1.** *Пусть  $\xi \in W^N$  и при этом  $\xi(x, y)$  отличается от тождественного нуля. Тогда имеет место следующее неравенство:*

$$\iint_D \left[ \left( \frac{\partial \xi}{\partial x} \right)^2 + \left( \frac{\partial \xi}{\partial y} \right)^2 \right] dx dy > 0.$$

**Доказательство.** Достаточно показать, очевидно, что в предположениях леммы  $\frac{\partial \xi}{\partial x}$  или  $\frac{\partial \xi}{\partial y}$  отлично от нуля хотя бы в одной точке  $(x, y)$  хотя бы одного из треугольников  $D_k, k = 1, 2, \dots, K$ , разбиения  $\overline{D} = \bigcup \overline{D}_k$ . Для этого покажем, что если

$$\frac{\partial \xi}{\partial x} \equiv 0, \quad \frac{\partial \xi}{\partial y} \equiv 0, \quad (x, y) \in D_k, \quad k = 1, 2, \dots, K,$$

то вопреки предположению леммы  $\xi(x, y) = 0$  всюду на  $D$ . Действительно в силу предполагаемых тождеств функция  $\xi(x, y)$  постоянна на каждом  $D_k$ ,  $\xi(x, y) = C_k, k = 1, 2, \dots, K$ . Ввиду непрерывности функции  $\xi \in W^N$  всюду на замкнутой области  $\overline{D} = D \bigcup \Gamma$  все эти постоянные равны друг другу:  $C_1 = C_2 = \dots = C_K = \text{const}$ ,

$$\xi(x, y) = \text{const}, \quad (x, y) \in \overline{D}.$$

Отсюда в силу условия  $\xi_\Gamma = 0$  следует, что  $\xi(x, y) \equiv 0$ , и лемма 1 доказана.  $\square$

**Лемма 2.** Пусть  $u(x, y)$  — решение задачи (1),(2), а  $\xi \in W^N$ . Тогда имеет место равенство

$$\iint_D \left[ \frac{\partial}{\partial x} \left( \xi \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( \xi \frac{\partial u}{\partial y} \right) \right] dx dy = 0. \quad (3)$$

Доказательство. Заметим, что

$$\begin{aligned} \iint_D \left[ \frac{\partial}{\partial x} \left( \xi \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( \xi \frac{\partial u}{\partial y} \right) \right] dx dy = \\ = \sum_{k=1}^K \iint_{D_k} \left[ \frac{\partial}{\partial x} \left( \xi \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( \xi \frac{\partial u}{\partial y} \right) \right] dx dy = \sum_{k=1}^K \int_{\Gamma_k} \xi \frac{\partial u}{\partial n} ds, \end{aligned}$$

где  $\partial u / \partial n$  — производная от  $u$  по направлению внешней по отношению к  $D_k$  нормали, а  $s$  — длина дуги  $\Gamma_k$ , отсчитываемая при обходе области  $D_k$  в направлении против часовой стрелки. В самом деле, достаточно воспользоваться равенством

$$\iint_{D_k} \left[ \frac{\partial}{\partial x} \left( \xi \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( \xi \frac{\partial u}{\partial y} \right) \right] dx dy = \int_{\Gamma_k} \xi \frac{\partial u}{\partial n} ds,$$

которое следует из того, что интеграл от дивергенции

$$\frac{\partial}{\partial x} \left( \xi \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( \xi \frac{\partial u}{\partial y} \right) = \operatorname{div}(\xi \operatorname{grad} u)$$

непрерывного на  $D_k$  (в силу предположений о решении  $u(x, y)$  задачи (1), (2) и о функции  $\xi \in W^N$ ) векторного поля  $\xi \operatorname{grad} u$  равен потоку этого поля через границу  $\Gamma_k = \partial D_k$ .

Далее имеем

$$\sum_{k=1}^K \int_{\Gamma_k} \xi \frac{\partial u}{\partial n} ds = \int_{\Gamma} \xi \frac{\partial u}{\partial n} ds. \quad (4)$$

В самом деле, интегралы по отрезкам  $P_i P_j$ , являющимся сторонами какого-либо треугольника разбиения, входят в сумму  $\sum_{k=1}^K \int_{\Gamma_k} \xi \frac{\partial u}{\partial n} ds$  дважды и при этом с противоположными знаками, поскольку отрезок  $P_i P_j$  является стороной сразу двух граничащих между собой треугольников разбиения. В силу этого в сумме (4) остается только сумма интегралов по тем сторонам криволинейных треугольников, которые являются отрезками границы  $\Gamma = \partial D$ , т.е. интеграл по  $\Gamma$ .

Для доказательства (3) остается заметить, что  $\int_{\Gamma} \xi \frac{\partial u}{\partial n} ds = 0$ , так как  $\xi(x, y) \in W^N$  обращается в нуль на границе  $\Gamma$ . Лемма 2 доказана.  $\square$

**Теорема.** Среди всех функций  $w \in W^N$  решение  $u(x, y)$  задачи (1), (2) придает выражению

$$I(w) = \iint_D \left[ \left( \frac{\partial w}{\partial x} \right)^2 + \left( \frac{\partial w}{\partial y} \right)^2 + 2fw \right] dx dy \quad (5)$$

наименьшее значение.

**Доказательство.** Пусть  $w(x, y) \in W^N$  — некоторая фиксированная функция. Введем обозначение

$$\xi(x, y) = w(x, y) - u(x, y).$$

Очевидно, что  $u(x, y) \in W^N$ , а значит,  $\xi \in W^N$ .

Докажем равенство

$$I(w) = I(u + \xi) = I(u) + \iint_D \left[ \left( \frac{\partial \xi}{\partial x} \right)^2 + \left( \frac{\partial \xi}{\partial y} \right)^2 \right] dx dy. \quad (6)$$

Непосредственно проверяется

$$\begin{aligned} I(u + \xi) &= I(u) + \iint_D \left[ \left( \frac{\partial \xi}{\partial x} \right)^2 + \left( \frac{\partial \xi}{\partial y} \right)^2 \right] dx dy + \\ &\quad + 2 \iint_D \left( \frac{\partial u}{\partial x} \frac{\partial \xi}{\partial x} + \frac{\partial u}{\partial y} \frac{\partial \xi}{\partial y} + \xi f \right) dx dy. \end{aligned}$$

Остается проверить, что третье слагаемое в правой части обращается в нуль. Действительно, из очевидных тождеств

$$\begin{aligned} \frac{\partial u}{\partial x} \frac{\partial \xi}{\partial x} + \frac{\partial u}{\partial y} \frac{\partial \xi}{\partial y} &\equiv \frac{\partial}{\partial x} \left( \xi \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( \xi \frac{\partial u}{\partial y} \right) - \xi \left( \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) \equiv \\ &\equiv \frac{\partial}{\partial x} \left( \xi \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( \xi \frac{\partial u}{\partial y} \right) - \xi f \end{aligned}$$

и из леммы 2 следует

$$\begin{aligned} \iint_D \left( \frac{\partial u}{\partial x} \frac{\partial \xi}{\partial x} + \frac{\partial u}{\partial y} \frac{\partial \xi}{\partial y} + \xi f \right) dx dy &= \\ &= \iint_D \left[ \frac{\partial}{\partial x} \left( \xi \frac{\partial u}{\partial x} \right) + \frac{\partial}{\partial y} \left( \xi \frac{\partial u}{\partial y} \right) \right] dx dy = 0. \end{aligned}$$

Но если  $\xi \in W^N$  не есть тождественный нуль, то в силу формулы (6) и леммы 1 очевидно, что  $I(w) > I(u)$ . Теорема доказана.  $\square$

Таким образом, задача (1), (2) допускает следующую вариационную постановку: среди всех функций  $w \in W^N$  найти ту, которая придает наименьшее значение выражению  $I(w)$ , определенному формулой (5).

Подчеркнем, что в силу доказанной теоремы выбор пространства  $W^N$ , фигурирующего в вариационной постановке задачи (1), (2), допускает произвол: выбор разбиения  $\{D_k\}$  остается в руках вычислиеля.

Этим можно пользоваться при дискретизации вариационной постановки задачи (1), (2) для ее численного решения.

**2. Вариационный метод Ритца.** Естественно ожидать и можно строго показать (см., напр., [6], с. 338 и далее), что приближением для решения  $u$  задачи (1), (2) могут служить те функции  $w$  из числа допустимых  $w \in W^N$ , на которых функционал  $I(w)$  принимает значения, близкие к минимальному значению  $I(u)$ .

Формальная схема отыскания приближенного решения по методу Ритца состоит в следующем. Зададим натуральное  $N$  и фиксируем какие-нибудь  $N$  линейно независимых функций

$$\omega_1(x, y), \omega_2(x, y), \dots, \omega_N(x, y); \quad \omega_n \in W^N. \quad (7)$$

Будем искать теперь такую линейную комбинацию функций (7)

$$\sum_{n=1}^N a_n \omega_n \equiv w_N(x, y, a_1, a_2, \dots, a_N), \quad (8)$$

где  $a_1, \dots, a_N$  — вещественные числа, которая придает наименьшее значение функционалу  $I(w)$  на множестве всех функций вида (8). Эту функцию вида (8) примем за приближенное решение при сделанном выборе базисных функций (7). Задача об отыскании функции  $w_N(x, y, a_1, \dots, a_N)$  несравненно проще задачи отыскания точного решения  $u(x, y)$ . В самом деле, речь идет об отыскании  $N$  чисел  $a_1, a_2, \dots, a_N$ , придающих наименьшее значение следующей квадратичной функции:

$$I(w_N) = \int_D \left[ \left( \frac{\partial}{\partial x} \sum_n a_n \omega_n \right)^2 + \left( \frac{\partial}{\partial y} \sum_n a_n \omega_n \right)^2 \right] dx dy + \\ + 2 \int_D f \sum_n a_n \omega_n dx dy. \quad (9)$$

Первое слагаемое в правой части (9) является квадратичной формой от  $a_1, \dots, a_N$ . Ввиду линейной независимости системы функций (7) и леммы 1 форма (9) положительно определена, а в силу этого функция (9) имеет единственный минимум. Этот минимум достигается при тех значениях  $a_n, n = 1, 2, \dots, N$ , которые удовлетворяют системе линейных алгебраических уравнений

$$\frac{\partial I[w_N(x, y, a_1, \dots, a_N)]}{\partial a_n} = 0, \quad n = 1, 2, \dots, N. \quad (10)$$

Подробно система (10) может быть записана в виде

$$\sum_{i=1}^N a_i \int_D \left( \frac{\partial \omega_i}{\partial x} \frac{\partial \omega_n}{\partial x} + \frac{\partial \omega_i}{\partial y} \frac{\partial \omega_n}{\partial y} \right) dx dy = - \int_D f \omega_n dx dy, \quad (11)$$

$$n = 1, 2, \dots, N.$$

Компоненты решения  $(a_1, a_2, \dots, a_N)$  системы (11) и служат коэффициентами в формуле (8) для функции  $w_N$ , которую мы приняли за приближенное решение задачи (1), (2).

**3. Метод конечных элементов на базе вариационного подхода.** Описанная выше вариационная схема приближенного вычисления решения задачи (8) при фиксированном выборе пространства  $W^N$

содержит произвол в выборе системы линейно независимых функций (7).

Вообще говоря, система линейных уравнений (11) имеет заполненную матрицу, т. е. все ее элементы могут отличаться от нуля. Таким образом, в каждое уравнение системы (11) могут входить значения искомого приближенного решения задачи (1), (2) во всех узлах  $P_1, P_2, \dots, P_N$ , так что система (11) не обладает удобством разностных схем, каждое уравнение которых связывает искомые значения только в нескольких узлах сетки, на которой подлежит вычислению искомое приближенное решение задачи (1), (2).

Однако можно так выбрать систему линейно независимых функций (7), чтобы возникшая система (11) связывала значения искомой функции только в соседних узлах сетки. При этом два узла  $P_i$  и  $P_j$  мы считаем соседними, если точки  $P_i$  и  $P_j$  являются вершинами одного и того же треугольника  $D_k$  разбиения области  $D$ .

Вариационные схемы, соответствующие такому выбору системы линейно независимых функций (7), называют *вариационно-разностными*. Вариационно-разностные схемы называют также схемами метода конечных элементов на базе вариационного подхода. Происхождение термина «конечный элемент» мы поясним в конце этого пункта.

Приведем простейший пример такого выбора системы функций (7), который приводит к вариационно-разностной схеме.

Определим функцию  $\omega_n(x, y)$  сначала в точках  $P_1, P_2, \dots, P_N$  и  $Q_1, Q_2, \dots, Q_M$ , положив

$$\omega_n(P_j) = \delta_n^j, \quad j, n = 1, 2, \dots, N;$$

$$\omega_n(Q_m) = 0, \quad m = 1, 2, \dots, M, \quad n = 1, 2, \dots, N.$$

Далее на пересечении треугольника  $D_k$  с многоугольником  $Q_1Q_2Q_3\dots, Q_MQ_1$  доопределим функцию  $\omega_n$  линейно по ее значениям в вершинах этого треугольника, а в точках  $(x, y) \in D$ , не принадлежащих этому многоугольнику, положим  $\omega_n(x, y) = 0$ .

Заметим, что значение искомого решения (8) в точке  $P_n$  совпадает со значением коэффициента  $a_n$ . При этом в уравнение (11) при фиксированном  $n$  войдут искомые величины  $u(P_n)$ , а также значения  $u(P_i)$  только в точках, соседних с  $P_n$ . В самом деле, пусть  $u(P_i)$  не является вершиной никакого треугольника  $D_k$  разбиения, для которого  $P_n$  является вершиной. При сделанном выборе функций (7) очевидно, что в этом случае носители (области отличия от нуля) функций  $\omega_n(x, y)$  и  $\omega_i(x, y)$  не пересекаются, так что коэффициент при  $a_i$  в уравнении (11) обращается в нуль.

В заключение заметим, что совокупность треугольников  $D_k$  разбиения  $D$ , имеющих общую вершину  $P_n$ , дополненную соответствующей функцией  $\omega_n(x, y)$ , называют иногда *конечным элементом*, что и дало название методу конечных элементов.

**4. Проекционный метод Галёркина.** Б.Г. Галёркин в 1916 г. предложил численный метод решения краевых задач, не требующий знания их вариационной постановки. Изложим общую схему этого метода применительно к задаче (1), (2), а затем тот ее вариант, который называют проекционно-разностным методом, или методом конечных элементов.

