



ИПМ им.М.В.Келдыша РАН • Электронная библиотека

Препринты ИПМ • Препринт № 74 за 2019 г.



ISSN 2071-2898 (Print)
ISSN 2071-2901 (Online)

Бахвалов П.А., Сурначёв М.Д.

Линейные схемы с
несколькими степенями
свободы для многомерного
уравнения переноса

Рекомендуемая форма библиографической ссылки: Бахвалов П.А., Сурначёв М.Д.
Линейные схемы с несколькими степенями свободы для многомерного уравнения переноса //
Препринты ИПМ им. М.В.Келдыша. 2019. № 74. 44 с. doi:[10.20948/prepr-2019-74](https://doi.org/10.20948/prepr-2019-74)
URL: <http://library.keldysh.ru/preprint.asp?id=2019-74>

О р д е н а Л е н и н а
ИНСТИТУТ ПРИКЛАДНОЙ МАТЕМАТИКИ
имени М.В.КЕЛДЫША
Р о с с и й с к о й а к а д е м и и н а у к

П. А. Бахвалов, М. Д. Сурначёв

Л и н е й н ы е с х е м ы
с несколькими степенями свободы
для многомерного уравнения переноса

Москва — 2019

Бахвалов П. А., Сурначёв М. Д.

Линейные схемы с несколькими степенями свободы для многомерного уравнения переноса

Рассматриваются линейные схемы с несколькими степенями свободы на одну ячейку для многомерного уравнения переноса. Решение по такой схеме может обладать ошибкой $O(h^p + th^q)$, причём p совпадает с порядком аппроксимации или превосходит его на единицу, а $q \geq p$. Доказывается, что существует такое отображение гладких функций на сеточное пространство, отличающееся от обычного (например, L_2 -проекции) на величину порядка h^p , в смысле которого схема будет обладать q -м порядком аппроксимации. В отличие от одномерного случая, локальное отображение с требуемыми свойствами может не существовать. Приводятся достаточные условия его существования.

Ключевые слова: аппроксимация и точность, суперсходимость

Pavel Alexeevich Bakhvalov, Mikhail Dmitrievich Surnachev

Linear schemes with several degrees of freedom for the multidimensional transport equation

We consider linear schemes with several degrees of freedom for the transport equation. Solution error of these schemes can have estimate $O(h^p + th^q)$, where p is equal to or greater by one than approximation order and $q \geq p$. We prove the existence of a mapping of smooth functions on the mesh space providing the q -th order of the truncation error and deviating from the standard mapping (L_2 -projection for example) by the order h^p . In contrast with 1D case local mapping with such properties generally does not exist. We prove sufficient existence conditions.

Key words: consistency and accuracy, superconvergence

1. Введение

Настоящая работа продолжает исследование линейных полудискретных схем с несколькими степенями свободы на одну ячейку для уравнения переноса с постоянным коэффициентом, начатое в [1–4]. Численная ошибка решения таких схем может обладать оценкой $O(h^p + th^q)$, причём p совпадает с порядком аппроксимации или превосходит его на единицу, а $q \geq p$. Например, для метода Галёркина с разрывными базисными функциями эта оценка наблюдается при $p = k + 1$, $q = 2k + 1$, где k – порядок используемых полиномов. В одномерном случае эта оценка доказана в [5], в двумерном случае на декартовой сетке – в [6]. На трансляционно-инвариантной треугольной сетке эта оценка также наблюдается, однако строгого доказательства этого факта авторам неизвестно.

Как ошибка аппроксимации, так и ошибка решения зависят не только от уравнений, определяющих разностную схему, но и от отображения функций из некоторого пространства решений дифференциального уравнения в пространство сеточных функций. Для того чтобы ошибка решения при использовании некоторого отображения Π_h (например, L_2 -проекции) имела оценку $O(h^p + th^q)$, достаточно, чтобы существовало такое отображение $\tilde{\Pi}_h$, что:

- 1) $\|\tilde{\Pi}_h - \Pi_h\| = O(h^p)$;
- 2) в смысле $\tilde{\Pi}_h$ имеет место q -й порядок аппроксимации.

В [3] мы показали, что в одномерном случае существование такого отображения является не только достаточным, но и необходимым. При этом получение оптимальных значений p и q сводится к поиску отображения $\tilde{\Pi}_h$ определённого вида методом неопределённых коэффициентов.

Настоящая работа посвящена многомерному случаю. Как и в [3], в предположении, что ошибка решения обладает оценкой $O(h^p + th^q)$, доказывается существование отображения гладких функций на сеточное пространство, удовлетворяющего условиям 1) и 2). Однако в общем случае локального отображения, обладающего такими свойствами, не существует; приводятся достаточные условия его существования. Также изучается вопрос зависимости формального порядка точности схемы и порядка точности в длительном счёте от направления волнового вектора и направления скорости переноса.

2. Постановка задачи

Пусть $\mathbf{a}_j, j = 1, \dots, d$, – линейно независимая система векторов, интерпретируемых как смещения сеточных блоков друг относительно друга при $h = 1$. Пусть $\mathbf{a}_j^*, j = 1, \dots, d$, – взаимный базис к $\{\mathbf{a}_j\}$, то есть $\mathbf{a}_i^* \cdot \mathbf{a}_j = \delta_{ij}$. Обозначим через T линейный оператор в пространстве \mathbb{R}^d , сопоставляющий вектору $\boldsymbol{\eta}$ вектор $T\boldsymbol{\eta}$ с компонентами

$$T\boldsymbol{\eta} = \sum_{j=1}^d \eta_j \mathbf{a}_j. \quad (2.1)$$

Если $\boldsymbol{\eta} \in \mathbb{Z}^d$, то $T\boldsymbol{\eta}$ есть смещение сеточного блока $\boldsymbol{\eta}$ относительно нулевого. Обозначим $U = (T^*)^{-1}$. Имеем

$$(T^* \mathbf{x})_j = \mathbf{x} \cdot \mathbf{a}_j, \quad U\mathbf{y} = (T^*)^{-1}\mathbf{y} = \sum_{j=1}^d y_j \mathbf{a}_j^*.$$

Если \mathbf{a}_j совпадают с векторами стандартного базиса, T и U тождественные.

Введём пространство $L_{2,per}(\mathbb{R}^d)$ как множество функций $f \in L_{2,loc}(\mathbb{R}^d)$, таких что при некотором $N_0 \in \mathbb{N}$ для всех $j = 1, \dots, d$ и почти всех $\mathbf{r} \in \mathbb{R}^d$ выполняется $f(\mathbf{r} + N_0 \mathbf{a}_j) = f(\mathbf{r})$. При этом будем называть N_0 *периодом* функции f . Для $f \in L_{2,per}(\mathbb{R}^d)$ определим норму

$$\|f\|^2 = \frac{1}{|\square|} \int_{\square} |f(\mathbf{r})|^2 dV,$$

где \square – параллелепипед, образованный векторами $N_0 \mathbf{a}_1, \dots, N_0 \mathbf{a}_d$.

Для $q \in \mathbb{N} \cup \{0\}$ определим пространство $H_{per}^q(\mathbb{R}^d) = L_{2,per}(\mathbb{R}^d) \cap W_{2,loc}^q(\mathbb{R}^d)$. Для $w \in H_{per}^q(\mathbb{R}^d)$, $q \in \mathbb{N} \cup \{0\}$, $r = 0, \dots, q$, введём обозначение

$$\|\nabla^r w\|^2 = \sum_{|\mathbf{m}|=r} \frac{r!}{\mathbf{m}!} \|D^{\mathbf{m}} w\|^2, \quad D^{\mathbf{m}} = \frac{\partial^{|\mathbf{m}|}}{\partial x_1^{m_1} \dots \partial x_d^{m_d}}. \quad (2.2)$$

Здесь $\mathbf{m} = (m_1, \dots, m_d)$ – мультииндекс: $m_i \geq 0$, $|\mathbf{m}| = m_1 + \dots + m_d$, $\mathbf{m}! = m_1! \dots m_d!$. Введём обозначение $\mathbf{r}^{\mathbf{m}} = x_1^{m_1} \dots x_d^{m_d}$ для $\mathbf{r} = (x_1, \dots, x_d)$.

На $H_{per}^q(\mathbb{R}^d)$ будем использовать семейство норм $\|f\|_{(q,h)}^2 = \sum_{r=0}^q h^{2r} \|\nabla^r f\|^2$.

Также для $q \in \mathbb{N} \cup \{0\}$ введём $C_{per}^q(\mathbb{R}^d) = C^q(\mathbb{R}^d) \cup L_{2,per}(\mathbb{R}^d)$. На $C_{per}(\mathbb{R}^d) \equiv C_{per}^0(\mathbb{R}^d)$ будем использовать норму $\|f\|_{\infty} = \sup_{\mathbf{r}} |f(\mathbf{r})|$. На $C^q(\mathbb{R}^d)$ будем использовать семейство норм

$$\|f\|_{(\infty,q,h)}^2 = \sum_{r=0}^q h^{2r} \|\nabla^r f\|_{\infty}^2, \quad \|\nabla^r f\|_{\infty}^2 = \sum_{|\mathbf{m}|=r} \frac{r!}{\mathbf{m}!} \|D^{\mathbf{m}} w\|_{\infty}^2.$$

Такое введение норм обеспечивает выполнение условия $\|f\|_{(q,h)} \leq \|f\|_{(\infty,q,h)}$.

В настоящей работе рассматривается начальная задача для линейного уравнения переноса

$$\frac{\partial v}{\partial t} + \boldsymbol{\omega} \cdot \nabla v = 0, \quad \mathbf{r} \in \mathbb{R}^d, \quad (2.3)$$

$$v(0, \mathbf{r}) = v_0(\mathbf{r}) \in L_{2,per}(\mathbb{R}^d). \quad (2.4)$$

Скорость переноса $\boldsymbol{\omega}$ полагается постоянной во времени и пространстве.

Введём следующие обозначения. M^0 – конечное множество степеней свободы в одном сеточном блоке. $M = \mathbb{Z}^d \times M^0$ – общее множество степеней свободы. Если $f \in \mathbb{C}^M$, то $f_\eta \in \mathbb{C}^{M^0}$ – часть вектора f в блоке $\eta \in \mathbb{Z}^d$. V_{per}^N – множество последовательностей с периодом N :

$$V_{per}^N = \{f \in \mathbb{C}^M : \forall \eta, \zeta \in \mathbb{Z}^d \ f_{\eta+N\zeta} = f_\eta\}.$$

$V_{per} = \bigcup_{N \in \mathbb{N}} V_{per}^N$ – множество периодических последовательностей со скалярным произведением, определяемым для $f \in V_{per}^{N(f)}$, $g \in V_{per}^{N(g)}$ формулой

$$(f, g) = \frac{1}{N^d} \sum_{\eta=(0,\dots,N-1)^d} (f_\eta, g_\eta), \quad N = N(f)N(g).$$

Здесь (f_η, g_η) – некоторое скалярное произведение на \mathbb{C}^{M^0} . Если функция f имеет период $N(f)$, то она, очевидно, также имеет период $nN(f)$ для любого $n \in \mathbb{N}$, но на значение скалярного произведения замена $N(f)$ на $nN(f)$ не влияет. На \mathbb{C}^{M^0} и V_{per} будем использовать нормы, порождённые этими скалярными произведениями: $\|f_\eta\|^2 = (f_\eta, f_\eta)$; $\|f\|^2 = (f, f)$.

Для аппроксимации (2.3) рассмотрим полудискретные схемы вида

$$\sum_{\zeta \in \mathcal{S}} Z_\zeta \frac{du_{\eta+\zeta}}{dt}(t) + \frac{1}{h} \sum_{\zeta \in \mathcal{S}} L_\zeta u_{\eta+\zeta}(t) = 0, \quad \eta \in \mathbb{Z}^d, \quad u_\eta \in \mathbb{C}^{M^0}, \quad (2.5)$$

где $\mathcal{S} \subset \mathbb{Z}^d$ – конечное множество, которое мы будем называть шаблоном схемы, а Z_ζ и L_ζ – действительнoзначные матрицы. Для $\zeta \notin \mathcal{S}$ будем полагать $Z_\zeta = L_\zeta = 0$.

Введём операторы $Z : \mathbb{C}^M \rightarrow \mathbb{C}^M$ и $L : \mathbb{C}^M \rightarrow \mathbb{C}^M$ равенствами

$$(Zu)_\eta = \sum_{\zeta \in \mathcal{S}} Z_\zeta u_{\eta+\zeta}, \quad (Lu)_\eta = \sum_{\zeta \in \mathcal{S}} L_\zeta u_{\eta+\zeta}. \quad (2.6)$$

Очевидно, $ZV_{per}^N \subseteq V_{per}^N$, $LV_{per}^N \subseteq V_{per}^N$, откуда $ZV_{per} \subseteq V_{per}$, $LV_{per} \subseteq V_{per}$.

Будем предполагать, что схема (2.5) является устойчивой, то есть существует $K > 0$, такая что для всех $u \in C^\infty([0, \infty), V_{per})$, являющихся решением (2.5), при всех $t \geq 0$ выполняется $\|u(t)\| \leq K \|u(0)\|$.

Всюду далее будем считать, что $1/h \in \mathbb{N}$. Для отображения данных на пространство сеточных функций используются операторы $\Pi_h, \mathcal{P}_h : L_{2,loc}(\mathbb{R}^d) \rightarrow \mathbb{C}^M$, задаваемые формулами

$$(\Pi_h f)_{\eta,\xi} = \int_{\mathbf{r} \in G} \mu_\xi(\mathbf{r}) f(h(\mathbf{r} + T\boldsymbol{\eta})) d\mathbf{r}, \quad (\mathcal{P}_h f)_{\eta,\xi} = \int_{\mathbf{r} \in G} \hat{\mu}_\xi(\mathbf{r}) f(h(\mathbf{r} + T\boldsymbol{\eta})) d\mathbf{r}, \quad (2.7)$$

где $\boldsymbol{\eta} \in \mathbb{Z}^d$ – индекс блока, $\xi \in M^0$ – индекс переменной внутри блока, $G \subset \mathbb{R}^d$ – некоторая ограниченная область, $\mu_\xi, \hat{\mu}_\xi \in L_2(G)$, $\int_G \mu_\xi(\mathbf{r}) = 1$. Если $f \in L_{2,per}(\mathbb{R}^d)$ имеет период N_0 , то $\Pi_h f, \mathcal{P}_h f \in V_{per}^{(N_0/h)}$. Операторы Π_h и \mathcal{P}_h как операторы из $L_{2,per}(\mathbb{R}^d)$ в V_{per} являются ограниченными равномерно по h .

