

УЧРЕЖДЕНИЕ  
РОССИЙСКОЙ АКАДЕМИИ НАУК  
ИНСТИТУТ ЯДЕРНОЙ ФИЗИКИ  
им. Г.И. Будкера СО РАН  
СИБИРСКОГО ОТДЕЛЕНИЯ РАН  
(ИЯФ СО РАН)

Э.А. Бибердорф, Н.И. Попова

КОНТРОЛЬ ТОЧНОСТИ  
РЕШЕНИЯ КРАЕВОЙ ЗАДАЧИ  
МЕТОДОМ ОРТОГОНАЛЬНОЙ  
ПРОГОНКИ

ИЯФ 2009-1

НОВОСИБИРСК  
2009



# Контроль точности решения краевой задачи методом ортогональной прогонки

*Э.А. Бибердорф*

Институт математики им. С.Л. Соболева  
630090, Новосибирск, Россия

*Н.И. Попова*

Институт ядерной физики им. Г.И.Будкера  
630090, Новосибирск, Россия

## Аннотация

Цель нашей работы заключается в развитии и популяризации современных алгоритмов нового поколения, позволяющих проводить вычисления с гарантированной оценкой точности. В данной работе мы впервые переходим от задач линейной алгебры к гораздо более сложной проблеме – оценке точности выполнения алгоритма решения задачи математической физики. В качестве объекта исследования рассмотрен алгоритм ортогональной прогонки решения краевой задачи для системы обыкновенных дифференциальных уравнений, который включает в себя такие элементы, как дискретизация, интерполяция, численное интегрирование и решение серий задач линейной алгебры. Здесь мы используем полученные нами ранее гарантированные оценки точности ортогонализации и решения линейных алгебраических систем [3]. В результате впервые практически получена гарантированная оценка погрешности (как алгоритмической, так и погрешности машинных округлений) решения краевой задачи. Создана новая версия **GALA-2.1** пакета программ **GALA-2.0** (Guaranteed Accuracy in Linear Algebra), описанного в монографии [3]. В новую версию входят программы реализующие такие вычислительные методы, как ”ортогональная прогонка”, ”регуляризация”, QR-разложение матриц с контролем оценки точности результата. Приведены примеры.

# Accuracy control of solving the boundary-value problem by the orthogonal sweep method

*E.A. Biberdorf*

Sobolev Institute of Mathematics  
630090, Novosibirsk, Russia

*N.I. Popova*

Budker Institute of Nuclear Physics  
630090, Novosibirsk, Russia

## Abstract

The aim of our paper is to develop and popularize the modern algorithms permitting one to carry out computations with the guaranteed accuracy estimation. In this paper we first pass from the linear algebra problem to the more complicated problem – the accuracy estimation of fulfilling the algorithm of solving the mathematical physics problem. The subject of our study is the orthogonal sweep algorithm for solving the boundary-value problem of the system of ordinary differential equations. The algorithm contains such elements as discretization, interpolation, numerical integration and solution of series of linear algebra problems. Here we used the guaranteed accuracy estimations of the orthogonalization and the solution of linear algebra systems obtained in [3]. As a result we first practically obtained the guaranteed estimation of error (both the algorithmic and round-off errors) for solving the boundary-value problem. The new version **GALA-2.1** of the program package **GALA-2.0** (Guaranteed Accuracy in Linear Algebra) described in the monograph [3] is created. The programs realizing such numerical methods as the orthogonal sweep method, regularization, QR-decomposition of matrixes with the guaranteed estimation of the result accuracy enter into new version. The numerical examples are given.

---

# 1 Введение

Данная работа продолжает начатую в 1999 году (см. [1, 2, 3]) деятельность по программной реализации алгоритмов, позволяющих проводить вычисления с гарантированной оценкой точности. Цель нашей работы заключается в развитии и популяризации этих современных алгоритмов.

Все вычисления на компьютере проводятся приближенно из-за округлений при выполнении элементарных арифметических операций (погрешности появляются уже при вводе данных в память компьютера). Поэтому при реализации на компьютере какого-либо численного алгоритма возникают вычислительные погрешности. В результате накопления машинных погрешностей вычисленное решение может сильно отличаться от истинного (см. вычислительные парадоксы в [8]) и важно уметь оценивать границы интервала, в котором находится точное решение.

В начале 50-х годов XX века Гивенс ввел метод, называемый методом обратного анализа погрешностей [14], который позволяет контролировать накопление ошибок округлений в процессе выполнения вычислительного алгоритма на компьютере. Суть метода заключается в том, что арифметические погрешности, возникающие в процессе вычислений трактуются как возмущения исходных данных задачи, что позволяет оценить погрешность решения, используя соответствующую теорию возмущений. Оказалось, однако, что для подавляющего большинства алгоритмов метод обратного анализа практически не применим из-за его громоздкости.

В работах [8, 11, 9, 13] опубликованы результаты исследования ряда алгоритмов линейной алгебры с точки зрения возможности применения метода обратного анализа для вычисления погрешности. Продолжая эту линию, мы получили формулы гарантированных оценок погрешностей выполнения алгоритмов решения основных задач линейной алгебры [1, 2, 3]. Эти алгоритмы были реализованы в виде программного пакета **GALA-2.0** (Guaranteed Accuracy in Linear Algebra), (Свидетельство об официальной регистрации программы для ЭВМ № 2007614804 Федеральной службы по интеллектуальной собственности, патентам и товар-

ным знакам). Основная особенность процедур этого пакета в том, что вычисленный результат сопровождается гарантированной оценкой его точности. Таким образом, результат приобретает характер не "гипотезы" (предположительного значения), а "теоремы" (приводятся границы интервала, в котором лежит точное значение решения ).

В данной работе мы впервые переходим от задач линейной алгебры к гораздо более сложной проблеме – оценке точности выполнения алгоритма решения задачи математической физики. Эти алгоритмы имеют комплексный характер и включают в себя такие элементы, как дискретизация, интерполяция, численное интегрирование и решение серий задач линейной алгебры. Следовательно, при получении гарантированной оценки решения необходимо учитывать погрешности всех этих операций.

Так как это **первая** работа данного направления, то мы в качестве объекта исследования выбрали алгоритм ортогональной прогонки решения краевой задачи для системы обыкновенных дифференциальных уравнений. Такой выбор обусловлен следующими причинами. Первая, это – простота алгоритма. Весь алгоритм состоит из цепочки однотипных шагов (ортогонализация + интегрирование задачи Коши) и решения цепочки линейных алгебраических систем уравнений. Такая структура облегчает нашу задачу, тем более, что мы можем использовать полученные нами ранее гарантированные оценки точности ортогонализации и решения линейных алгебраических систем [3]. Вторая, это – свобода в выборе алгоритмов для промежуточных шагов. Как уже говорилось, для решения алгебраических задач мы выбрали алгоритмы, с которыми мы работали ранее, а для решения задач Коши на данном начальном этапе был выбран простейший алгоритм Эйлера. Интерполяцию начальных данных мы тоже выбираем самую простую – линейную. Третья причина, это – надежность алгоритма. За основу нами была взята работа [9]. В ней доказано, что результат ортогональной прогонки устойчив к ошибкам округления. Фактически, в этой работе показано, что к ортогональной прогонке применим метод обратного анализа погрешностей и итоговая погрешность решения оценивается пропорционально вносимой погрешности. Порядок же коэффициента пропорциональности не оценивается. Таким образом, наша задача свелась к тому, чтобы для конкретных алгоритмов внутренних шагов оценить вносимую погрешность и вычислить коэффициент пропорциональности. Результат этой деятельности мы представляем ниже.

Во втором разделе дана постановка краевой задачи и краткое описание метода ортогональной прогонки. В третьем разделе приведены факты из теории возмущений, на которых базируются итоговые оценки. В

четвертом разделе мы последовательно разбираем шаги ортогональной прогонки на предмет накопления арифметической погрешности. И наконец, пятый раздел посвящен численным примерам.

Очевидно, что поскольку мы вычисляем **гарантированные оценки** погрешности, то они должны были получиться завышенными. Забегая вперед, отметим, что действительность превзошла наши ожидания. Полученные оценки оказались достаточно грубыми. Основную причину такого положения мы видим в несовершенстве подхода, предложенного в [9], который мы взяли за основу. Дальнейшими шагами по преодолению сложившейся ситуации должны стать как модификация способа применения метода обратного анализа к ортогональной прогонке, так и совершенствование оценок погрешностей промежуточных шагов.

Тем не менее еще раз подчеркнем, что в данной работе впервые получены гарантированные оценки погрешности решения задачи математической физики. Это говорит о том, что несмотря на сложность соответствующих алгоритмов, решение задач математической физики **может** сопровождаться гарантированной оценкой.

Данная работа связана с участием авторов в Проекте № 46 (Конкурса междисциплинарных, интеграционных проектов фундаментальных исследований СО РАН) "Исследование и моделирование физиологических, молекулярногенетических и биофизических механизмов формирования артериальной гипертензии с целью создания оптимальных программ ранней диагностики, прогнозирования осложнений и их профилактики". Впервые предпринята попытка компьютерного моделирования сердечно-сосудистой системы человека как численного решения системы гемодинамики методом прогонки. Работа завершена изданием монографии [4] совместно с биологами, медиками, физиками и математиками.

## 2 Метод ортогональной прогонки

### 2.1 Предварительные замечания

Метод ортогональной прогонки предназначен для численного решения краевой задачи (1). Он зарекомендовал себя на практике как эффективное и надежное средство численного решения краевой задачи. Метод обладает повышенной устойчивостью к погрешностям округлений при реализации метода на компьютере в предположении, что решаемая краевая задача имеет матрицы Грина  $G(x, s)$ ,  $G_L(x)$ ,  $G_R(x)$  ограниченные не слишком большой постоянной  $K$  (11).

Опишем коротко метод численного решения краевой задачи

$$\begin{aligned} \frac{du}{dx} &= A(x)u(x) + f(x), \\ Lu(x_0) &= \varphi, \quad Ru(x_m) = \psi, \\ x_0 &\leq x \leq x_m, \quad d = x_m - x_0, \end{aligned} \tag{1}$$

где  $A$  – матрица размера  $n \times n$ ,  $f$  – вектор-функция размера  $n$ ,  $L, R$  – прямоугольные матрицы размера  $k \times n$  и  $p \times n$  соответственно:

$$L = \begin{pmatrix} l_{11} & l_{12} & \dots & l_{1n} \\ l_{21} & & & \\ \vdots & & & \\ l_{k1} & l_{k2} & \dots & l_{kn} \end{pmatrix}, \quad R = \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ r_{21} & & & \\ \vdots & & & \\ r_{p1} & r_{p2} & \dots & r_{pn} \end{pmatrix},$$

причем  $k + p = n$  и ранги матриц  $L$  и  $R$  соответственно равны  $k$  и  $p$ . Элементы матрицы  $A(x)$  и вектор-функции  $f(x) = (f_1(x), \dots, f_n(x))^T$  предполагаются непрерывными на отрезке  $[x_0, x_m]$ . Метод основан на представлении решения задачи (1) через решения серий задач Коши.

Пусть  $z_1, z_2, \dots, z_k$  – полная ортонормированная система линейно независимых векторов удовлетворяющих краевому условию  $Lz(x_0) = 0$ . Пусть вектор  $z_f$  удовлетворяет левому условию  $Lz_f(x_0) = \varphi$  и ортогонален всем  $z_j$  ( $j = 1, 2, \dots, k$ ):  $(z_f, z_j) = 0$ . Очевидно, что если решение задачи (1) существует, то оно представимо в виде:

$$u(x) = y_f(x) + \sum_{j=1}^p \alpha_j y_j(x),$$

где вектор-функции  $y_f(x), y_1(x), \dots, y_p(x)$  – решения следующих задач Коши:

$$\begin{aligned} \frac{dy_f(x)}{dx} &= A(x)y_f(x) + f(x), \quad y_f(x_0) = z_f(x_0), \\ \frac{dy_j(x)}{dx} &= A(x)y_j(x), \quad y_j(x_0) = z_j(x_0), \quad j = 1, \dots, p, \end{aligned} \tag{2}$$

а коэффициенты  $\alpha_j$   $j = 1, \dots, p$  определяются из правого граничного условия  $Ru(x_m) = \psi$  как решения систем линейных уравнений

$$\sum_{j=1}^p \alpha_j R y_j(x_m) = \psi - R y_f(x_m).$$



Но во многих случаях описанная процедура практически невыполнима из-за ряда причин (например, "сплющивания" системы векторов  $y_f(x), y_1(x), \dots, y_p(x)$ ). Для преодоления явления "сплющивания" базиса, можно применить несколько раз процедуру его ортогонализации. Что и было сделано С. К. Годуновым в [5] (ортогонализация Грама-Шмита), затем С. В. Кузнецовым (ортогональные отражения Хаусхолдера) [9]. При этом весь отрезок  $[x_0, x_m]$  разбивается на  $m$  интервалов. На каждом интервале решается серия задач Коши (4) и проводится ортогонализация базисов векторов. Такой метод решения краевой задачи называется методом ортогональной прогонки (левое граничное условие как бы "перегоняется" с одного интервала на следующий за ним справа).

Здесь в методе ортогональной прогонки мы используем преобразования отражения Хаусхолдера, что дает следующие преимущества:

1) такой способ ортогонализации характеризуется повышенной устойчивостью к ошибкам округления по сравнению с методом Грама-Шмита;

2) точность получения ортогональной системы векторов при использовании преобразований отражения не зависит от числа обусловленности прямоугольной матрицы, составленной из векторов, подлежащих ортогонализации;

3) позволяет осуществлять более крупные шаги ортогонализации;

4) использование метода отражений для ортогонализации систем векторов позволяет реально провести оценку точности полученного численного решения, то есть оценить близость численного решения к истинному.