Вновь выберем систему базисных функций (7), но будем считать (временно), что функции  $w_n(x, y)$  имеют непрерывные вторые производные всюду в  $D$ . Вновь будем искать приближенное решение  $w_N(x, y)$  в виде (8). Подставим выражение (8) в левую часть уравнения (1). Получим

$$\frac{\partial^2}{\partial x^2}[w_N(x, y, a_1, \dots, a_N)] + \frac{\partial^2}{\partial y^2}[w_N(x, y, a_1, \dots, a_N)] - f(x, y) = \\ = \delta_N(x, y, a_1, \dots, a_N).$$

где  $\delta_N(x, y, a_1, \dots, a_N)$  — возникающая невязка. Если бы  $\delta_N$  оказалась ортогональной ко всем функциям  $w(x, y) \in W^N$  в смысле скалярного умножения  $(w'(x, y), w''(x, y)) = \int_D w' w'' dx dy$ , то невязка  $\delta_N(x, y, a_1, \dots, a_N)$  была бы точным нулем, а функция  $w_N$  — точным решением. Поэтому для построения приближенного решения (8) подберем параметры  $a_1, a_2, \dots, a_N$  так, чтобы проекции невязки  $\delta_N(x, y, a_1, \dots, a_N)$  на каждую из базисных  $w_n$  были равны нулю, то есть чтобы невязка была ортогональна ко всем базисным функциям (7):

$$(\delta_N, w_n) = 0, \quad n = 1, 2, \dots, N. \quad (12)$$

В развернутом виде эта система запишется так:

$$\int_D \left( \frac{\partial^2 w_N}{\partial x^2} + \frac{\partial^2 w_N}{\partial y^2} \right) \omega_n dx dy = \int_D f \omega_n dx dy, \quad n = 1, 2, \dots, N. \quad (13)$$

Преобразуем левую часть равенства (13). Интегрируя по частям, видим, что благодаря  $\omega_n|_{\Gamma} = 0$  имеет место равенство

$$\int_D \left( \frac{\partial^2 w_N}{\partial x^2} + \frac{\partial^2 w_N}{\partial y^2} \right) \omega_n dx dy = - \int_D \left( \frac{\partial w_N}{\partial x} \frac{\partial \omega_n}{\partial x} + \frac{\partial w_N}{\partial y} \frac{\partial \omega_n}{\partial y} \right) dx dy = \\ = - \sum_{i=1}^N a_i \int_D \left( \frac{\partial \omega_i}{\partial x} \frac{\partial \omega_n}{\partial x} + \frac{\partial \omega_i}{\partial y} \frac{\partial \omega_n}{\partial y} \right) dx dy.$$

Поэтому система (13) в точности совпадает с системой (11) метода Ритца.

От дополнительного предположения о наличии вторых производных у базисных функций, сделанного при выводе уравнений системы (13), можно отказаться.

Строгое обоснование последнего утверждения можно получить, аппроксимируя функции  $\omega \in W^N$  функциями, имеющими непрерывные вторые производные всюду в  $\bar{D}$ .

**5. Метод конечных элементов на базе проекционного метода Галёркина.** Используя ту же сетку (рис. 40) точек  $Q_m, m = 1, 2, \dots, M; P_n, n = 1, 2, \dots, N$ , и те же кусочно линейные базисные функции, что и при построении вариационного варианта метода конечных элементов, мы получим такую конкретизацию проекционного метода, которая называется проекционно-разностной схемой. Эту схему называют также схемой метода конечных элементов на базе метода Галёркина.

В нашем примере схемы метода конечных элементов на базе метода Ритца и на базе метода Галёркина совпадают. Однако схема на базе метода Галёркина допускает более широкие приложения, так как не требует существования удобной вариационной постановки исходной краевой задачи.

**6. Библиографический комментарий.** Методу конечных элементов и его приложениям посвящена обширная литература. Имеется большое количество компьютерных программ, реализованных для решения задач теории упругости, динамики вязкой жидкости и др.: см. [25–27] и Интернет.

### § 3. Вычисление решений сеточных аналогов краевых задач

Изложим и сопоставим некоторые употребительные способы вычисления решений больших систем линейных алгебраических уравнений, возникающих при дискретизации эллиптических краевых задач с помощью простейших разностных схем на регулярных сетках или схем метода конечных элементов, допускающих существенно нерегулярные сетки.

**1. Точное решение разностного аналога задачи Дирихле для уравнения Пуассона в квадратной области.** Разностная схема, о которой идет речь, изложена в § 1. В случае однородных краевых условий, т. е.  $u|_G = 0$ , решение можно выписать в виде конечного ряда Фурье. Соответствующие формулы приведены в § 7 из гл. 4. В случае, если число  $N$  ( $h = 1/N$ ) точек сетки является степенью двух:  $N = 2^k$ ,  $k = \log_2 N$ , вычисления можно организовать так, чтобы потребовалось всего  $O(N^2 \ln N)$  арифметических операций. Обычно в математическом обеспечении компьютера имеется программа, реализующая соответствующий алгоритм, который называют *быстрым преобразованием Фурье* (БПФ).

В случае  $\psi(s) \neq 0$  также можно воспользоваться методом Фурье. Для этого достаточно заметить, что значения решения  $u_{mn}$  во внутренних точках  $(mh, nh) \in D_h^0$  совпадают с решением разностной задачи Дирихле с нулевыми условиями на границе и с правой частью

$\tilde{\varphi}_{mn}((mh, nh) \in D_h^0)$ , которая совпадает с  $\varphi_{mn}$  в тех точках, которые более чем на  $h$  отстоят от границы, а в приграничных точках  $\tilde{\varphi}_{mn}$  определяется формулой

$$\tilde{\varphi}_{mn} = \varphi_{mn} - \frac{1}{h^2} [\tilde{\psi}(s_{m+1,n}) + \tilde{\psi}(s_{m-1,n}) + \tilde{\psi}(s_{m,n+1}) + \tilde{\psi}(s_{m,n-1})],$$

$$\tilde{\psi}(s_{i,j}) = \begin{cases} 0, & (i, j) \in D_h^0, \\ \psi(s_{i,j}), & (i, j) \in \Gamma_h. \end{cases}$$

**2. Методы простых итераций, трехслойный метод Чебышёва и метод сопряженных градиентов.** В гл. 5 мы рассмотрели названные методы для решения абстрактных систем линейных уравнений вида

$$Ax = f; \quad A: R^n \rightarrow R^n, \quad A = A^* > 0.$$

Пусть  $\mu(A)$  — число обусловленности оператора  $A$ . Были установлены формулы, показывающие, что для уменьшения погрешности  $\varepsilon_0$  нулевого приближения  $x_0$  в заданное число  $\varepsilon > 0$  раз требуется найти приближение  $x_p$  с номером  $p$ , который выражается соответственно формулами вида

$$p_1 = p_1(\mu, \varepsilon) = O(\mu \ln \varepsilon),$$

$$p_2 = p_2(\mu, \varepsilon) = O(\sqrt{\mu} \ln \varepsilon),$$

$$p_3 = p_3(\mu, \varepsilon) = O(\sqrt{\mu} \ln \varepsilon).$$

При этом первые два метода вычислительно устойчивы, а метод сопряженных градиентов при больших значениях числа  $\mu = \mu(A)$  может оказаться неустойчивым.

Для ускорения сходимости можно воспользоваться эквивалентными по спектру, или энергетически эквивалентными, операторами. Оператор  $B = B^* > 0$  энергетически эквивалентен оператору  $A$  с постоянными  $\gamma_1 > 0$ ,  $\gamma_2 > 0$ , если

$$\gamma_1(Bx, x) \leq (Ax, x) \leq \gamma_2(Bx, x).$$

Если решение задачи вида  $Bx = \varphi$  легко вычисляется, а  $\gamma_2/\gamma_1 \ll \mu(A)$ , то имеет смысл перейти от системы  $Ax = f$  к системе  $Cx = \psi$ , где  $C = B^{-1}A$ ,  $\psi = B^{-1}f$ , и решать эту последнюю итерациями. Ускорение сходимости будет достигнуто за счет того, что  $\mu_B(C) \leq \gamma_2/\gamma_1 \ll \mu$ .

Все эти соображения и факты можно использовать при вычислении решений систем, возникающих при дискретизации эллиптических задач. Свойство  $A = A^* > 0$  для них часто имеет место.

Для системы, возникающей в методе конечных элементов, это свойство вытекает из вариационного происхождения этой системы, так как матрица системы есть матрица Грама  $\omega^N$  из § 2 системы линейно независимых элементов.

Ускорение сходимости за счет выбора оператора  $B$  легко получить для задачи Дирихле в случае эллиптического уравнения с переменными коэффициентами в квадратной области, используя в качестве  $B$

оператор для задачи с постоянными коэффициентами, для обращения которого можно воспользоваться алгоритмом БПФ.

В случае более сложных областей, краевых условий или нерегулярных сеток ускорение за счет построения оператора  $B$ , который был бы легко обратим и для которого  $\mu_B(C) \ll \mu(A)$ ,  $C = B^{-1}A$ , затруднительно.

Подробное изложение итерационных методов сеточных аналогов краевых задач см. в [7, 20].

В настоящее время широкое распространение и развитие (особенно в связи с уравнениями метода конечных элементов на нерегулярных сетках) получил многосеточный метод Федоренко, которому мы посвятим следующий параграф.

## § 4. Многосеточный метод Федоренко

В работе Р.П. Федоренко (Релаксационный метод решения разностных эллиптических уравнений // ЖВМ и МФ. — 1961. — Т. 1, № 5. — С. 922–927) предложен метод итерационного решения разностных эллиптических задач, названный *релаксационным*. Для уменьшения нормы первоначальной погрешности вдвое этот метод требует всего  $CN$  арифметических действий, где  $C = \text{const}$ , а  $N$  — число точек в области. Очевидно, что по порядку числа арифметических действий этот метод неулучшаем, так как число расчетных точек также  $O(N)$ .

Границы применимости метода Федоренко почти такие же, как у метода простых итераций. Регулярность сетки не предполагается. Дополнительным ограничением является требование «плавности», «гладкости» первых собственных функций оператора задачи, которое для эллиптических задач обычно выполнено.

Упомянутая работа Федоренко была развита в ряде последующих работ многих авторов.

Обоснования многосеточного метода Федоренко довольно громоздки, поэтому мы ограничимся качественным описанием идеи метода и самого алгоритма Федоренко в простом случае, отсылая за доказательствами и общими конструкциями к оригинальным работам и [24].

**1. Идея метода.** При решении итерациями задачи

$$\begin{aligned} \Delta_h u_{mn} - \varphi_{mn} &= 0, \quad m, n = 1, 2, \dots, M-1, \\ u_{mn}|_{\Gamma_h} &= \psi(s_{mn}) \end{aligned} \tag{1}$$

будем отправляться от метода простых итераций

$$\begin{aligned} u_{mn}^{p+1} &= u_{mn}^p + \tau(\Delta_h u_{mn}^p - \varphi_{mn}), \quad m, n = 1, 2, \dots, M-1, \\ p &= 0, 1, \dots, \\ u_{mn}^{p+1}|_{\Gamma_h} &= \psi(s_{mn}), \quad [u_{mn}^0] \text{ задано}, \end{aligned} \tag{2}$$

который в целом сходится очень медленно, но неравномерно на различных гармониках. Погрешность  $\varepsilon^p = u^p - u$  в соответствии с § 7 гл. 4 записывается в виде конечного ряда Фурье

$$\varepsilon^p = \sum [\lambda_{rs}(\tau)]^p c_{rs}^{(0)} \psi^{(r,s)}, \quad (3)$$

где  $c_{rs}^{(0)}$  — коэффициенты разложения погрешности  $\varepsilon^0 = u - u^0$  нулевого приближения,  $\lambda_{rs} = 1 - 4\tau M^2 [\sin^2(\pi r/(2M)) + \sin^2(\pi s/(2M))]$ . Числа  $\lambda_{rs}$  лежат на отрезке  $\lambda_{лев} \leq \lambda \leq \lambda_{прав}$ , где

$$\lambda_{лев} = \lambda_{M-1,M-1} \approx 1 - 8\tau M^2, \quad \lambda_{прав} = \lambda_{11} \approx 1 - 2\pi^2\tau.$$

Положим

$$\tau = \frac{3}{16M^2}. \quad (4)$$

Если при этом условии хотя бы одно из чисел  $r$  или  $s$  больше, чем  $M/2$ , то

$$|\lambda_{rs}| < 0,6. \quad (5)$$

Поэтому вклад высокочастотных гармоник  $\psi^{(r,s)}$  ( $r > M/2$  или  $s > M/2$ ) в погрешность (3) за один шаг итерационного процесса уменьшается почти вдвое и вскоре становится малым. После нескольких итераций по формуле (2) погрешность будет в основном состоять только из гладких компонент (гармоники  $\psi^{(r,s)}$ ,  $r < M/2$ ,  $s < M/2$ ), потому что низкочастотные гармоники  $\psi^{(r,s)}$  умножаются на числа  $\lambda_{rs}^p$ , которые ближе к единице. Очень медленно гасится вклад первой гармоники  $\psi^{(1,1)}$  при сделанном выборе  $\tau$ , так как

$$\lambda_{11} \approx 1 - \frac{3\pi^2}{8M^2} (\approx 1). \quad (6)$$

Обозначим полученное в процессе итераций (2) приближение  $u^p$  через  $U$ , а погрешность  $\varepsilon^p = u^p - u = U - u$  через  $v$ . Если бы мы знали погрешность  $v$ , то нашли бы искомое решение  $u = U - v$ . Однако относительно  $v$  мы знаем только, что она удовлетворяет уравнению

$$\Delta_h v = \xi, \quad v|_{\Gamma_h} = 0, \quad (7)$$

где известная сеточная функция  $\xi$  — это невязка, возникающая при подстановке  $u^p = U$  в уравнение (1):

$$\Delta_h u^p - \Delta_h u = \Delta_h U - \varphi.$$

Задача (7) для определения поправки  $v$  проще исходной задачи (1) лишь в том отношении, что заранее известно, что  $v$  — гладкая сеточная функция. Поэтому для определения  $v$  вместо задачи (7) можно приближенно рассматривать такую же задачу на вдвое более крупной сетке, которая (при четном  $M$ ) является подсеткой исходной сетки:

$$\Delta_{2h} v^* = \xi^*, \quad v^*|_{\Gamma_{2h}} = 0. \quad (1^*)$$

Звездочкой мы обозначили величины на укрупненной сетке. Задачу (1\*) решаем итерациями по формулам

$$(v_{mn}^*)^{p+1} = (v_{mn}^*)^p + \tau^* [\Delta_{2h}(v_{mn}^*)^p - \xi_{mn}^*],$$

$$m, n = 1, 2, \dots, M^* - 1,$$

$$(v_{mn}^*)^{p+1} \Big|_{\Gamma_{2h}^*} = 0,$$
(2\*)

приняв за нулевое приближение  $(v_{mn}^*)^0 \equiv 0$ . Здесь

$$M^* = \frac{M}{2}, \quad \tau^* = 4\tau.$$

Каждый шаг итераций (2\*) вчетверо менее трудоемок, чем шаг итераций (2), потому что расчетных точек вчетверо меньше. Кроме того, благодаря  $\tau^* = 4\tau$  быстрее происходит погашение самой медленно убывающей компоненты погрешности. В соответствии с (6)

$$\lambda_{11}^* = 1 - \frac{3\pi^2}{8(M^*)^2} = 1 - 4 \frac{3\pi^2}{8M^2} < \lambda_{11},$$

и для уменьшения вклада  $(\psi^*)^{(1,1)}$  в заданное число раз нужно вчетверо меньше итераций. Результат итераций по формуле (2\*) обозначим  $V^*$ . Проинтерполируем  $V^*$  на исходную сетку (линейно). Гладкие компоненты будут получены почти правильно. Возникающая при интерполяции погрешность будет мала относительно интерполируемой гладкой функции, но ее разложение Фурье будет содержать все гармоники (погрешность интерполяции негладкая из-за изломов в узлах интерполяции). Кроме того, негладкая компонента, не имеющая отношения к искомой поправке, при интерполяции тоже дает случайный вклад в негладкую составляющую полученной при интерполяции функции  $V$ . Итак, гладкая компонента разности  $U - V$  близка к гладкой компоненте искомого решения  $u = U - v$ , но негладкая компонента не очень мала и носит случайный характер.