Также будем рассматривать оператор $\tilde{\Pi}_h^{(p,q)} : W_{2,loc}^q(\mathbb{R}^d) \rightarrow \mathbb{C}^M$, задаваемый формулой

$$(\tilde{\Pi}_h^{(p,q)} f)_\eta = (\Pi_h f)_\eta + \sum_{0 < |\mathbf{m}| \leq q} h^{|\mathbf{m}|} \mathfrak{C}^{(\mathbf{m})} (\mathcal{P}_h D^{\mathbf{m}} f)_\eta, \quad (2.8)$$

где $\mathfrak{C}^{(\mathbf{m})}$ – некоторые диагональные матрицы размера $|M^0| \times |M^0|$.

Будем также рассматривать операторы $\mathring{\Pi}_h, \mathring{\mathcal{P}}_h : C_{loc} \rightarrow \mathbb{C}^M$, задаваемые формулами

$$(\mathring{\Pi}_h f)_{\eta,\xi} = f(h(\boldsymbol{\rho}_\xi + T\boldsymbol{\eta})), \quad (\mathring{\mathcal{P}}_h f)_{\eta,\xi} = f(h(\hat{\boldsymbol{\rho}}_\xi + T\boldsymbol{\eta})), \quad (2.9)$$

и оператор $\mathring{\tilde{\Pi}}_h^{(p,q)} : C_{loc}^q \rightarrow \mathbb{C}^M$, определяемый как

$$\left(\mathring{\tilde{\Pi}}_h^{(p,q)} f \right)_\eta = \left(\mathring{\Pi}_h f \right)_\eta + \sum_{p \leq |\mathbf{m}| \leq q} h^{|\mathbf{m}|} \mathfrak{C}^{(\mathbf{m})} (\mathring{\mathcal{P}}_h (D^{\mathbf{m}} f))_\eta. \quad (2.10)$$

Оператор Π_h вида (2.7) является частным случаем (2.8) при $p = q = 0$, поэтому все утверждения, которые будут доказаны для $\tilde{\Pi}_h$, справедливы и для Π_h . Аналогично $\mathring{\Pi}_h$ вида (2.9) является частным случаем (2.10). Иногда для краткости верхние индексы p и q будем опускать.

Определение 1. Пусть $\tilde{\Pi}_h$ – некоторый оператор вида (2.8) или (2.10). Ошибкой решения в смысле $\tilde{\Pi}_h$ с начальными данными v_0 будем называть величину

$$\varepsilon_h(t, v_0, \tilde{\Pi}_h) = u(t) - \tilde{\Pi}_h v(t, \cdot), \quad (2.11)$$

где $u(t)$ – решение (2.5) с условием $u(0) = \tilde{\Pi}_h v_0$, а $v(t, \mathbf{r}) = v_0(\mathbf{r} - \boldsymbol{\omega}t)$.

Определение 2. Пусть Π – оператор из $H_{loc}^q(\mathbb{R}^d)$ или $C^q(\mathbb{R}^d)$ в \mathbb{C}^M . Ошибкой аппроксимации на функции f в смысле Π будем называть величину $\epsilon_h(f, \Pi) \in \mathbb{C}^M$, определяемую формулой

$$\epsilon_h(f, \Pi) = -Z\Pi(\omega \cdot \nabla f) + \frac{1}{h}L\Pi f. \quad (2.12)$$

Если $\tilde{\Pi}_h^{(q)}$ – оператор вида (2.8), то для функций $f \in H_{per}^{q+1}(\mathbb{R}^d)$, имеющих период N_0 , выполняется $\epsilon_h(f, \tilde{\Pi}_h^{(p,q)}) \in V_{per}^{(N_0/h)}$.

Определение порядка точности в длительном счёте, данное в [3], переносится на многомерный случай.

Определение 3. Предположим, что существуют такие константы C_1 и C_2 , что при всех начальных условиях $v_0 \in H_{per}^r(\mathbb{R}^d)$, $r \geq Q + 1$, и при всех t и h схема обладает оценкой ошибки

$$\|\epsilon_h(t, v_0, \tilde{\Pi}_h)\| \leq C_1 h^P \|\nabla^P v_0\| + C_2 (t + h) h^Q \|\nabla^{Q+1} v_0\|, \quad (2.13)$$

где $+\infty \geq Q \geq P > 0$. Тогда будем говорить, что схема обладает формальным порядком точности P и порядком точности в длительном счёте Q в смысле $\tilde{\Pi}_h$ на $H_{per}^r(\mathbb{R}^d)$. Если выполняется оценка $\|\epsilon_h(t, v_0, \tilde{\Pi}_h)\| \leq C_1 h^P \|\nabla^P v_0\|$ при $r \geq P > 0$, будем говорить, что $Q = \infty$, а если $\epsilon_h(t, v_0, \tilde{\Pi}_h) \equiv 0$, будем говорить, что $P = Q = \infty$. Если для всех $P' > P$, $Q' \geq P'$ и для всех $P' = P$, $Q' > Q$ не существует оценки вида (2.13), то такие параметры P и Q будем называть оптимальными.

Определение порядков на $C_{per}^r(\mathbb{R}^d)$ отличается от определения 3 заменой $H_{per}^r(\mathbb{R}^d)$ на $C_{per}^r(\mathbb{R}^d)$ и норм $\|\cdot\|$ в правой части (2.13) на $\|\cdot\|_\infty$.

3. Спектральное представление схемы

Для $\phi \in \mathbb{C}^d$ введём матрицы

$$Z(\phi) = \sum_{\eta \in S \subset \mathbb{Z}^d} Z_\eta \exp(i\phi \cdot \eta), \quad L(\phi) = \sum_{\eta \in S \subset \mathbb{Z}^d} L_\eta \exp(i\phi \cdot \eta), \quad (3.1)$$

$$A(\phi) = -Z^{-1}(\phi)L(\phi) + i(T^{-1}\omega) \cdot \phi I. \quad (3.2)$$

Здесь $\phi \cdot \eta = \phi_1 \eta_1 + \dots + \phi_d \eta_d$. Для $\phi \in \mathbb{R}^d$ детерминант $Z(\phi)$ отделён от нуля. Схема (2.5) устойчива с константой K тогда и только тогда, когда при всех $\phi \in \mathbb{R}^d$ и $\nu > 0$ выполняется

$$\|\exp(A(\phi)\nu)\| \leq K. \quad (3.3)$$

В частности, если $A(\phi) = -A^*(\phi)$, то схема устойчива с $K = 1$.

Обозначим через ϵ вектор размерности M^0 , состоящий из единиц. Будем отождествлять его с вектором $\epsilon \in V_{per}$, все блочные компоненты которого равны ϵ . Всюду далее будем предполагать, что схема точна на константе, то есть $L(0)\epsilon = 0$.

Функции аппроксимационной ошибки и ошибки решения в образах Фурье имеют вид

$$\hat{\epsilon}(\phi, \tilde{\Pi}_h) = A(\phi)(\tilde{\Pi}_1 e^{i(U\phi)\cdot r})_0, \quad (3.4)$$

$$\hat{\epsilon}(\phi, \nu, \tilde{\Pi}_h) = \left(e^{\nu A(\phi)} - I \right) (\tilde{\Pi}_1 e^{i(U\phi)\cdot r})_0. \quad (3.5)$$

Здесь и ниже индекс 1 означает подстановку $h = 1$, а индекс 0 – взятие блочной компоненты в блоке $\eta = 0$.

Для $\mathbf{k} \in \mathbb{Z}^d$ определим волновой вектор

$$\alpha(\mathbf{k}) = \frac{2\pi}{Nh} U \mathbf{k} = \frac{2\pi}{Nh} \sum_{j=1}^d k_j \mathbf{a}_j^*. \quad (3.6)$$

Для всех $\mathbf{k} \in \mathbb{Z}^d$ выполняется

$$\|\epsilon_h(t, e^{i\alpha(\mathbf{k})\cdot r}, \tilde{\Pi}_h)\| = \|\hat{\epsilon}(T^* \alpha(\mathbf{k})h, t/h, \tilde{\Pi}_h)\|. \quad (3.7)$$

Следующие три утверждения позволяют переносить оценки функций (3.4) и (3.5) на произвольные достаточно гладкие периодические функции.

Утверждение 3.1. Пусть $\tilde{\Pi}_h^{(p,q)}$ – оператор вида (2.8). Пусть $P_A \in \mathbb{N} \cup \{0\}$. Тогда следующие утверждения эквивалентны:

- для некоторого $C > 0$ в некоторой окрестности $\phi = 0$ справедливо $\|\hat{\epsilon}(\phi, \tilde{\Pi}_h^{(p,q)})\| \leq C|\phi|^{P_A+1}$;
- существуют такие $C_1, C_2 \geq 0$, что для всех $v_0 \in H_{per}^r(\mathbb{R}^d)$, $r = \max\{P_A, q\} + 1$, справедливо

$$\|\epsilon_h(v_0, \tilde{\Pi}_h^{(p,q)})\| \leq C_1 \|\nabla^{P_A+1} v_0\| h^{P_A} + C_2 \|\nabla^r v_0\| h^{r-1}; \quad (3.8)$$

- для всех мультииндексов \mathbf{m} , таких что $|\mathbf{m}| \leq P_A$, выполняется

$$\epsilon_1(\mathbf{r}^{\mathbf{m}} / \mathbf{m}!, \tilde{\Pi}_1^{(p,q)}) = 0. \quad (3.9)$$

Утверждение 3.2. Пусть $\tilde{\Pi}_h^{(p,q)}$ – некоторый оператор вида (2.8). Пусть $P, Q \in \mathbb{N}$, $Q \geq \max\{P, q\}$. Тогда следующие утверждения эквивалентны:

- для некоторых $C_1, C_2 \geq 0$ в некоторой окрестности $\phi = 0$ справедливо

$$\|\hat{\varepsilon}(\phi, \nu, \tilde{\Pi}_h^{(p,q)})\| \leq C_1|\phi|^P + C_2\nu|\phi|^{Q+1}; \quad (3.10)$$

- для некоторых $C_1, C_2 \geq 0$, при любых начальных данных $v_0 \in H_{per}^{Q+1}(\mathbb{R}^d)$ для ошибки решения справедлива оценка (2.13), то есть схема обладает порядком точности P и порядком точности в длительном счёте Q .

Утверждение 3.3. Пусть $\overset{\circ}{\Pi}_h^{(p,q)}$ – некоторый оператор вида (2.10). Пусть $P, Q \in \mathbb{N}$, $Q \geq \max\{P, q\}$. Тогда следующие утверждения эквивалентны:

- для некоторых $C_1, C_2 \geq 0$ в некоторой окрестности $\phi = 0$ справедливо

$$\|\hat{\varepsilon}(\phi, \nu, \overset{\circ}{\Pi}_h^{(p,q)})\| \leq C_1|\phi|^P + C_2\nu|\phi|^{Q+1}; \quad (3.11)$$

- схема обладает порядком точности P и порядком точности в длительном счёте Q на $C_{per}^{Q+1}(\mathbb{R}^d)$.

Пусть $\mathbb{R}^{n \times n}$ и $\mathbb{C}^{n \times n}$ – пространства действительно- и комплекснозначных матриц размера n . Через $\mathcal{A}(\cdot, \mathbb{R}^{n \times n})$ будем обозначать множество функций из \mathbb{C}^d в $\mathbb{C}^{n \times n}$, аналитических при $\phi = 0$, таких что при $i\phi \in \mathbb{R}^d$ выполняется $A(\phi) \in \mathbb{R}^{n \times n}$. Проводимый ниже анализ опирается на следующие две теоремы, доказательство которых основано на теории возмущений (см. [7]) и приведено в [2].

Теорема 3.4. Пусть $A \in \mathcal{A}(\cdot, \mathbb{R}^{n \times n})$. Пусть $A(0)$ имеет нулевое собственное значение кратности n_0 , $0 < n_0 < n$. Тогда в некоторой окрестности $\phi = 0$ справедливо представление $A(\phi) = S(\phi)M(\phi)S^{-1}(\phi)$, где $S, M, S^{-1} \in \mathcal{A}(\cdot, \mathbb{R}^{n \times n})$,

$$M(\phi) = \begin{pmatrix} M^{(*)}(\phi) & 0 \\ 0 & M^{(0)}(\phi) \end{pmatrix}, \quad (3.12)$$

причём $M^{(0)}(0)$ – нильпотентная матрица размера n_0 , а $M^{(*)}(0)$ – невырожденная матрица размера $n - n_0$.

Теорема 3.5. Пусть $d = 1$ и $A \in \mathcal{A}(\cdot, \mathbb{R}^{n \times n})$. Пусть существует $K > 0$, такое что для всех $\phi \in \mathbb{R}$ и всех $\nu > 0$ справедливо $\|\exp(\nu A(\phi))\| \leq K$. Тогда в окрестности $\phi = 0$ матрица $A(\phi)$ представима в блочно-диагональном виде

$$A(\phi) = S(\phi)M(\phi)S^{-1}(\phi), \quad (3.13)$$

$$M(\phi) = \begin{pmatrix} M_0(\phi) & 0 & \dots & 0 & 0 \\ 0 & \phi M_1(\phi) & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \phi^m M_m(\phi) & 0 \\ 0 & 0 & \dots & 0 & M_\infty(\phi) \equiv 0 \end{pmatrix}, \quad (3.14)$$

где $S, M, S^{-1} \in \mathcal{A}(\cdot, \mathbb{R}^{n \times n})$, блоки $M_k(\phi)$, $k \in \mathbb{N} \cup \{0, \infty\}$ квадратные, невырожденные при $\phi = 0$ (за исключением $k = \infty$), и некоторые из них могут отсутствовать (иметь нулевой размер).