Ниже дан алгоритм численного решения краевой задачи (1) методом ортогональной прогонки.

## 2.2 Алгоритм

**Дано:** краевая задача (1).

**Разбивка отрезка**  $[x_0, x_m]$ .

Разобьем отрезок  $[x_0, x_m]$  на  $m$  участков точками

$$x_0 < x_1 < x_2 < \dots < x_{s-1} < x_s < \dots < x_m$$

так, что

$$|x_s - x_{s-1}| \leq \frac{C}{\max_{x \in [x_0, x_m]} \|A(x)\|}, \quad C \approx 1 \div 3. \quad (3)$$

**Определение базиса решений однородного левого граничного условия.**

Определяем  $z_j(x_0)$ ,  $j = 1, \dots, p$  – полную ортонормированную систему векторов

$$(z_j, z_i) = \begin{cases} 1, & \text{если } j = i, \\ 0, & \text{если } j \neq i, \end{cases}$$

удовлетворяющих условию

$$Lz_j(x_0) = 0.$$

**Вычисление решения неоднородного левого граничного условия.**

Находим вектор  $z_f(x_0)$ , ортогональный ко всем  $z_j(x_0)$  (т.е.  $(z_f, z_j) = 0$ ), удовлетворяющий уравнению

$$Lz_f(x_0) = \varphi.$$

**Прямая прогонка.**

**Цикл по  $s = 1, \dots, m$ , отрезок  $[x_{s-1}, x_s]$ .**

**Решение задач Коши.**

Путем численного интегрирования следующей серии задач Коши находим векторы  $y_f(x_s)$ ,  $y_j(x_s)$ :

$$\begin{aligned} \frac{dy_f(x)}{dx} &= A(x)y_f(x) + f(x), & y_f(x_{s-1}) &= z_f(x_{s-1}), \\ \frac{dy_j(x)}{dx} &= A(x)y_j(x), & y_j(x_{s-1}) &= z_j(x_{s-1}), \quad j = 1, \dots, p. \end{aligned} \tag{4}$$

**Ортогонализация.**

В результате ортогонализации набора векторов  $y_1(x_s), \dots, y_p(x_s), y_f(x_s)$ , используя, например, QR-разложение, получаем векторы  $z_1(x_s), \dots, z_p(x_s), z_f(x_s)$  и матрицу  $\Omega_s$  размера  $(p+1) \times (p+1)$  такую, что

$$Y_s = Z_s \Omega_s, \tag{5}$$

где

$$Y_s = [y_1(x_s), \dots, y_p(x_s), y_f(x_s)], \quad Z_s = [z_1(x_s), \dots, z_p(x_s), z_f(x_s)],$$

причем  $Z_s^* Z_s = I_n$ , а  $\Omega_s$  – верхнетреугольная матрица вида

$$\Omega_s = \left( \begin{array}{cccc|c} & & & & \theta_1^{(s)} \\ & & & & \vdots \\ & \Theta_s & & & \theta_p^{(s)} \\ 0 & 0 & \dots & 0 & 1 \end{array} \right).$$

**Конец цикла.**

**Обратная прогонка.**

**Правое граничное условие.** Вектор-коэффициент  $\alpha$  однозначно определяется из граничного условия

$$Ru(x_m) = RZ_m \alpha = \psi,$$

т. е. из решения системы линейных уравнений

$$\sum_{j=1}^p \alpha_j Rz_j(x_m) = \psi - Rz_f(x_m). \quad (6)$$

*Замечание.* Если решение системы единственно, то можно утверждать, что решение задачи (1) единственно и наоборот.

**Вычисление решения  $u(x)$ .**

Находим значения искомой вектор-функции  $u(x)$  в точках

$$x_m, x_{m-1}, \dots, x_s, \dots, x_1, x_0$$

следующим образом. Пусть

$$\beta_s = \begin{pmatrix} \beta_s^{(1)} \\ \beta_s^{(2)} \\ \vdots \\ \beta_s^{(p)} \\ 1 \end{pmatrix}, \quad s = m, m-1, \dots, 2, 1.$$

Положим  $\beta_m^{(j)} = \alpha_j$ ,  $j = 1, \dots, p$ . Тогда, учитывая, что на каждом шаге  $s$  при прямой прогонке уже вычислены  $Z_s$  и  $\Omega_s$  (5), находим

$$u(x_s) = Z_s \beta_s, \quad (7)$$

где  $\beta_s$  определяется через  $\beta_{s+1}$  путем решения системы уравнений с хорошо обусловленной верхнетреугольной матрицей  $\Omega_{s+1}$ :

$$\Omega_{s+1}\beta_s = \beta_{s+1}. \quad (8)$$

### Оценка вычислительных погрешностей.

Определить погрешность результата по формуле (75).

**Результат** выполнения алгоритма – вычисленное решение краевой задачи (1) с оценкой точности результата.

Схематично этот алгоритм представлен на рис. 1.

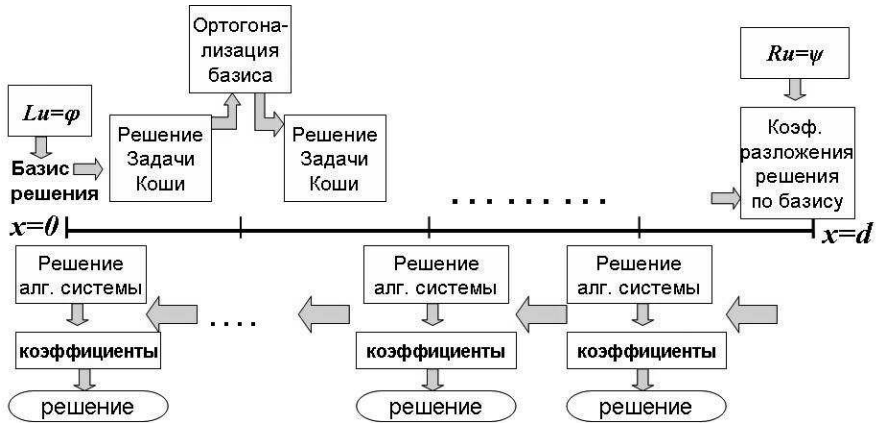


Рис. 1. Схема метода ортогональной прогонки.

## 3 Теория возмущений

Пусть краевая задача

$$\begin{aligned} \frac{du(x)}{dx} &= A(x)u(x) + f(x), \quad 0 \leq x \leq d \\ Lu(0) &= \varphi, \quad Ru(d) = \psi, \end{aligned} \quad (9)$$

правильно поставлена. (Здесь  $A(x)$  – матрица размера  $n \times n$ ,  $L, R$  – прямоугольные матрицы размера  $k \times n$  и  $p \times n$  соответственно, причем  $k+p = n$  (ранги матриц  $L, R$  равны  $k$  и  $p$ ), вектор-функция  $f(x)$  – размера  $n$ ).

Это означает, что краевая задача (9) имеет единственное решение для любых  $\varphi, \psi$  и любой непрерывной на отрезке  $[0, d]$  вектор-функции  $f(x)$ . Тогда решение задачи представляется с помощью матриц Грина  $G(x, s), G_L(x), G_R(x)$  (см. [6])

$$u(x) = G_L(x)\varphi + \int_0^d G(x, s)f(s)ds + G_R(x)\psi, \quad (10)$$

где  $G_L(x)$  – матрица-функция размера  $n \times k$ , решение задачи

$$\begin{cases} \frac{dG_L(x)}{dx} = A(x)G_L(x), \\ LG_L(0) = I_k, \quad RG_L(d) = O_{p \times k}, \end{cases}$$

матрица-функция  $G_R(x)$  размера  $n \times p$ , решение задачи

$$\begin{cases} \frac{dG_R(x)}{dx} = A(x)G_R(x), \\ LG_R(0) = O_{k \times p}, \quad RG_R(d) = I_p, \end{cases}$$

и матрица-функция  $G(x, s)$  размера  $n \times n$ , решение задачи

$$\begin{cases} \frac{dG(x, s)}{dx} = A(x)G(x, s), \\ LG(0, s) = O_{k \times n}, \quad RG(d, s) = O_{p \times n}, \\ G(s + 0, s) - G(s - 0, s) = I_n, \end{cases}$$

где  $O_{l \times m}$  – нулевая матрица размера  $l \times m$ .

Будем говорить, что задача (9) хорошо обусловлена, если для любых  $x$  и  $s$  из интервала  $[0, d]$  верны оценки

$$\|G_L(x)\| \leq K, \quad \|G(x, s)\| \leq K, \quad \|G_R(x)\| \leq K. \quad (11)$$

Из (10) следует, что для решения хорошо обусловленной задачи (9) справедлива оценка

$$\max_{x \in [0, d]} \|u(x)\| \leq K(\|\varphi\| + \|\psi\| + d \max_{x \in [0, d]} \|f(x)\|). \quad (12)$$

Это позволяет утверждать, что решение задачи (9) устойчиво по отношению к возмущениям.

Предположим, что мы одновременно рассматриваем две "близкие" краевые задачи на отрезке  $0 \leq x \leq d$ :

$$\left\{ \begin{array}{l} \frac{du(x)}{dx} = A(x)u(x) + f(x), \\ Lu(0) = \varphi, \quad Ru(d) = \psi; \end{array} \right. \quad \left\{ \begin{array}{l} \frac{d\tilde{u}(x)}{dx} = \tilde{A}(x)\tilde{u}(x) + \tilde{f}(x), \\ \tilde{L}\tilde{u}(0) = \tilde{\varphi}, \quad \tilde{R}\tilde{u}(d) = \tilde{\psi}. \end{array} \right. \quad (13)$$

Матрицы  $\tilde{L}, \tilde{R}$ , входящие в граничные условия, предполагаются имеющими то же число строк и столбцов, что и  $L, R$  соответственно. Утверждение, что первая и вторая задачи "близкие", можно понимать как утверждение, что имеют место неравенства

$$\begin{aligned} \|\tilde{A} - A\| &\leq \epsilon_A, & \|\tilde{f} - f\| &\leq \epsilon_f, \\ \|\tilde{\varphi} - \varphi\| &\leq \epsilon_\varphi, & \|\tilde{\psi} - \psi\| &\leq \epsilon_\psi, \\ \|\tilde{L} - L\| &\leq \epsilon_L, & \|\tilde{R} - R\| &\leq \epsilon_R, \end{aligned} \quad (14)$$

где

$$\max\{\epsilon_A, \epsilon_f, \epsilon_\varphi, \epsilon_\psi, \epsilon_L, \epsilon_R\} = \varepsilon - \text{достаточно маленькое число.}$$

Всюду  $f(x), \tilde{f}(x)$  – непрерывные вектор-функции. Оказывается, что из разрешимости первой задачи при не слишком большом  $\varepsilon$  вытекает разрешимость второй.

**Теорема 1** Пусть  $u(x)$  – решение правильно поставленной задачи (9), для функций Грина, которой справедливы оценки (11). Пусть наряду с краевой задачей (9) имеется близкая к ней краевая задача (13), причем для элементов, определяющих эту краевую задачу, справедливы оценки (14). Тогда при

$$\varepsilon < \min\{1/[2K(2+d)], 1/(2+d)\} \quad (15)$$

существует и единственно  $\tilde{u}(x)$  – решение задачи (13) и для него верна оценка

$$\max_{x \in [0, d]} \|\tilde{u}(x) - u(x)\| \leq \varepsilon \mu, \quad (16)$$

причем  $\mu$  – число обусловленности задачи равно

$$\mu = K(2+d)(1 + \max_{x \in [0, d]} \|\tilde{u}(x)\|) \leq K(2+d)(1 + 2KF), \quad (17)$$

где

$$F = 1 + \|\varphi\| + \|\psi\| + d \max_{x \in [0, d]} \|f(x)\|.$$

(Доказательство – в [9].)

Оценка означает, что решение  $u(x)$  краевой задачи непрерывно зависит от коэффициентов  $A, L, R$  и правых частей  $\varphi, \psi, f(x)$ . В оценку непрерывности входит длина  $d$  отрезка, на котором ищется решение, и оценка  $K$  норм матриц Грина. Видно, что первое выражение для числа обусловленности  $\mu$  в (17) удобно использовать, когда уже есть вычисленное решение  $\tilde{u}(x)$ , второе – при прогнозировании погрешности, когда численное решение еще не известно.

**Вспомогательные леммы.** Нам понадобятся следующие леммы.

**Лемма 1** Пусть  $Z$  и  $\Delta$  – матрицы размера  $n \times p$  ( $n > p$ ), обладающие свойствами

$$Z^*Z = I_p, \quad \|\Delta\| \leq \varepsilon,$$

где  $\varepsilon$  – достаточно мало. Тогда для любой квадратной  $(n \times n)$ -матрицы  $M$  существует квадратная  $(n \times n)$ -матрица  $\tilde{M}$  такая, что

$$MZ + \Delta = \tilde{M}Z$$

и верна оценка

$$\|\tilde{M} - M\| \leq p\varepsilon. \quad (18)$$

**Лемма 2** Пусть даны матрицы  $Z$  и  $\Delta$  размера  $n \times p$  ( $n > p$ ), такие, что

$$Z^*Z = I_p, \quad \|\Delta\| \leq \varepsilon,$$

где  $\varepsilon$  – достаточно мало. Матрица  $U(x_s)$  размера  $n \times p$  – решение задачи Коши:

$$\frac{dU(x)}{dx} = A(x)U(x), \quad U(x_{s-1}) = Z,$$

причем

$$|x_s - x_{s-1}| \leq \frac{C}{\max_{x \in [x_{s-1}, x_s]} \|A(x)\|}.$$

Тогда матрица  $\tilde{U}(x_s) = U(x_s) + \Delta$  может быть получена, как решение задачи Коши:

$$\frac{d\tilde{U}(x)}{dx} = \tilde{A}(x)\tilde{U}(x), \quad \tilde{U}(x_{s-1}) = Z,$$

причем, при  $\varepsilon \leq \frac{1}{2pe^C}$  справедлива оценка

$$\max_{x \in [x_{s-1}, x_s]} \|\tilde{A}(x) - A(x)\| \leq 2e^C(C+1)p\varepsilon/(x_s - x_{s-1}) \leq \frac{1}{2SK}, \quad (19)$$

где  $S = \max\{2, d\}$ . (Доказательство – в [9].)