Поэтому следует проделать еще несколько шагов исходного итерационного процесса (2), приняв  $U - V$  за начальное приближение. Это погасит привнесенную при интерполяции негладкую составляющую погрешности, которую итерации (2) гасят почти вдвое за один шаг.

**2. Описание алгоритма.** При большом  $M$  (мелкой сетке) задача (1\*) на укрупненной сетке все еще трудоемка. Поэтому при ее решении в свою очередь целесообразно проделать еще одно укрупнение сетки вдвое, а при решении задачи на вчетверо укрупненной сетке снова использовать процесс удвоения сетки и увеличения  $h$  и т. д. Будем считать для простоты, что  $M = 2^k$ , т. е.  $M$  является некоторой степенью двух.

На исходной сетке делаем несколько шагов итераций (2), чтобы «сгладить» погрешность. Погрешность нам неизвестна, поэтому можно следить за этим по невязке  $\Delta_h u^p - \varphi$ , которая тоже сглаживается.

Результат вычислений  $U = u^p$  запоминаем. Затем для поправки  $v$  рассматриваем задачу на укрупненной сетке, делаем несколько итераций (2\*), чтобы сгладить «поправку к поправке», и результат  $V^*$  запоминаем (он занимает вчетверо меньше места в памяти, чем  $U$ ). Для вычисления поправки к  $V^*$  рассматриваем задачу на еще вдвое укрупненной сетке, делаем несколько итераций с шагом  $\tau^{**} = 4\tau^* = 16\tau$  и запоминаем результат  $V^{**}$ . Этот процесс вычисления поправок к поправкам на вдвое укрупненных сетках продолжается  $k$  раз, пока не дойдет до самой крупной сетки и поправки  $\tilde{V}^{(k)*}$ .

Затем начинаем возвращение на мелкую сетку. Сначала с самой крупной сетки интерполируем полученную там последнюю поправку  $\tilde{V}^{(k)*}$  на вдвое более мелкую сетку, вносим эту проинтерполированную поправку в  $V^{(k-1)*}$  и делаем несколько итераций, чтобы погасить привнесенную при интерполяции погрешность. Результат этих итераций интерполируем на еще вдвое более мелкую сетку, уточняем с его помощью хранящуюся для этой сетки поправку  $V^{(k-2)*}$ , делаем несколько итераций и производим следующую интерполяцию. На предпоследнем шаге после внесения в  $\tilde{V}^*$  поправки и итераций получим поправку  $V^*$ , которую интерполируем на исходную сетку. Проделав несколько итераций (2) над  $u - V$ , получим результат.

## ЧАСТЬ IV

# МЕТОДЫ ГРАНИЧНЫХ УРАВНЕНИЙ ДЛЯ ЧИСЛЕННОГО РЕШЕНИЯ КРАЕВЫХ ЗАДАЧ

---

Разностные методы численного решения краевых задач имеют весьма широкую область применимости. Тем не менее существенные трудности для использования разностных методов возникают при разностной аппроксимации краевых условий общего вида (в частности, нелокальных) в областях произвольной формы, поскольку в результате этой аппроксимации должна получаться устойчивая и удобная для численного решения разностная схема. Существенную трудность для применения разностных методов представляют также внешние задачи с условиями на бесконечности.

Эти трудности в ряде случаев удается преодолеть, используя редукцию исходной краевой задачи относительно неизвестной функции в области к некоторым уравнениям относительно других неизвестных функций, определенных на границе этой области. К тому же при такой редукции на единицу уменьшается геометрическая размерность задачи (граница двумерной области — линия). Мы познакомим с двумя способами такой редукции, а также и с соответствующими методами численного решения уравнений на границе.

Первый метод редукции исходной задачи к уравнениям на границе есть метод граничных интегральных уравнений (ГИУ) классической теории потенциала, идущей от работ Фредгольма начала XX века. Дискретизация и численное решение возникающих граничных интегральных уравнений осуществляются с помощью специальных квадратурных формул, или так называемого метода граничных элементов (МГЭ).

Второй метод редукции исходной задачи приводит к таким уравнениям на границе, в структуру которых входит оператор проектирования (он был предложен Кальдероном в 1963 г.). Мы излагаем некоторую модификацию граничных уравнений с проекторами (ГУРП) Кальдерона, предпринятую для того, чтобы ГУРП стали удобны для дискретизации и численного решения методом разностных потенциалов (МРП), предложенным в 1969 г. Рябеньким [17].

ГУРП и МРП можно применять во многих случаях, когда ГИУ и МГЭ неприменимы. Мы ограничимся изложением идей, лежащих в основе названных методов, и описанием границ применимости этих методов, а также укажем литературу, где можно найти более полное изложение.

## ГЛАВА 12

**ГРАНИЧНЫЕ ИНТЕГРАЛЬНЫЕ УРАВНЕНИЯ  
И МЕТОД ГРАНИЧНЫХ ЭЛЕМЕНТОВ  
ДЛЯ ИХ ЧИСЛЕННОГО РЕШЕНИЯ**

Познакомимся со способами перехода от краевых задач к интегральным уравнениям на границе, с приемом, лежащим в основе дискретизации граничных интегральных уравнений, а также укажем трудности, ограничивающие класс задач, для которого применение метода достаточно естественно.

**§ 1. Способы редукции краевых задач к ГИУ**

Для выяснения основных идей ограничимся рассмотрением внутренних и внешних задач Дирихле и Неймана для уравнения Лапласа

$$\Delta u \equiv \frac{\partial^2 u}{\partial x_1^2} + \frac{\partial^2 u}{\partial x_2^2} + \frac{\partial^2 u}{\partial x_3^2} = 0.$$

Пусть  $D$  — ограниченная область в трехмерном пространстве с достаточно гладкой границей  $\Gamma$ ,  $\varphi(x)|_{\Gamma}$  — заданная функция. Рассмотрим следующие четыре задачи:

$$\Delta u = 0, \quad x \in D, \quad u|_{\Gamma} = \varphi(x)|_{\Gamma}; \quad (1)$$

$$\Delta u = 0, \quad x \in D, \quad \partial u / \partial n|_{\Gamma} = \varphi(x)|_{\Gamma}; \quad (2)$$

$$\Delta u = 0, \quad x \notin D, \quad u|_{\Gamma} = \varphi(x)|_{\Gamma}; \quad (3)$$

$$\Delta u = 0, \quad x \notin D, \quad \partial u / \partial n|_{\Gamma} = \varphi(x)|_{\Gamma}. \quad (4)$$

Задачи (1), (3) суть соответственно внутренняя и внешняя задачи Дирихле, а задачи (2) и (4) — соответственно внутренняя и внешняя задачи Неймана.

Решение внешних задач будем искать в классе функций  $u(x)$ , стремящихся к нулю при  $r \rightarrow \infty$ , где  $r = (x_1^2 + x_2^2 + x_3^2)^{1/2}$ .

В математической физике доказывается, что внутренняя и внешняя задачи Дирихле, а также внешняя задача Неймана всегда имеют одно и только одно решение. Внутренняя задача Неймана имеет решение (которое определено с точностью до произвольного постоянного слагаемого) только в том случае, если интеграл от функции  $\varphi$  по поверхности  $\Gamma$  обращается в нуль:

$$\int_{\Gamma} \varphi \, ds = 0, \quad (5)$$

где  $ds$  — элемент площади на поверхности  $\Gamma$ .

Для перехода от задач (1)–(4) к интегральным уравнениям используется функция

$$G(x) = -\frac{1}{4\pi} \cdot \frac{1}{r}. \quad (6)$$

Эта функция является фундаментальным решением уравнения Лапласа в трехмерном пространстве. Это означает, что стремящееся к нулю при  $r \rightarrow \infty$  решение уравнения

$$\Delta u = f(x)$$

с произвольной достаточно гладкой правой частью  $f(x)$ , обращающейся в нуль вне какой-либо сферы, можно записать формулой

$$u(x) = \int \int \int G(x - y) f(y) dy_1 dy_2 dy_3, \quad (7)$$

где интегрирование ведется по всему трехмерному пространству.

Для перехода от задач (1)–(4) к интегральным уравнениям на границе используются потенциал простого слоя

$$V(x) = \int_{\Gamma} G(x - y) \rho(y) ds_y \quad (8)$$

и потенциал двойного слоя

$$W(x) = \int_{\Gamma} \frac{\partial G(x - y)}{\partial n_y} \sigma(y) ds_y. \quad (9)$$

Здесь  $ds_y$  — элемент площади на поверхности  $\Gamma$ , ограничивающей область  $D$ ,  $n_y$  — направление внешней по отношению к  $D$  нормали в точке  $y \in \Gamma$ . Функции  $\rho(y)$  и  $\sigma(y)$  называются соответственно плотностями потенциалов простого и двойного слоев.

Известно (и легко видеть), что  $G(x)$  (а значит, и потенциалы (8), (9) при  $x \notin \Gamma$ ) удовлетворяет уравнению Лапласа при  $x \neq 0$ . Таким образом, потенциалы (8), (9) суть семейства гармонических функций, зависящих от произвольных плотностей  $\rho(y)$ ,  $\sigma(y)$ ,  $y \in \Gamma$ .

Будем искать решения задач (1), (3) в виде потенциалов двойного слоя (9), а решения задач (2), (4) в виде потенциалов простого слоя (8). В курсах уравнений математической физики показано, что возникают следующие интегральные уравнения:

для внутренней задачи Дирихле (1)

$$\sigma(x) + \frac{1}{2\pi} \int_{\Gamma} \sigma(y) \frac{\partial G}{\partial n_y} dS_y = -\frac{1}{2\pi} \varphi(x), \quad x \in \Gamma; \quad (10)$$

для внешней задачи Неймана (4)

$$\rho(x) - \frac{1}{2\pi} \int_{\Gamma} \rho(y) \frac{\partial G}{\partial n_y} dS_y = \frac{1}{2\pi} \varphi(x), \quad x \in \Gamma; \quad (11)$$

для внутренней задачи Неймана (2)

$$\rho(x) + \frac{1}{2\pi} \int_{\Gamma} \rho(y) \frac{\partial G}{\partial n_y} dS_y = -\frac{1}{2\pi} \varphi(x), \quad x \in \Gamma; \quad (12)$$

для внешней задачи Дирихле (3)

$$\sigma(x) - \frac{1}{2\pi} \int_{\Gamma} \sigma(y) \frac{\partial G}{\partial n_y} dS_y = \frac{1}{2\pi} \varphi(x), \quad x \in \Gamma. \quad (13)$$

В теории потенциала устанавливается, что уравнения (10), (11) имеют, и притом единственны, решения  $\sigma(y)$ ,  $\rho(y)$  при произвольной функции  $\varphi(x)$  ( $x \in \Gamma$ ). Уравнение (12) имеет решение лишь в случае выполнения условия (5). Это связано с существом исходной задачи (2), которая разрешима лишь при этом условии. Однако уравнение (13) также разрешимо не для всякой функции  $\varphi(x)$ , хотя исходная внешняя задача Дирихле (3) разрешима всегда. Переход к интегральному уравнению (13) в этом случае оказался непригоден.

В теории потенциала показывается, что потеря разрешимости уравнения (13) для внешней задачи Дирихле связана с тем, что внутренняя задача Неймана для уравнения Лапласа в случае  $\varphi(x) \equiv 0$  имеет нетривиальное (ненулевое) решение. Это обстоятельство называют *резонансом внутренней области*. В данном случае вместо (13) легко построить другое интегральное уравнение, которое равносильно исходной задаче, но мы здесь этого делать не будем.

Пусть вместо уравнения Лапласа рассматривается уравнение Гельмгольца

$$\Delta u + \mu^2 u = 0, \quad \mu > 0,$$

с условием излучения на бесконечности. Соответственно

$$G(x) = -\frac{1}{4\pi} \cdot \frac{e^{i\mu r}}{r}.$$

Тогда вид потенциалов простого и двойного слоев, а также интегральных уравнений (10)–(13) будет тот же (только  $G(x)$  другая). Внешняя задача Дирихле и внешняя задача Неймана будут иметь единственныесolutions при любой  $\varphi(x)$  ( $x \in \Gamma$ ) независимо от  $\mu$ .

Однако интегральное уравнение вида (13) для внешней задачи Дирихле будет иметь решение при произвольной  $\varphi(x)$  лишь в том случае, если при данном  $\mu$  внутренняя задача Неймана при  $\varphi(x) \equiv 0$  имеет только тривиальное решение, что имеет место не всегда. Аналогично обстоит дело и с разрешимостью интегрального уравнения вида (11) для внешней задачи Неймана: оно разрешимо при любой  $\varphi(x)$  лишь в том случае, если внутренняя задача Дирихле при  $\varphi(x) \equiv 0$  и данном  $\mu$  имеет только тривиальное решение.

Указанная трудность редукции к интегральному уравнению, связанная с наличием резонанса дополнительной области для уравнения Гельмгольца, существенно затрудняет расчеты. Мы не будем останавливаться на этом подробнее.

Интегральные уравнения теории потенциала конструктивно построены в настоящее время для основных краевых задач теории упругости и некоторых других, для которых известны удобные аналитические представления фундаментальных решений.

Для построения граничных интегральных уравнений наряду с теми или иными потенциалами часто используют связь между значениями решения  $u|_{\Gamma}$  и значениями  $\partial u / \partial n|_{\Gamma}$  на границе  $\Gamma$  области  $D$ . Эта связь получается из формулы Грина

$$u(x) = \int_{\Gamma} \left( G(x, y) \frac{\partial u}{\partial n_y} - \frac{\partial G}{\partial n_y} u \right) dS_y, \quad (14)$$

представляющей решение  $u(x)$  во внутренней точке  $x$  области  $D$  через значения  $u|_{\Gamma}$ ,  $\partial u / \partial n|_{\Gamma}$ .