4. Общие замечания

Утверждение 4.1. *Существование такого K , что $\|e^{\nu A}\| \leq K$ для всех $\nu > 0$, равносильно одновременному выполнению двух условий:*

- *все собственные числа λ матрицы A таковы, что $\operatorname{Re} \lambda \leq 0$;*
- *собственные числа λ матрицы A , такие что $\operatorname{Re} \lambda = 0$, полупростые.*

Действительно, представим матрицу A в виде $A = SJS^{-1}$, где J – её жорданова нормальная форма. Тогда $e^{\nu A} = Se^{\nu J}S^{-1}$. Доказываемое утверждение становится очевидным, если привести явный вид для $e^{\nu J}$.

Утверждение 4.2. *Пусть матрица $A(\phi)$ определена (3.2). Тогда собственное значение $\lambda = 0$ у матрицы $A(0)$ полупростое, то есть не может быть жордановых клеток размера больше 1, соответствующих $\lambda = 0$.*

Действительно, в силу точности на константе $A(0)\epsilon = 0$, поэтому $\lambda = 0$ – собственное значение $A(0)$. Из устойчивости схемы следует выполнение (3.3), в частности, для $\phi = 0$. В силу утверждения 4.1 собственные значения $\lambda = 0$ полупростое.

Утверждение 4.3. *Пусть \check{L} – ограничение L на V_{per} . Пусть $\operatorname{Ker} \check{L} = \operatorname{span}\{\epsilon\}$. Тогда собственное значение $\lambda = 0$ матрицы $A(0)$ простое.*

По определению (3.1)–(3.2) имеем $A(0) = -Z^{-1}(0)L(0)$. Рассмотрим произвольный вектор $w_0 \in \mathbb{C}^{M^0}$, такой что $w_0 \notin \operatorname{span}\{\epsilon\}$. Пусть $w \in V_{per}$ – вектор, все блочные компоненты которого одинаковы и равны w_0 . Тогда $\check{L}w \neq 0$. Но поскольку у $\check{L}w$ все блочные компоненты одинаковы и равны $L(0)w_0$, справедливо $L(0)w_0 \neq 0$. Отсюда $A(0)w_0 = -Z^{-1}(0)L(0)w_0 \neq 0$. Значит, $w_0 = \epsilon$ является единственным (с точностью до множителя) решением $A(0)w_0 = 0$. Отсюда с учётом утверждения 4.2 следует, что $\lambda = 0$ является простым собственным значением $A(0)$.

Обратное утверждение, вообще говоря, неверно: если собственное значение $A(0)$ простое, то это не гарантирует отсутствия неконстантных решений уравнения $\check{L}w = 0$. Действительно, рассмотрим одномерное уравнение переноса с $\omega = 1$ и схему с направленной разностью, выписанной независимо для двух переменных внутри блока. После этого в каждом втором блоке компоненты вектора внутри блока поменяем местами. Получаем схему с блочными компонентами

$$Z_0 = I, \quad L_0 = I, \quad L_{-1} = \begin{pmatrix} 0 & -1 \\ -1 & 0 \end{pmatrix}.$$

У матрицы $L(0)$ есть два простых собственных значения: $\lambda = 0$ и $\lambda = 2$. Однако вектор $w \in V_{per}$, имеющий период $N = 2$, с компонентами $w_\eta = (-1)^\eta(1, -1)^T$ лежит в ядре \check{L} , так же как и ϵ .

Утверждение 4.4. Пусть $A \neq 0$, $\|A\| \leq 1$ и

$$f(A) = \sum_{k=1}^{\infty} \frac{A^{k-1}}{k!}. \quad (4.1)$$

Тогда $\|(f(A))^{-1}\| \leq 4$. Если дополнительно A не вырождена, то

$$\|(e^A - I)^{-1}\| \leq 4\|A^{-1}\|. \quad (4.2)$$

Доказательство см. в [3].

По определению $\hat{\epsilon}(\phi, \nu, \tilde{\Pi}_h)$ имеет вид (3.5). В силу теоремы 3.4 в некоторой окрестности $\phi = 0$ справедливо представление $A(\phi) = S(\phi)M(\phi)S^{-1}(\phi)$, где $M(\phi)$ имеет вид (3.12), откуда

$$\hat{\epsilon}(\phi, \nu, \tilde{\Pi}_h) = S(\phi) \begin{bmatrix} e^{M^{(*)}(\phi)\nu} - I & 0 \\ 0 & e^{M^{(0)}(\phi)\nu} - I \end{bmatrix} S^{-1}(\phi) (\tilde{\Pi}_1 e^{i(U\phi)\cdot r})_0, \quad (4.3)$$

причём $M^{(*)}(\phi)$ и $M^{(0)}(\phi)$ – аналитические функции ϕ , $M^{(*)}(0)$ невырожденная, а $M^{(0)}(0)$ – нильпотентная.

Введём векторы $v^{(*)}(\phi, \tilde{\Pi}_h)$ и $v^{(0)}(\phi, \tilde{\Pi}_h)$, размерности которых соответствуют блокам в представлении (4.3), так чтобы

$$\begin{pmatrix} v^{(*)}(\phi, \tilde{\Pi}_h) \\ v^{(0)}(\phi, \tilde{\Pi}_h) \end{pmatrix} = S^{-1}(\phi) (\tilde{\Pi}_1 e^{i(U\phi)\cdot r})_0. \quad (4.4)$$

Обозначим

$$\begin{aligned}\hat{\varepsilon}^{(1)}(\phi, \nu, \tilde{\Pi}_h) &= S(\phi) \begin{pmatrix} [\exp(M^{(*)}(\phi)\nu) - I] v^{(*)}(\phi, \tilde{\Pi}_h) \\ 0 \end{pmatrix}, \\ \hat{\varepsilon}^{(2)}(\phi, \nu, \tilde{\Pi}_h) &= S(\phi) \begin{pmatrix} 0 \\ [\exp(M^{(0)}(\phi)\nu) - I] v^{(0)}(\phi, \tilde{\Pi}_h) \end{pmatrix}.\end{aligned}\quad (4.5)$$

Очевидно, выполняется

$$\hat{\varepsilon}(\phi, \nu, \tilde{\Pi}_h) = \hat{\varepsilon}^{(1)}(\phi, \nu, \tilde{\Pi}_h) + \hat{\varepsilon}^{(2)}(\phi, \nu, \tilde{\Pi}_h). \quad (4.6)$$

Спектр матрицы $M^{(0)}(0)$ состоит из одной точки $\lambda = 0$, откуда в силу утверждения 4.2 следует, что $M^{(0)}(0) = 0$, и $M^{(0)}(\phi) = \sum_j M_j^{(0)}(\phi)\phi_j$. Поэтому

$$M^{(0)}(T^*\alpha h) \frac{t}{h} = \sum_{j=1}^d M_j^{(0)}(T^*\alpha h) (T^*\alpha)_j t$$

является аналитической функцией t , h и α . Следовательно, таковой является и $\hat{\varepsilon}^{(2)}(T^*\alpha h, t/h, \tilde{\Pi}_h)$.

В частности, если $L(0) = 0$, то $A(0) = 0$ и $\hat{\varepsilon}^{(1)}(\phi, \nu, \tilde{\Pi}_h) = 0$. Таким образом, имеем следующий результат.

Утверждение 4.5. *Если $L(0) = 0$, то $\hat{\varepsilon}(\alpha h, t/h, \tilde{\Pi}_h)$ является аналитической функцией t , h и α .*

Утверждение 4.6. *Пусть схема (2.5) обладает оценкой (2.13) в смысле Π_h . Тогда в некоторой окрестности $\phi = 0$ при всех $\nu > 0$ выполняется*

$$\|v^{(*)}(\phi, \Pi_h)\| \leq C|\phi|^P, \quad \|\hat{\varepsilon}^{(1)}(\phi, \nu, \Pi_h)\| \leq (K+1)C|\phi|^P. \quad (4.7)$$

По утверждению 3.2 из оценки (2.13) следует (3.10). В силу аналитичности матрицы $S^{-1}(\phi)$ в некоторой окрестности $\phi = 0$, отсюда получаем

$$\|\hat{\varepsilon}^{(m)}(\phi, \nu, \Pi_h)\| \leq C_1|\phi|^P + C_2\nu|\phi|^{Q+1},$$

для $m = 1, 2$. То есть для

$$E = \left[\exp\left(M^{(*)}(\phi)\nu\right) - I \right] v^{(*)}(\phi)$$

имеем

$$\|E\| \leq C'_1|\phi|^P + C'_2\nu|\phi|^{Q+1}.$$

Матрица $M^{(*)}(0)$ невырождена, поэтому для достаточно малых ϕ можно считать, что $\|(M^{(*)}(\phi))^{-1}\| \leq 2\|(M^{(*)}(0))^{-1}\|$ и $\|M^{(*)}(\phi)\| \leq 2\|M^{(*)}(0)\|$. Положим $\nu = 1/(2\|M^{(*)}(0)\|)$. Тогда в силу (4.2) имеем

$$\begin{aligned} \|v^{(*)}(\phi)\| &\leq \left\| \left[\exp \left(M^{(*)}(\phi)\nu \right) - I \right]^{-1} \right\| \|E\| \leq \\ &\leq 8\|(M^{(*)}(0))^{-1}\| \left(C_1 |\phi|^P + C_2 \frac{1}{2\|M^{(*)}(0)\|} |\phi|^{Q+1} \right). \end{aligned}$$

Отсюда получаем первое неравенство в (4.7). Второе неравенство следует из него в силу устойчивости.

Утверждение 4.7. Пусть схема (2.5) обладает оценкой (2.13) в смысле Π_h . Тогда в некоторой окрестности $\phi = 0$ выполняется

$$\|M^{(0)}(\phi)v^{(0)}(\phi, \Pi_h)\| \leq C|\phi|^{P+1}. \quad (4.8)$$

Аналогично предыдущему утверждению, для

$$E = \left[\exp \left(M^{(0)}(\phi)\nu \right) - I \right] v^{(0)}(\phi, \Pi_h)$$

имеем

$$\|E\| \leq C'_1 |\phi|^P + C'_2 \nu |\phi|^{Q+1}.$$

Рассмотрим некоторое ϕ . Если $M^{(0)}(\phi) = 0$, то искомое утверждение очевидно. Далее будем считать, что $M^{(0)}(\phi) \neq 0$. Преобразуем E следующим образом:

$$E = f \left(M^{(0)}(\phi)\nu \right) M^{(0)}(\phi)\nu v^{(0)}(\phi, \Pi_h),$$

где f определено (4.1). Положим $\nu = 1/\|M^{(0)}(\phi)\|$. Тогда в силу утверждения 4.4 имеем

$$\left\| M^{(0)}(\phi)\nu v^{(0)}(\phi, \Pi_h) \right\| = \left\| \left[f \left(M^{(0)}(\phi)\nu \right) \right]^{-1} E \right\| \leq 4\|E\|$$

и, следовательно,

$$\left\| M^{(0)}(\phi)v^{(0)}(\phi, \Pi_h) \right\| \leq C'_1 \frac{1}{\nu} |\phi|^P + C'_2 |\phi|^{Q+1}.$$

Поскольку $1/\nu = \|M^{(0)}(\phi)\|$, а $M^{(0)}(0) = 0$, имеем $1/\nu \leq \tilde{C}|\phi|$. Таким образом, получаем (4.8).

Утверждение 4.8. Пусть схема (2.5) обладает оценкой (2.13) в смысле Π_h . Тогда для любого $q \in \mathbb{N} \cup \{0\}$ существует оператор $\tilde{\Pi}_h^{(P,q)}$ вида (2.8), такой что величины $v^{(*)}(\phi, \tilde{\Pi}_h^{(P,q)})$ и $v^{(0)}(\phi, \tilde{\Pi}_h^{(P,q)})$, определённые (4.4), удовлетворяют оценкам

$$\begin{aligned} \|v^{(*)}(\phi, \tilde{\Pi}_h^{(P,q)})\| &\leq C|\phi|^{q+1}, \\ \|v^{(0)}(\phi, \tilde{\Pi}_h^{(P,q)}) - v^{(0)}(\phi, \Pi_h)\| &\leq C|\phi|^{q+1}. \end{aligned} \quad (4.9)$$

При $q < P$ выполняется $\tilde{\Pi}_h^{(P,q)} \equiv \Pi_h$, и искомая оценка следует из утверждения 4.6. Пусть $q \geq P$. Введём в окрестности $\phi = 0$ аналитическую функцию

$$V(\phi) = S(\phi) \begin{pmatrix} 0 & 0 \\ 0 & I \end{pmatrix} S^{-1}(\phi) (\Pi_1 e^{i(U\phi)\cdot r})_0. \quad (4.10)$$

По построению

$$\|V(\phi) - (\Pi_1 e^{i(U\phi)\cdot r})_0\| \leq C\|v^{(*)}(\phi)\| \leq \tilde{C}|\phi|^P.$$

Определим диагональную матрицу $\mathfrak{C}(\phi)$ равенством

$$V(\phi) - (\Pi_1 e^{i(U\phi)\cdot r})_0 = \mathfrak{C}(\phi) (\mathcal{P}_1 e^{i(U\phi)\cdot r})_0 \quad (4.11)$$

и введём диагональные матрицы

$$\mathfrak{C}^{(m)} = \frac{1}{m!} \frac{1}{i^{|m|}} \left. \frac{\partial^{|m|} \mathfrak{C}(T^*\psi)}{\partial \psi_1^{m_1} \dots \partial \psi_d^{m_d}} \right|_{\psi=0}. \quad (4.12)$$

Покажем, что $\mathfrak{C}^{(m)} = 0$ при $|m| < P$. Поскольку матрица $\mathfrak{C}(\phi)$ диагональная, а $(\mathcal{P}_1 e^{i(U\phi)\cdot r})_0$ стремится к ϵ при $\phi \rightarrow 0$, порядок малости $\mathfrak{C}(\phi)$ по ϕ совпадает с порядком малости правой части (4.11). Ввиду невырожденности $S(0)$ он совпадает с порядком малости вектора $v^{(*)}(\phi)$, который имеет величину не больше $O(|\phi|^P)$. Следовательно, $\mathfrak{C}(\phi) = O(|\phi|^P)$, а это и означает, что её первые $P - 1$ производные нулевые.