## 4 Детали алгоритма и его погрешности

Численное решение задачи (1), полученное по Алгоритму 2.2, обозначим  $u^{[c]}(x)$  ("с" от английского "computation"). Вычисленное решение  $u^{[c]}(x)$  можно рассматривать как решение краевой задачи

$$\begin{aligned} \frac{du^{[c]}(x)}{dx} &= \tilde{A}(x)u^{[c]}(x) + \tilde{f}(x), \\ Lu^{[c]}(x_0) &= \tilde{\varphi}, \quad Ru^{[c]}(x_m) = \tilde{\psi}, \end{aligned} \quad (20)$$

где возмущения

$$\|\tilde{\varphi} - \varphi\| = \|\Delta\varphi\| \leq \epsilon_\varphi, \quad \|\tilde{\psi} - \psi\| = \|\Delta\psi\| \leq \epsilon_\psi,$$

$$\max_{x \in [x_0, x_m]} \|A - \tilde{A}\| = \|\Delta A\| \leq \epsilon_A, \quad \max_{x \in [x_0, x_m]} \|f(x) - \tilde{f}(x)\| = \|\Delta f_1 + \Delta f_2\| \leq \epsilon_f,$$

невелики. Следовательно, краевая задача (20) может считаться достаточно близкой к исходной (1) и применение Теоремы 1 дает оценку решения

$$\max_{x \in [x_0, x_m]} \|u^{[c]}(x) - u(x)\| \leq \epsilon\mu, \quad (21)$$

причем число обусловленности  $\mu$  равно

$$\mu = K(2 + d)(1 + \max_{x \in [x_0, x_m]} \|\tilde{u}(x)\|).$$

Видно, что для вычисления погрешности решения (21) необходимо знать  $\epsilon$  и  $K$ . Для оценки  $K$  (11) существуют алгоритмы вычисления матриц Грина методом ортогональной встречной прогонки (см. [6]). Чтобы оценить погрешности в неравенствах (14) и выбрать  $\epsilon = \max\{\epsilon_\varphi, \epsilon_\psi, \epsilon_L, \epsilon_R, \epsilon_A, \epsilon_f\}$  нужно повторить процесс решения краевой задачи по Алгоритму 2.2, учитывая погрешности, которые возникают в промежуточных вычислениях.

Забегая вперед скажем, что методика описанная в [9] позволяет рассматривать вычислительные погрешности, возникающие при прямой прогонке, как возмущения матрицы  $\Delta A$  и правой части  $\Delta f_1$ , погрешности, возникающие при обратной прогонке – как возмущение правой части  $\Delta f_2$ , в граничных же точках их можно представить как возмущения  $\Delta\varphi$ ,  $\Delta\psi$  и, применяя метод обратного анализа погрешностей, оценить. Мы используем эту методику только при прямой прогонке. Погрешности обратной прогонки оцениваются напрямую.



## 4.1 Ввод данных

В зависимости от исследуемой задачи природа погрешности входных данных для алгоритма прогонки может быть самой разной. Если входные данные вычисляются на основе точных данных, то погрешность этих вычислений порядка  $\varepsilon_1$  ( $\varepsilon_1$  – параметр разрядной сетки компьютера и при вычислениях с двойной точностью  $\varepsilon_1 \sim 10^{-16}$ ). А так как эта погрешность вносится в решение один раз, то не нужно следить за ее накоплением, т. е. мы пренебрегаем погрешностями  $\|\tilde{L} - L\| \leq \varepsilon_L$  и  $\|\tilde{R} - R\| \leq \varepsilon_R$ .

## 4.2 Левое граничное условие

Сначала коротко изложим способ определения  $z_j$  и  $z_f$  из левого граничного условия, затем оценим погрешности.

Допустим, необходимо решить систему

$$LZ = 0, \quad (22)$$

где матрица  $L$  размера  $k \times n$ , ( $k + p = n$ ), вектор  $Z$  размера  $n$ .

С помощью ортогональных преобразований отражения получаем ортогональную матрицу  $Q$  такую, что матрица  $L^*$  представима в виде произведения двух матриц – верхнетреугольной  $\Omega$  размера  $n \times k$  и квадратной матрицы  $Q$  размера  $n \times n$ . Для построения отражений воспользуемся рекомендациями из Главы 3 монографии [3]. В результате получим верхнетреугольную матрицу  $\Omega_0$

$$\Omega_0 = Q_0 L^* \quad \text{или} \quad L^* = Q_0^* \Omega_0.$$

При этом система (22) запишется в виде

$$\Omega_0^* Q_0 Z_0 = 0 \quad \text{или} \quad \Omega_0^* Y_0 = 0, \quad \text{где} \quad Y_0 = Q_0 Z_0. \quad (23)$$

Для системы (23) полный ортогональный базис подпространства векторов, удовлетворяющих  $\Omega_0^* Y_0 = 0$ , составляют векторы  $e_{k+1}, \dots, e_n$ , где  $e_i$  – вектор, у которого все компоненты, кроме  $i$ -ой, равны 0, а  $i$ -я компонента равна 1.

Отсюда получаем решение в (22)  $Z_0 = Q_0^* Y_0$ . Столбцы  $z_j(x_0)$ ,  $j = 1, \dots, p$ , матрицы  $Z_0$  являются искомым начальным базисом.

Вектор  $z_f(x_0)$  находится из решения недоопределенной системы линейных уравнений

$$Lz_f = \varphi. \quad (24)$$

В результате вычисления на компьютере вместо векторов  $z_j(x_0)$  и  $z_f(x_0)$  будут найдены близкие к ним  $z_j^{[c]}(x_0) = z_j(x_0) + \zeta_j^{(0)}$  и  $z_f^{[c]}(x_0) = z_f(x_0) + \zeta_f^{(0)}$ .

Так как векторы  $z_j$  получаются из столбцов матрицы  $Q_0$ , то, очевидно, погрешности  $\zeta_j^{(0)}$  связаны с погрешностью вычисления матрицы  $Q_0$ :

$$\|\zeta_j^{(0)}\| \leq \|Q_0^{[c]} - Q_0\| \leq \epsilon_z, \quad (25)$$

где  $\epsilon_z$  вычисляется по следующим формулам (см.[3], стр. 66,96):

$$\begin{aligned} \epsilon_z &= (p+2)\sqrt{n} \tilde{\epsilon}_P(n), \\ \tilde{\epsilon}_P(n) &= \epsilon_1(5n+2\gamma+10). \end{aligned} \quad (26)$$

Здесь параметры компьютера:  $\gamma$  – основание машинной арифметики,  $\epsilon_1$  – относительная погрешность единицы. Обычно для персональных компьютеров  $\gamma = 2$  и  $\epsilon_1 \approx 10^{-16}$ .

Для оценки точности решения недоопределенной системы (24) воспользуемся упрощенным вариантом оценки (8.71) из [3]:

$$\begin{aligned} \|z_f - z_f^{[c]}\| &\leq [\tilde{\epsilon}_P(n)(3(k+1)\sqrt{6k}+1)\mu(L) + \tilde{\epsilon}_P(n)(n+1)] \|z_f^{[c]}\| \leq \\ &\leq [\tilde{\epsilon}_P(n)(3(k+1)\sqrt{6k}+1)\mu(L) + \tilde{\epsilon}_P(n)(n+1)] \|u^{[c]}\| = \epsilon_{z_f}^{(0)}. \end{aligned}$$

Здесь  $\mu(L)$  – число обусловленности матрицы  $L$ .

В итоге погрешность вычисления начального базиса оценивается величиной

$$\max\{\epsilon_{z_f}^{(0)}, \epsilon_z\} = \epsilon_0, \quad (27)$$

т. е. вычислительный процесс при прямой прогонки (интегрирование системы (4) ) начинается с неточных векторов  $z_j^{[c]}(x_0)$  и  $z_f^{[c]}(x_0)$ , определенных с погрешностью (27).

### 4.3 Оценка погрешности ортогонализации

В результате интегрирования системы (4) на отрезке  $[x_{s-1}, x_s]$  мы получим приближенные значения  $y_1^{[c]}(x_s), \dots, y_p^{[c]}(x_s), y_f^{[c]}(x_s)$  в точке  $x_s$ , которые рассматриваем как точные решения системы (62) и следующим шагом в Алгоритме 2.2 является их ортогонализация.

Для ортогонализации базисных векторов мы использовали метод ортогональных отражений, подробно описанный в Главах 3, 4 монографии [3]. Согласно этому методу, если

$$Y_s^{[c]} = [y_1^{[c]}(x_s), \dots, y_p^{[c]}(x_s), y_f^{[c]}(x_s)],$$

то точное  $QR$ -разложение

$$\bar{R}_s = \bar{Q}_s^* Y_s^{[c]} \quad \text{или} \quad Y_s^{[c]} = \bar{Q}_s \bar{R}_s,$$

где

$$\bar{Q}_s = P_{p+1} P_p \dots P_1,$$

$P_j$  - преобразования отражения, аннулирующие элементы  $j$ -го столбца матрицы  $P_{j-1} \dots P_1 \bar{Y}_s$ . Отсюда, если  $\bar{q}_1, \bar{q}_2, \dots, \bar{q}_{p+1}$  - первые  $p+1$  столбцов матрицы  $\bar{Q}_s$ , то стартовые значения для интегрирования на следующем шаге берутся в виде

$$\bar{z}_i(x_s) = q_i, \quad \bar{z}_f(x_s) = \bar{r}_s q_{p+1},$$

а матрица  $\bar{\Omega}_s$ , используемая при обратной прогонке, совпадает с верхней квадратной частью матрицы  $\bar{R}_s$  за исключением последнего элемента на диагонали:  $\bar{\omega}_{p+1, p+1}^{(s)} = 1$ . В действительности при реализации этого подхода на компьютере возникают арифметические погрешности. В результате вместо векторов  $\bar{z}_i, \bar{z}_f$  и матрицы  $\bar{\Omega}_s$  будут вычислены  $z_i^{[c]}(x_s), z_f^{[c]}(x_s)$  и  $\Omega_s^{[c]}$ .

Для оценок погрешностей вычисления  $\bar{Q}_s$  и  $\bar{\Omega}_s$  мы как и при выводе оценки (25) воспользуемся неравенствами (3.18) и (4.18) из монографии [3]:

$$\|\bar{z}_i(x_s) - z_i^{[c]}(x_s)\| \leq \|Q_s^{[c]} - \bar{Q}_s\| \leq \epsilon_z, \quad (28)$$

где  $\epsilon_z$  определяется по формуле (26).

Чтобы подготовить оценку точности для  $z_f^{[c]}(x_s)$  воспользуемся неравенством треугольника

$$\begin{aligned} \|\bar{z}_f - z_f^{[c]}\| &= \|\bar{q}_{p+1} \cdot \bar{r}_s \pm \bar{q}_{p+1} \cdot r_s^{[c]} - q_{p+1}^{[c]} \cdot r_s^{[c]}\| \leq \\ &\leq \|\bar{q}_{p+1}\| |\bar{r}_s - r_s^{[c]}| + \|\bar{q}_{p+1} - q_{p+1}^{[c]}\| |r_s^{[c]}| \end{aligned} \quad (29)$$

Далее оцениваем отдельно каждое слагаемое.

Напомним, что точность  $QR$ -разложения оценивается по формуле (4.18) из [3]:

$$\|R_s^{[c]} - \bar{R}_s\| \leq (p+2)\sqrt{p+1} \tilde{\epsilon}_P(n) \|\bar{R}_s\| \leq \epsilon_z \|\bar{R}_s\|. \quad (30)$$

Также воспользуемся оценками

$$\|z_f(x)\| \leq \|u(x)\|, \quad \|z_f^{[c]}(x)\| \leq \|u^{[c]}(x)\|.$$

По формуле (2.35) из [3] (оценка точности умножения вектора на скаляр)

$$\|z_f^{[c]}\| \leq (1+\alpha) \|q_{p+1}^{[c]} \cdot r^{[c]}\|,$$

где  $|\alpha| \leq \varepsilon_1$

$$(1+\alpha)(1-\epsilon_z) \leq \|z_f^{[c]} \cdot \frac{1}{r_s^{[c]}}\| = (1+\alpha) \|q_{p+1}^{[c]}\| \leq (1+\alpha)(1+\epsilon_z).$$

Следовательно, для всех точек  $x_s$  верна оценка

$$|r_s^{[c]}| \leq \frac{\|z_f^{[c]}\|}{(1-\varepsilon_1)(1-\epsilon_z)} \leq \frac{\|u^{[c]}\|}{(1-\varepsilon_1)(1-\epsilon_z)}.$$

Для верхнетреугольной матрицы верна цепочка неравенств

$$\|\bar{R}_s\| = \|Y_s^{[c]}\| \leq \|\bar{Y}_s\| + \epsilon_I \leq e^C \|Z_{s-1}^{[c]}\| + \epsilon_I \leq e^C (1 + \epsilon_z + |r_{s-1}^{[c]}|) + \epsilon_I,$$

где  $\epsilon_I$  – погрешность интегрирования задачи Коши (см. формулу (44)).