Переходя к пределу при  $x \rightarrow x_0 \in \Gamma$  и учитывая скачок потенциала двойного слоя, т. е. скачок во втором слагаемом правой части (14), получаем интегральное соотношение, о котором идет речь. В случае задачи Неймана можно подставить значение  $\partial u / \partial n|_{\Gamma} = \varphi(x)$  и получить разрешимое уравнение Фредгольма второго рода относительно  $\partial u / \partial n|_{\Gamma}$ . Удобство по сравнению с (10) состоит в том, что вычисляется  $\partial u / \partial n|_{\Gamma}$ , а не вспомогательная плотность  $\sigma(y)$ .

## § 2. Границные элементы и дискретизация ГИУ

Дискретизация ГИУ, построенных в § 1, осложняется тем, что  $G(x, y)$  и  $\partial G / \partial n_y$  имеют особенности. Проиллюстрируем построение квадратурных формул, приближенно заменяющих интегралы, на примере уравнения (10) из § 1.

Поверхность  $\Gamma$ , являющуюся границей области  $D$ , разобьем на конечное число пересекающихся только, быть может, по своим границам криволинейных треугольников — граничных элементов. В каждом из граничных элементов искомую плотность будем считать постоянной. Перенумеруем граничные элементы и обозначим приближенное значение плотности внутри граничного элемента  $\Gamma_j$  через  $\sigma_j$ . Заменим интеграл, входящий в уравнение (10) из § 1, суммой:

$$\int_{\Gamma} \sigma(y) \frac{\partial G}{\partial n_y} dS_y \approx \sum_j \sigma_j \int_{\Gamma_j} \frac{\partial G}{\partial n_y} dS_y. \quad (1)$$

Интегралы, входящие в правую часть (1), зависят только от  $\Gamma_j$ , а также от  $x$  как от параметра.

Обозначим  $x_k$  какую-либо точку граничного элемента  $\Gamma_k$ , а значения интегралов обозначим  $a_{kj}$ :

$$a_{kj} = \int_{\Gamma_j} \frac{\partial G(x_k - y)}{\partial n_y} dS_y. \quad (2)$$

Перейдем теперь от уравнения (10) из § 1 к системе линейных алгебраических уравнений

$$2\pi\sigma_k + \sum_j a_{kj}\sigma_j = -\varphi(x_k), \quad k, j = 1, 2, \dots, J, \quad (3)$$

где  $J$  — число граничных элементов.

При стремлении числа граничных элементов к бесконечности, а диаметра граничных элементов  $\Gamma_j$  к нулю решение системы (3) стремится к искомой плотности  $\sigma(y)$  потенциала двойного слоя, которая является решением уравнения (10) из § 1. Само искомое решение внутренней задачи Дирихле представляется с помощью потенциала двойного слоя, который по найденной приближенно из (3) плотности можно вычислить. Вычисление коэффициентов  $a_{kj}$  осуществляется точно или приближенно. В качестве граничных элементов выбираются не обязательно криволинейные треугольники. Искомая функция внутри граничного элемента не обязательно заменяется постоянной: можно принять, что это многочлен заданного вида с неопределенными коэффициентами. Соответствующие изменения должны претерпеть уравнения (3).

Существует развитая техника построения граничных элементов, дискретизации граничных интегральных уравнений и точного или итерационного решения возникающих алгебраических систем. Мы ограничимся сказанным.

### § 3. Область применимости ГИУ для численного решения краевых задач

В настоящее время существуют пакеты программ для численного решения многих краевых задач: для уравнений Лапласа и Гельмгольца; для системы Ламе, описывающей упругие деформации в однородном или кусочно однородном материале; для системы Стокса, описывающей медленные течения вязкой несжимаемой жидкости, и многих других (см., напр., [13, 14]).

Очевидная выгода использования ГИУ по сравнению с методом конечных разностей состоит в понижении на единицу геометрической размерности задачи за счет перехода к уравнениям на границе, а также в автоматическом учете граничных условий и условий на бесконечности. Кроме того, переход к интегральным уравнениям позволяет иногда построить алгоритмы, автоматически учитывающие гладкость входных данных и искомых функций, т. е. алгоритмы без насыщения (см. [1], а также статью: Белых В.Н.// ДАН СССР. — 1989. — Т. 304, № 3. — С. 629–631).

Главное ограничение для непосредственного использования изложенной схемы вычисления решений краевых задач с помощью ГИУ состоит в том, что для построения интегральных уравнений (а главное, для их дискретизации) требуется располагать удобным представлением

ядер этих уравнений. Ядра выражаются через фундаментальные решения, которые допускают простое представление в виде формул лишь для некоторых уравнений с постоянными коэффициентами. Впрочем, метод граничных элементов во взаимодействии с различными итерационными процедурами применяют даже для некоторых нелинейных краевых задач.

Отметим еще одну трудность. Пусть фундаментальное решение известно и записано в виде простой формулы. Редукция краевой задачи к равносильному интегральному уравнению может оказаться непростым делом даже при этом условии, как мы объяснили на примере внутренних резонансов в случае задач Дирихле и Неймана для уравнений Лапласа и Гельмгольца.

## ГЛАВА 13

### МЕТОД РАЗНОСТНЫХ ПОТЕНЦИАЛОВ

Метод разностных потенциалов (МРП) предназначен для дискретного моделирования некоторых задач математической физики. В частности, одним из главных приложений МРП является численное решение внутренних и внешних краевых задач для линейных уравнений с частными производными. Наиболее известными в этой традиционной области приложений вычислительной математики являются *метод конечных разностей* (МКР) и *метод граничных элементов* (МГЭ).

МРП в случае его применимости объединяет достоинства МКР и МГЭ, предоставляет принципиально новые возможности и обходит некоторые трудности, присущие названным методам.

Именно, метод конечных разностей особенно удобен для решения задач в простых областях (квадрат, куб, круг, шар, тор и т. п.) при простых краевых условиях, но встречает дополнительные трудности при разностной аппроксимации сложных краевых условий, особенно если расчетная область имеет сложную форму.

Метод граничных интегральных уравнений в случае его применимости автоматически учитывает граничные условия и реализует понижение размерности задачи путем редукции задачи на границу расчетной области; однако этот метод требует знания удобного аналитического представления фундаментального решения уравнения задачи, что весьма ограничивает область применимости.

Метод разностных потенциалов объединяет указанные выше достоинства метода конечных разностей и метода граничных интегральных уравнений, но, в отличие от метода конечных разностей, не требует разностной аппроксимации граничных условий и, в отличие от метода

граничных интегральных уравнений, не требует знания фундаментального решения.

Метод разностных потенциалов уже нашел применения в вычислительных задачах газовой динамики, упругости, акустики, электродинамики, в математическом моделировании активного подавления шума.

Мы дадим общее представление о конструкциях и идеях применения МРП, используя для этих целей различные модельные задачи, связанные с уравнением Пуассона

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = f(x, y). \quad (\text{I})$$

Мы называем эти задачи модельными, чтобы подчеркнуть, что они допускают широкие обобщения, некоторые из которых изложены в [17].

Для дискретизации модельных задач мы будем использовать простейший пятиточечный разностный аналог уравнения Пуассона (I):

$$\sum_{n \in N_m} a_{mn} u_n = f_m, \quad (\text{II})$$

где суммирование ведется по точкам  $n$ , принадлежащим пятиточечному шаблону  $N_m$ . Именно, мы используем квадратную сетку  $m = (m_1 h, m_2 h)$  с шагом  $h$ ,  $m_1, m_2 = 0, \pm 1, \pm 2, \dots$ , и шаблон  $N_m$ , состоящий из пяти точек

$$N_m = \{(m_1, m_2), (m_1 \pm 1, m_2), (m_1, m_2 \pm 1)\}.$$

Для краткости мы пишем  $m = (m_1, m_2)$  и  $n = (n_1, n_2)$  вместо  $m = (m_1 h, m_2 h)$ ,  $n = (n_1 h, n_2 h)$  соответственно. Коэффициенты  $a_{mn}$  задаются формулами

$$a_{mn} = \begin{cases} -4h^{-2}, & \text{если } n = m, \\ h^{-2}, & \text{если } n = (m_1 \pm 1, m_2) \text{ или } n = (m_1, m_2 \pm 1). \end{cases}$$

В §1 мы сформулируем модельные задачи, в §2 построим разностные потенциалы и рассмотрим их свойства. В §3 мы используем развитый в §2 аппарат МРП для решения модельных задач из §1 с целью дать иллюстрации некоторых общих подходов и возможностей МРП.

## § 1. Постановка модельных задач

Мы поставим здесь следующие шесть задач.

**1. Внутренняя краевая задача.** Вычислить приближенно решение задачи Дирихле

$$\Delta u = 0, \quad (x, y) \in D^+, \\ u|_{\Gamma} = \varphi(s), \quad (1)$$

где  $D^+$  — ограниченная область с границей  $\Gamma = \partial D^+$ , а  $\varphi(s)$  — заданная функция длины дуги  $s$  вдоль границы  $\Gamma$ .

Будем считать для определенности, что  $D^+$  вместе с границей  $\Gamma = \partial D^+$  лежит внутри единичного квадрата  $D^0$ ,  $0 \leq x, y \leq 1$ .

**2. Внешняя краевая задача.** Вычислить приближенно значения ограниченного решения  $u(x, y)$  задачи

$$\Delta u = 0, \quad (x, y) \notin D^+, \quad u|_{\Gamma} = \varphi(s). \quad (2)$$

Воспользуемся вместо поставленной следующей модельной для нее задачей:

$$\Delta u = 0, \quad (x, y) \in D^-, \quad D^- = D^0 \setminus D^+, \quad u|_{\Gamma} = \varphi(s), \quad u|_{\partial D_0} = 0. \quad (3)$$

Задача (3) моделирует задачу (2) в любой фиксированной окрестности области  $D = D^+$  тем точнее, чем больше размеры квадрата  $D^0$  и чем дальше граница  $\partial D^0$  отстоит от границы  $\Gamma = \partial D$ , вне которой определено решение внешней задачи (2). Мы будем рассматривать здесь задачу (3) вместо задачи (2).

**3. Задача о построении искусственных граничных условий.** Пусть в квадрате  $D^0$  поставлена краевая задача

$$Lu = f(x, y), \quad (x, y) \in D^0, \quad (4)$$

$$u|_{\partial D^0} = 0, \quad (5)$$

для дифференциального уравнения (4). О задаче (4), (5) предполагается, что она имеет, и притом единственное, решение при любой правой части  $f(x, y)$ . Предположим далее, что решение этой задачи интересует нас не всюду в  $D^0$ , а только в некоторой (малой) подобласти  $D^+ \subset D^0$ , расположенной в окрестности центра квадрата  $D^0$ . Далее, пусть уравнение (4) вне этой подобласти  $D^+$  принимает вид

$$\Delta u = 0, \quad (x, y) \in D^- = D^0 \setminus D^+. \quad (6)$$

Введем искусственно границу  $\Gamma$ , которой не было в исходной задаче (4), (5), приняв за  $\Gamma$  границу  $\partial D^+$  расчетной подобласти  $D^+$ . Поставим задачу построить такое условие  $ru|_{\Gamma} = 0$  на искусственной границе  $\Gamma$ , чтобы решение задачи

$$Lu = f(x, y), \quad (x, y) \in D^+, \quad (7)$$

$$ru|_{\Gamma} = 0 \quad (8)$$

при любой  $f(x, y)$  совпадало на  $D^+$  с решением исходной задачи (4), (5), (6).

Условие (8) будем называть искусственным граничным условием (ИГУ).

Можно сказать, что условие (8) должно равносильно заменять уравнение Лапласа (6) вне расчетной подобласти  $D^+$ , а также граничное условие (5) на (удаленной) границе  $\partial D^0$  исходной области. В нашей модельной задаче условие (5) на удаленной границе  $\partial D^0$  используется взамен условия ограниченности решения уравнения (4) на бесконечности.

Можно сказать также, что ИГУ получаются путем переноса условия (5) с удаленной границы области  $D^0$  на (искусственную) границу  $\partial D^+$ , возникшую при выделении расчетной подобласти  $D^+$ .

**4. Задача о вычислении вклада каждой из двух заряженных подобластей в значения потенциала на границе между ними.** Будем интерпретировать решение  $u(x, y)$  задачи

$$\Delta u = f(x, y), \quad (x, y) \in D^0 = D^+ \cup D^-, \quad (9)$$

$$u|_{\partial D^0} = 0$$

как потенциал, индуцированный зарядами, которые распределены с плотностью  $f(x, y)$ .

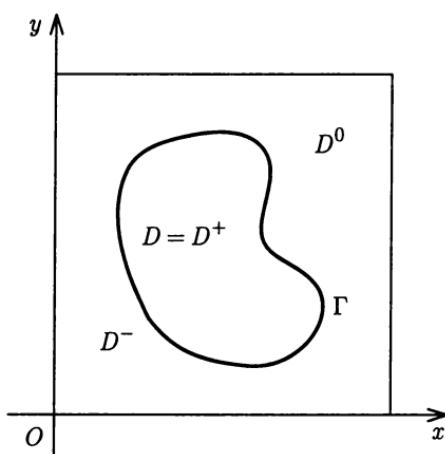


Рис. 41

Пусть функция  $f(x, y)$  неизвестна, но известен (например, может быть измерен) потенциал  $u(x, y)$  в окрестности границы  $\Gamma$  между подобластями  $D^+$  и  $D^-$ .

Задача состоит в том, чтобы, зная в окрестности  $\Gamma$  сумму

$$u = u^+ + u^-$$

решений задачи

$$\Delta u^+ = \theta(D^-)f(x, y), \quad u^+|_{\partial D^0} = 0 \quad (10)$$

и задачи

$$\Delta u^- = \theta(D^+)f(x, y), \quad u^-|_{\partial D^0} = 0, \quad (11)$$

где

$$\theta(\Omega) = \begin{cases} 1, & \text{если } (x, y) \in \Omega, \\ 0, & \text{если } (x, y) \notin \Omega, \end{cases}$$

определить каждое слагаемое в отдельности.

Более точная постановка задачи: зная сумму

$$\left( \begin{array}{c} u \\ \frac{\partial u}{\partial n} \end{array} \right) \Big|_{\Gamma} = \left( \begin{array}{c} u^+ \\ \frac{\partial u^+}{\partial n} \end{array} \right) \Big|_{\Gamma} + \left( \begin{array}{c} u^- \\ \frac{\partial u^-}{\partial n} \end{array} \right) \Big|_{\Gamma}, \quad (12)$$

определить каждое слагаемое в отдельности.