Введём оператор $\tilde{\Pi}_h^{(P,q)}$ равенством (2.8). Вычисляя $(\tilde{\Pi}_1^{(P,q)} e^{i(U\phi)\cdot r})_0$ напрямую из определения (2.8), получаем

$$(\tilde{\Pi}_1^{(P,q)} e^{i(U\phi)\cdot r})_0 = (\Pi_1 e^{i(U\phi)\cdot r})_0 + \left[\sum_{m: P \leq |m| \leq q} \mathfrak{C}^{(m)} (iU\phi)^m \right] (\mathcal{P}_1 e^{i(U\phi)\cdot r})_0.$$

Так как $\mathfrak{C}^{(m)} = 0$ при $|m| < P$, выражение в квадратных скобках совпадает с тейлоровским разложением функции $\mathfrak{C}(T^*\psi)$ по степеням ψ до порядка q (см. (4.12)) при $\psi = U\phi$. Так как $T^*U\phi = \phi$, получаем

$$(\tilde{\Pi}_1^{(P,q)} e^{i(U\phi)\cdot r})_0 = (\Pi_1 e^{i(U\phi)\cdot r})_0 + \mathfrak{C}(\phi) (\mathcal{P}_1 e^{i(U\phi)\cdot r})_0 + O(|\phi|^{q+1}).$$

В силу (4.11) отсюда следует

$$(\tilde{\Pi}_1^{(P,q)} e^{i(U\phi)\cdot r})_0 = V(\phi) + O(|\phi|^{q+1}).$$

Домножая это равенство слева на $S^{-1}(\phi)$, получаем искомые оценки.

5. Определение формального порядка

В [3] был предложен алгоритм установления оптимальных значений формального порядка точности P и порядка точности в длительном счёте Q в одномерном случае. Покажем, что его часть, связанная с определением P , переносится на многомерный случай.

Алгоритм 1 (определения формального порядка точности).

1. Определяем $P_A = \max\{m \in \mathbb{N} \cup \{0\} : (\epsilon_1(\mathbf{r}^m, \Pi_1))_0 = 0 \ \forall m : |\mathbf{m}| \leq m\}$.
2. Вычисляем $f^m = -(\epsilon_1(\mathbf{r}^m/m!, \Pi_1))_0$ для всех $\mathbf{m} : |\mathbf{m}| = P_A + 1$.
3. Если для всех $\mathbf{m} : |\mathbf{m}| = P_A + 1$ выполняется $f^m \in \text{Im}L(0)$, то полагаем $P' = P_A + 1$, иначе полагаем $P' = P_A$.

Утверждение 5.1. Значение P' , найденное алгоритмом 1, совпадает с оптимальным значением порядка точности P .

Если $P' = P_A + 1$, то системы $L(0)\mathfrak{C}^{(m)}\mathfrak{e} = f^{(m)}$, $|\mathbf{m}| = P_A + 1$, совместны. Определим отображение $\tilde{\Pi}_h^{(P_A+1, P_A+1)}$, в которое подставим решения этих систем. В силу устойчивости и неравенства треугольника имеем оценку (2.13) при $P = Q = P_A + 1$ (доказательство повторяет одномерный случай).

Если $P = P_A + 1$, то по утверждению 4.8 найдутся такие коэффициенты $\mathfrak{C}^{(m)}$, $|\mathbf{m}| = P$, что для отображения $\tilde{\Pi}_h^{(P,P)}$, определённого (2.8), будут выполняться неравенства (4.9) для $q = P$. Первое из них даёт

$$\|M^{(*)}(\phi)v^{(*)}(\phi, \tilde{\Pi}_h^{(P,P)})\| \leq C|\phi|^{P+1}.$$

Второе с учётом утверждения 4.7 даёт

$$\|M^{(0)}(\phi)v^{(0)}(\phi, \tilde{\Pi}_h^{(P,P)})\| \leq C|\phi|^{P+1}.$$

Таким образом, имеем

$$\|A(\phi)(\tilde{\Pi}_h^{(P,P)} e^{i(U\phi)\cdot r})_0\| \leq C|\phi|^{P+1},$$

и по утверждению 3.2, в смысле $\tilde{\Pi}_h^{(P,P)}$ схема обладает порядком аппроксимации P . Следовательно, коэффициенты $\mathfrak{C}^{(m)}$, $|\mathbf{m}| = P$, этого оператора удовлетворяют системе

$$L(0)\mathfrak{C}^{(m)}\mathfrak{e} = -(\epsilon_1(\mathbf{r}^m/m!, \tilde{\Pi}_1))_0,$$

поэтому правая часть этого равенства для всех \mathbf{m} , $|\mathbf{m}| = P$, лежит в образе $L(0)$, и алгоритм выдаст значение $P' = P_A + 1$.

Следствие 5.2. *Если $L(0) = 0$, то порядок точности схемы не может превосходить её порядок аппроксимации.*

Действительно, если схема обладает порядком аппроксимации P_A , то для какого-то m : $|m| = P_A$ выполняется $f^m \neq 0$. Поскольку $L(0) = 0$, то $f^m \notin \text{Im}L(0)$, и алгоритм выдаст значение $P' = P_A$. По утверждению получаем $P = P_A$.

6. Существование отображения

Перейдём к рассмотрению вопроса, переносятся ли на многомерный случай результаты, полученные в [3], применительно к порядку точности в длительном счёте.

Утверждение 6.1. *Пусть $A \in \mathbb{C}^{n \times n}$ и $v \in \mathbb{C}^n$. Пусть при всех $\nu \geq 0$ справедливы неравенства $\|e^{\nu A}\| \leq K$ и*

$$\|(e^{\nu A} - I)v\| \leq (\tilde{C}_1 + \tilde{C}_2\nu)\|v\|. \quad (6.1)$$

Тогда для $w = (A^*A + \varepsilon^2)^{-1}\varepsilon^2v$, где $\varepsilon = \tilde{C}_2/\tilde{C}_1$ и при $\tilde{C}_1 = 0$ или $\tilde{C}_2 = 0$ под w понимается соответствующее предельное значение, справедливы оценки

$$\|v - w\| \leq \delta\tilde{C}_1\|v\|, \quad \|Aw\| \leq \delta\tilde{C}_2\|v\|, \quad (6.2)$$

где δ зависит только от n и K .

Это утверждение доказано в [4].

Теорема 6.2. *Пусть Π_h определён (2.7). Пусть схема (2.5) в смысле Π_h на $H_{per}^{Q+1}(\mathbb{R}^d)$ обладает формальным порядком точности P и порядком точности в длительном счёте Q . Тогда существует оператор $\tilde{\Pi}_h : L_{2,per}(\mathbb{R}^d) \rightarrow V_{per}$, такой что для некоторого C выполняются условия*

$$\|\tilde{\Pi}_h f - \Pi_h f\| \leq Ch^P\|f\|, \quad (6.3)$$

$$\|\epsilon_h(f, \tilde{\Pi}_h)\| \leq Ch^Q\|f\| \quad (6.4)$$

для всех h и $f \in H_{per}^{Q+1}(\mathbb{R}^d)$.

Обозначим $v(\phi) = (\Pi_1 e^{i(U\phi) \cdot r})_0$. Поскольку норма экспоненты равна единице, имеем $\|v(\phi)\| \leq \|\Pi_1\|$. В силу утверждения (3.10) имеем

$$\left\| \left(e^{\nu A(\phi)} - I \right) v(\phi) \right\| \leq C_1 |\phi|^P + C_2 \nu |\phi|^{Q+1}.$$

Это неравенство совпадает с (6.1) при подстановке $A = A(\phi)$, $v = v(\phi)$, $\tilde{C}_1 = C_1|\phi|^P$, $\tilde{C}_2 = C_2|\phi|^{Q+1}$. В силу утверждения 6.1 существует $w(\phi)$, такой что выполняется

$$\|v(\phi) - w(\phi)\| \leq \delta C_1 |\phi|^P \|\Pi_1\|, \quad \|A(\phi)w(\phi)\| \leq \delta C_2 |\phi|^{Q+1} \|\Pi_1\|.$$

Введём $\omega'(\phi) = \omega(\phi)$ при $|\phi| \leq 1$ и $\omega'(\phi) = 0$ при $\|\phi\| > 1$. Тогда получаем

$$\|v(\phi) - \omega'(\phi)\| \leq \max\{1, \delta C_1\} |\phi|^P \|\Pi_1\|, \quad \|A(\phi)\omega'(\phi)\| \leq \delta C_2 |\phi|^{Q+1} \|\Pi_1\|.$$

Определим вначале действие $\tilde{\Pi}_h$ на функцию вида $f(\mathbf{r}) = \exp(i\alpha(\mathbf{k}) \cdot \mathbf{r})$, где $\mathbf{k} \in \mathbb{Z}^d$, а $\alpha(\mathbf{k})$ определена (3.6). Положим

$$\left(\tilde{\Pi}_h f\right)_\eta = \exp(i\phi \cdot \boldsymbol{\eta}/h) \omega'(\phi). \quad (6.5)$$

Тогда будет обеспечиваться

$$\|\tilde{\Pi}_h f - \Pi_h f\| = \|(\tilde{\Pi}_h f)_0 - (\Pi_h f)_0\| = \|\omega'(\phi) - v(\phi)\| \leq \max\{1, \delta C_1\} |\phi|^P \|\Pi_1\|, \quad (6.6)$$

$$\begin{aligned} \|\epsilon_h(f, \tilde{\Pi}_h)\| &= \|(\epsilon_h(f, \tilde{\Pi}_h))_0\| = \frac{1}{h} \left\| \sum_{\boldsymbol{\eta} \in \mathcal{S}} (-ih\boldsymbol{\omega} \cdot \boldsymbol{\alpha}(\mathbf{k}) Z_\eta + L_\eta) (\tilde{\Pi}_h f)_\eta \right\| = \\ &= \frac{1}{h} \|Z(\phi) A \omega'(\phi)\| \leq \frac{1}{h} \|Z(\phi)\| \delta C_2 |T^* \boldsymbol{\alpha}(\mathbf{k}) h|^{Q+1} \leq \tilde{C} h^Q |\boldsymbol{\alpha}(\mathbf{k})|^{Q+1}. \end{aligned} \quad (6.7)$$

Далее, если $f \in L_{2,per}(\mathbb{R}^d)$, то она может быть представлена своим сходящимся в L_2 рядом Фурье:

$$f = \sum_{\mathbf{k} \in \mathbb{Z}^d} U_{\mathbf{k}} \exp(i\alpha(\mathbf{k}) \cdot \mathbf{r}),$$

причём $\|f\|^2 = \sum |U_{\mathbf{k}}|^2 |\mathbf{k}|^2$. Действие $\tilde{\Pi}_h$ на каждую из гармоник определено (6.5), а действие $\tilde{\Pi}_h$ на f определим исходя из линейности $\tilde{\Pi}_h$. Покажем, что $\tilde{\Pi}_h$ является ограниченным. Действительно, в силу ортогональности образов $\exp(i\alpha(\mathbf{k}) \cdot \mathbf{r})$ для разных $\mathbf{k} \in I_N^d$ под действием $\tilde{\Pi}_h$ имеем

$$\begin{aligned} \|\tilde{\Pi}_h f\|^2 &= \sum_{\mathbf{k} \in \mathbb{Z}^d} |U_{\mathbf{k}}|^2 |w'(2\pi\mathbf{k}/N)|^2 \leq \|f\|^2 \sup_{\phi \in \mathbb{R}^d} \|w(\phi)\|^2 \leq \\ &\leq \|f\|^2 \sup_{\phi \in \mathbb{R}^d} (\|v(\phi)\| + \|w'(\phi) - v(\phi)\|)^2 \leq \|f\|^2 (2 + \delta C_1)^2 \|\Pi_h\|^2. \end{aligned}$$

По теореме 15 в [1] из ограниченности $\tilde{\Pi}_h$ из неравенств (6.6) и (6.7) следует (6.3) и (6.4) для произвольной функции $f \in H_{per}^{Q+1}(\mathbb{R}^d)$.

В одномерном случае $w(\phi)$, даваемое утверждением 6.1, является аналитической функцией ϕ . Эту функцию можно приблизить суммой первых $Q + 1$ членов ряда Тейлора. Зная коэффициенты при ϕ^k , $k = P, \dots, Q$, легко построить такое отображение $\tilde{\Pi}_h^{(P,Q)}$ вида (2.10), в смысле которого будет иметь место Q -й порядок аппроксимации (подробно построение этого отображения описано в [1]). В многомерном случае функция $w(\phi)$, как правило, не является аналитической, и оператор $\tilde{\Pi}_h^{(P,Q)}$ вида (2.10), доставляющий Q -й порядок аппроксимации, вообще говоря, не существует. Соответствующий пример будет приведён ниже, а пока рассмотрим два частных случая, когда оператор $\tilde{\Pi}_h^{(P,Q)}$ искомого вида существует и, следовательно, оптимальное значение порядка точности в длительном счёте удаётся установить алгоритмически.

7. Простой случай

Начнём со случая, когда собственное значение $\lambda = 0$ матрицы $A(0)$ является простым. Будем называть такой случай *простым*.

Утверждение 7.1. Пусть схема (2.5) обладает оценкой (2.13) в смысле Π_h . Пусть размер матрицы $M^{(0)}(\phi)$ равен 1. Тогда в некоторой окрестности $\phi = 0$ выполняется

$$\|M^{(0)}(\phi)\| \leq C|\phi|^{Q+1}. \quad (7.1)$$

Обозначим единственный элемент матрицы $M^{(0)}(\phi)$ через $\lambda(\phi)$. Он является аналитической функцией ϕ . Пусть j_0 – минимальный порядок малости по ϕ членов в разложении $\lambda(\phi)$ по степеням ϕ_1, \dots, ϕ_d . Если $\lambda(\phi) \equiv 0$, будем считать $j_0 = \infty$.