Подставим сюда оценку (30):

$$|r_s^{[c]} - \bar{r}_s| \leq \|R_s^{[c]} - \bar{R}_s\| \leq \epsilon_z \|\bar{R}_s\| \leq \epsilon_z \left[ e^C \left( 1 + \epsilon_z + \frac{\|u^{[c]}(x_{s-1})\|}{(1-\varepsilon_1)(1-\epsilon_z)} \right) + \epsilon_I \right]$$

Для получения окончательного результата используем также факт

$$\|\bar{q}_{p+1} - q_{p+1}^{[c]}\| \leq \epsilon_z$$

и подставляем найденные оценки в (29):

$$\begin{aligned} \|z_f^{[c]}(x_s) - \bar{z}_f(x_s)\| &\leq \|\bar{q}_{p+1}\| |\bar{r} - r^{[c]}| + \|\bar{q}_{p+1} - q_{p+1}^{[c]}\| |r^{[c]}| \leq \\ &\leq \epsilon_z \left[ e^C \left( 1 + \epsilon_z + \frac{\|u^{[c]}(x_{s-1})\|}{(1-\varepsilon_1)(1-\epsilon_z)} \right) + \epsilon_I \right] + \epsilon_z \frac{\|u^{[c]}(x_s)\|}{(1-\varepsilon_1)(1-\epsilon_z)} = \end{aligned}$$

$$= \epsilon_z \left[ e^C(1 + \epsilon_z) + \epsilon_I + (1 + e^C) \frac{K_{sol}}{(1 - \epsilon_1)(1 - \epsilon_z)} \right] = \epsilon_{zf}^{(s)} \quad (31)$$

Для обратной прогонки понадобится также оценка точности матрицы  $\Theta_s$  – подматрицы матрицы  $\Omega_s$  (см. Алгоритм 2.2). Вновь используя формулу (4.18) из [3], получаем неравенство:

$$\|\bar{\Theta}_s - \Theta_s^{[c]}\| \leq \epsilon_z \|\bar{\Theta}_s\|. \quad (32)$$

#### 4.4 Кусочно-линейное приближение матрицы $A(x)$ и правой части $f(x)$

Часто алгоритм ортогональной прогонки применяется к задачам, для которых матричная и векторная функции  $A(x)$  и  $f(x)$  заданы только в отдельных точках, т. е. матрицу и правую часть мы сможем вычислить только в точках  $x_s$ :  $x_s - x_{s-1} = d/m$ ,  $0 = x_0 < x_1 < \dots < x_s < \dots < x_m = d$ . Для того, чтобы приблизить  $A(x)$  и  $f(x)$  непрерывными функциями, применим простейший метод интерполяции при помощи кусочно-линейных функций. Таким образом, задача (1) заменяется на следующую

$$\begin{aligned} \frac{du_{lin}(x)}{dx} &= A_{lin}(x)u_{lin} + f_{lin}(x), \\ Lu_{lin}(x_0) &= \varphi, \quad Ru_{lin}(x_N) = \psi, \end{aligned} \quad (33)$$

где на отрезке  $[x_{s-1}, x_s]$  функции  $A_{lin}(x)$  и  $f_{lin}(x)$  представляются в виде

$$A_{lin}(x) = A_1^{(s)}x + A_0^{(s)}, \quad f_{lin}(x) = f_1^{(s)}x + f_0^{(s)},$$

так что

$$\begin{aligned} A_{lin}(x_{s-1}) &= A(x_{s-1}), & A_{lin}(x_s) &= A(x_s) \\ f_{lin}(x_{s-1}) &= f(x_{s-1}), & f_{lin}(x_s) &= f(x_s) \end{aligned}$$

Для того, чтобы оценить погрешности  $\|f(x) - f_{lin}(x)\|$  и  $\|A(x) - A_{lin}(x)\|$ , используем теорему о полиномиальной интерполяции.

**Теорема 2** Пусть  $f(x)$  имеет  $n + 1$  непрерывную производную, тогда для фиксированных узлов  $x_0 < x_1 < \dots < x_n$  существует единственный полином  $p(x)$  степени не больше  $n$  такой, что  $p(x_j) = f(x_j)$ . Причем для любой точки  $x$ ,  $x_0 < x < x_n$  имеется представление

$$f(x) - p(x) = \frac{(x - x_0)(x - x_1) \dots (x - x_n)}{(n + 1)!} f^{(n+1)}(z), \quad \text{где } z \in [x_0, x_n].$$

В нашем случае интерполяция производится полиномами первого порядка. Поэтому введем величину  $K_{dif}$ , которая оценивает вторые производные элементов матричной и векторной функций  $A(x)$  и  $f(x)$  так, что

$$\|A''(x)\| \leq K_{dif} \quad \text{и} \quad \|f''(x)\| \leq K_{dif} \quad (34)$$

при  $x \in [x_0, x_m]$ . Тогда по теореме о полиномиальной интерполяции получаем требуемые оценки

$$\|f(x) - f_{lin}(x)\| \leq K_{dif} \frac{h^2}{2} = \epsilon_{lin}, \quad \|A(x) - A_{lin}(x)\| \leq K_{dif} \frac{h^2}{2} = \epsilon_{lin}. \quad (35)$$

На основании этих неравенств по Теореме 1 получаем оценку для разности решений исходной задачи (1) и задачи с кусочно-линейными коэффициентами

$$\max_{x \in [x_0, x_m]} \|u - u_{lin}\| \leq \frac{h^2}{2} K_{dif} \mu, \quad (36)$$

где

$$\mu = K_{lin}(2 + d)(1 + \max_{x \in [x_0, x_m]} \|u_{lin}\|).$$

Здесь  $K_{lin}$  означает оценку матричной функции Грина задачи (33).

Далее в тех случаях, когда кусочная линейность функций несущественна, мы будем опускать обозначение "lin".

## 4.5 Погрешности интегрирования

### Метод интегрирования

Для интегрирования мы используем самый простой метод – метод Эйлера. У этого метода есть недостаток – уменьшение его погрешности (т. е. соответствующее уменьшение шага интегрирования) требует больших вычислительных затрат. Несмотря на это мы используем этот метод по следующим соображениям. Во-первых, во многих задачах интегрируются не сами функции, а их приближения. В данном случае используется также наиболее простое кусочно-линейное приближение (см. Раздел 4.4). Следовательно, уже на этапе интерполяции коэффициентов задачи вносится определенная погрешность. И при выборе метода интегрирования достаточно следить за тем, чтобы погрешность интегрирования не превышала погрешности интерполяции. В предлагаемом варианте это условие будет соблюдено.

*Замечание.* Из этих рассуждений также следует, что для того, чтобы повысить общую точность решения задачи, нужно параллельно повышать порядок интерполяции и точность интегрирования.

Вторая причина использования именно метода Эйлера – это его простота, которая позволяет относительно легко проследить за накоплением и распространением погрешностей в процессе интегрирования.

При этом метод Эйлера можно все же несколько видоизменить, если учесть специальный кусочно-линейный вид матрицы коэффициентов и правой части с тем, чтобы минимизировать погрешности интегрирования.

### Погрешность метода интегрирования на одном шаге

Рассмотрим отрезок  $[z_0, z_1]$ , где  $z_1 - z_0 = \Delta$ . Пусть  $v(z)$  есть решение задачи Коши

$$\frac{dv}{dz} = A_{lin}(z)v + f_{lin}(z), \quad v(z_0) = v_0.$$

тогда

$$v(z_1) = v_0 + \int_{z_0}^{z_1} (A_{lin}(z)v(z) + f_{lin}(z))dz$$

Сравним это выражение со следующим:

$$w(z_1) = v_0 + (A_{mid}v_0 + f_{mid})\Delta.$$

Здесь введено обозначение

$$A_{mid} = \frac{A_{lin}(z_1) - A_{lin}(z_0)}{2}, \quad f_{mid} = \frac{f_{lin}(z_1) - f_{lin}(z_0)}{2}.$$

Так как

$$A_{mid} = (A_0 + A_1 z_1)(z_1 - z_0) - A_1 \frac{(z_1 - z_0)^2}{2}$$

и

$$\int_{z_0}^{z_1} f_{lin}(z)dz = f_{mid},$$

то

$$v(z_1) - w_{z_1} = \int_{z_0}^{z_1} A_{lin}(z)v(z)dz - A_{lin}(z_1)v_0\Delta + A_1 v_0 \frac{\Delta^2}{2}.$$

Перегруппируем слагаемые и введем для них обозначения

$$v(z_1) - w(z_1) = A_0 \left( \int_{z_0}^{z_1} v(z)dz - v_0\Delta \right) +$$

$$+A_1 \left( \int_{z_0}^{z_1} v(z)z dz - z_1 v_0 \Delta + v_0 \frac{\Delta}{2} \right) = A_0 S_0 + A_1 S_1.$$

Оценка первого слагаемого дает следующий результат

$$\|S_0\| = \left\| \int_{z_0}^{z_1} (v(z) - v_0) dz \right\| \leq \max |v'(z)| \frac{\Delta^2}{2}.$$

А так как

$$z_1 = \frac{z_1 - z_0}{2} + \frac{z_1 + z_0}{2} \quad \text{и} \quad z_1 \Delta - \Delta/2 = \frac{z_1^2 - z_0^2}{2},$$

то

$$\begin{aligned} \|S_1\| &= \left\| \int_{z_0}^{z_1} (v(z) - v_0) z dz \right\| = \\ &= \left\| \int_{z_0}^{z_1} v'(\xi(z))(z - z_0) z dz \right\| \leq \max |v'(z)| \left( z_0 \frac{\Delta^2}{2} + \frac{\Delta^3}{3} \right). \end{aligned}$$

Учтем также, что для средней точки  $z_{mid}$  отрезка  $[z_0, z_1]$  имеют место выражения:

$$z_{mid} = \frac{\Delta}{2} + z_0, \quad z_0 \frac{\Delta^2}{2} + \frac{\Delta^3}{3} = \frac{2}{3} \Delta^2 \left( \frac{3}{4} z_0 + \frac{\Delta}{2} \right) \leq \frac{2\Delta^2}{3} z_{mid}.$$

Тогда общая погрешность одного шага представленной модификации метода Эйлера оценивается по формуле

$$\|v(z_1) - w(z_1)\| \leq \frac{2}{3} \Delta^2 K_{sol} \|A_{mid}\|, \quad (37)$$

где

$$\begin{aligned} \|A_{mid}\| &= \|A_0\| + |z_{mid}| \|A_1\|, \\ \max |v'(z)| &\leq \|v\|_{C^1} = \max |v(z)| + \max |v'(z)| = K_{sol}, \end{aligned} \quad (38)$$

где  $z_{mid} = (z_1 - z_0)/2$ .

### Арифметическая погрешность на одном шаге интегрирования

Сравним результат точного выражения

$$w(z_1) = v_0 + (A_{mid} v_0 + f_{mid}) \Delta$$

с результатом, полученным при вычислении на компьютере

$$w^{[c]}(z_1) = v_0 \oplus (A_{mid}^{[c]} \otimes v_0 \oplus f_{mid}^{[c]}) \otimes \Delta,$$



где

$$A_{mid}^{[c]} = (A_{lin}(z_1) \ominus A_{lin}(z_0)) \oslash 2, \quad f_{mid}^{[c]} = (f_{lin}(z_1) \ominus f_{lin}(z_0)) \oslash 2.$$

Здесь и далее в кружочках обозначены машинные (компьютерные) бинарные арифметические операции над машинными числами. Так как моделирование погрешностей машинного вычитания происходит по правилу

$$a \ominus b = (1 + \alpha)(a - b) + \beta, \quad |\alpha| \leq \varepsilon_1, \quad |\beta| \leq \varepsilon_0,$$

(см. [3]), а деление на 2 происходит с абсолютной погрешностью не более  $\varepsilon_0$ , то

$$\begin{aligned} \|\Omega_A\| &= \|A_{mid}^{[c]} - A_{mid}\| \leq \varepsilon_0(1 + n/2) + \varepsilon_1\sqrt{n}\|A_{mid}\|, \\ \|\omega_f\| &= \|f_{mid}^{[c]} - f_{mid}\| \leq \varepsilon_0(1 + \sqrt{n}/2) + \varepsilon_1\|f_{mid}\|, \end{aligned}$$

где  $n \times n$  – размер матрицы  $A$ . Следовательно,

$$w^{[c]}(z_1) = v_0 \oplus [(A_{mid} + \Omega_A) \otimes v_0 \oplus (f_{mid} + \omega_f)] \otimes \Delta.$$

При дополнительных условиях

$$\|A_{mid}\| \geq \frac{\varepsilon_0(1 + n/2)}{\varepsilon_1}$$

и

$$\|f_{mid}\| \geq \frac{\varepsilon_0(1 + \sqrt{n}/2)}{\varepsilon_1}$$

оценки имеют относительный характер

$$\|\Omega_A\| \leq \varepsilon_1(\sqrt{n} + 1)\|A_{mid}\|, \quad \|\omega_f\| \leq 2\varepsilon_1\|f_{mid}\|.$$

Именно их и имеет смысл использовать, так как дополнительные условия нарушаются только в исключительных случаях для практически нулевой матрицы  $A$  и правой части  $f$ .