**5. Задача вычисления влияния одной заряженной подобласти на другую по известному суммарному потенциалу на их общей границе.** Пусть известна левая часть  $\left( \begin{array}{c} u \\ \frac{\partial u}{\partial n} \end{array} \right) \Big|_{\Gamma}$  равенства (12). Определить функции

$$u^+(x, y), \quad (x, y) \in D^+, \quad (13)$$

$$u^-(x, y), \quad (x, y) \in D^-, \quad (14)$$

не зная  $f(x, y)$ .

**6. Задача активного экранирования.** Пусть функция  $f(x, y)$  не известна исследователю, но заданы (например, измерены) значения потенциала  $u|_{\Gamma}$  и его нормальной производной  $\frac{\partial u}{\partial n}|_{\Gamma}$  на контуре  $\Gamma$ , разделяющем  $D^+$  и  $D^-$ .

Введем в правую часть (9) функцию-параметр — слагаемое  $\delta f(x, y)$ , которое будем называть активным управлением. Задача (9) перейдет в задачу

$$\Delta w = f(x, y) + \delta f(x, y), \quad (x, y) \in D^0, \quad (15)$$

$$w|_{\partial D^0} = 0.$$

Задача активного экранирования потенциала в подобласти  $D^+ \subset D^0$  от влияния зарядов, расположенных в подобласти  $D^- \subset D^0$ , состоит в следующем: надо построить все активные управления  $\delta f(x, y)$ , при включении которых решение  $w(x, y)$  задачи (15) на экранируемой подобласти  $D^+$  совпадает с решением  $v(x, y)$  задачи

$$\Delta v = \begin{cases} f(x, y), & (x, y) \in D^+, \\ 0, & (x, y) \in D^-, \end{cases} \quad (16)$$

$$v|_{\partial D^0} = 0. \quad (17)$$

**З а м е ч а н и е.** Подчеркнем, что решение  $u(x, y)$  задачи (9), решение  $w(x, y)$  задачи (15) и решение  $v(x, y)$  задачи (16), (17) не только не интересуют, но и не могут быть найдены по заданным входным данным. Наша цель состоит лишь в том, чтобы найти все те  $\delta f(x, y)$ , воздействие которых на потенциал  $u(x, y)$  в подобласти  $D^+$  равносильно удалению всех зарядов, расположенных в подобласти  $D^-$ , т. е. вне  $D^+$ .

Очевидно, что функция

$$\delta f(x, y) \equiv \begin{cases} 0, & (x, y) \in D^+, \\ -f(x, y), & (x, y) \in D^-, \end{cases} \quad (18)$$

является активным экранирующим управлением.

Но этим тривиальным экранирующим управлением нельзя воспользоваться, так как  $f(x, y)$  неизвестна. Однако даже если бы  $f(x, y)$  была известной функцией, тривиальное экранирующее управление (18) может оказаться неудобным.

## § 2. Разностные потенциалы

Мы построим здесь один из основных вариантов разностных потенциалов — разностные потенциалы типа Коши для простейшего пятиточечного аналога (II) уравнения Пуассона (I).

Развитый при этом аппарат будет использован затем для демонстрации подходов МРП на примерах решения модельных задач, поставленных в § 1.

Во всех конструкциях потенциалов, которые мы построим, используется та или иная вспомогательная разностная задача, к вычислению решения которой сводится вычисление потенциалов.

**1. Вспомогательная разностная задача.** Пусть  $\tilde{D}^0$  — некоторая область на плоскости  $Oxy$ . Совокупность точек  $m = (m_1 h, m_2 h)$ , попавших в эту область, обозначим  $M^0$ . Рассмотрим разностное уравнение (II) на множестве  $M^0$ :

$$\sum_{n \in N_m} a_{mn} u_n = f_m, \quad m \in M^0. \quad (1)$$

Левая часть этого уравнения имеет смысл, очевидно, для функций  $u_{N^0} = \{u_n\}$ ,  $n \in N^0$ , сеточная область определения  $N^0$  которых есть

$$N^0 = \bigcup N_m, \quad m \in M^0.$$

Дополним уравнение (1) тем или иным линейным однородным граничным условием, обозначив это условие следующим образом:

$$lu = 0.$$

Совокупность всех функций  $u_{N^0} = \{u_n\}$ ,  $n \in N^0$ , удовлетворяющих этому условию  $lu = 0$ , образует некоторое линейное пространство, которое мы обозначим  $U_{N^0}$ . Выполнение условия  $lu = 0$  в таком случае можно переписать в форме следующего включения:

$$u_{N^0} \in U_{N^0}. \quad (2)$$

Выбор дополнительного условия  $lu = 0$  и соответствующего пространства  $U_{N^0}$  подчиним только одному условию: задача (1), (2) должна иметь одно и только одно решение при произвольной функции  $f_{M^0} = \{f_m\}$ ,  $m \in M^0$ , стоящей в правой части уравнения (1). При

выполнении этого условия задачу (1), (2) будем называть вспомогательной разностной задачей.

Мы видим, что при построении вспомогательных разностных задач вида (1), (2) существует большая свобода в выборе области  $\tilde{D}^0$  и соответствующих сеточных областей  $M^0$  и  $N^0 = \bigcup N_m$ ,  $m \in M^0$ , а также в выборе дополнительного условия  $lu = 0$ , которое определяет подпространство  $U_{N^0}$  всех функций, определенных на  $N^0$ . В зависимости от нашего выбора  $M^0$  и  $U_{N^0}$  вспомогательная разностная задача (1), (2) может оказаться удобной или неудобной для вычисления ее решения.

Условимся считать для определенности, что  $\tilde{D}_0 = D^0$  есть единичный квадрат  $0 \leq x, y \leq 1$  со сторонами, лежащими на линиях сетки  $\{x = kh, y = lh\}$ , где  $k$  и  $l$  — какие-нибудь целые числа, а  $h$  — шаг квадратной сетки. Множества  $M^0$  и  $N^0$  в этом случае изображены на рис. 42.

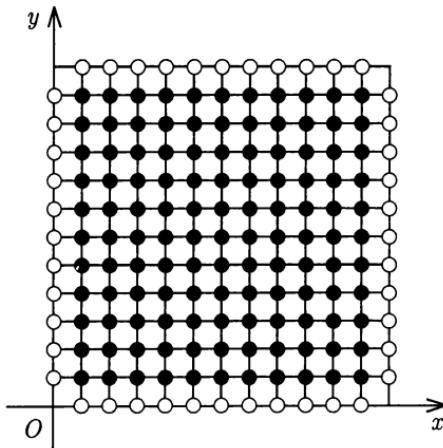


Рис. 42

Множество  $M^0$  состоит из точек, обозначенных черными кружками, а множество  $N^0 = \bigcup N_m$ ,  $m \in M^0$ , возникает при дополнении множества  $M^0$  теми точками сетки, которые обозначены белыми кружками.

Пространство  $U_{N^0}$  определим как пространство всех функций  $u_{N^0} = \{u_n\}$ ,  $n \in N^0$ , которые обращаются в нуль в точках  $n \in N^0$ , лежащих на сторонах квадрата.

В этом случае задача (1), (2) есть разностная задача Дирихле для уравнения Пуассона с нулевыми условиями на границе. Хорошо известно, что она имеет единственное решение при произвольной правой части  $f_{M^0} = \{f_m\}$  и что ее решение может быть вычислено, например, с помощью метода разделения переменных, или метода Фурье.

Если  $L/h = 2^k$ , где  $L$  — длина стороны квадрата, а  $k$  — некоторое целое число, то известен вариант метода Фурье (быстрое преобразование Фурье), требующий  $O(h^{-2}|\ln h|)$  арифметических операций для вычисления решения задачи (1), (2).

**2. Разностный потенциал  $P^+v_\gamma$ .** Обозначим через  $M^+$  множество точек  $m$ , лежащих внутри или на границе  $D^+$ , и рассмотрим уравнение

$$\sum a_{mn}u_n = f_m, \quad m \in M^+. \quad (3)$$

Левая часть этого уравнения имеет смысл для тех функций  $\{u_n\}$ , которые определены на множестве  $N^+ = \bigcup N_m$ ,  $m \in M^+$ .

Всюду дальше будем считать, что шаг  $h$  настолько мал, что все точки  $n \in N^+$  лежат внутри квадрата  $D_0$ , т. е.  $N^+ \subset M^0$ .

Мы собираемся определить разностный потенциал для решений уравнения (3). Рассмотрим какой-нибудь квадрат  $D^0$ , содержащий область  $D^+$  вместе с ее границей, и введем вспомогательную разностную задачу (1), (2). Обозначим

$$M^- = \{m \mid m \in D^-\} = M^0 \setminus M^+.$$

Рассмотрим систему (II) на множестве  $M^-$ :

$$\sum a_{mn}u_n = f_m, \quad m \in M^-. \quad (4)$$

Левая часть системы (4) имеет смысл для тех функций  $\{u_n\}$ , которые определены на множестве

$$N^- = \bigcup N_m, \quad m \in M^-.$$

Таким образом, система (1) распалась на две подсистемы (3) и (4), решения которых определены на  $N^+$  и  $N^-$  соответственно.

Определим границу  $\gamma$  между сеточными областями  $N^+$  и  $N^-$ , положив (рис. 43)

$$\gamma = N^+ \cap N^-.$$

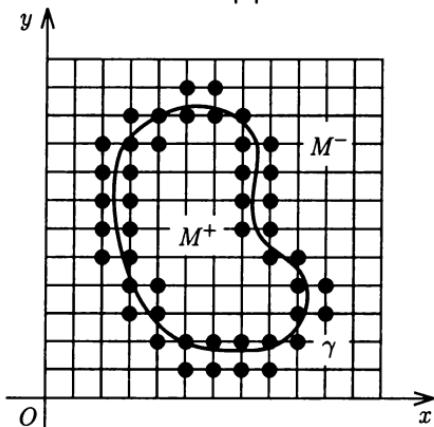


Рис. 43

Введем линейное пространство  $V_\gamma$  всех тех функций  $v_\gamma$ , которые определены на  $\gamma$ .

Функции  $v_\gamma$  из пространства  $V_\gamma$  будем называть *плотностями* и определим потенциал  $u^+ = P^+v_\gamma$  с плотностью  $v_\gamma$  из пространства плотностей  $V_\gamma$ .

**Определение 1.** Разностным потенциалом  $u^+ = P^+v_\gamma$ , с плотностью  $v_\gamma$  будем называть функцию  $u^+ = u_{N^+} = \{u_n\}$ , определенную на  $N^+$ , которая совпадает на  $N^+$ ,  $N^+ \subset N^0$ , с решением  $u_{N^0} = \{u_n\}, n \in N^0$ , вспомогательной разностной задачи (1), (2) с правой частью  $f_m$ , определенной равенствами

$$f_m = \begin{cases} 0, & \text{если } m \in M^+, \\ \sum_{n \in N_m} a_{mn} v_n, & \text{если } m \in M^-, \end{cases} \quad (5)$$

где

$$v_n = \begin{cases} v_\gamma|_n, & \text{если } n \in \gamma, \\ 0, & \text{если } n \notin \gamma. \end{cases} \quad (6)$$

Подчеркнем, что вычисление разностного потенциала  $u^+ = P^+v_\gamma$  с заданной плотностью  $v_\gamma \in V_\gamma$  сводится к вычислению решения вспомогательной разностной задачи (1), (2) с правой частью (5), (6).

**Теорема 1.** Разностный потенциал  $u^+ = u_{N^+} = P^+v_\gamma$  обладает следующими свойствами:

1) разностный потенциал  $u^+ = P^+v_\gamma$  удовлетворяет на  $M^+$  однородному уравнению

$$\sum a_{mn} u_n = 0, \quad m \in M^+; \quad (7)$$

2) если  $v_\gamma$  совпадает на  $\gamma$  с некоторым решением  $v_{N^+} = \{v_n\}$  однородного уравнения (7), которое может быть доопределено всюду на  $N^0$  до некоторой функции  $v_{N^0} \in U_{N^0}$ , то такое решение  $v_{N^+}$  единственно и совпадает с разностным потенциалом  $u^+ = P^+v_\gamma$ , т. е. в этом случае решение  $v_{N^+} = \{v_n\}, n \in N^+$ , восстанавливается по своим значениям на границе  $\gamma$  по формуле

$$v_{N^+} = P^+v_\gamma. \quad (8)$$

**Доказательство.** Утверждение 1) справедливо в силу равенства (5), из которого следует, что правая часть уравнения (1) задачи (1), (2), определяющей потенциал  $P^+v_\gamma$ , обращается в нуль для  $m \in M^+$ .

Докажем второе утверждение теоремы.

Пусть  $\{v_n\}, n \in N^+$  — решение однородного уравнения (7), сужение которого на  $\gamma$  совпадает с  $v_\gamma$ .

Рассмотрим функцию  $w_n \equiv \theta_{N^0}(N^+)v_{N^0}|_n, n \in N^0$ , которая тоже принадлежит пространству  $U_{N^0}, w_{N^0} \in U_{N^0}$ .

Очевидно, что функция  $w_{N^0} = \{w_n\} \in U_{N^0}$  удовлетворяет уравнению (1) с правой частью (5). Ввиду единственности решения задачи (1), (2) решение  $w_{N^0}$  этой задачи с правой частью (5) совпадает на  $N^+$  с разностным потенциалом  $v_{N^+} = P^+v_\gamma$ , что одновременно доказывает как единственность функции  $v_{N^+}$ , так и ее представимость по формуле (8).  $\square$

Обозначим  $P_\gamma^+ : V_\gamma \rightarrow V_\gamma$  оператор, которой определен равенством

$$P_\gamma^+ v_\gamma = P^+ v_\gamma|_\gamma,$$

т.е. сопоставляет каждой функции  $v_\gamma \in V_\gamma$  сужение  $v_\gamma$  разностного потенциала  $u^+ = P^+ v_\gamma$  со всего множества  $N^+$  на  $\gamma \subset N^+$ .

**Теорема 2.** Функция  $v_\gamma \in V_\gamma$  является сужением на  $\gamma$  (следом) некоторого решения  $v_{N^+} = \{v_n\}$  однородного уравнения (7) в том и только том случае, если выполняется равенство

$$P_\gamma^+ v_\gamma = v_\gamma. \quad (9)$$

**Доказательство.** Пусть  $v_\gamma$  допускает доопределение до некоторого  $v_{N^+}$ , о котором говорится в условии теоремы. Тогда в силу второго утверждения теоремы 1 это доопределение единственно и согласно формуле (8) восстанавливается по своему сужению  $v_\gamma$ . Рассматривая равенство (8) не на всем  $N^+$ , а только на  $\gamma \subset N^+$ , получаем (9).