В некоторой окрестности $\phi = 0$ имеем $\|v^{(*)}(\phi)\| \leq C|\phi|^P$. Следовательно, $v^{(0)}(0) \neq 0$, иначе было бы $v(0) = 0$, откуда из (4.4) следовало бы $(\Pi_1 1)_0 = 0$. Напомним, что $v^{(0)}(\phi)$ – это однокомпонентная величина. Отсюда для достаточно малых $|\phi|$ можно записать

$$\left| e^{\nu\lambda(\phi)} - 1 \right| \leq \frac{2}{|v^{(0)}(0)|} (C_1|\phi|^P + C_2|\phi|^{Q+1}(\nu + 1)).$$

Положим $\nu = |\phi|^{-Q-1/2}$, тогда правая часть будет стремиться к нулю при $\phi \rightarrow 0$. Следовательно, $\lambda(\phi)|\phi|^{-Q-1/2} \rightarrow 0$. Но поскольку $\lambda(\phi)$ – аналитическая функция ϕ при $\phi = 0$, имеем $j_0 \geq Q + 1$, то есть $\lambda(\phi) \leq C|\phi|^{Q+1}$.

Утверждение 7.2. Пусть Π_h и \mathcal{P}_h определены (2.7). Пусть схема (2.5) в смысле Π_h на $H_{per}^{Q+1}(\mathbb{R}^d)$ обладает формальным порядком точности P и порядком

точности в длительном счёте Q . Пусть собственное значение $\lambda = 0$ матрицы $A(0)$, определённой (3.2), является простым. Тогда существуют такие диагональные матрицы $\mathfrak{C}^{(m)} \in \mathbb{R}^{M^0}$, что схема обладает порядком аппроксимации Q в смысле модифицированного отображения $\tilde{\Pi}_h^{(P,Q)} : H_{per}^{Q+1}(\mathbb{R}^d) \rightarrow V_{per}$, определяемого (2.8).

Будем пользоваться введёнными выше обозначениями. Поскольку по условию $\lambda = 0$ является простым, то размер блока $M^{(0)}(\phi)$ равен 1. Обозначим единственный элемент этой матрицы через $\lambda(\phi)$. По утверждению 4.8 существуют такие коэффициенты $\mathfrak{C}^{(m)}$, $P \leq |m| \leq Q$, что выполняется неравенство

$$\|v^{(*)}(\phi, \tilde{\Pi}_h^{(P,Q)})\| \leq C|\phi|^{Q+1}.$$

Отсюда получаем

$$\|M^{(*)}(\phi)v^{(*)}(\phi, \tilde{\Pi}_h^{(P,Q)})\| \leq C|\phi|^{Q+1}.$$

В силу утверждения 7.1 имеем $\|M^{(0)}(\phi)\| \leq C|\phi|^{Q+1}$, откуда

$$\|M^{(0)}(\phi)v^{(0)}(\phi, \tilde{\Pi}_h^{(P,Q)})\| \leq C|\phi|^{Q+1}.$$

Таким образом,

$$\|A(\phi)(\tilde{\Pi}_h^{(P,Q)} e^{i(U\phi)\cdot r})_0\| \leq C|\phi|^{Q+1},$$

и по утверждению 3.2 в смысле $\tilde{\Pi}_h^{(P,Q)}$ схема обладает порядком аппроксимации Q .

Утверждение 7.3. Пусть $\mathring{\Pi}_h$ и $\mathring{\mathcal{P}}_h$ – операторы вида (2.9). Пусть схема (2.5) в смысле $\mathring{\Pi}_h$ на $C_{per}^{Q+1}(\mathbb{R}^d)$ обладает формальным порядком точности P и порядком точности в длительном счёте Q . Пусть собственное значение $\lambda = 0$ матрицы $A(0)$, определённой (3.2), является простым. Тогда существуют такие диагональные матрицы $\mathfrak{C}^{(m)} \in \mathbb{R}^{M^0}$, что схема обладает порядком аппроксимации Q в смысле модифицированного отображения $\mathring{\Pi}_h : C_{per}^{Q+1}(\mathbb{R}^d) \rightarrow V_{per}$, определяемого (2.10).

Это утверждение следует из предыдущего в силу возможности приближения операторов $\mathring{\Pi}_h$ и $\mathring{\mathcal{P}}_h$ операторами вида (2.7) с любым наперёд заданным порядком точности. Строгое доказательство повторяет доказательство для одномерного случая, проведённое в [3].

Утверждения 7.2 и 7.3 позволяют обобщить алгоритм 2 определения оптимальных значений P и Q , приведённый в [3], на многомерный случай. Уравнения на нахождение коэффициентов $\mathfrak{C}^{(m)}$ записываются в виде

$$L(0)\mathfrak{C}^{(m)}\mathfrak{e} = f^m, \quad f^m = - \left(\epsilon_1 \left(\frac{\mathbf{r}^m}{\mathbf{m}!}, \tilde{\Pi}_1^{(p,|\mathbf{m}|-1)} \right) \right)_0, \quad (7.2)$$

где $\tilde{\Pi}_h^{(p,m-1)}$ – оператор, определённый равенством

$$\left(\tilde{\Pi}_h^{(p,m-1)} f \right)_\eta = (\Pi_h f)_\eta + \sum_{p \leq |\mathbf{n}| \leq m-1} h^{|\mathbf{n}|} \mathfrak{C}^{(\mathbf{n})} (\mathcal{P}_h(D^\mathbf{n} f))_\eta. \quad (7.3)$$

Алгоритм 2 (определения формального порядка точности и порядка точности в длительном счёте для простого случая).

1. Определяем $P_A = \max\{m \in \mathbb{N} \cup \{0\} : (\epsilon_1(\mathbf{r}^m, \Pi_1))_0 = 0 \ \forall m : |\mathbf{m}| \leq m\}$.
2. Полагаем $t = P_A + 1$.
3. Вычисляем $f^m = -(\epsilon_1(\mathbf{r}^m/\mathbf{m}!, \tilde{\Pi}_1^{(P_A+1,m-1)}))_0$ для всех $\mathbf{m} : |\mathbf{m}| = m$, подставляя в $\tilde{\Pi}_h^{(P_A+1,m-1)}$ ранее найденные коэффициенты $\mathfrak{C}^{(\mathbf{n})}$, $P_A + 1 \leq |\mathbf{n}| \leq m - 1$.
4. Если для всех $\mathbf{m} : |\mathbf{m}| = t$ выполняется $f^m \in \text{Im}L(0)$, то:
 - находим диагональные матрицы $\mathfrak{C}^{(\mathbf{m})}$ из систем $L(0)\mathfrak{C}^{(\mathbf{m})}\mathfrak{e} = f^m$ (с точностью до $c_m I$, $c_m \in \mathbb{R}$);
 - увеличиваем значение t на единицу;
 - возвращаемся к п. 3.
5. Полагаем $Q' = t - 1$.
6. Если $Q' = P_A$, полагаем $P' = P_A$, иначе полагаем $P' = P_A + 1$.

Доказательство корректности этого алгоритма полностью повторяет доказательство корректности алгоритма 2 в [3].

8. Квазиодномерный случай

Будем называть схему квазиодномерной, если матрица $A(\phi)$ зависит только от проекции ϕ на некоторое выделенное направление. Примером такой схемы является метод Галёркина с разрывными базисными функциями на двумерной регулярной сетке (то есть либо состоящей из одинаковых параллелепипедов, либо полученной из неё разбиением этих параллелепипедов диагоналями однородным образом) в случае, когда скорость переноса направлена вдоль одной из сеточных линий.

Утверждение 8.1. Пусть Π_h и \mathcal{P}_h – отображения вида (2.7). Пусть схема (2.5) в смысле Π_h обладает формальным порядком точности P и порядком

точности в длительном счёте Q . Пусть шаблон схемы \mathcal{S} лежит в некотором одномерном подмножестве \mathbb{Z}^d , то есть существует такой $\eta \in \mathbb{Z}^d$, что $\mathcal{S} \subset \{m\eta, m = -M, \dots, M\}$. Тогда существуют такие диагональные матрицы $\mathfrak{C}^{(m)} \in \mathbb{R}^{M^0}$, что схема обладает порядком аппроксимации Q в смысле модифицированного отображения $\tilde{\Pi}_h : H_{per}^{Q+1}(\mathbb{R}^d) \rightarrow V_{per}$, определяемого (2.8).

Доказательство этого утверждения в основном повторяет одномерный случай (см. [3]), поэтому приведём только его схему. Без ограничения общности будем считать, что выделенное направление совпадает с осью x .

Зависимость матрицы только от ϕ_x позволяет получить для неё представление (3.13)–(3.14). По аналогии с (4.4) введём вектор $v(\phi)$ следующим образом:

$$v(\phi) = S^{-1}(\phi_x) \left(\Pi_1 e^{i(U\phi) \cdot r} \right)_0, \quad (8.1)$$

и через $v_j(\phi)$ будем обозначать его компоненты, размерность которых соответствует размерам блоков матриц $M_j(\phi_x)$ в представлении (3.14). Введём целые числа p_j условиями

$$c_j \psi^{p_j} \leq \sup_{\phi: |\phi|=\psi} \|v_j(\phi)\| \leq 2c_j \psi^{p_j}, \quad (8.2)$$

где $c_j \neq 0$, и $p_j = \infty$, если $v_j(\phi) \equiv 0$. Повторяя доказательство для одномерного случая, получаем $p_j \geq \min\{P, Q + 1 - j\}$. Определим в окрестности $\phi = 0$ аналитическую функцию

$$V(\phi) = S(\phi_x) \begin{pmatrix} \delta_0 I & 0 & \dots & 0 & 0 \\ 0 & \delta_1 I & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & \delta_m I & 0 \\ 0 & 0 & \dots & 0 & \delta_\infty I \end{pmatrix} S^{-1}(\phi_x) (\Pi_1 e^{i(U\phi) \cdot r})_0, \quad (8.3)$$

где размеры блоков соответствуют представлению (3.14) матрицы M , а $\delta_j = 0$, если $p_j \geq P$, и $\delta_j = 1$ при $p_j < P$. Определим коэффициенты $\mathfrak{C}^{(m)}$ равенствами (4.11), (4.12). Как и в доказательстве утверждения 7.2, получаем $(\tilde{\Pi}_1^{(P,Q)} e^{i(U\phi) \cdot r})_0 = V(\phi) + O(|\phi|^{Q+1})$ и $A(\phi) (\tilde{\Pi}_1^{(P,Q)} e^{i(U\phi) \cdot r})_0 = O(|\phi|^{Q+1})$, что в силу утверждения 3.1 влечёт Q -й порядок аппроксимации в смысле $\tilde{\Pi}_h^{(P,Q)}$.

Таким образом, в квазиодномерном случае анализ структуры ошибки решения по схеме (2.5) и получение оптимальной оценки точности может быть проведено путём поиска модифицированного отображения $\tilde{\Pi}_h$ в виде (2.8).

9. Контрпример

Условия, накладываемые в утверждениях 7.2 и 8.1, являются существенными. Обоснуем это примером схемы, которая обладает оценкой ошибки $\|\varepsilon_h(v_0, t, \mathring{\Pi}_h)\| \leq C_1 h \|\nabla v_0\|_\infty + C_2 h^2 (h + t) \|\nabla^3 v_0\|_\infty$, но для которой не существует отображение $\mathring{\Pi}_h$ вида (2.10), в смысле которого схема обладает вторым порядком аппроксимации.

Рассмотрим уравнение переноса (2.3)–(2.4) в \mathbb{R}^2 со скоростью переноса $\omega = 0$. Пусть $\mathbf{a}_1 = (1, 0)^T$, $\mathbf{a}_2 = (0, 1)^T$, так что $T = U = I$. Рассмотрим схему вида (2.5) при $|M^0| = 3$, ненулевые коэффициенты которой равны $Z_{(0,0)} = I$,

$$\begin{aligned} L_{(-1,-1)} &= W, & L_{(0,-1)} &= -F, & L_{(1,-1)} &= -W, \\ L_{(-1,0)} &= -E, & L_{(0,0)} &= 0, & L_{(1,0)} &= E, \\ L_{(-1,1)} &= -W, & L_{(0,1)} &= F, & L_{(1,1)} &= W, \end{aligned}$$

где

$$W = \begin{pmatrix} 0 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & -1 & 0 \end{pmatrix}, \quad F = \begin{pmatrix} 1 & -2 & 1 \\ -2 & 1 & 1 \\ 1 & 1 & -2 \end{pmatrix}, \quad E = \begin{pmatrix} 1 & 0 & -1 \\ 0 & -1 & 1 \\ -1 & 1 & 0 \end{pmatrix}.$$

Поскольку $L_\eta = -L_{-\eta}^*$, выполняется $A(\phi)^* = -A(\phi)$ и схема является устойчивой. Будем считать, что $\mathring{\Pi}_h$ определяется как $(\mathring{\Pi}_h f)_{\eta, \xi} = f(\eta_1 h, \eta_2 h)$.

Вначале проведём анализ аппроксимационной ошибки, вычисляя величины $(\varepsilon_1(f, \mathring{\Pi}_1))_0 = (L\mathring{\Pi}_1 f)_0 = \sum_\eta L_\eta (\mathring{\Pi}_1 f)_\eta$ для $f = \mathbf{r}^m$. Имеем $(\mathring{\Pi}_1 \mathbf{r}^m)_\eta = \eta^m \mathbf{e}$,

$$\sum_\eta L_\eta = 0, \quad \sum_\eta \eta_1 L_\eta = 2E, \quad \sum_\eta \eta_2 L_\eta = 2F, \quad (9.1)$$

$$\sum_\eta \eta_1^2 L_\eta = 0, \quad \sum_\eta \eta_1 \eta_2 L_\eta = 4W, \quad \sum_\eta \eta_2^2 L_\eta = 0. \quad (9.2)$$

Поскольку $E\mathbf{e} = F\mathbf{e} = 0$, схема обладает первым порядком аппроксимации, однако, так как $W\mathbf{e} \neq 0$, вторым порядком аппроксимации схема не обладает.