Оценка точности машинного умножения матрицы на вектор имеет вид

$$\|A \otimes v - Av\| \leq (n + 1)\sqrt{n}\varepsilon_1\|A\|\|v\|$$

при условии

$$\begin{aligned} \sqrt{2n\varepsilon_0/\varepsilon_1} &\leq \|v\| \leq \sqrt{\varepsilon_\infty/2}, \\ \sqrt{2n\varepsilon_0/\varepsilon_1} &\leq \|a_i\| \leq \sqrt{\varepsilon_\infty/2}, \end{aligned}$$

где  $a_i$  –  $i$ -я строка матрицы  $A$ . Эти условия не являются слишком обременительными и, как правило, выполняются. Следовательно,

$$w^{[c]}(z_1) = v_0 \oplus [(A_{mid} + \Omega_A)v_0 + \nu] \oplus (f_{mid} + \omega_f) \otimes \Delta,$$

где

$$\|\nu\| \leq (n+1)\sqrt{n}\varepsilon_1(1 + \varepsilon_1(\sqrt{n} + 1))\|A_{mid}\|\|v_0\|.$$

Для моделирования погрешностей машинного сложения и умножения на скаляр будем использовать формулы (см. Главу 2 в [3]):

$$a \oplus b = (1 + \alpha)(a + b), \quad |\alpha| \leq \varepsilon_1,$$

$$a \otimes v = (1 + \alpha)(av), \quad |\alpha| \leq \varepsilon_1.$$

*Замечание.* Мы не учитываем абсолютные погрешности порядка  $\varepsilon_0$  не только потому, что они малы по сравнению с другими погрешностями, но и потому, что они возникают только в единичных случаях, когда результат операции близок к нулю, и, следовательно, не накапливаются.

В результате применения этих формул получим

$$w^{[c]}(z_1) = v_0 + [(A_{mid} + \Omega_A)v_0 + \nu + f_{mid} + \omega_f + \xi] \Delta + \zeta + \eta,$$

где

$$\|\xi\| \leq \varepsilon_1\|(A_{mid} + \Omega_A)v_0 + \nu + f_{mid} + \omega_f\|,$$

$$\|\zeta\| \leq \varepsilon_1\Delta\|(A_{mid} + \Omega_A)v_0 + \nu + f_{mid} + \omega_f + \xi\|,$$

$$\|\eta\| \leq \varepsilon_1\|v_0 + [(A_{mid} + \Omega_A)v_0 + \nu + f_{mid} + \omega_f + \xi] \Delta + \zeta\|.$$

Группируя все оценки погрешностей, получим

$$w^{[c]}(z_1) = w(z_1) + \zeta + \eta + [\Omega_A v_0 + \nu + \omega_f + \xi] \Delta = w(z_1) + \chi,$$

а для величины  $\chi$ , которая оценивает арифметические погрешности на одном шаге интегрирования, нетрудно получить неравенство.

$$\|w(z_1) - w^{[c]}(z_1)\| = \|\chi\| \leq \varepsilon_1 \left[ (1 + \Delta[4 + c_1])\|A_{mid}\|\|v_0\| + 8\|f_{mid}\| \right], \quad (39)$$

где  $c_1 = (1 + \varepsilon_1\sqrt{n})(n+1)(\sqrt{n} + 1) \approx n^{3/2}$ ,  $n$  – размер вектора.

### Накопление погрешностей интегрирования

Рассмотрим отрезок  $[z_0, z_1]$ , где  $z_1 - z_0 = \Delta$ . Предположим теперь, что значение функции  $v(z)$  в начальной точке  $z_0$  задано с погрешностью:  $\tilde{v}_0 = v(z_0) + \vartheta_0$ . Вычислим, как это скажется на погрешности в точке  $z_1$ . Для этого будем сравнивать величины

$$v(z_1) = v_0 + \int_{z_0}^{z_1} (A_{lin}(z)v(z) + f_{lin}(z))dz \quad (40)$$

и

$$\tilde{w}^{[c]}(z_1) = \tilde{v}_0 \oplus (A_{mid}^{[c]} \otimes \tilde{v}_0 \oplus f_{mid}^{[c]}) \otimes \Delta.$$

Применим очевидное неравенство

$$\|v(z_1) - \tilde{w}^{[c]}(z_1)\| \leq \|v(z_1) - w(z_1)\| + \|w(z_1) - \tilde{w}(z_1)\| + \|\tilde{w}(z_1) - \tilde{w}^{[c]}(z_1)\|. \quad (41)$$

Первое слагаемое в правой части этого неравенства было оценено ранее (37). Для того, чтобы оценить последнее слагаемое, модифицируем (39):

$$\begin{aligned} \|\tilde{w}(z_1) - \tilde{w}^{[c]}(z_1)\| &\leq \varepsilon_1 \left[ (1 + \Delta[4 + c_1]\|A_{mid}\|)\|\tilde{v}_0\| + 8\|f_{mid}\| \right] \leq \\ &\leq \varepsilon_1 \left[ (1 + \Delta[4 + c_1]\|A_{mid}\|)\|v_0\| + 8\|f_{mid}\| \right] + \varepsilon_1 \left[ (1 + \Delta[4 + c_1]\|A_{mid}\|)\|\vartheta_0\| \right]. \end{aligned}$$

Легко также вывести оценку для среднего слагаемого в (41):

$$\|w(z_1) - \tilde{w}(z_1)\| = \|(I + \Delta A_{mid})\vartheta_0\| \leq (1 + \Delta\|A_{mid}\|)\|\vartheta_0\|.$$

Следовательно, оценка имеет вид

$$\begin{aligned} \|v(z_1) - \tilde{w}^{[c]}(z_1)\| &\leq (1 + (\varepsilon_1 + \Delta(1 + \varepsilon_1(4 + c_1))\|A_{mid}\|)\|\vartheta_0\| + \\ &+ \left(\frac{2}{3}\Delta^2\|A_{mid}\| + \varepsilon_1 + \varepsilon_1\Delta(4 + c_1)\|A_{mid}\|)K_{sol} + \varepsilon_1 8\|f_{mid}\|. \end{aligned}$$

Усиливая неравенство, получим

$$\|v(z_1) - \tilde{w}^{[c]}(z_1)\| = \|\vartheta_1\| \leq (1 + r)\|\vartheta_0\| + q,$$

где

$$r = \varepsilon_1 + \Delta(1 + \varepsilon_1(4 + c_1))K_A\varepsilon_1, \quad K_A = \max\|A(z)\|,$$

$$q = (\varepsilon_1 + \left(\frac{2}{3}\Delta^2 + \varepsilon_1\Delta(4 + c_1)\right)K_A)K_{sol} + \varepsilon_1 8K_f, \quad K_f = \max\|f(z)\|,$$

а  $K_{sol}$  определяется в (38).

Таким образом, мы получили погрешность  $\|\vartheta_1\|$  вычисления значения функции  $v(z)$  по формуле (40) с учетом погрешности метода Эйлера, машинных погрешностей и начальной погрешности  $\|\vartheta_0\|$ .

Перейдем теперь к исследованию вопроса о накоплении алгоритмической и арифметической погрешностей в процессе интегрирования по отрезку. Итак, отрезок  $[x_0, x_m]$  разбит на  $m$  равных частей  $[x_{s-1}, x_s]$ ,  $s = 1, 2, \dots, m+1$ , длина которых составляет  $x_s - x_{s-1} = (x_m - x_0)/m = h$ . Рассмотрим отдельно некоторый отрезок  $[x_{s-1}, x_s]$ , по которому мы будем интегрировать векторную функцию  $y(x)$ , являющуюся решением задачи Коши

$$\frac{dy}{dx} = A_{lin}(x)y + f_{lin}(x), \quad y(x_{s-1}) = y_{s-1}.$$

С этой целью мы произведем дополнительное разбиение отрезка  $[x_{s-1}, x_s]$ , поделив его на  $N$  равных частей с шагом интегрирования  $\Delta = h/N = x^{(k)} - x^{(k-1)}$ :

$$x_{s-1} = x^{(0)} < x^{(1)} < x^{(2)} < \dots < x^{(N)} = x_s.$$

Из предыдущих рассуждений получаем

$$\begin{aligned} \|\vartheta_N\| &\leq (1+r)\|\vartheta_{N-1}\| + q \leq (1+r)^2\|\vartheta_{N-2}\| + (1+r)q + q \leq \dots \\ &\leq (1+r)^N\|\vartheta_0\| + q \sum_{i=0}^{N-1} (1+r)^i, \end{aligned}$$

где

$$\vartheta_i = y(x^{(i)}) - \tilde{y}^{[cl]}(x^{(i)}).$$

Так как

$$(1+r)^N \leq \frac{1}{1-Nr} \quad \text{и} \quad \sum_{i=0}^{N-1} (1+r)^i = \frac{(1+r)^N - 1}{r} \leq \frac{N}{1 - (N-1)r/2},$$

то в итоге получаем

$$\|\vartheta_N\| \leq \frac{1}{1-Nr}\|\vartheta_0\| + \frac{N}{1 - (N-1)r/2}q.$$

Заметим, что  $rN \approx \Delta \varepsilon_1 K_A N = \varepsilon_1 h K_A$ ,  $q \approx 2/3 \Delta^2 K_A K_{sol}$ , тогда можно прояснить порядок накопленной погрешности

$$\|\vartheta_N\| \lesssim \frac{\|\vartheta_0\|}{1 - \varepsilon_1 h K_A} + \frac{2}{3} \frac{\Delta h}{1 - \varepsilon_1 h K_A / 2} K_A K_{sol} = k_b \|\vartheta_0\| + \epsilon_I = \epsilon_C. \quad (42)$$

Таким образом, мы получили погрешность вычисления векторной функции  $y(x)$  в точке  $x_s$  с учетом накопления всех погрешностей, возникающих при интегрировании на отрезке  $[x_{s-1}, x_s]$ : начальной в точке  $x_{s-1}$ , машинной и алгоритмической. В оценке (42) мы используем следующие обозначения, которые понадобятся в дальнейшем:

$$k_b = 1/(1 - \varepsilon_1 h K_A) \quad (43)$$

– коэффициент влияния погрешности начальных данных на результат интегрирования задачи Коши ("b" от слова "begin"),

$$\varepsilon_I = \frac{2}{3} \frac{\Delta h}{1 - \varepsilon_1 h K_A / 2} K_A \quad (44)$$

– алгоритмическая и арифметическая погрешность интегрирования задачи Коши модифицированным методом Эйлера (индекс "I" означает "Integral"),  $\varepsilon_C$  – общая погрешность интегрирования задачи Коши ("C" по имени "Cauchy").

## 4.6 Обратный анализ погрешностей прямой прогонки

Рассмотрим погрешности, возникающие при прямой прогонке на шаге  $s$  (Алгоритм 2.2). Если бы вычисления выполнялись точно, то формулы  $s$  шага прямой прогонки выглядели бы следующим образом

$$y_f(x_s) = X(x_s, x_{s-1}) z_f(x_{s-1}) + X(x_s, x_{s-1}) \int_{x_{s-1}}^{x_s} [X(\xi, x_{s-1})]^{-1} f(\xi) d\xi, \quad (45)$$

и

$$y_j(x_s) = X(x_s, x_{s-1}) z_j(x_{s-1}), \quad (46)$$

где  $X(x, x_s)$  – матрицант, т. е. матрица-функция размера  $n \times n$ , решение задачи Коши

$$\frac{dX(x, x_s)}{dx} = A(x)X(x, x_s), \quad X(x_s, x_s) = I_n, \quad (47)$$

но из левого граничного условия начальный базис определяется с погрешностью (27). Таким образом, дальнейший вычислительный процесс начинается с этих неточных векторов. Далее при прямой прогонке путем численного интегрирования на отрезке  $[x_{s-1}, x_s]$  серии задач Коши (4)

с учетом погрешностей ортогонализации и интегрирования в точке  $x_s$  будут вычислены значения векторов  $y_f^{[c]}(x_s)$ ,  $y_j^{[c]}(x_s)$ :

$$y_f^{[c]}(x_s) = X(x_s, x_{s-1})\bar{z}_f(x_{s-1}) + X(x_s, x_{s-1}) \int_{x_{s-1}}^{x_s} [X(\xi, x_{s-1})]^{-1} f(\xi) d\xi + \eta_f^{(s)}, \quad (48)$$

и

$$y_j^{[c]}(x_s) = X(x_s, x_{s-1})\bar{z}_j(x_{s-1}) + \eta_j^{(s)}, \quad (49)$$

где

$$\|\eta_f^{(s)}\| \leq k_b \epsilon_{z_f}^{(s)} + \epsilon_I, \quad \|\eta_j^{(s)}\| \leq k_b \epsilon_z + \epsilon_I,$$

$\epsilon_{z_f}^{(s)}$ ,  $\epsilon_z$ ,  $k_b$ ,  $\epsilon_I$  определяются по формулам (31), (26), (43), (44). Обозначим

$$\max\{\|\eta_f^{(s)}\|, \|\eta_j^{(s)}\|\} = \delta. \quad (50)$$

Рассмотрим уравнение (49). Из утверждения Леммы 1 следует, что существует матрица  $\tilde{X}(x_s, x_{s-1})$  для которой верно

$$X(x_s, x_{s-1})\bar{z}_j(x_{s-1}) + \eta_j^{(s)} = \tilde{X}(x_s, x_{s-1})\bar{z}_j(x_{s-1}) \quad (51)$$

и

$$\|\tilde{X}(x_s, x_{s-1}) - X(x_s, x_{s-1})\| \leq p\delta.$$

Это означает, что существует такая система дифференциальных уравнений

$$\frac{du(x)}{dx} = \tilde{A}(x)u(x), \quad (52)$$

что векторы  $y_1^{[c]}(x_s), y_2^{[c]}(x_s), \dots, y_p^{[c]}(x_s)$  есть решения следующих задач Коши

$$\frac{dy_j^{[c]}(x)}{dx} = \tilde{A}(x)y_j^{[c]}(x), \quad y_j^{[c]}(x_{s-1}) = \bar{z}_j(x_{s-1})$$

и

$$y_j^{[c]}(x_s) = \tilde{X}(x_s, x_{s-1})\bar{z}_j(x_{s-1}),$$

где  $\tilde{X}(x_s, x_{s-1})$  – матрицант системы (52).