Обратно, пусть (9) выполнено. Рассмотрим разностный потенциал  $u^+ = P^+ v_\gamma$  с плотностью  $v_\gamma$ . В силу теоремы 1 этот потенциал есть решение уравнения (7). С другой стороны,  $v_\gamma$  представляет собой граничные значения этого потенциала в силу определения  $P_\gamma^+ v_\gamma$  и (9):

$$u^+|_\gamma = P^+ v_\gamma|_\gamma = P_\gamma^+ v_\gamma = v_\gamma. \quad \square$$

**3. Разностный потенциал  $u^- = P^- v_\gamma$ .** Разностный потенциал  $P^- v_\gamma$  с плотностью  $v_\gamma$  есть функция, определенная на  $N^-$  совершенно аналогично тому, как был определен разностный потенциал  $u^+ = P^+ v_\gamma$ .

Для получения разностного потенциала  $P^- v_\gamma$  и его свойств достаточно в определении 1, теоремах 1 и 2 и их доказательствах поменять везде знак «+» на знак «-».

**4. Разностный потенциал  $w^\pm = P^\pm v_\gamma$  со скачком, или разностный потенциал типа Коши.**

**Определение 2.** Разностным потенциалом  $w^\pm = P^\pm v_\gamma$  типа Коши с плотностью  $v_\gamma$  назовем функцию  $w^\pm = P^\pm v_\gamma$ , определенную в точках  $n \in N^0$  равенствами

$$w_n^\pm = \begin{cases} w_n^+ = P^+ v_\gamma|_n, & \text{если } n \in N^+, \\ w_n^- = -P^- v_\gamma|_n, & \text{если } n \in N^-. \end{cases} \quad (10)$$

Очевидно, что разностный потенциал типа Коши с плотностью  $v_\gamma$  есть двузначная на  $\gamma$  функция, так как каждая точка  $n \in \gamma$  принадлежит одновременно как множеству  $N^+$ , так и  $N^-$ . Дадим еще и другое определение разностного потенциала  $w^\pm = P^\pm v_\gamma$  со скачком  $v_\gamma$ , которое окажется равносильным определению 2 разностного потенциала типа Коши с плотностью  $v_\gamma$ . Это новое определение будет полезно для получения свойств разностного потенциала  $w^\pm = P^\pm v_\gamma$ , а также для понимания глубокой аналогии между разностным потенциалом типа Коши и классическим интегралом типа Коши из теории аналитических функций.

Кроме того, новое определение разностного потенциала  $w^\pm = P^\pm v_\gamma$ , позволит вычислять все три разностных потенциала

$$u^+ = P^+ v_\gamma, \quad u^- = P^- v_\gamma, \quad w^\pm = P^\pm v_\gamma$$

одновременно, решая один раз вспомогательную задачу (1),(2) при некоторой специальной правой части  $f_m, m \in M^0$ .

Введем термины, которые понадобятся для нового определения разностного потенциала  $w^\pm = P^\pm v_\gamma$ .

Функции  $u_{N^0} \in U_{N^0}$  из пространства  $U_{N^0}$  будем называть *регулярными*. Пусть  $u_n^+$  и  $u_n^-$  — две какие-нибудь регулярные функции. Определим *кусочно регулярную* функцию  $u_n^\pm, n \in N^0$ , положив

$$u_n^\pm = \begin{cases} u_n^+, & \text{если } n \in N^+, \\ u_n^-, & \text{если } n \in N^-. \end{cases} \quad (11)$$

Введем линейное пространство  $U^\pm$  всех кусочно регулярных функций вида (11). Функция (11) в каждой точке  $n \in \gamma$  границы  $\gamma$  принимает два значения  $u_n^+$  и  $u_n^-$ . Скачком  $v_\gamma = [u^\pm]$  назовем однозначную функцию  $v_\gamma$ , определенную в точках  $n \in \gamma$  формулой

$$v_\gamma|_n = [u^\pm]_n = u_n^+ - u_n^-, \quad n \in \gamma. \quad (12)$$

Заметим, что пространство  $U_{N^0}$  регулярных функций можно рассматривать как подпространство пространства  $U^\pm$ , состоящее из всех  $u^\pm \in U^\pm$ , имеющих нулевой скачок  $v_\gamma = 0_\gamma$ .

Кусочно регулярную функцию (11) будем называть *кусочно регулярным решением* задачи

$$\sum a_{mn} u_n^\pm = 0, \quad (13)$$

$$u^\pm \in U^\pm, \quad (14)$$

если функции  $u_n^+, n \in N^+$ , и  $u_n^-, n \in N^-$ , удовлетворяют однородным уравнениям

$$\sum a_{mn} u_n^+ = 0, \quad m \in M^+, \quad (15)$$

$$\sum a_{mn} u_n^- = 0, \quad m \in M^-, \quad (16)$$

**Теорема 3.** Пусть  $v_\gamma$  — произвольная функция из  $V_\gamma$ . Тогда существует одно и только одно кусочно регулярное решение задачи (13), (14), имеющее скачок  $v_\gamma$ . Это решение определяется формулой

$$u^\pm = v^\pm - u, \quad (17)$$

где уменьшаемое есть произвольная кусочно регулярная функция  $v^\pm \in U^\pm$ , имеющая заданный скачок  $v_\gamma$ , а вычитаемое есть регулярное решение разностной вспомогательной задачи вида (1), (2) с правой частью

$$f_m = \begin{cases} \sum a_{mn} v_n^+, & m \in M^+, \\ \sum a_{mn} v_n^-, & m \in M^-. \end{cases} \quad (18)$$

**Доказательство.** Прежде всего заметим, что существуют кусочно регулярные функции  $v^\pm \in U^\pm$  с заданным скачком  $v_\gamma$ . Одну из таких функций, очевидно, задает формула

$$v_n^\pm = \begin{cases} v_n^+, & n \in N^+, \\ v_n^-, & n \in N^-, \end{cases} \quad (19)$$

где

$$v_n^+ = \begin{cases} v_\gamma|_n, & \text{если } n \in \gamma \subset N^+, \\ 0, & \text{если } n \in N^+ \setminus \gamma, \end{cases} \quad (20)$$

$$v_n^- \equiv 0, \quad \text{если } n \in N^-. \quad (21)$$

Вычитаемое в формуле (17) существует, так как задача (1), (2) имеет одно и только одно решение при произвольной правой части  $f_m$ ,  $m \in M^0$ , в частности, заданной формулой (18).

В силу очевидной формулы

$$[u^\pm] = [v^\pm] - [u] = v_\gamma - [u]$$

и равенства

$$[u]|_n = 0, \quad n \in \gamma,$$

имеющего место для любой регулярной функции, можно утверждать, что функция  $u^\pm$ , задаваемая формулой (17), есть кусочно регулярная функция со скачком  $v_\gamma$ . Покажем, что функция  $u^\pm$ , задаваемая равенством (17), есть кусочно регулярное решение. Для этого надо проверить, что функции

$$u_n^+ = v_n^+ - u_n, \quad n \in N^+, \quad (22)$$

$$u_n^- = v_n^- - u_n, \quad n \in N^-, \quad (23)$$

удовлетворяют однородным уравнениям (15) и (16) соответственно. Проверим наше утверждение, относящееся к функции (22). Подставим  $u_n^+$  в уравнение (15):

$$\sum a_{mn} u_n^+ = \sum a_{mn} v_n^+ - \sum a_{mn} u_n, \quad m \in M^+. \quad (24)$$

Но  $\{u_n\}$  — решение разностной вспомогательной задачи с правой частью (18). Поэтому

$$\sum a_{mn} u_n = f_m = \sum a_{mn} v_n^+, \quad m \in M^+.$$

Таким образом, уменьшаемое и вычитаемое в правой части (24) совпадают, так что (15) выполнено. Аналогично проверяется утверждение о том, что функция (23) удовлетворяет (16).

Осталось показать, что существует только одно решение задачи (13), (14), имеющее заданный скачок  $v_\gamma$ . Предположим, что таких решений два. Тогда разность этих решений есть некоторое кусочно регулярное решение с нулевым скачком. Тогда эта разность является регулярной функцией, т. е. принадлежит  $U_{N^0}$ . Но в этом случае задача (13), (14) совпадает со вспомогательной задачей (1), (2), имеющей

тождественно нулевую правую часть. Однако ввиду единственности решения вспомогательной разностной задачи (1), (2) соответствующее решение тождественно обращается в нуль, т. е. кусочно регулярные решения, имеющие один и тот же скачок, совпадают между собой.  $\square$

Доказанная теорема 3 делает корректным следующее определение.

**Определение 3.** *Разностным потенциалом  $u^\pm = P^\pm v_\gamma$  со скачком  $v_\gamma$  будем называть кусочно регулярное решение  $u^\pm$  задачи (13), (14), имеющее заданный скачок  $v_\gamma \in V_\gamma$ .*

**Теорема 4.** *Определение 2 разностного потенциала  $w^\pm = P^\pm v_\gamma$  типа Коши с плотностью  $v_\gamma \in V_\gamma$  и определение 3 разностного потенциала  $u^\pm = P^\pm v_\gamma$  со скачком  $v_\gamma$  равносильны, т. е.  $w^\pm = u^\pm$ .*

**Доказательство.** Очевидно, что формула (10) задает некоторое кусочно регулярное решение задачи (13), (14). Ввиду доказанной единственности решения задачи (13), (14) с заданным скачком  $v_\gamma$  для доказательства теоремы остается показать, что кусочно регулярное решение (10) имеет скачок  $v_\gamma$ , т. е. что

$$[w^\pm] = u_\gamma^+ - u_\gamma^- = P_\gamma^+ v_\gamma - (-P_\gamma^- v_\gamma) = P_\gamma^+ v_\gamma + P_\gamma^- v_\gamma = v_\gamma.$$

В этой цепочке равенств неочевидно только последнее, т. е. равенство

$$P_\gamma^+ v_\gamma + P_\gamma^- v_\gamma = v_\gamma, \quad v_\gamma \in V_\gamma, \quad (25)$$

которое и осталось доказать.

Напомним, что  $P_\gamma^+ v_\gamma$  совпадает на  $\gamma$  с разностным потенциалом  $u^+ = P^+ v_\gamma$ , который в свою очередь совпадает на  $N^+$  с решением вспомогательной разностной задачи

$$\sum a_{mn} u_n^+ = f_m^+, \quad m \in M^0, \quad u_{N^0}^+ \in U_{N^0} \quad (26)$$

с правой частью, определяемой равенствами (5), (6):

$$f_m^+ = \begin{cases} 0, & \text{если } m \in M^+, \\ \sum a_{mn} v_n, & \text{если } m \in M^-. \end{cases} \quad (27)$$

Совершенно аналогично, выражение  $P_\gamma^- v_\gamma$  совпадает на  $\gamma$  с разностным потенциалом  $u^- = P^- v_\gamma$ , который в свою очередь совпадает на  $N^- \subset N^0$  с решением вспомогательной разностной задачи

$$\sum a_{mn} u_n^- = f_m^-, \quad m \in M^0, \quad u_{N^0}^- \in U_{N^0} \quad (28)$$

с правой частью, которая определяется следующими равенствами:

$$f_m^- = \begin{cases} \sum a_{mn} v_n, & m \in M^+, \\ 0, & m \in M^-. \end{cases} \quad (29)$$

где  $v_n$  по-прежнему определяется по  $v_\gamma$  с помощью равенства (6).

Равенство (29) аналогично (5), но используется при определении потенциала  $u^- = P^- v_\gamma$ .

Сложим почленно равенства (26) и (28). Предварительно заметим, что

$$f_m^+ + f_m^- = \sum a_{mn} v_n, \quad m \in M^0. \quad (30)$$

Тогда очевидно, что

$$\sum a_{mn} (u_n^+ + u_n^-) = \sum a_{mn} v_n, \quad m \in M^0. \quad (31)$$

Функция  $v_{N^0}$ , определяемая по  $v_\gamma$  равенством (6), принадлежит пространству  $U_{N^0}$ . Функции  $u_{N^0}^+$  и  $u_{N^0}^-$ , будучи решениями соответственно задач (26) и (28) вида (1), (2), также принадлежат  $U_{N^0}$ . Поэтому сумма  $z_{N^0} \equiv u_{N^0}^+ + u_{N^0}^-$  тоже принадлежит  $U_{N^0}$ :

$$u_{N^0}^+ + u_{N^0}^- \in U_{N^0}, \quad (32)$$

и является решением следующей задачи вида (1), (2):

$$\sum a_{mn} z_n = \sum a_{mn} v_n, \quad m \in M^0, \quad z_{N^0} \in U_{N^0},$$

относительно функции  $z_{N^0}$ . Очевидно, что функция  $v_{N^0}$  удовлетворяет этому уравнению, а поэтому ввиду единственности решения  $z_{N^0} = v_{N^0}$ , или  $z_n = u_n^+ + u_n^- = v_n$ ,  $n \in N^0$ . Но тогда можно написать равенство

$$u_n^+ + u_n^- = P^+ v_\gamma|_n + P^- v_\gamma|_n = v_\gamma|_n, \quad n \in \gamma \subset N^0,$$

которое совпадает с доказываемым равенством (25).  $\square$

**Замечание.** Пусть потенциал  $u^\pm = P^\pm v_\gamma$  для заданного  $v_\gamma \in V_\gamma$  найден по формуле (17), требующей решения некоторой задачи вида (1), (2) с правой частью (18). Тем самым найдены функции  $u^+$  и  $u^-$ :

$$u_n^\pm = \begin{cases} u_n^+, & n \in N^+, \\ u_n^-, & n \in N^-. \end{cases} \quad (33)$$

В силу теоремы 3 найдены также  $P^+ v_\gamma$  и  $P^- v_\gamma$ , а именно:

$$P^+ v_\gamma|_n = u_n^+, \quad \text{если } n \in N^+;$$

$$P^- v_\gamma|_n = -u_n^-, \quad \text{если } n \in N^-.$$

**5. Аналогия между разностным потенциалом типа Коши и классическим интегралом типа Коши.** Пусть  $\Gamma$  — некоторый несамопересекающийся замкнутый контур, разбивающий комплексную плоскость  $z = x + iy$  на ограниченную  $D^+$  и дополнительную к ней неограниченную  $D^-$  части. Классический интеграл типа Коши

$$u^\pm(z) = \frac{1}{2\pi i} \oint_{\Gamma} \frac{v_\Gamma(\xi)}{\xi - z} d\xi, \quad (34)$$

где контур  $\Gamma$  обходится в направлении против часовой стрелки, можно определить как кусочно аналитическую стремящуюся к нулю на бесконечности функцию, претерпевающую скачок  $v_\Gamma = [u^\pm]_\Gamma$  на контуре  $\Gamma$ .

Здесь  $u^+(z)$  и  $u^-(z)$  – значения интеграла (34) для  $z \in D^+$  и  $z \in D^-$  соответственно.