Покажем, что выполняется оценка

$$\|\varepsilon_h(t, v_0, \mathring{\Pi}_h)\| \leq C_1 h \|\nabla v_0\|_\infty + C_2 h^2 (t + h) \|\nabla^3 v_0\|_\infty. \quad (9.3)$$

Достаточно показать эту оценку для функций вида $v_0 = e^{i\alpha(\mathbf{k}) \cdot \mathbf{r}}$, $\mathbf{k} \in \mathbb{Z}^2$. Действительно, если для $v_0 = e^{i\alpha(\mathbf{k}) \cdot \mathbf{r}}$ выполняется (9.3), то для $\phi = h\alpha(\mathbf{k})$ имеем

$$\left\| \hat{\varepsilon}(\nu, \phi, \mathring{\Pi}_h) \right\| = \left\| \varepsilon_h(\nu h, e^{i\phi \cdot \mathbf{r}/h}, \mathring{\Pi}_h) \right\| \leq C_1 |\phi| + C_2 (\nu + 1) |\phi|^3,$$

откуда в силу утверждения 3.3 следует (9.3) для всех $v_0 \in C_{per}^3(\mathbb{R}^2)$.

Показать (9.3) для $v_0 = e^{i\alpha(\mathbf{k})\cdot\mathbf{r}}$, $\mathbf{k} \in \mathbb{Z}^2$, можно двумя способами. Первым способом является введение модифицированного отображения $\mathring{\mathbb{P}}_h$. Действительно, зафиксируем направление $\mathbf{e} = \alpha(\mathbf{k})/|\alpha(\mathbf{k})|$. Будем искать отображение $\mathring{\mathbb{P}}_h$ в виде

$$\mathring{\mathbb{P}}_{h,e}f = \mathring{\mathbb{P}}_h f + h\mathfrak{C}_e\mathring{\mathbb{P}}_h \frac{\partial f}{\partial \mathbf{e}} \quad (9.4)$$

исходя из условия 2-го порядка аппроксимации рассматриваемой схемы на функциях вида $v(\mathbf{r}) = \hat{v}(\mathbf{e} \cdot \mathbf{r})$. Поскольку схема обладает 1-м порядком аппроксимации на любых гладких функциях, достаточно рассмотреть функцию $v(\mathbf{r}) = (\mathbf{e} \cdot \mathbf{r})^2/2$. Имеем аппроксимационную ошибку

$$\begin{aligned} \left(\epsilon_1(v, \mathring{\mathbb{P}}_{1,e}) \right)_0 &= \sum_{\boldsymbol{\eta}} L_{\boldsymbol{\eta}} \left(\frac{(\mathbf{e} \cdot \boldsymbol{\eta})^2}{2} \boldsymbol{\epsilon} + (\mathbf{e} \cdot \boldsymbol{\eta}) \mathfrak{C}_e \boldsymbol{\epsilon} \right) = \\ &= 4e_x e_y W \boldsymbol{\epsilon} + 2(e_x E + e_y F) \mathfrak{C}_e \boldsymbol{\epsilon}. \end{aligned} \quad (9.5)$$

Подставляя теперь выражения для E , F и W и приравнявая $(\epsilon_1(v, \mathring{\mathbb{P}}_{1,e}))_0$ к нулю, получаем

$$\begin{pmatrix} e_x + e_y & -2e_y & -e_x + e_y \\ -2e_y & -e_x + e_y & e_x + e_y \\ -e_x + e_y & e_x + e_y & -2e_y \end{pmatrix} \mathfrak{C}_e \boldsymbol{\epsilon} = -2e_x e_y \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}. \quad (9.6)$$

Матрица $e_x E + e_y F$ симметрическая. Поскольку сумма столбцов равна нулю, у неё есть собственное значение $\lambda = 0$, которому соответствует собственный вектор $\boldsymbol{\epsilon}$. С другой стороны, главный минор размера 2 равен $(e_x + e_y)(-e_x + e_y) - 4e_y^2 = -3e_y^2 - e_x^2$ и отличен от нуля при $\mathbf{e} \neq 0$. Следовательно, критерием совместности этой системы является ортогональность правой части вектору $\boldsymbol{\epsilon}$, и он, очевидно, выполняется. Можно показать, что решением этой системы является диагональная матрица

$$\mathfrak{C}_e = \frac{2e_x e_y}{e_x^2 + 3e_y^2} \begin{pmatrix} e_y & 0 & 0 \\ 0 & -e_x & 0 \\ 0 & 0 & -e_y \end{pmatrix} + c_e I. \quad (9.7)$$

Положим $c_e = 0$. Таким образом, для любого направления \mathbf{e} существует отображение $\mathring{\mathbb{P}}_{h,e}$ вида (9.4), в смысле которого схема обладает вторым порядком аппроксимации на плоских волнах, таких что волновой вектор направлен по \mathbf{e} .

По неравенству треугольника и в силу устойчивости с $K = 1$ имеем

$$\begin{aligned} \|\epsilon_h(t, e^{i\alpha(\mathbf{k})\cdot\mathbf{r}}, \mathring{\mathbb{P}}_h)\| &\leq \|\epsilon_h(t, e^{i\alpha(\mathbf{k})\cdot\mathbf{r}}, \mathring{\mathbb{P}}_{h,e})\| + 2 \left\| (\mathring{\mathbb{P}}_{h,e} - \mathring{\mathbb{P}}_h) e^{i\alpha(\mathbf{k})\cdot\mathbf{r}} \right\| \leq \\ &\leq t \left\| \epsilon_h(e^{i\alpha(\mathbf{k})\cdot\mathbf{r}}, \mathring{\mathbb{P}}_{h,e}) \right\| + 2h \|\mathfrak{C}_e\| |\alpha(\mathbf{k})|. \end{aligned} \quad (9.8)$$

Оценим первое слагаемое в правой части. Очевидно,

$$\left\| \epsilon_h(e^{i\alpha(\mathbf{k}) \cdot \mathbf{r}}, \overset{\circ}{\tilde{\Pi}}_{h,e}) \right\| = \left\| \left(\epsilon_h(e^{i\alpha(\mathbf{k}) \cdot \mathbf{r}}, \overset{\circ}{\tilde{\Pi}}_{h,e}) \right)_0 \right\|.$$

Далее, по определению

$$\left(\epsilon_h(e^{i\alpha(\mathbf{k}) \cdot \mathbf{r}}, \overset{\circ}{\tilde{\Pi}}_{h,e}) \right)_0 = - \sum_{\eta \in \mathcal{S}} L_\eta \left(\overset{\circ}{\tilde{\Pi}}_h e^{i\alpha(\mathbf{k}) \cdot \mathbf{r}} \right)_\eta.$$

Введём $g(x) = 1 + ix - x^2/2$. В силу точности на многочленах 2-го порядка от $(\alpha(\mathbf{k}) \cdot \mathbf{r})$ получаем

$$\begin{aligned} \left(\epsilon_h(e^{i\alpha(\mathbf{k}) \cdot \mathbf{r}}, \overset{\circ}{\tilde{\Pi}}_{h,e}) \right)_0 &= -\frac{1}{h} \sum_{\eta \in \mathcal{S}} L_\eta \left(\overset{\circ}{\tilde{\Pi}}_{h,e} \left[e^{i\alpha(\mathbf{k}) \cdot \mathbf{r}} - g(\alpha(\mathbf{k}) \cdot \mathbf{r}) \right] \right)_\eta = \\ &= -\frac{1}{h} \sum_{\eta \in \mathcal{S}} L_\eta \left[e^{i\alpha(\mathbf{k}) \cdot h\eta} - g(h\alpha(\mathbf{k}) \cdot \eta) \right] \epsilon - \\ &- \sum_{\eta \in \mathcal{S}} L_\eta(i\alpha(\mathbf{k}) \cdot \mathbf{e}) \left[e^{i\alpha(\mathbf{k}) \cdot h\eta} - (1 + ih\alpha(\mathbf{k}) \cdot \eta) \right] \mathfrak{C}_e \epsilon. \end{aligned}$$

Видно, что эта величина имеет оценку $O(|\alpha(\mathbf{k})|^3 h^2)$ с константой, не зависящей от направления α . Подставляя эту оценку в (9.8), получаем (9.3) с подстановкой $e^{i\alpha(\mathbf{k}) \cdot \mathbf{r}}$ вместо v_0 , причём константы C_1 и C_2 не зависят от $\alpha(\mathbf{k})$, что и требовалось доказать.

Оценку (9.3) можно показать и при помощи спектрального анализа. Мы не будем проводить рассуждение строго, ограничившись общими соображениями. Имеем

$$\begin{aligned} A(\phi) = -L(\phi) &= W \left(e^{i(\phi_1+\phi_2)} + e^{-i(\phi_1+\phi_2)} - e^{i(\phi_1-\phi_2)} - e^{i(-\phi_1+\phi_2)} \right) + \\ &+ E \left(e^{i\phi_1} - e^{-i\phi_1} \right) + F \left(e^{i\phi_1} - e^{-i\phi_1} \right) = \\ &= -4W \sin(\phi_1) \sin(\phi_2) + 2iE \sin(\phi_1) + 2iF \sin(\phi_2). \end{aligned}$$

Введём обозначения $u = \sin(\phi_1)$, $v = \sin(\phi_2)$. Тогда

$$A(\phi) = 2i \begin{pmatrix} v + u & -2v + 2iuv & v - u \\ -2v - 2iuv & v - u & v + u + 2iuv \\ v - u & v + u - 2iuv & -2v \end{pmatrix}.$$

Характеристическое уравнение для матрицы $A(\phi)/2i$ имеет вид

$$\lambda^3 - (9v^2 + 3u^2 + 8v^2u^2)\lambda + 4v^2u^2(v - u) = 0.$$

После замены переменных $\lambda = \mu(9v^2 + 3u^2 + 8v^2u^2)^{1/2}$ получаем

$$\mu^3 - \mu + \omega = 0,$$

где

$$\omega = 4v^2u^2(v - u)(9v^2 + 3u^2 + 8v^2u^2)^{-3/2}.$$

При $(u, v) \rightarrow (0, 0)$ выполняется $\omega \rightarrow 0$. При $\omega = 0$ имеем три корня: $\mu = 0$, $\mu = \pm 1$; им соответствует корень $\lambda = 0$ кратности 3. При малых ω имеем $\mu_1(\omega) \approx \omega$, что соответствует $|\lambda_1(\phi)| \sim |\phi|^3$, а остальные собственные значения $|\lambda_2(\phi)|, |\lambda_3(\phi)| \sim |\phi|$.

Пусть $\Lambda(\phi) = \text{diag}\{\lambda_1(\phi), \lambda_2(\phi), \lambda_3(\phi)\}$, $A(\phi) = S(\phi)\Lambda(\phi)S^{-1}(\phi)$. Компоненты вектора $v(\phi) = S^{-1}(\phi)(\mathring{\Pi}_1 e^{i\phi \cdot r})_0$, соответствующие второму и третьему собственным значениям, имеют величину не больше чем порядка $O(|\phi|)$, поскольку в противном случае не может быть сеточной сходимости. Значит, компонента, соответствующая первому собственному значению, имеет величину порядка $O(1)$, поскольку хотя бы одна компонента должна иметь величину порядка $O(1)$. Отсюда

$$\hat{\varepsilon}(\phi, \nu, \mathring{\Pi}_h) = S(\phi) \left(e^{\nu\Lambda(\phi)} - 1 \right) v(\phi) = O(|\phi| + \nu|\phi|^3).$$

В силу утверждения 3.2 для ошибки решения по рассматриваемой схеме имеет место оценка $O(h + h^2t)$.

Наконец, продемонстрируем эту оценку в численном счёте. Положим начальные данные равными $v_0(\mathbf{r}) = e^{i\alpha \cdot \mathbf{r}}$ при $\alpha = (2\pi, -2\pi)$. Расчёты будем проводить на сетках с $h = 1/60$, $h = 1/120$ и $h = 1/240$. Зависимость нормы ошибки решения от времени приведена на рис. 1. Видно, что при малых временах ошибка имеет величину порядка $O(h)$, тогда как растущая со временем компонента ошибки – величину порядка $O(h^2)$. Интегрирование по времени осуществлялось явным 3-стадийным методом Рунге–Кутты 3-го порядка. Шаг по времени был выбран равным $h/20$; его дальнейшее уменьшение визуально не отражалось на величине ошибки.