При этом погрешности вычислений интерпретируются как малые возмущения кусочно-линейного приближения матрицы  $A$  и при  $\delta \leq 1/2pe^C$  для любого интервала  $[x_{s-1}, x_s]$  верна оценка (19) из Леммы 2:

$$\max_{x \in [x_{s-1}, x_s]} \|\tilde{A}(x) - A_{lin}(x)\| \leq \delta p 2e^C (C + 1)/h = \epsilon_A. \quad (53)$$

Теперь рассмотрим уравнение (48) и получим оценку возмущения правой части  $f$ . Положим для остальных  $x \in [x_{s-1}, x_s]$

$$\tilde{X}(x, x_{s-1}) = X(x, x_{s-1}) + \frac{x - x_{s-1}}{x_s - x_{s-1}} \left( \tilde{X}(x_s, x_{s-1}) - X(x_s, x_{s-1}) \right). \quad (54)$$

Перепишем выражение (48) следующим образом

$$\begin{aligned} y_f^{[c]}(x_s) &= \tilde{X}(x_s, x_{s-1}) \bar{z}_f(x_{s-1}) + \eta_f^{(s)} + \left( X(x_s, x_{s-1}) - \tilde{X}(x_s, x_{s-1}) \right) \bar{z}_f(x_{s-1}) \\ &\quad + \tilde{X}(x_s, x_{s-1}) \int_{x_{s-1}}^{x_s} [\tilde{X}(\xi, x_{s-1})]^{-1} f_1(\xi) d\xi, \end{aligned} \quad (55)$$

где

$$f_1(x) = \tilde{X}(x, x_{s-1}) [\tilde{X}(x_s, x_{s-1})]^{-1} X(x_s, x_{s-1}) [X(x, x_{s-1})]^{-1} f(x).$$

С целью оценить разность  $\|f(x) - f_1(x)\|$  выполним следующие преобразования:

$$\begin{aligned} f(x) - f_1(x) &= \left( I - \tilde{X}(x, x_{s-1}) [\tilde{X}(x_s, x_{s-1})]^{-1} X(x_s, x_{s-1}) [X(x, x_{s-1})]^{-1} \right) f(x) \\ &\quad \doteq \left( X(x_{s-1}, x) [X(x_{s-1}, x_s)]^{-1} - \tilde{X}(x, x_{s-1}) [\tilde{X}(x_s, x_{s-1})]^{-1} \right) \\ &\quad \times X(x_s, x_{s-1}) [X(x, x_{s-1})]^{-1} f(x) \\ &= \left( X(x, x_{s-1}) [X(x_s, x_{s-1})]^{-1} - X(x, x_{s-1}) [\tilde{X}(x_s, x_{s-1})]^{-1} \right) \\ &\quad - \frac{x - x_{s-1}}{x_s - x_{s-1}} \left( \tilde{X}(x_s, x_{s-1}) - X(x_s, x_{s-1}) \right) [\tilde{X}(x_s, x_{s-1})]^{-1} \\ &\quad \times X(x_s, x_{s-1}) [X(x, x_{s-1})]^{-1} f(x) \\ &= \left( X(x, x_{s-1}) \left( [X(x_s, x_{s-1})]^{-1} - [\tilde{X}(x_s, x_{s-1})]^{-1} \right) \right. \\ &\quad \left. - \frac{x - x_{s-1}}{x_s - x_{s-1}} \left( \tilde{X}(x_s, x_{s-1}) - X(x_s, x_{s-1}) \right) [\tilde{X}(x_s, x_{s-1})]^{-1} \right) \\ &\quad \times X(x_s, x_{s-1}) [X(x, x_{s-1})]^{-1} f(x). \end{aligned}$$

Обозначим для краткости  $X = X(x_s, x_{s-1})$ ,  $\tilde{X} = \tilde{X}(x_s, x_{s-1})$ ,  $\mathcal{X} = X - \tilde{X}$ , тогда нетрудно получить следующие оценки

$$\|\tilde{X}^{-1}\| = \|(X - \mathcal{X})^{-1}\| = \|X^{-1}(I - \mathcal{X}X^{-1})\|.$$

Если выполнено условие  $\|\mathcal{X}\| \cdot \|X^{-1}\| < 1$ , а в нашем случае это означает  $p\delta e^C < 1$ , то можно применить ряд Неймана и получить оценку

$$\begin{aligned} \|\tilde{X}^{-1}\| &= \|X^{-1} \sum_{j=0}^{\infty} (\mathcal{X}X^{-1})^j\| \leq \|X^{-1}\| \sum_{j=0}^{\infty} (\|\mathcal{X}\| \cdot \|X^{-1}\|)^j = \\ &= \frac{\| [X(x_s, x_{s-1}) ]^{-1} \|}{1 - \|\mathcal{X}\| \| [X(x_s, x_{s-1}) ]^{-1} \|}. \end{aligned}$$

Аналогичным образом получаем

$$\begin{aligned} \|\tilde{X}^{-1} - X^{-1}\| &= \|(X - \mathcal{X})^{-1} - X^{-1}\| = \|X^{-1} \left( \sum_{j=0}^{\infty} (\mathcal{X}X^{-1})^j - I \right)\| \leq \\ &\leq \frac{\| [X(x_s, x_{s-1}) ]^{-1} \|^2 \|\mathcal{X}\|}{1 - \|\mathcal{X}\| \| [X(x_s, x_{s-1}) ]^{-1} \|}. \end{aligned}$$

Поскольку в нашем случае

$$\begin{aligned} \|\tilde{X}(x_s, x_{s-1}) - X(x_s, x_{s-1})\| &= \|\mathcal{X}\| \leq p\delta, \\ \|X(x_s, x_{s-1})\| &\leq e^C, \quad \| [X(x_s, x_{s-1}) ]^{-1} \| \leq e^C, \end{aligned}$$

то

$$\begin{aligned} \| [ \tilde{X}(x_s, x_{s-1}) ]^{-1} \| &\leq \frac{e^C}{1 - p\delta e^C}, \\ \| [ \tilde{X}(x_s, x_{s-1}) ]^{-1} - [ X(x_s, x_{s-1}) ]^{-1} \| &\leq p\delta \frac{e^{2C}}{1 - p\delta e^C}. \end{aligned}$$

Объединяя полученные оценки, выводим итоговое неравенство

$$\begin{aligned} &\max_{x \in [x_{s-1}, x_s]} \|f(x) - f_1(x)\| \leq \\ &\leq \left( p\delta \frac{e^{3C}}{1 - p\delta e^C} + p\delta \frac{e^C}{1 - p\delta e^C} \right) e^{2C} \max_{x \in [x_{s-1}, x_s]} \|f(x)\| = \\ &= p\delta \frac{e^{3C}}{1 - p\delta e^C} (e^{2C} + 1) \max_{x \in [x_{s-1}, x_s]} \|f(x)\|. \end{aligned} \tag{56}$$

Введем обозначение

$$\chi_s = \eta_f^{(s)} + \left( X(x_s, x_{s-1}) - \tilde{X}(x_s, x_{s-1}) \right) \bar{z}_f(x_{s-1}).$$



Очевидно, что

$$\|\chi_s\| \leq \delta + p\delta \|\tilde{z}_f(x_{s-1})\| \leq \delta(1 + pK_{sol}^{[c]}),$$

так как

$$\|\tilde{z}_f(x_{s-1})\| \leq \max_{x \in [x_0, x_m]} \|u^{[c]}(x)\| \leq K_{sol}^{[c]}. \quad (57)$$

Рассмотрим на интервале  $[x_{s-1}, x_s]$  функцию

$$f_2(x) = \tilde{X}(x, x_{s-1})[\tilde{X}(x_s, x_{s-1})]^{-1}\chi_s/(x_s - x_{s-1}).$$

На основе проведенных выше оценок, получаем неравенство

$$\|f_2(x)\| \leq (e^C + p\delta) \frac{e^C}{1 - p\delta e^C} \frac{\delta(1 + pK_{sol}^{[c]})}{h}. \quad (58)$$

Тогда выражение (55) можно переписать в виде

$$y_f^{[c]}(x_s) = \tilde{X}(x_s, x_{s-1})\tilde{z}_f(x_{s-1}) + \tilde{X}(x_s, x_{s-1}) \int_{x_{s-1}}^{x_s} [\tilde{X}(\xi, x_{s-1})]^{-1} \tilde{f}(\xi) d\xi,$$

где

$$\tilde{f}(\xi) = f_1(\xi) + f_2(\xi).$$

Это означает, что вектор  $y_f^{[c]}(x_s)$  является значением решения задачи Коши:

$$\frac{dy_f^{[c]}(x)}{dx} = \tilde{A}(x)\tilde{y}_f(x) + \tilde{f}(x), \quad \tilde{y}_f(x_{s-1}) = \tilde{z}_f(x_{s-1})$$

в точке  $x_s$ . Отсюда, исходя из (56), (58), (35) имеем оценку

$$\begin{aligned} & \max_{x \in [x_{s-1}, x_s]} \|f_{lin}(x) - \tilde{f}(x)\| \leq \max_{x \in [x_{s-1}, x_s]} \|f_{lin}(x) - f_1(x)\| + \max_{x \in [x_{s-1}, x_s]} \|f_2(x)\| \leq \\ & \leq \frac{\delta}{1 - p\delta e^C} \left( pe^{3C}(e^{2C} + 1) \max_{x \in [x_{s-1}, x_s]} \|f(x)\| + \frac{1}{h} e^C (e^C + p\delta)(1 + pK_{sol}^{[c]}) \right) = \epsilon_F, \end{aligned} \quad (59)$$

где  $\delta$  определяется по формуле (50). Очевидно, последняя оценка может быть сделана только при условии

$$p\delta e^C < 1. \quad (60)$$

Таким образом, обратный анализ позволяет нам рассматривать вычисленные векторы как точные решения серии задач Коши со слегка возмущенными матрицей  $\tilde{A}$  (53) и правой частью  $f$  (59).

**Резюме.** В результате прямой прогонки в каждой точке  $x_s$  будут вычислены матрицы  $Z_s^{[c]}$  и  $\Omega_s^{[c]}$ . Согласно Алгоритму 2.2 решение задачи в каждой точке  $x_s$  представимо в виде линейной комбинации столбцов матрицы  $Z_s^{[c]} = [\hat{Z}_s, \bar{z}_f(x_s)]$ , где  $\hat{Z}_s = [z_1(x_s), \dots, z_p(x_s)]$ .

Вычисленные верхне-треугольные матрицы  $\Omega_s^{[c]}$  имеют специальный вид

$$\Omega_s^{[c]} = \begin{bmatrix} \Theta_s^{[c]} & \theta_s^{[c]} \\ 0 & 1 \end{bmatrix}. \quad (61)$$

Если из правого условия и при обратной прогонке будут найдены векторы  $\bar{\alpha}_s$  (черта сверху означает, что величины не являются точными) такие, что

$$[R\bar{Z}_m]\alpha_m = \psi - R\bar{z}_f(x_s), \quad \bar{\Theta}_s\bar{\alpha}_{s-1} = \alpha_s - \bar{\theta}_s,$$

то функция  $\tilde{u}(x_s) = \bar{Z}_s\bar{\alpha}_s + \bar{z}_f(x_s)$  в точках  $x_s$  совпадает с решением системы

$$\begin{cases} \frac{d\tilde{u}}{dx} = \tilde{A}(x)\tilde{u}(x) + \tilde{f}(x), \\ L\tilde{u}(x_0) = \varphi, \quad R\tilde{u}(x_m) = \psi, \end{cases} \quad (62)$$

которое согласно Теореме 1 отличается от  $u_{lin}$  на величину

$$\|u_{lin}(x_s) - \tilde{u}(x_s)\| \leq \max\{\epsilon_A, \epsilon_F\}K_{lin}(2+d)(1 + \max\|\tilde{u}(x_s)\|). \quad (63)$$

## 4.7 Правое граничное условие

Мы уже оценили влияние вычислительных и алгоритмических погрешностей прямой прогонки на решение, используя метод обратного анализа погрешностей. Влияние погрешностей обработки правого граничного условия и обратной прогонки мы в отличие от [9] оценим напрямую.

Обратная прогонка начинается с того, что  $u(x_m)$  представляется в виде

$$u(x_m) = \bar{z}_f(x_m) + \bar{Z}(x_m)\bar{\alpha}_m,$$

а вектор-коэффициент  $\bar{\alpha}_m$  находится из правого граничного условия (6):

$$[R\bar{Z}(x_m)]\bar{\alpha}_m = \psi - R\bar{z}_f(x_m).$$

В действительности же, с учетом погрешностей, решение ищется в виде

$$u^{[c]}(x_m) = z_f^{[c]}(x_m) + \hat{Z}^{[c]}(x_m)\alpha_m^{[c]} \quad (64)$$

так, что

$$[R\hat{Z}^{[c]}(x_m)]\alpha_m^{[c]} = \psi - Rz_f^{[c]}(x_m). \quad (65)$$

Очевидно, что погрешность решения в точке  $x_m$  есть

$$u(x_m) - u^{[c]}(x_m) = \bar{z}_f(x_m) + \bar{Z}(x_m)\bar{\alpha}_m - z_f^{[c]}(x_m) - \hat{Z}^{[c]}(x_m)\alpha_m^{[c]}.$$

Преобразуем скобки в последнем выражении к удобному виду

$$\begin{aligned} u(x_m) - u^{[c]}(x_m) &= \left( \bar{z}_f(x_m) - z_f^{[c]}(x_m) \right) + \bar{Z}(x_m)(\bar{\alpha}_m - \alpha_m^{[c]}) + \\ &+ \left( \bar{Z}(x_m) - \hat{Z}^{[c]}(x_m) \right) \alpha_m^{[c]}. \end{aligned} \quad (66)$$

Следовательно

$$\begin{aligned} &\|u(x_m) - u^{[c]}(x_m)\| \leq \\ &\leq \|\bar{z}_f(x_m) - z_f^{[c]}(x_m)\| + \|\bar{Z}(x_m)\| \|\bar{\alpha}_m - \alpha_m^{[c]}\| + \|\bar{Z}(x_m) - \hat{Z}^{[c]}(x_m)\| \|\alpha_m^{[c]}\|. \end{aligned} \quad (67)$$

Оценим каждое слагаемое в отдельности.