Интеграл типа Коши можно интерпретировать в качестве некоторого потенциала для решений системы Коши–Римана

$$\frac{\partial a}{\partial x} = \frac{\partial b}{\partial y}, \quad \frac{\partial b}{\partial x} = -\frac{\partial a}{\partial y}, \quad (35)$$

которая связывает вещественную и мнимую части аналитической функции.

Таким образом, разностный потенциал типа Коши

$$w^\pm = P^\pm v_\gamma \quad (36)$$

играет для решений уравнения вида (1) ту же роль, которую интеграл типа Коши играет для решений системы Коши–Римана.

### § 3. Решение модельных задач

Мы используем разностные потенциалы, построенные в § 2, для решения модельных задач, поставленных в § 1.

**1. Внутренняя краевая задача.** Мы рассмотрим задачу 1, § 1:

$$\Delta u = 0, \quad (x, y) \in D^+, \quad (1)$$

$$u|_{\Gamma} = \varphi(s), \quad \Gamma = \partial D^+. \quad (2)$$

Будем использовать вспомогательную разностную задачу вида (1), (2) из § 2:

$$\sum a_{mn} u_n = f_m, \quad m \in M^0, \quad (3)$$

$$u_{N^0}|_n = 0, \quad n \in \partial D^0. \quad (4)$$

Мы построим два алгоритма для численного решения задачи (1), (2). В обоих этих алгоритмах уравнение (1) заменяется разностным уравнением

$$\sum a_{mn} u_n = 0, \quad m \in M^+ \subset D^+. \quad (5)$$

Определяем, следуя § 2, множества  $N^+$ ,  $N^-$ ,  $\gamma = N^+ \cap N^-$ , пространство  $V_\gamma$  всех функций, определенных на  $\gamma$ , а также потенциал

$$u_N^+ = P^+ v_\gamma. \quad (6)$$

а) *Алгоритм с использованием разностной аппроксимации граничного условия.* В этом алгоритме условие Дирихле аппроксимируется тем или иным линейным разностным граничным условием, которое мы запишем символическим равенством

$$l u_\gamma = \varphi^{(h)}. \quad (7)$$

В силу теоремы 2 из § 2 условие

$$u_\gamma = P_\gamma^+ u_\gamma \quad (8)$$

необходимо и достаточно для того, чтобы  $u_\gamma \in V_\gamma$  было сужением решения уравнения (5) на  $\gamma \subset N^0$ . След  $u_\gamma$  решения задачи (5), (7), очевидно, совпадает с решением системы (7), (8):

$$u_\gamma - P_\gamma^+ u_\gamma = 0, \quad l u_\gamma = \varphi^{(h)}. \quad (9)$$

После того, как решение  $u_\gamma$  системы (9) вычислено, искомое решение  $u_{N^+}$  уравнения (5) восстанавливается по своему сужению  $u_\gamma$  на границу  $\gamma$  в силу теоремы 1 из § 2 по формуле

$$u_{N^+} = P^+ u_\gamma. \quad (10)$$

Относительно системы (9) заметим, что для отыскания ее решения можно воспользоваться возможностью вычислять  $v_\gamma - P_\gamma^+ v_\gamma$  для любого  $v_\gamma \in V_\gamma$  с помощью решения вспомогательной разностной задачи (3), (4).

Подчеркнем, что редукция задачи (5), (7) к задаче (9) резко уменьшает число неизвестных: вместо  $u_N$  надо искать  $u_\gamma$ .

б) *Спектральный подход к аппроксимации граничных условий.* Изложим подход, не требующий разностной аппроксимации граничного условия.

Выбираем систему базисных функций  $\psi_1(s), \psi_2(s), \dots$ , заданных на границе  $\Gamma = \partial D$  области  $D$ . Будем предполагать, что нормальная производная  $\frac{\partial u}{\partial n} \Big|_\Gamma$  искомого решения задачи (1), (2) может с любой точностью быть представлена в виде

$$\frac{\partial u}{\partial n} \Big|_\Gamma \approx \sum_{k=1}^N c_k \psi_k(s), \quad (11)$$

где  $c_k$  — некоторые числа, а  $N$  — достаточно большое число слагаемых.

Следует подчеркнуть, что количество членов в представлении  $\frac{\partial u}{\partial n}$  с какой-либо заранее требуемой точностью может зависеть от базиса  $\psi_1(s), \dots, \psi_n(s), \dots$  и от класса функций, которому принадлежит функция  $\frac{\partial u}{\partial n} \Big|_\Gamma$  аргумента  $s$ .

По заданной функции  $u|_\Gamma = \varphi(s)$  и по  $\frac{\partial u}{\partial n} \Big|_\Gamma$  вида (11) построим некоторую функцию  $v_\gamma$ , определив значение  $v_\nu$ ,  $\nu \in \gamma$ , с помощью формулы Тейлора

$$v_\nu = v_\nu(c_1, \dots, c_N) = \varphi(s_\nu) + \rho_\nu \sum_{k=1}^N c_k \psi_k(s_\nu). \quad (12)$$

В формуле (12) через  $s_\nu$  обозначена точка, которая является основанием перпендикуляра, опущенного из  $\nu$  на  $\Gamma$ . Число  $\rho_\nu$  есть расстояние от точки  $\nu$  до точки  $s_\nu$  на  $\Gamma$ , взятое со знаком «+», если  $\nu$  лежит вне  $D^+$ , и со знаком «-», если  $\nu$  лежит внутри  $D^+$ .

Очевидно, желательно, чтобы функция  $v_\gamma = v_\gamma(\mathbf{c})$  удовлетворяла условию (9):

$$v_\gamma(\mathbf{c}) - P_\gamma^+ v_\gamma(\mathbf{c}) = 0. \quad (13)$$

Линейная система (13) относительно неизвестных  $c_1, c_2, \dots, c_N$  содержит столько скалярных уравнений, каково число  $|\gamma|$  точек границы  $\gamma$ . При фиксированном  $N$  и достаточно мелкой сетке (т. е. при большом  $|\gamma|$ ) эта система переопределена. Решая эту систему в смысле метода наименьших квадратов, можно найти числа  $c_1, \dots, c_N$ . После этого находим функцию  $v_\gamma(\mathbf{c})$  по формуле (12), а затем функцию

$$u_{N+} \approx P^+ v_\gamma, \quad (14)$$

которую принимаем за приближенное решение.

Подходящий выбор скалярного умножения в методе наименьших квадратов, экономный алгоритм вычисления обобщенного решения  $c_1, c_2, \dots, c_N$  системы (13), а также сходимость приближенного решения (14) к точному обсуждаются в [17], ч. I.

**2. Внешняя краевая задача.** Дифференциальное уравнение (2) из § 1 заменим разностным уравнением

$$\sum a_{mn} u_n = 0, \quad m \in M^- = M^0 \setminus M^+. \quad (15)$$

Для разностной аппроксимации граничного условия  $u|_\Gamma = \varphi(s)$  воспользуемся той же формулой (11). Условие на  $u_\gamma$ , при котором эта функция является следом решения задачи (15) с условием (4) на  $\partial D^0$ , в силу результатов § 2 примет вид

$$u_\gamma - P_\gamma^- v_\gamma = 0. \quad (16)$$

Для отыскания  $v_\gamma$  надо решить систему (16), а затем восстановить решение  $u_{N-}$  внешней разностной задачи по его сужению  $u_\gamma$  на границе, воспользовавшись приближенным равенством  $u_\gamma \approx v_\gamma$  и формулой из § 2:

$$u_{N-} = P^- u_\gamma \approx P^- v_\gamma. \quad (17)$$

Мы видим, что алгоритм без разностной аппроксимации граничного условия в случае внешней задачи (3) из § 1 отличается от изложенного алгоритма для внутренней задачи тем, что вместо переопределенного уравнения (13) относительно коэффициентов  $c_1, c_2, \dots, c_N$ , входящих в (12), нужно воспользоваться переопределенным уравнением

$$v_\gamma(\mathbf{c}) - P^- v_\gamma(\mathbf{c}) = 0. \quad (18)$$

**3. Задача о построении ИГУ.** В разностной постановке эта задача приобретает следующий вид.

Рассматривается задача

$$\sum a_{mn} u_n = \begin{cases} f_m, & m \in M^+, \\ 0, & m \in M^-, \end{cases} \quad (19)$$

$$u_{N^0} \in U_{N^0}. \quad (20)$$

Используются те же  $M^0, M^+, M^-, N^+, N^-, \gamma, V_\gamma, U_{N^0}$ , что и в предыдущей задаче. Требуется построить граничное условие  $lu_\gamma = 0$  на границе  $\gamma = N^+ \cap N^-$  расчетной подобласти  $N^+$  таким образом, чтобы решение задачи

$$\sum a_{mn} u_n = f_m, \quad m \in M^+, \quad (21)$$

$$lu_\gamma = 0 \quad (22)$$

совпадало на  $N^+$  с решением  $u_{N^0}$  задачи (19), (20). Мы называем в этом случае множество  $\gamma$  искусственной границей, а условие (22) искусственным граничным условием, так как их не было в первоначальной постановке задачи (19), (20). Они возникли только из-за того, что мы решили ограничиться расчетом в области  $N^+$ , перенеся условие  $u_n = 0$ , если  $n \in \partial D^0$ , с границы  $\partial D^0$  на  $\gamma$ . В качестве (22) можно использовать условие (16), так как это условие необходимо и достаточно, чтобы решение  $u_{N^+}$  задачи (21), (22) допускало доопределение всюду на  $N^-$  с выполнением условий

$$\sum a_{mn} u_n = 0, \quad m \in M^-, \quad (23)$$

$$u_n = 0, \quad \text{если } n \in \partial D^0.$$

**4. Задача о вычислении вклада в потенциал.** В разностной постановке эта задача принимает следующий вид.

Пусть  $u_{N^0}$  удовлетворяет уравнению

$$\begin{aligned} \sum a_{mn} u_n &= f_m, \quad m \in M^0, \\ u_n &= 0, \quad \text{если } n \in \partial D^0. \end{aligned} \quad (24)$$

Решение этой задачи есть сумма решений  $u_{N^0}^+$  и  $u_{N^0}^-$  следующих двух задач (25) и (26):

$$\begin{aligned} \sum a_{mn} u_n^+ &= \theta_{M^0}(M^-, m) f_m, \quad m \in M^0, \\ u_{N^0}^+ &\in U_{N^0}; \end{aligned} \quad (25)$$

$$\begin{aligned} \sum a_{mn} u_n^- &= \theta_{M^0}(M^+, m) f_m, \quad m \in M^0, \\ u_{N^0}^- &\in U_{N^0}, \end{aligned} \quad (26)$$

где

$$\theta_{M^0}(X, m) = \begin{cases} 1, & \text{если } m \in X \subset M^0, \\ 0, & \text{если } m \in M^0 \setminus X. \end{cases}$$

Допустим, нам известны значения решения задачи (24) на  $\gamma$ , т. е. известна  $u_\gamma$ .

Требуется найти каждое из слагаемых в сумме

$$u_\gamma = u_\gamma^+ + u_\gamma^-.$$

Подчеркнем, что значения  $f_m$ ,  $m \in M^0$ , нам неизвестны. Решение поставленной задачи 4 из § 1 дают формулы

$$u_\gamma^+ = P_\gamma^+ u_\gamma, \quad (27)$$

$$u_\gamma^- = P_\gamma^- u_\gamma, \quad (28)$$

где

$$P_\gamma^+ u_\gamma = P^+ u_\gamma |_\gamma; \quad P_\gamma^- u_\gamma = P^- u_\gamma |_\gamma.$$

Мы обоснуем формулы (27), (28) немного ниже.

**5. Задача о влиянии подобластей.** Выпишем решение разностного аналога задачи 5 из § 1, т. е. две функции

$$u_n^+, \quad n \in N^+; \quad u_n^-, \quad n \in N^-, \quad (29)$$

являющиеся сужениями решений задач (25) и (26) на  $N^+$  и  $N^-$  соответственно. Именно, оказывается, что

$$u_n^+ = P^+ u_\gamma - n, \quad n \in N^+; \quad (30)$$

$$u_n^- = P^- u_\gamma - n, \quad n \in N^-. \quad (31)$$

Для обоснования формул (30), (31) введем функцию

$$w_n^\pm = \begin{cases} u_n^+, & n \in N^+, \\ -u_n^-, & n \in N^-. \end{cases} \quad (32)$$

Очевидно, что  $w_n^\pm$ ,  $n \in N^0$ , является кусочно регулярным решением задачи (13), (14) со скачком  $u_\gamma$  на  $\gamma$ :

$$[w^\pm]_n = w_n^+ - w_n^- = u^+ + u_n^- = u_n, \quad n \in \gamma. \quad (33)$$

Но в таком случае в силу определения разностного потенциала  $w^\pm = P^\pm u_\gamma$  со скачком  $u_\gamma$  функция (32) как раз и есть этот потенциал. Тогда в силу эквивалентности определений потенциала типа Коши с плотностью  $v_\gamma$  и потенциала со скачком  $v_\gamma$  можно написать

$$w_n^\pm = \begin{cases} w_n^+ = P^+ u_\gamma, & n \in N^+, \\ w_n^- = -P^- u_\gamma, & n \in N^-. \end{cases} \quad (34)$$

Сравнивая формулу (34) с формулой (32), получим формулы (31). Рассматривая равенства (31) только в точках  $n \in \gamma$ , получим формулы (27) и (28), т.е. решение задачи 4.

**6. Задача активного экранирования.** Задачу об активном экранировании рассмотрим в следующей разностной постановке.

Рассмотрим разностную краевую задачу

$$\sum a_{mn} u_n = f_m, \quad m \in M^0, \quad (35)$$

$$u_{N^0} \in U_{N^0}. \quad (36)$$

При этом ни  $f_m$ ,  $m \in M^0$ , ни  $u_n$ ,  $n \in N^0$ , не заданы; известен только след  $u_\gamma$  решения задачи (35) на сеточной границе  $\gamma$ .

Введем в правую часть уравнения (35) функцию-параметр — слагаемое  $\delta f_m$ , которое будем называть активным управлением. Задача (35), (36) перейдет в задачу

$$\sum a_{mn} w_n = f_m + \delta f_m, \quad m \in M^0, \quad (37)$$

$$w_{N^0} \in U_{N^0}. \quad (38)$$

Разностная задача активного экранирования решения в подобласти  $N^+ \subset N^0$  от влияния зарядов  $f_m$ , расположенных в подобласти  $M^- \subset M^0$ , состоит в следующем.