Теперь, покажем, что не существует отображения $\mathring{\Pi}_h$ вида

$$\begin{aligned} \left(\mathring{\Pi}_h f \right)_\eta &= \left(\mathring{\Pi}_h f \right)_\eta + h \left(\mathfrak{e}^{(x)} \left(\mathring{\Pi}_h f_x \right)_\eta + \mathfrak{e}^{(y)} \left(\mathring{\Pi}_h f_y \right)_\eta \right) + \\ &+ h^2 \left(\mathfrak{e}^{(xx)} \left(\mathring{\Pi}_h f_{xx} \right)_\eta + \mathfrak{e}^{(xy)} \left(\mathring{\Pi}_h f_{xy} \right)_\eta + \mathfrak{e}^{(yy)} \left(\mathring{\Pi}_h f_{yy} \right)_\eta \right), \end{aligned} \quad (9.9)$$

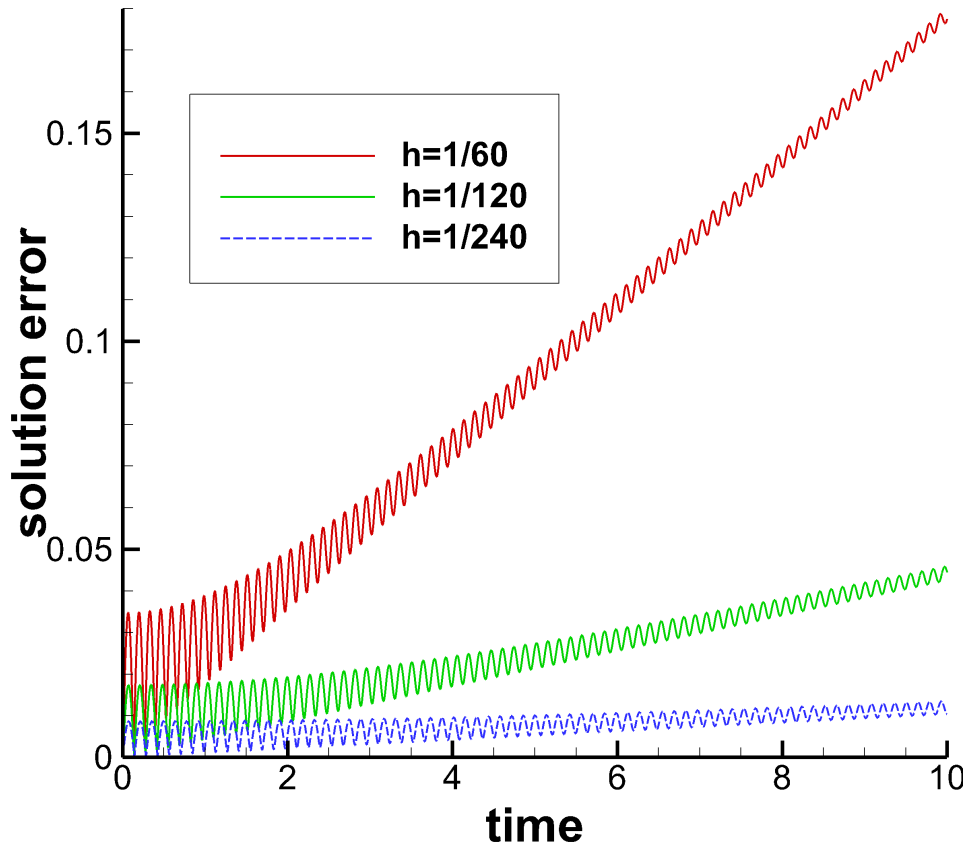


Рис. 1. Зависимость ошибки решения от времени

в смысле которого схема обладает 2-м порядком аппроксимации. Действительно, поскольку $A(0) = 0$, аппроксимационная ошибка в смысле отображения (9.9) на квадратичном полиноме не зависит от коэффициентов $\mathfrak{E}^{(xx)}$, $\mathfrak{E}^{(xy)}$, $\mathfrak{E}^{(yy)}$. Поэтому отображение, доставляющее 2-й порядок аппроксимации, нужно искать в виде

$$\left(\overset{\circ}{\Pi}_h f\right)_\eta = \left(\dot{\Pi}_h f\right)_\eta + h \left(\mathfrak{E}^{(x)} \left(\dot{\Pi}_h f_x\right)_\eta + \mathfrak{E}^{(y)} \left(\dot{\Pi}_h f_y\right)_\eta \right).$$

Аппроксимационная ошибка должна быть равна нулю на любой квадратичной функции; достаточно проверять функции $f = x^2$, $f = y^2$, $f = xy$. На функции $f = x^2$ образ под действием этого проектора совпадает с (9.4) при $e = (1, 0)^T$, где $\mathfrak{E}_e \equiv \mathfrak{E}^{(x)}$. Система (9.6) на нахождение этих коэффициентов при подстановке $\alpha_1 \neq 0$, $\alpha_2 = 0$ имеет решение $\mathfrak{E}^{(x)} = \kappa_x I$, где κ_x – любое. Аналогичным образом получаем $\mathfrak{E}^{(y)} = \kappa_y I$, где κ_y – любое. Из двух уравнений, таким образом, мы однозначно (с точностью до κI , не влияющих на аппроксимационную ошибку) определили коэффициенты отображения. Однако легко удостовериться, что аппроксимационная ошибка на функции xy не зависит от κ_x и κ_y и в силу (9.2) равна $4Wh$. Таким образом, искомого отображения не существует, что и требовалось доказать.

10. Особые направления волнового вектора

С формальной точки зрения, мы можем рассматривать только волны с волновыми векторами $\alpha(\mathbf{k})$, $\mathbf{k} \in \mathbb{Z}^d$. То есть число направлений, по которым могут распространяться плоские волны, счётно. Более того, при разных N векторы $\varepsilon_h(t, e^{i\alpha(\mathbf{k})\cdot\mathbf{r}}, \Pi_h)$ имеют разные размерности. Но ввиду равенства (3.7) величина $\|\varepsilon_h(t, e^{i\alpha(\mathbf{k})\cdot\mathbf{r}}, \Pi_h)\|^2$ является аналитической функцией $\alpha(\mathbf{k})$, и это позволяет по непрерывности доопределить её на некоторую окрестность нуля $\alpha \in \tilde{G} \subset \mathbb{R}^d$.

Пусть Ω – единичная сфера в \mathbb{R}^d . Будем называть направление $e \in \Omega$ *минус-особым* по аппроксимации (точности, точности в длительном счёте), если существует последовательность $e_n \rightarrow e$, такая что при всех n порядок аппроксимации (формальный порядок точности, порядок точности в длительном счёте) для волн, распространяющихся в направлении e_n , равен A , но для волн, распространяющихся в направлении e , он равен $\tilde{A} < A$. Аналогично будем называть направление *плюс-особым*, если $\tilde{A} > A$. Рассмотрим вопрос, какими могут быть особые направления.

Пример схемы, для которой существует плюс-особое направление, построить легко. Рассмотрим конечно-разностную схему 1-го порядка на равномерной сетке для аппроксимации уравнения переноса $\partial u / \partial t + \partial u / \partial x = 0$ в двумерном случае ($\mathbf{r} = (x, y)^T \in \mathbb{R}^2$, $\boldsymbol{\omega} = (1, 0)^T$):

$$\frac{du_{i,j}}{dt} + \frac{u_{i,j} - u_{i-1,j}}{h} = 0.$$

Рассмотрим ошибку аппроксимации в смысле некоторого локального отображения Π_h вида (2.7). Для решений, зависящих только от y , схема будет точна, то есть порядки аппроксимации, точности и точности в длительном счёте формально будут бесконечными. Для всех остальных решений схема будет обладать первым порядком аппроксимации, первым порядком точности и первым порядком точности в длительном счёте.

Теперь рассмотрим вопрос, возможны ли минус-особые направления.

Утверждение 10.1. *Минус-особых направлений по аппроксимации не существует.*

Порядок аппроксимации P_A в смысле отображений $\hat{\Pi}_h = \Pi_h$ или $\hat{\Pi}_h = \dot{\Pi}_h$ для плоских волн, распространяющихся по некоторому направлению \tilde{e} , равносильно выполнению равенств $(\epsilon_1((\mathbf{r} \cdot \tilde{e})^m, \hat{\Pi}_1))_0 = 0$ для $m = 1, \dots, P_A$, где $\epsilon_h(f, \hat{\Pi}_h)$ определено (2.12). Выражения $(\epsilon_1((\mathbf{r} \cdot \tilde{e})^m, \hat{\Pi}_1))_0$ являются аналитическими функциями \tilde{e} , поэтому если они равны нулю для некоторой последовательности $e_n \rightarrow e$, то они равны нулю и для e . Таким образом, минус-особых направлений по аппроксимации не существует.

Утверждение 10.2. *Минус-особых направлений по точности не существует.*

Аналогичное утверждение можно получить для порядка точности. Если для волн, распространяющихся по направлениям e_n , имеет место порядок точности P , то величины $\hat{\varepsilon}^{(1)}(T^* \alpha h, t/h, \Pi_h)$ и $\hat{\varepsilon}^{(2)}(T^* \alpha h, t/h, \Pi_h)$ являются величинами P -го порядка малости по h . Поскольку $\hat{\varepsilon}^{(2)}(T^* \alpha h, t/h, \Pi_h)$ – аналитическая по h , это означает, что члены в её разложении слагаемые при h^p , $p < P$, отсутствуют. Поскольку эта функция также аналитическая по α , свойство отсутствия членов h^p , $p < P$, переносится с направлений e_n на направление e . Величина $\hat{\varepsilon}^{(1)}(T^* \alpha h, t/h, \Pi_h)$ имеет представление (4.5), причём верно (4.7). Следовательно, $\|v^{(*)}(e_n | \alpha | h)\| \leq \hat{C} h^P |\alpha|^P$. Но $v^{(*)}(\alpha h)$ по построению является аналитической функцией α и h , поэтому $\|v^{(*)}(e | \alpha | h)\| \leq \hat{C} h^P |\alpha|^P$, и в силу устойчивости получаем, что $\hat{\varepsilon}^{(1)}(T^* \alpha h, t/h, \Pi_h)$ имеет P -й порядок малости по h . Таким образом, минус-особых направлений по точности не существует.

Утверждение 10.3. *Пусть собственное значение $\lambda = 0$ матрицы $A(0)$ является простым. Тогда минус-особых направлений по точности в длительном счёте не существует.*

Действительно, пусть схема имеет некий формальный порядок P . Порядок точности в длительном счёте Q определяется только его “физическим” собственным значением $\lambda(\phi)$. Напомним, что $\lambda(\phi)$ является аналитической функцией в окрестности $\phi = 0$. Порядок точности в длительном счёте Q на $v_0 \in H_{per}^{Q+1}(\mathbb{R}^d)$ имеется тогда и только тогда, когда в некоторой окрестности нуля выполняется $|\lambda(\phi)| \leq C |\phi|^{Q+1}$. Аналогично, порядок точности в длительном счёте Q при начальных данных вида $v_0(r) = f(r \cdot e)$, имеется тогда и только тогда, когда $|\lambda(e\phi)| \leq C |\phi|^{Q+1}$. Отсюда видно, что в простом случае минус-особых направлений по точности в длительном счёте не бывает.

Условие простоты собственного значения $\lambda = 0$ матрицы $A(0)$ является существенным. Пример схемы, обладающей минус-особым направлением по точности в длительном счёте, будет приведён в следующей работе.

Замечание

В [4] доказано следующее утверждение.

Утверждение 10.4. Пусть $A \in \mathbb{C}^{n \times n}$ удовлетворяет (3.3), $v \in \mathbb{C}^n$. Пусть существуют такие $C_1, C_2 \geq 0$, что при всех $\nu \geq 0$ справедлива оценка

$$\|(e^{\nu A} - I)v\| \leq (C_1 + C_2\nu)\|v\|. \quad (10.1)$$

Тогда

$$v^* \left(A^* A + \frac{C_2^2}{C_1^2} \right)^{-1} A^* Av \leq \tilde{\delta}^2 C_1^2 \|v\|^2, \quad (10.2)$$

где $\tilde{\delta}$ зависит только от n и K . Обратно, если выполняется (10.2), то

$$\|(e^{\nu A} - I)v\| \leq \tilde{\delta}(K + 1)(C_1 + C_2\nu)\|v\|. \quad (10.3)$$

Покажем, что из него следует 1) сходимость с порядком аппроксимации (устойчивость схемы предполагается) и 2) утверждение 5.2. Действительно, пусть схема обладает порядком аппроксимации P_A . Тогда по утверждению 3.1 имеем $\|Av\| = O(|\phi|^{P_A+1})$. Значит, $v^* A^* Av = O(|\phi|^{2(P_A+1)})$, то есть

$$v^* (|\phi|^2)^{-1} A^* Av = O(|\phi|^{2P_A}).$$

Добавлением положительно полуопределённой эрмитовой матрицы под минус первую степень мы не можем увеличить значение выражения, поэтому

$$v^* (A^* A + |\phi|^2)^{-1} A^* Av = O(|\phi|^{2P_A}),$$

что по утверждению 10.4 влечёт неравенство (10.3) с $C_1 = O(|\phi|^{P_A})$ и $C_2 = O(|\phi|^{P_A+1})$.

Пусть теперь $A(0) = 0$ и схема обладает порядком точности P . Тогда $A(\phi)$ выражается в виде произведения $|\phi|$ на некоторую функцию (не обязательно аналитическую), ограниченную в нуле, и, следовательно, $A^* A = |\phi|^2 F(\phi)$, где F ограничено в нуле. Так как схема обладает порядком точности P , по утверждению 10.4 выполняется 10.2 с подстановкой $C_1 = c_1 |\phi|^P$, $C_2 = c_2 |\phi|^{P+1}$, то есть

$$v^* (|\phi|^2(I + F(\phi)))^{-1} A^* Av = O(|\phi|^{2P}).$$

Ввиду ограниченности F в нуле отсюда следует

$$v^* A^* Av = O(|\phi|^{2P+2}),$$

то есть $\|Av\| = O(|\phi|^{P+1})$. В силу утверждения 3.1 схема обладает P -м порядком аппроксимации.

Минус-особое направление по точности в длительном счёте

Выше было показано, что минус-особых направлений по аппроксимации и точности не бывает. Приведём пример схемы, для которой есть минус-особые направления по точности в длительном счёте. А именно, обладающей 1-м порядком точности в длительном счёте на волнах, распространяющихся в горизонтальном направлении, и 2-м порядком точности на остальных плоских волнах.

Рассмотрим уравнение (2.3) с $\omega = 0$, и пусть $\mathbf{a}_1 = (1,0)^T$, $\mathbf{a}_2 = (0,1)^T$. Положим $M^0 = \{L, R\}$ и определим матрицы

$$W = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}, \quad E = \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}.$$

Ненулевые коэффициенты схемы (2.5) выберем равными $Z_{0,0} = I$,

$$L_{0,0} = -2W, \quad L_{\pm 1,0} = \pm E, \quad L_{0,\pm 1} = W.$$

Оператор $\mathring{\Pi}_h$ зададим равенствами $(\mathring{\Pi}_h f)_{\eta,L} = (\mathring{\Pi}_h f)_{\eta,R} = f(\eta_1 h, \eta_2 h)$. Матрица $A(\phi) = -L(\phi)$ имеет вид

$$A(\phi_1, \phi_2) = \begin{pmatrix} 2i \sin \phi_1 & -2i \sin \phi_1 - 4 \sin^2(\phi_2/2) \\ -2i \sin \phi_1 + 4 \sin^2(\phi_2/2) & 2i \sin \phi_1 \end{pmatrix}. \quad (11.4)$$

Поскольку $(A(\phi))^* = -(A(\phi))$, схема устойчива.