Погрешности  $\|\bar{z}_f(x_m) - z_f^{[c]}(x_m)\|$  и  $\|\bar{Z}(x_m) - \hat{Z}^{[c]}(x_m)\|$  возникают в процессе ортогонализации и оцениваются по формулам (28) и (31).

Далее, из представления (64) следует, что

$$\|\alpha_m^{[c]}\| \leq \|u^{[c]}(x_m)\|.$$

Осталось оценить погрешность вычисления  $\alpha^{[c]}$ . Так как вектор  $\alpha_m^{[c]}$  является решением системы алгебраических уравнений (65), то можно применить следующую оценку. Если есть две системы линейных алгебраических уравнений: исходная  $Ax = f$  и возмущенная  $(A + \Delta A)(x + \Delta x) = f + \Delta f$ , то возмущение решения оценивается по формуле

$$\begin{aligned} \|\Delta x\| &\leq (\|\Delta f\| + \|\Delta A\| \|x\|) \|(A + \Delta A)^{-1}\| \leq \\ &\leq (\|\Delta f\| + \|\Delta A\| \|x\|) \frac{\|A^{-1}\|}{1 - \|\Delta A\| \|A^{-1}\|}. \end{aligned} \quad (68)$$

Рассмотрим вспомогательную краевую задачу

$$\begin{cases} \frac{du^R}{dx} = A(x)u^R, \\ Lu^R(x_0) = 0, \quad Ru^R(x_m) = \psi \end{cases}.$$

Структура решения этой задачи  $u^R(x_m) = \hat{Z}^R(x_m)\alpha_m^R$ , так как компонента  $z_f^R = 0$  ввиду отсутствия правых частей в уравнении и в левом краевом условии. При этом имеет место оценка  $\|u^R(x)\| \leq K\|\psi\|$ , где  $K$  как и ранее означает оценку нормы матрицы Грина. Так как вследствие ортогональности матрицы  $\hat{Z}^R(x_m)$  имеет место неравенство

$$\|u^R(x_m)\| = \|\hat{Z}^R(x_m)\alpha_m^R\| = \|\alpha_m^R\| \leq K\|\psi\|,$$

выполненное при любом векторе  $\psi$ , то  $\|[R\hat{Z}^R(x_m)]^{-1}\| \leq K$ . А по построению решение задачи с произвольной правой частью

$$\begin{cases} \frac{du}{dx} = A(x)u + f(x), \\ Lu^R(x_0) = \varphi, \quad Ru^R(x_m) = \psi \end{cases}$$

в точке  $x_m$  имеет структуру  $u(x_m) = \hat{Z}(x_m)\alpha_m + z_f(x_m)$ , причем  $\hat{Z}(x_m) = \hat{Z}^R(x_m)$ , то  $\|[R\hat{Z}(x_m)]^{-1}\| \leq K$ .

Рассмотрим теперь возмущенную задачу

$$\begin{cases} \frac{d\tilde{u}^R}{dx} = \tilde{A}(x)\tilde{u}^R, \\ L\tilde{u}^R(x_0) = 0, \quad R\tilde{u}^R(x_m) = \psi. \end{cases}$$

Воспользуемся фактом, что если  $\mathcal{L}$  – некий обратимый линейный оператор, то

$$\|(\mathcal{L} + \Delta\mathcal{L})^{-1}\| \leq \frac{\|\mathcal{L}^{-1}\|}{1 - \|\mathcal{L}^{-1}\|\|\Delta\mathcal{L}\|}$$

при достаточно малом возмущении  $\Delta\mathcal{L}$ :  $\|\mathcal{L}^{-1}\|\|\Delta\mathcal{L}\| < 1$ . В нашей ситуации  $\mathcal{L}$  это дифференциальный оператор задачи и норма обратного оператора оценивается через норму функции Грина  $\|\mathcal{L}^{-1}\| \leq K(2+d)$ . Для возмущения оператора используем оценку  $\|\Delta\mathcal{L}\| = \max\|A(x) - \tilde{A}(x)\| \leq \epsilon_A$

и после подстановки получаем

$$\|(\mathcal{L} + \Delta\mathcal{L})^{-1}\| \leq \frac{K(2+d)}{1 - \epsilon_A K(2+d)}.$$

Возвращаясь к решению возмущенной задачи и используя полученный результат, приходим к оценке

$$\|\tilde{u}^R\| \leq \frac{K(2+d)}{1 - \epsilon_A K(2+d)} \|\psi\|,$$

а следовательно,

$$\|(R\tilde{Z}(x_m))^{-1}\| \leq \frac{K(2+d)}{1 - \epsilon_A K(2+d)}, \quad \|\bar{\alpha}_m\| \leq \frac{K(2+d)}{1 - \epsilon_A K(2+d)} \|\psi\|.$$

Так как имеет место неравенство

$$\|R\hat{Z}^{[c]}(x_m) - R\tilde{Z}(x_m)\| = \|R(\hat{Z}^{[c]}(x_m) - \tilde{Z}(x_m))\| \leq p\epsilon_z \|R\| = p\epsilon_z \|R\tilde{Z}(x_m)\|,$$

то на основании оценки (4.7) получаем в точке  $x_m$

$$\|\bar{\alpha}_m - \alpha_m^{(c)}\| \leq (\epsilon_F + p\epsilon_z \|\bar{\alpha}_m\|) \frac{K(2+d)\|R\|\|\psi\|}{1 - K(2+d)(\epsilon_A + p\epsilon_z\|R\|\|\psi\|)} = \epsilon_\alpha \quad (69)$$

при условии

$$K(2+d)(\epsilon_A + p\epsilon_z\|R\|\|\psi\|) < 1.$$

При выводе оценки (69) мы не учли погрешность с которой  $\alpha_m^{[c]}$  определяется при численном решении системы (65) на компьютере, а также погрешность перемножения матриц  $R$  и  $\hat{Z}^{[c]}(x_m)$  размера  $p \times n$ ,  $n \times p$  и погрешность вычисления правой части. Чтобы сделать это нужно сравнить решения следующих систем

$$[R \otimes \hat{Z}^{[c]}(x_m)] \otimes \alpha_m^{[c]} = \psi \ominus R \otimes z_f^{[c]}(x_m),$$

$$[R \otimes \hat{Z}^{[c]}(x_m)] \cdot \hat{\alpha}_m = \psi \ominus R \otimes z_f^{[c]}(x_m),$$

$$[R\hat{Z}^{[c]}(x_m)] \cdot \tilde{\alpha}_m = \psi - Rz_f^{[c]}(x_m),$$

$$[R\tilde{Z}(x_m)] \cdot \bar{\alpha}_m = \psi - R\bar{z}_f(x_m)$$

(крюжочками обозначены машинные операции, см. [3]) и оценить выражение

$$\|\alpha_m^{[c]} - \bar{\alpha}_m\| \leq \|\alpha_m^{[c]} - \hat{\alpha}_m\| + \|\hat{\alpha}_m - \tilde{\alpha}_m\| + \|\tilde{\alpha}_m - \bar{\alpha}_m\|. \quad (70)$$

Сначала рассмотрим второе слагаемое. Для простоты введем обозначения

$$W = R\hat{Z}^{[c]}(x_m), \quad x = \tilde{\alpha}_m, \quad w = \psi - Rz_f^{[c]}(x_m)$$

и запишем систему (65) в виде

$$Wx = w. \quad (71)$$

Погрешность машинного (компьютерного) вычисления правой части равна (см. формулы (2.37), (2.40) в [3])

$$\begin{aligned} \|\Delta w\| &= \|[\psi \ominus R \otimes z_f^{[c]}(x_m)] - [\psi - Rz_f^{[c]}(x_m)]\| \\ &\leq 2\varepsilon_1 \|\psi\| + [2\varepsilon_1(1 + \delta_0) + \delta_0] \|R\| K_{sol}^{[c]}, \end{aligned}$$

где  $\delta_0 = (n + 1)\sqrt{p}\varepsilon_1$  и  $\|\Delta w\|/\|w\| \leq \delta_1$ . Погрешность машинного перемножения двух матриц  $A$ ,  $B$  размера  $p \times n$ ,  $n \times p$  определяем используя формулу (2.40) из [3]:

$$\|(A \otimes B) - AB\| \leq (n + 1)p\varepsilon_1 \|A\| \|B\|.$$

Таким образом, матрица  $W$  вычислена с погрешностью

$$\frac{\|W_{\text{маш}} - W\|}{\|W\|} \leq (n + 1)p\varepsilon_1 = \delta_2$$

и в действительности решается система

$$[W + \Delta W](x + \Delta x) = w + \Delta w,$$

из которой  $x$  определяется с погрешностью (см. формулу (8.33) в [3]):

$$\frac{\|\Delta x\|}{\|x\|} \leq (\delta_2 + \delta_1) \frac{\mu(W)}{1 - \delta_2 \mu(W)} = \delta_3,$$

где  $\mu(W)$  – число обусловленности матрицы  $W$ . Отсюда для второго слагаемого в (70) имеем

$$\frac{\|\hat{\alpha}_m - \tilde{\alpha}_m\|}{\|\tilde{\alpha}_m\|} \leq \delta_3 \quad \text{или} \quad \|\hat{\alpha}_m - \tilde{\alpha}_m\| \leq \frac{\delta_3}{1 + \delta_3} \|\hat{\alpha}_m\|.$$

Первое слагаемое в (70) оценивает погрешность решения на компьютере системы с квадратной матрицей  $W$  размера  $n$  (см. формулу (8.7) и таблицу 8.1 в [3]), которая при упрощении равна

$$\frac{\|\alpha_m^{[c]} - \hat{\alpha}_m\|}{\|\hat{\alpha}_m\|} \leq \tilde{\epsilon}_P(p)\mu(W)[3(p + 1)\sqrt{p} + 1] + \tilde{\epsilon}_P(p)(p + 1) = \delta_4,$$

где  $\tilde{\epsilon}_P(p)$  определяется по формуле (26). Следовательно,

$$\|\alpha_m^{[c]} - \hat{\alpha}_m\| \leq \frac{\delta_4}{1 + \delta_4} \|\alpha_m^{[c]}\|.$$

Третье слагаемое в (70) есть  $\epsilon_\alpha$  (69). Складывая все оценки, окончательно получаем

$$\|\bar{\alpha}_m - \alpha_m^{[c]}\| \leq \frac{(1 + \delta_3)\delta_4 + \delta_3}{(1 + \delta_3)(1 + \delta_4)} \|\alpha_m^{[c]}\| + \epsilon_\alpha = \epsilon_\alpha(x_m).$$

Отсюда, используя (67), имеем оценку погрешности решения в точке  $x_m$ :

$$\|u(x_m) - u^{[c]}(x_m)\| \leq \epsilon_{zf}^{(m)} + \epsilon_\alpha(x_m) + \epsilon_z(1 + \epsilon_\alpha(x_m)) = \epsilon_{U_m}. \quad (72)$$

## 4.8 Обратный анализ погрешностей обратной прогонки

При обратной прогонке особое влияние оказывают погрешности вычисления матриц  $\Omega_s$  (см. (32)). На интервале  $[x_{s-1}, x_s]$  вместо решения систем (8)

$$\bar{\Omega}_s \bar{\beta}_{s-1} = \bar{\beta}_s$$

будет получено решение системы

$$\Omega_s^{[c]} \beta_{s-1}^{[c]} = \beta_s^{[c]}.$$

Оценку погрешностей удобнее проводить в терминах подматриц (см. (61))

$$\Theta_s^{[c]} \alpha_{s-1}^{[c]} = \alpha_s^{[c]} - \theta_s^{[c]}.$$

Тогда

$$\begin{aligned} & \|\alpha_{s-1}^{[c]} - \bar{\alpha}_{s-1}\| \leq \\ & \leq \left( \|\alpha_{s-1}^{[c]} - \bar{\alpha}_{s-1}\| + \|\theta_s^{[c]} - \bar{\theta}_s\| + \|\Theta_s^{[c]} - \bar{\Theta}_s\| \|\bar{\alpha}_{s-1}\| \right) \frac{\|\bar{\Theta}_s^{-1}\|}{1 - \|\Theta_s^{[c]} - \bar{\Theta}_s\| \|\bar{\Theta}_s^{-1}\|}. \end{aligned}$$

Так как

$$\|\Omega_s^{[c]} - \bar{\Omega}_s\| \leq \epsilon_z \|\bar{\Omega}_s\| \leq \epsilon_z \|\tilde{X}(x_{s-1}, x_s)\| \leq \epsilon_z (e^C + p\delta),$$

то

$$\begin{aligned}\|\Theta_s^{[c]} - \bar{\Theta}_s\| &\leq \epsilon_z(e^C + p\delta), \\ \|\theta_s^{[c]} - \bar{\theta}_s\| &\leq \epsilon_z(e^C + p\delta).\end{aligned}$$

С учетом неравенств

$$\|\bar{\Theta}_s^{-1}\| \leq \frac{e^C}{1 - p\delta e^C} \quad \text{и} \quad \|\bar{\alpha}_{s-1}\| \leq \|\tilde{u}(x_{s-1})\|$$

имеем

$$\begin{aligned}&\|\alpha_{s-1}^{[c]} - \bar{\alpha}_{s-1}\| \leq \\ &\leq \left( \|\alpha_s^{[c]} - \bar{\alpha}_s\| + \epsilon_z(e^C + p\delta(1 + K_{sol})) \right) \frac{e^C}{1 - (1 + \epsilon_z)p\delta - \epsilon_z e^C}.\end{aligned}$$

Отсюда по индукции получаем оценку

$$\|\alpha_{s-1}^{[c]} - \bar{\alpha}_{s-1}\| \leq \|\alpha_m^{[c]} - \bar{\alpha}_m\| C_\alpha^{m-s} + p\delta(1 + K_{sol}) \frac{C_\alpha^{m-s} - 1}{C_\alpha - 1} = \epsilon_\alpha(x_{s-1}), \quad (73)$$

где

$$C_\alpha = \frac{e^C}{1 - (1 + \epsilon_z)p\delta - \epsilon_z e^C}$$

при условии

$$(1 + \epsilon_z)p\delta - \epsilon_z e^C < 1.$$

Решение в точке  $x_s$  определяется по формуле

$$u^{[c]}(x_s) = z_f^{[c]}(x_s) + Z_s^{[c]} \alpha_s^{[c]}.$$

Используя имеющиеся оценки (28), (73), (31), нетрудно вывести неравенство

$$\begin{aligned}\|u^{[c]}(x_s) - \tilde{u}(x_s)\| &\leq \|Z_s^{[c]} - \bar{Z}_s\| \|\bar{\alpha}_s\| + \|\bar{Z}_s\| \|\alpha_s^{[c]} - \bar{\alpha}_s\| + \|\bar{z}_f(x_s) - z_f^{[c]}(x_s)\| \leq \\ &\leq \epsilon_z K_{sol} + \epsilon_\alpha(x_s) + \epsilon_{z_f}^{(s)}\end{aligned} \quad (74)$$

– погрешность вносимая в решение задачи при обратной прогонки.