Надо найти все те активные управления  $\delta f_m$ , при включении которых решение задачи (37), (38) на подобласти  $N^+ \subset N^0$  совпадает с решением  $v_{N^0}$  следующей задачи:

$$\sum a_{mn} v_n = \begin{cases} f_m, & m \in M^+, \\ 0, & m \in M^-, \end{cases} \quad (39)$$

$$v_{N^0} \in u_{N^0}. \quad (40)$$

Общее решение  $\delta f_m$  задачи активного экранирования дает следующая

**Теорема.** Решение  $w_{N^0}$  задачи (37), (38) совпадает на  $N^+ \subset N^0$  с решением  $v_{N^0}$  задачи (39), (40) в том и только том случае, если  $\delta f_m$  имеет вид

$$\delta f_m = \begin{cases} 0, & m \in M^+, \\ -\sum a_{mn} z_n, & m \in M^-, \end{cases} \quad (41)$$

где  $z_n$ ,  $n \in N^0$ , суть значения произвольной функции  $z_{N^0} \in U_{N^0}$ , которая на границе  $\gamma$  совпадает с решением  $u_{N^0}$  задачи (35), (36), т.е. для которой справедливо  $z_\gamma = u_\gamma$ .

**Доказательство.** Сначала покажем, что решение  $u_{N^0} \in U_{N^0}$  задачи вида

$$\begin{aligned} \sum a_{mn} u_n &= \varphi_m, \quad m \in M^0, \\ u_{N^0} &\in U_{N^0} \end{aligned} \quad (42)$$

обращается в нуль на  $N^+$  в том и только том случае, если  $\varphi_m$  имеет вид

$$\varphi_m = \begin{cases} 0, & m \in M^+, \\ \sum a_{mn} \xi_n, & m \in M^-, \end{cases} \quad (43)$$

где  $\xi_{N^0} \in U_{N^0}$  — произвольная функция, обращающаяся в нуль на  $\gamma$ .

В самом деле, произвольная функция  $u_{N^0} \in U_{N^0}$  является решением задачи (42) в том и только том случае, если функция  $\varphi_{M_0}$  определяется тождеством

$$\varphi_m \equiv \sum_{n \in N_m} a_{mn} u_n, \quad m \in M^0.$$

В предположении, что  $u_n = 0$  всюду на  $N^+$ , получаем, что  $\varphi_m = 0$  при  $m \in M^+$ , поскольку в этом случае имеют место включения  $n \in N_m \subset N^+$ . После этого замечания видно, что тождество, определяющее  $\varphi_m$ , совпадает с (43).

Теперь покажем, что формула (41) дает общее решение задачи экранирования, т. е. что для функций  $\delta f_m$  вида (41) и только для них решение задачи (37), (38) совпадает на  $N^+$  с решением  $v_{N^0} \in U_{N^0}$  задачи (39), (40). Вычтем из уравнения (37) уравнение (39) почленно. Обозначим  $z_{N^0} = w_{N^0} - v_{N^0}$ . Получим

$$\sum a_{mn} z_n = \begin{cases} \delta f_m, & m \in M^+, \\ \delta f_m + f_m, & m \in M^-, \end{cases} \quad (44)$$

причем в силу (35)

$$f_m = \sum a_{mn} u_n, \quad m \in M^-.$$

В силу доказанного выше  $z_n \equiv 0$ ,  $n \in N^+$ , в том и только том случае, если правая часть (44), т.е. функция

$$\psi_m = \begin{cases} \delta f_m, & m \in M^+, \\ \delta f_m + f_m = \delta f_m + \sum a_{mn} u_n, & m \in M^-. \end{cases}$$

имеет вид (43). Как легко видеть, это требование равносильно тому, что  $\delta f_m$  имеет вид (41).  $\square$

\* \* \*

Метод разностных потенциалов был предложен В.С. Рябеньким в докторской диссертации (1969) и интенсивно развивается многими авторами. В [17] отражен уровень развития метода, достигнутый к 2001 г.

## СПИСОК ЛИТЕРАТУРЫ

1. *Бabenko K.I.* Основы численного анализа. — М.: Наука, 1986.
2. *Бахвалов Н.С., Жидков Н.П., Кобельков Г.М.* Численные методы. — М.: Наука, 1987.
3. *Годунов С.К., Антонов А.Г. и др.* Гарантированная точность решения систем линейных уравнений в евклидовых пространствах. — Новосибирск: Наука, 1988.
4. *Годунов С.К., Забродин А.В. и др.* Численное решение многомерных задач газовой динамики. — М.: Наука, 1976.
5. *Годунов С.К., Рябенький В.С.* Введение в теорию разностных схем. — М.: Физматгиз, 1962.
6. *Годунов С.К., Рябенький В.С.* Разностные схемы. — М.: Наука, 1977.
7. *Дьяконов Е.Г.* Минимизация вычислительной работы. — М.: Наука, 1989.
8. *Жуков А.И.* Метод Фурье в вычислительной математике. — М.: Наука, 1992.
9. *Калиткин Н.Н.* Численные методы. — М.: Наука, 1978.
10. *Коробов Н.М.* Теоретико-числовые методы в прикладном анализе. — М.: Физматгиз, 1963.
11. *Локуциевский О.В., Гавриков М.Б.* Начала численного анализа. — М.: ТОО Янус, 1995.
12. *Магомедов К.М., Холодов А.С.* Сеточно-характеристический метод. — М.: Наука, 1988.
13. *Марчук Г.И.* Методы вычислительной математики. — М.: Наука, 1980.
14. *Михлин С.Г., Морозов Н.Ф., Паукштот М.В.* Граничные интегральные уравнения в теории упругости. — М.: Наука, 1993.
15. *Партон В.З., Перлин П.И.* Интегральные уравнения теории упругости. — М.: Наука, 1978.
16. *Петровский И.Г.* Лекции об уравнениях с частными производными. — М.: Гостехиздат, 1950.
17. *Рябенький В.С.* Метод разностных потенциалов и его приложения. — М.: ФИЗМАТЛИТ, 2002.
18. *Рябенький В.С., Филиппов А.Ф.* Об устойчивости разностных уравнений. — М.: Гостехиздат, 1956.
19. *Самарский А.А.* Теория разностных схем. — Изд. 2-е. — М.: Наука, 1983.
20. *Самарский А.А., Гулин А.В.* Численные методы. — М.: Наука, 1989.
21. *Самарский А.А., Николаев Е.С.* Методы решения сеточных уравнений. — М.: Наука, 1978.

22. Соболев С.Л. Уравнения математической физики. — Изд. 4-е. — М.: Наука, 1994.
23. Соболев С.Л. Введение в теорию кубатурных формул. — М.: Наука, 1973.
24. Соболь И.М. Численные методы Монте-Карло. — М.: Наука, 1973.
25. Стренг Г., Фикс Дж. Теория метода конечных элементов. — М.: Мир, 2007.
26. Темам Ф. Уравнения Навье–Стокса. — М.: Мир, 1981.
27. Шайдуров В.В. Многосеточные методы конечных элементов. — М.: Наука, 1989.
28. Федоренко Р.П. Введение в вычислительную физику. — М.: Изд-во МФТИ, 1994.
29. Иванов В.Д., Косарев В.И. и др. Лабораторный практикум «Основы вычислительной математики». — М.: МЗ Пресс, 2003.

# ОГЛАВЛЕНИЕ

Предисловие к третьему изданию . . . . .	3
Предисловие к первому изданию . . . . .	3
<b>Введение . . . . .</b>	<b>5</b>
§ 1. Дискретизация . . . . .	7
§ 2. Обусловленность . . . . .	9
§ 3. Погрешность . . . . .	10
§ 4. О методах вычисления . . . . .	15
 <b>ЧАСТЬ I. Табличное задание и интерполяция функций.</b>	
<b>Квадратуры</b>	
<b>ГЛАВА 1. Алгебраическая интерполяция . . . . .</b>	<b>20</b>
§ 1. Существование и единственность интерполяционного многочлена . . . . .	20
§ 2. Классическая кусочно многочленная интерполяция . . . . .	27
§ 3. Кусочно многочленная гладкая интерполяция (сплайны) . . . . .	34
§ 4. Интерполяция функций двух переменных . . . . .	42
<b>ГЛАВА 2. Тригонометрическая интерполяция . . . . .</b>	<b>44</b>
§ 1. Интерполяция периодических функций . . . . .	45
§ 2. Интерполяция функций на отрезке. Связь между алгебраической и тригонометрической интерполяциями . . . . .	53
<b>ГЛАВА 3. Вычисление определенных интегралов. Квадратуры . . . . .</b>	<b>59</b>
§ 1. Квадратурные формулы трапеций и Симпсона . . . . .	59
§ 2. Сочетание численных и аналитических методов при вычислении интегралов с особенностями . . . . .	67
§ 3. Кратные интегралы . . . . .	68
 <b>ЧАСТЬ II. Системы скалярных уравнений</b>	
<b>ГЛАВА 4. Системы линейных алгебраических уравнений. Методы отыскания точного решения . . . . .</b>	<b>72</b>
§ 1. Формы записи совместных СЛАУ . . . . .	73
§ 2. Нормы . . . . .	77

§ 3. Обусловленность СЛАУ . . . . .	82
§ 4. Методы исключения Гаусса . . . . .	88
§ 5. Связь между задачей на минимум квадратичной функции и СЛАУ . . . . .	96
§ 6. Метод сопряженных градиентов как метод точного решения СЛАУ . . . . .	97
§ 7. Конечные ряды Фурье и запись точного решения разностного аналога задачи Дирихле для уравнения Пуассона . . . . .	99
<b>ГЛАВА 5. Методы последовательных приближений (итерационные методы) решения систем линейных алгебраических уравнений . . . . .</b>	104
§ 1. Методы простых итераций . . . . .	105
§ 2. Метод Чебышёва и метод сопряженных градиентов . . . . .	118
<b>ГЛАВА 6. Переопределенные СЛАУ. Метод наименьших квадратов . . . . .</b>	120
§ 1. Примеры задач, приводящих к переопределенным СЛАУ . . . . .	120
§ 2. Переопределенные СЛАУ и обобщенные решения в общем случае . . . . .	122
<b>ГЛАВА 7. Численное решение нелинейных скалярных уравнений и систем уравнений . . . . .</b>	129
§ 1. Метод простых итераций . . . . .	130
§ 2. Метод линеаризации Ньютона . . . . .	134
 ЧАСТЬ III. Метод конечных разностей для численного решения дифференциальных уравнений	
<b>ГЛАВА 8. Численное решение задач для обыкновенных дифференциальных уравнений . . . . .</b>	138
§ 1. Примеры разностных схем. Сходимость . . . . .	138
§ 2. Аппроксимация дифференциальной краевой задачи разностной схемой . . . . .	145
§ 3. Определение устойчивости разностной схемы. Сходимость как следствие аппроксимации и устойчивости . . . . .	152
§ 4. Схемы Рунге–Кутты . . . . .	160
§ 5. Методы решения краевых задач . . . . .	164
<b>ГЛАВА 9. Разностные схемы для уравнений с частными производными . . . . .</b>	167
§ 1. Основные определения и их иллюстрация . . . . .	168
§ 2. Некоторые приемы построения аппроксимирующих разностных схем . . . . .	181
§ 3. Спектральный признак устойчивости разностной задачи Коши . . . . .	196

§ 4. Принцип замороженных коэффициентов . . . . .	206
§ 5. Явные и неявные разностные схемы для уравнения теплопроводности . . . . .	216
<b>ГЛАВА 10. Понятие о разрывных решениях и способах их вычисления . . . . .</b>	<b>218</b>
§ 1. Дифференциальная формулировка интегрального закона сохранения . . . . .	219
§ 2. Построение разностных схем . . . . .	226
<b>ГЛАВА 11. Разностные методы для эллиптических задач . . . . .</b>	<b>232</b>
§ 1. Аппроксимация и устойчивость простейшей разностной схемы . . . . .	233
§ 2. Понятие о методе конечных элементов . . . . .	239
§ 3. Вычисление решений сеточных аналогов краевых задач . . . . .	247
§ 4. Многосеточный метод Федоренко . . . . .	249
 <b>ЧАСТЬ IV. Методы граничных уравнений для численного решения краевых задач</b>	
<b>ГЛАВА 12. Граничные интегральные уравнения и метод граничных элементов для их численного решения . . . . .</b>	<b>254</b>
§ 1. Способы редукции краевых задач к ГИУ . . . . .	254
§ 2. Граничные элементы и дискретизация ГИУ . . . . .	257
§ 3. Область применимости ГИУ для численного решения краевых задач . . . . .	258
<b>ГЛАВА 13. Метод разностных потенциалов . . . . .</b>	<b>259</b>
§ 1. Постановка модельных задач . . . . .	260
§ 2. Разностные потенциалы . . . . .	264
§ 3. Решение модельных задач . . . . .	273
<b>Список литературы . . . . .</b>	<b>280</b>

Учебное издание

*РЯБЕНЬКИЙ Виктор Соломонович*

**ВВЕДЕНИЕ В ВЫЧИСЛИТЕЛЬНУЮ МАТЕМАТИКУ**

Редактор *В.С. Аролович*  
Оригинал-макет: *В.В. Затекин*  
Оформление переплета: *Н.В. Гришина*

Подписано в печать 24.12.07. Формат 60×90/16. Бумага офсетная.  
Печать офсетная. Усл. печ. л. 18. Уч.-изд. л. 19.8. Тираж 1500 экз.  
Заказ № 2937.

Издательская фирма «Физико-математическая литература»  
МАИК «Наука/Интерпериодика»  
117997, Москва, ул. Профсоюзная, 90  
E-mail: fizmat@maik.ru, fmlsale@maik.ru;  
<http://www.fml.ru>

Отпечатано с готовых диапозитивов  
в ОАО «Ивановская областная типография»  
153008, г. Иваново, ул. Типографская, 6  
E-mail: 091-018@adminet.ivanovo.ru

ISBN 978-5-9221-0926-0



9 785922 109260

**Издательская фирма  
«Физико-математическая литература»  
МАИК «Наука/Интерпериодика»  
117997 Москва, Профсоюзная ул., 90**

---

**В издательстве «Физматлит» вышли из печати  
следующие книги:**

Кондратьев А.С., Райгородский П.А.  
Задачи по термодинамике, статистической физике  
и кинетической теории

Агошков В.П., Дубовский П.Б., Шутляев В.П.  
Методы решения задач математической физики

Будак Б.М., Самарский А.А., Тихонов А.Н.  
Сборник задач по математической физике

Владимиров В.С., Жаринов В.В.  
Уравнения математической физики. Учебник для вузов

Кулиев В.Д.  
Сингулярные краевые задачи  
  
и другие книги

Наиболее полную информацию о книгах Вы можете найти  
в Интернете по адресу <http://www.fml.ru>

**По вопросам приобретения книг обращаться:  
Издательская фирма  
«Физико-математическая литература»  
117997 Москва, Профсоюзная ул., 90  
тел./факс (495) 334-7421, e-mail: [fizmat@maik.ru](mailto:fizmat@maik.ru)**

В.С. РЯБЕНЬКИЙ ВВЕДЕНИЕ В ВЫЧИСЛИТЕЛЬНУЮ МАТЕМАТИКУ