## 4.9 Оценка погрешности решения

Таким образом, в результате вычислений по Алгоритму 2.2 в каждой точке  $x_s$  будет получено решение  $u^{[c]}(x_s)$  краевой задачи близкой к исходной (1). Для определения погрешности решения необходимо собрать все полученные оценки (74), (63), (36). При этом нужно, чтобы были выполнены условия (например, типа (60)) при которых они выводились.

Окончательная оценка погрешности решения в точке  $x_s$  имеет вид

$$\begin{aligned} \|u^{[c]}(x_s) - u(x_s)\| &\leq \|u^{[c]}(x_s) - \tilde{u}(x_s)\| + \|\tilde{u}(x_s) - u_{lin}(x_s)\| + \|u_{lin}(x_s) - u(x_s)\| \leq \\ &\leq \epsilon_z K_{sol}^{[c]} + \epsilon_\alpha(x_s) + \epsilon_{z_f}^{(s)} + \epsilon\mu, \end{aligned} \quad (75)$$

где

$$\epsilon = \max\{\epsilon_A, \epsilon_F\} + \frac{h^2}{2} K_{dif} \quad \mu = K_{lin}(2+d)(1 + K_{sol}^{[c]}).$$

Здесь  $K_{lin}$  – оценка норм матриц Грина задачи (33),  $\epsilon_A$  находится из (53),  $\epsilon_f$  – из (59),  $\epsilon_z$  – из (28),  $\epsilon_{z_f}^{(s)}$  – из (31),  $\epsilon_\alpha(x_s)$  – из (73),  $K_{dif}$  – из (34),  $h = (x_0 - x_m)/m = x_s - x_{s-1}$  – длина интервалов на которые разбит отрезок  $[x_0, x_m]$ , а  $K_{sol}^{[c]}$  оценивается из (57) как

$$\max_{x \in [x_0, x_m]} \|u^{[c]}(x)\| \leq K_{sol}^{[c]}.$$

## 5 Примеры

### 5.1 Пример 1

Приведем простой модельный пример использования описанного выше метода численного решения краевой задачи. Рассмотрим следующую краевую задачу

$$\begin{aligned} \frac{du}{dx} &= A(x)u(x) + f(x), \\ Lu(x_0) &= \varphi, \quad Ru(x_m) = \psi, \\ x_0 \leq x \leq x_m, \quad x_0 &= 0.0, \quad x_m = 1.0, \quad d = x_m - x_0 = 1.0, \end{aligned} \quad (76)$$

где

$$\begin{aligned} A &= \begin{pmatrix} 0 & 1 \\ 2 & 0 \end{pmatrix}, \quad f(x) = \begin{pmatrix} 0 \\ -2x \end{pmatrix}, \\ L &= (0 \ 1), \quad R = (1 \ 0), \quad \varphi = 1, \quad \psi = 1. \end{aligned}$$

Сравним точное решение

$$u = \begin{pmatrix} x \\ 1 \end{pmatrix}$$

с численным решением задачи методом ортогональной прогонки, полученным с помощью программы Sweep из пакета **GALA-2.1**.

**Результаты вычислений.** В результате вычислений мы получили, что реальная погрешность вектора решения равна

$$\max_{x \in [x_0, x_m]} \|u^{[c]}(x) - u(x)\| \leq 3.28 \cdot 10^{-12}. \quad (77)$$

Оценка погрешности решения  $u^{[c]}(x_s)$  в точке  $x_s$  определяется по формуле (75), где число обусловленности  $\mu = K(2 + d)(1 + K_{sol}^{[c]}) = 23.18$  для данной задачи величина постоянная (оценка  $K$  для матриц Грина вычислена программой EstimGreenMat из пакета **GALA-2.1**), а остальные величины зависят от шага  $h = d/m$  разбижки отрезка  $[x_0, x_m]$  и от шага интегрирования  $\Delta = h/N$  на каждом интервале  $[x_s, x_{s+1}]$ . Видно, что для получения приемлемой оценки необходимо, чтобы выполнялось и условие (15) Теоремы 1. Для этого пришлось значительно уменьшать  $\Delta$ .

Например, оценка погрешности вычисленного решения равна

$$\max_{x \in [x_0, x_m]} \|u^{[c]}(x) - u(x)\| \leq 7.5 \cdot 10^{-2} \quad (78)$$

при  $m = 8$  и  $N = 500$ .

Видно, что даже в таком простом примере гарантированная погрешность (78) решения задачи оказалась значительно завышенной по сравнению с реальной (77). Это говорит о том, что оценка погрешности решения краевой задачи, основанная на Теореме 1 и на методике, описанной в работе [9], не является удовлетворительной с практической точки зрения. Видимо необходимо разрабатывать другой способ оценки погрешности решения краевой задачи (1), что мы и предполагаем сделать в дальнейшем.

## 5.2 Пример 2

Следующий пример – это компьютерное моделирование сердечно-сосудистой системы человека. Описанный алгоритм решения краевой задачи использовался для решения системы гемодинамики (одномерной

модели кровотока в сосуде) [4]. Мы рассматривали упрощенную гидродинамическую модель течения крови в сосудах. Ее математические свойства (гиперболичность) позволяют формулировать для нее корректные краевые задачи [4]. Одномерная модель выводится с помощью интегрирования уравнений Навье-Стокса по произвольному осевому сечению  $S$  (усреднением этих уравнений по поперечному направлению) и представляет собой систему из двух дифференциальных уравнений в частных производных (см. [4, 16, 15]):

$$\begin{aligned} \frac{\partial A}{\partial t} + \frac{\partial Q}{\partial z} &= 0, \\ \frac{\partial Q}{\partial t} + \alpha \frac{\partial}{\partial z} \left( \frac{Q^2}{A} \right) + \frac{A}{\rho} \frac{\partial p}{\partial p} + K_r \left( \frac{Q}{A} \right) &= 0, \end{aligned} \quad (79)$$

где  $A$ ,  $Q$  – неизвестные. Здесь  $A = A(t, z)$  – площадь поперечного сечения сосуда,  $Q = Q(t, z)$  – поток массы,  $p = p(t, z)$  – давление,  $\alpha = \text{const}$  – коэффициент Кориолиса,  $\rho = \text{const}$  – плотность,  $K_r = \text{const}$  – коэффициент трения,  $t$  – время,  $z$  – осевая координата.

Методом прямых задача сводится к решению системы обыкновенных дифференциальных уравнений на каждом временном слое, которая затем решается описанным методом ортогональной прогонки. Эта математическая модель *благодаря* своей простоте позволяет моделировать кровоток во всем артериальном сосудистом русле (глобальное моделирование), т. е. для всех 55 артерий сосудистого дерева человека. При этом, *несмотря* на простоту модели она позволила провести моделирование основных параметров кровотока как в норме, так и в ряде патологий. **Результаты вычислений** подробно изложены в Главе 4 монографии [4] и дипломной работе [10].

Как уже отмечалось выше, метод ортогональной прогонки, допускает *контроль за накоплением вычислительной погрешности*. В своей работе мы делаем на этом особый акцент. Связано это с тем, что увеличение объемов вычислений при современном уровне математического моделирования приводит к усилению роли вычислительной погрешности. Но многие алгоритмы не позволяют оценить ее вклад в решение, вследствие чего результат математического моделирования принимает характер *гипотезы*, которая в ряде случаев может достаточно сильно отличаться от точного решения исходной задачи.

В данном случае такой неопределенный характер численного решения неприемлем по двум основным причинам. Во-первых, пользователями данного алгоритма являются биологи, физиологи и другие специалисты,

далекие от вычислительной математики. В силу своей специальности они не знают и не обязаны знать особенности внутреннего устройства вычислительных программ. Но при этом пользователи должны иметь возможность делать обоснованные выводы из результатов математического моделирования. Во-вторых, в данном случае моделируются процессы, связанные с жизнью и здоровьем человека. Следовательно, используемые алгоритмы должны вести себя предсказуемым и контролируемым образом.

## Список литературы

- [1] *Бибердорф Э.А., Попова Н.И.* Решение линейных систем с гарантированной оценкой точности результатов (часть первая). Новосибирск, 1999 (Препринт ИЯФ СО РАН; 99-49).
- [2] *Бибердорф Э.А., Попова Н.И.* Вычисления с гарантированной оценкой точности (часть вторая). Решение спектральных задач. Новосибирск, 2001 (Препринт ИЯФ СО РАН; 2001-21).
- [3] *Бибердорф Э.А., Попова Н.И.* Гарантированная точность современных алгоритмов линейной алгебры. Новосибирск: Издательство СО РАН, 2006.
- [4] *Бибердорф Э.А., Блохин А.М., Попова Н.И., Трахнин Ю.Л.* Система кровообращения и артериальная гипертония: биофизические и генетико-физиологические механизмы, математическое и компьютерное моделирование. Глава 4: Глобальное моделирование артериальной системы человека. Новосибирск: Издательство СО РАН, 2008 (в печати).
- [5] *Годунов С.К.* О численном решении краевых задач для систем линейных обыкновенных дифференциальных уравнений // Успехи мат. наук, 1961, т.16, №3, с.171-175.
- [6] *Годунов С.К.* Обыкновенные дифференциальные уравнения с постоянными коэффициентами. Новосибирск: Издательство НУ, 1994.
- [7] *Годунов С.К., Антонов А.Г., Кириллюк О.П., Костин В.И.* Гарантированная точность решения систем линейных уравнений в евклидовых пространствах. Новосибирск: Наука, 1992.
- [8] *Годунов С.К.* Лекции по современным аспектам линейной алгебры. Новосибирск: Научная книга, 2002.
- [9] *Кузнецов С.В.* Развитие метода ортогональной прогонки. Диссертация. Институт математики СО РАН, Новосибирск: 1988.
- [10] *Леонова Т.И.* Построение, численное моделирование и анализ одномерной модели гемодинамики. Дипломная работа, механико-математический факультет НГУ, Новосибирск: 2008.
- [11] *Уилкинсон Д.Х.* Алгебраическая проблема собственных значений/ Пер. с англ. М.: Наука, 1970.

- [12] *Форсайт Дж., Малькольм М., Молер К.* Машинные методы математических вычислений/ Пер. с англ. М.: Мир, 1980.
- [13] *Хаусхолдер А.С.* Основы численного анализа. М.: Изд-во иностр. лит., 1956.
- [14] *Givens W.* Numerical Computation of the Characteristic Values of a Real Symmetric Matrix. Report ORNL 1574, Oak Ridge, Tenn., Oak Ridge National Laboratory, 1954.
- [15] *Lamponi D.N.* One dimensional and multiscale models for blood flow circulation. Thesis, Lausanne, 2004.
- [16] *Quarteroni A., Formaggia L.* Mathematical modelling and numerical simulation of the cardiovascular system. In: *Ciarlet P.G. (ed.). Handbook of numerical analysis, v.12, special volume: Ayache N. (ed.). Computational Models for the Human Body.* Amsterdam, Elsevier Science & Technology, 2003, p.3-129.

*Э.А. Бибердорф, Н.И. Попова*

**Контроль точности решения краевой задачи  
методом ортогональной прогонки**

*E.A. Biberdorf, N.I. Popova*

**Accuracy control of solving the boundary-value problem  
by the orthogonal sweep method**

ИЯФ 2009-1

Ответственный за выпуск А.В. Васильев

Работа поступила 30.12.2008 г.

---

Сдано в набор 11.01.2009 г.

Подписано в печать 12.01.2009 г.

Формат бумаги 60×90 1/16 Объем 3.0 печ.л., 2.4 уч.-изд.л.

Тираж 90 экз. Бесплатно. Заказ № 1

---

Обработано на IBM PC и отпечатано на  
ротапринте ИЯФ им. Г.И. Будкера СО РАН  
*Новосибирск, 630090, пр. академика Лаврентьева, 11.*