Как и в предыдущем примере, будем искать оператор $\mathring{\Pi}_h$ вида (9.4), доставляющий второй порядок аппроксимации. Имеем аппроксимационную ошибку

$$\left(\epsilon_1(v, \mathring{\Pi}_{1,e}) \right)_0 = \sum_{\eta} L_{\eta} \left(\frac{(e \cdot \eta)^2}{2} \mathbf{e} + (e \cdot \eta) \mathfrak{C}_e \mathbf{e} \right). \quad (11.5)$$

Подставляя в (11.5) значения коэффициентов L_{η} , получаем

$$e_y^2 W \mathbf{e} + 2e_x E |\alpha| \mathfrak{C}_e \mathbf{e} = 0.$$

Подставляя далее коэффициенты W и E ,

$$2e_x \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} (\mathfrak{C}_e \mathbf{e}) + e_y^2 \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} = 0.$$

Видно, что при $e_x \neq 0$ система совместна и имеет решение

$$\mathfrak{C}_e = -\frac{e_y^2}{2e_x} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + c_e I,$$

тогда как при $e_x = 0$ система несовместна. Таким образом, для всех плоских волн с волновым вектором α , за исключением распространяющихся в горизонтальном направлении, схема обладает 2-м порядком точности в длительном счёте. Можно записать оценку

$$\|\varepsilon_h(t, e^{i\alpha \cdot r}, \mathring{\Pi}_h)\| \leq \min \left\{ \frac{\alpha_y^2}{2|\alpha_x|} h + C_2 |\alpha|^3 h^2 t, C_3 |\alpha|^2 h t \right\}.$$

Интересно также выписать для построенной схемы ошибку решения в явном виде, вычислив матричную экспоненту. Введём обозначения $a = 2 \sin \phi_1$, $b = 4 \sin^2(\phi_2/2)$, $m = \sqrt{a^2 + b^2}$, $v = a + bi$. Тогда матрица $A(\phi)$, заданная (11.4) примет вид

$$A(\phi) = \begin{pmatrix} ia & -ia - b \\ -ia + b & ia \end{pmatrix} = S \Lambda S^{-1}, \quad (11.6)$$

где

$$S = \begin{pmatrix} -1 & 1 \\ v/m & v/m \end{pmatrix}, \quad \Lambda = \begin{pmatrix} ia + im & 0 \\ 0 & ia - im \end{pmatrix}, \quad S^{-1} = \frac{1}{2} \begin{pmatrix} -1 & m/v \\ 1 & m/v \end{pmatrix}.$$

Матричная экспонента равна

$$\exp(\nu A(\phi)) = e^{i\nu a} \begin{pmatrix} \cos(\nu m) & -i \frac{m}{v} \sin(\nu m) \\ -i \frac{v}{m} \sin(\nu m) & \cos(\nu m) \end{pmatrix}. \quad (11.7)$$

Поскольку $A(\phi)$ была косоэрмитовой, матричная экспонента получилась унитарной матрицей (отметим, что $|v| = m$). Функция ошибки решения (3.5) имеет вид

$$\hat{\varepsilon}(\phi, \nu, \mathring{\Pi}_h) = \left[e^{\nu A(\phi)} - I \right] (\mathring{\Pi}_1 e^{i\phi \cdot r})_0. \quad (11.8)$$

Поскольку $A(0) = 0$, в силу утверждения 5.2 функция $\hat{\varepsilon}(\alpha h, t/h, \mathring{\Pi}_h)$ является аналитической функцией t, h, α . То есть она представима в виде ряда по степеням h . Опуская выкладки ввиду их громоздкости (для этого нужно вернуться к “физическим” переменным $h, t = h\nu, \alpha_x = \phi_1/h, \alpha_y = \phi_2/h$), приведём первый член этого ряда:

$$\hat{\varepsilon}(\alpha h, t/h, \mathring{\Pi}_h) = -\frac{\alpha_y^2}{4\alpha_x} (1 - \exp(-4t\alpha_x)) h + O(h^2). \quad (11.9)$$

При любом фиксированном $\alpha_x \neq 0$ слагаемое при h является ограниченной функцией времени и в длительном счёте имеет место второй порядок точности. Покажем теперь, что при $\alpha_x = 0$ второй порядок точности в длительном счёте не имеет места. Из неограниченности по времени члена при h^1 в (11.9)

это, вероятно, не следует. Поэтому рассмотрим выражение для матричной экспоненты (11.7) целиком. Подставляя в него $a = 0$, $v = ib$, $m = b$, получаем

$$\exp(\nu A(\phi)) = \begin{pmatrix} \cos(\nu b) & -\sin(\nu b) \\ \sin(\nu b) & \cos(\nu b) \end{pmatrix},$$

откуда

$$\hat{\varepsilon}(\alpha h, t/h, \mathring{\Pi}_h) = \begin{pmatrix} \cos(\nu b) - 1 - \sin(\nu b) \\ \cos(\nu b) - 1 + \sin(\nu b) \end{pmatrix},$$

где $\nu b = 4t \sin^2(\alpha_y h/2)/h$. Очевидно, что полученное выражение для ошибки решения оценкой $O(h + h^2 t)$ не обладает.

Вспомогательное утверждение

Утверждение 11.5. Пусть $m, n_j \in \mathbb{R}_+^d$, $j = 1, \dots, N$. Следующие два утверждения эквивалентны.

1. $\exists C \forall \phi \in [0,1]^d \phi^m \leq C \sum_j \phi^{n_j}$;
2. $\exists \delta_j \geq 0 \sum \delta_j = 1, \sum \delta_j n_j \leq m$.

Действительно, пусть выполняется условие 2. Тогда при $\phi \in [0,1]^d$ имеем

$$\phi^m \leq \phi^{\sum \delta_j n_j} = \prod_j (\phi^{n_j})^{\delta_j} \leq \sum_j \phi^{n_j}.$$

Последнее неравенство написано в силу неравенства Юнга. Таким образом, имеем условие 1.

Обратно, пусть условие 2 не выполняется. Пусть $M = \text{Conv}\{n_j\} + \mathbb{R}_+^d$, где плюсом обозначена сумма Минковского. Очевидно, M выпуклое и $m \notin M$. По теореме отделимости найдётся такое $q \in \mathbb{R}^d$, что для всех $a \in M$ выполняется $q \cdot a > C > q \cdot m$. В частности, подставив $a = n_j$ и взяв минимум по j , получаем

$$\min_j \{q \cdot n_j\} > C > q \cdot m. \quad (11.10)$$

Покажем, что $q \in \mathbb{R}_+^d$. Допустим противное, что какая-то компонента вектора q отрицательная, и обозначим её номер через k . Выберем любой $a \in M$ и рассмотрим последовательность

$$a_l = a + l(0, \dots, 1, \dots, 0)^T,$$

где единица стоит в k -й позиции. По построению такой вектор лежит в M , но $q \cdot a_l \rightarrow -\infty$ и поэтому не может превосходить C для всех l . Следовательно, $q \in \mathbb{R}_+^d$.

Допустим теперь, что выполняется условие 1. Подставим в него $\phi_1 = \varepsilon^{q_1}$, \dots , $\phi_d = \varepsilon^{q_d}$ и устремим ε к нулю. Сравнивая степени ε в левой и правой частях неравенства, получаем

$$\mathbf{q} \cdot \mathbf{m} \geq \min_j \{\mathbf{q} \cdot \mathbf{n}_j\},$$

что противоречит (11.10).

12. Странное особое направление

В рассмотренном выше примере для волн, распространяющихся вдоль особого направления, порядок точности в длительном счёте был на единицу меньше, чем для окрестных направлений. Но возможен и такой случай, когда направление e не является минус-особым, но при приближении к нему константы C_1 и/или C_2 в оценке вида (2.13) для $v_0(\mathbf{r}) = e^{i|\alpha|e_n \cdot \mathbf{r}}$ неограниченно растут.

Как и в предыдущих двух примерах, рассмотрим уравнение (2.3)–(2.4) в \mathbb{R}^2 со скоростью переноса $\boldsymbol{\omega} = 0$, и пусть \mathbf{a}_1 и \mathbf{a}_2 совпадают с единичными ортами. Положим $M^0 = \{1, 2, 3, 4, 5\}$ и определим матрицы

$$W = \begin{pmatrix} 0 & 1 & -1 & 0 & 0 \\ -1 & 0 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \quad E = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix},$$

$$R = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 \\ 0 & 0 & 0 & -1 & 1 \end{pmatrix}, \quad Q = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & -1 & 0 \end{pmatrix}.$$

Ненулевые коэффициенты схемы (2.5) выберем равными $Z_{0,0} = I$,

$$L_{0,0} = -2W + 6Q, \quad L_{\pm 1,0} = W - 4Q \mp (2E - R), \quad L_{\pm 2,0} = Q,$$

$$L_{\pm 1,1} = L_{\pm 1,-1} = \pm E.$$

Оператор $\mathring{\Pi}_h$ зададим равенствами

$$(\mathring{\Pi}_h f)_{\eta,\xi} = f \left(\begin{pmatrix} \eta_1 \\ \eta_2 \end{pmatrix} h \right), \quad \xi = 1, \dots, 4; \quad (\mathring{\Pi}_h f)_{\eta,5} = f \left(\begin{pmatrix} \eta_1 + 1 \\ \eta_2 \end{pmatrix} h \right).$$

Матрица $A(\boldsymbol{\phi}) = -L(\boldsymbol{\phi})$ имеет вид

$$A(\phi_1, \phi_2) = 4 \sin^2 \left(\frac{\phi_1}{2} \right) W + 8i \sin \phi_1 \sin^2 \left(\frac{\phi_2}{2} \right) E -$$

$$-2i \sin \left(\frac{\phi_1}{2} \right) R - 16 \sin^4 \left(\frac{\phi_1}{2} \right) Q. \quad (12.1)$$

Поскольку $(A(\phi))^* = -(A(\phi))$, схема устойчива.

Прежде всего, определим порядок аппроксимации. Имеем

$$\hat{\epsilon}(\phi, \mathring{\Pi}_h) = A(\phi_1, \phi_2) \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ e^{i\phi_1} \end{pmatrix} = -2i \sin\left(\frac{\phi_1}{2}\right) (1 - e^{i\phi_1}) \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ -1 \end{pmatrix} + O(|\phi|^3).$$

Таким образом, $\hat{\epsilon}(\phi, \mathring{\Pi}_h)$ имеет второй порядок малости по $|\phi|$ и, следовательно, схема имеет 1-й порядок аппроксимации. Далее, поскольку $A(0) = 0$, в силу следствия 5.2 схема обладает 1-м порядком точности и не обладает вторым.

Теперь займёмся исследованием этой схемы на порядок точности в длительном счёте. Уравнения для первых трёх компонент численного решения не связаны с уравнениями для последних двух компонент. Поэтому рассмотрим их в отдельности.

Вначале рассмотрим систему для первых трёх компонент сеточной функции. Соответствующее ограничение матрицы $A(\phi_1, \phi_2)$ имеет вид

$$A(\phi_1, \phi_2) = \begin{pmatrix} iy & x & -x \\ -x & 0 & x \\ x & -x & -iy \end{pmatrix}, \quad (12.2)$$

где $x = 4 \sin^2(\phi_1/2)$, $y = 8 \sin(\phi_1) \sin^2(\phi_2/2)$. Собственными числами матрицы $A(\phi_1, \phi_2)$ являются $\pm ir$ и 0 , где $r = \sqrt{3x^2 + y^2}$. Будем считать, что $\sin \phi_1 \geq 0$; случай $\sin \phi_1 < 0$ может быть рассмотрен аналогично. Также исключим случай $\phi_1 = \phi_2 = 0$, где матрица нулевая. Матрица правых собственных векторов имеет вид

$$S = \begin{pmatrix} 3(a-i)(1+b) & 3a(a+i) & a \\ 3a(1+b+3ia) & (1+b)(1-b-3ia) & a-ib \\ -3a(2-b) & -(b+1)(b+2) & a \end{pmatrix},$$

где $a = x/r$, $b = y/r$. Детерминант этой матрицы равен $\det S = -18(b+1)(2ia+b)$ и не обращается в 0 , поскольку $a \geq 0$, $b \geq 0$ и одновременно не обращаются в ноль. Обратная матрица имеет вид

$$S^{-1} = \frac{1}{\det S} \begin{pmatrix} -(1+b)(i-3a+2ib) & -3a(b+ia+1) & 3a(b+2ia) \\ 3a(2ib-3a-i) & -3a(ib-3a+i) & 3(b+1)(2ia+b) \\ -18a(b+1)(2ia+b) & 18(b+1)(ia^2+ab-i) & -18a(b+1)(2ia+b) \end{pmatrix}.$$

Если вычислить действие матрицы S^{-1} на образ начальных данных, получаем

$$S^{-1}\mathbf{e} = \frac{1}{\det S} \begin{pmatrix} -3ib(1+ia+b) \\ 3b(1+3ia+b) \\ -18(b+1)(ia^2+ab-i) \end{pmatrix}.$$

Последняя компонента не вносит вклад в ошибку, так как соответствующее собственное значение равно 0. Первые две компоненты ведут себя одинаковым образом, поэтому достаточно рассмотреть одну, например первую. Имеем

$$(S^{-1}\mathbf{e})_1 = \frac{ib(1+ia+b)}{6(b+1)(2ia+b)} = \frac{iy(r+ix+y)}{6(y+r)(2ix+y)} = \frac{iy}{6(2ix+y)} - \frac{xy}{6(y+r)(2ix+y)}.$$

Подставим выражения для x и y . Пренебрегая членами высоких порядков малости, получаем $x \approx \phi_1^2$ и $y \approx 2\phi_1\phi_2^2$. Таким образом,

$$(S^{-1}\mathbf{e})_1 \approx \frac{i\phi_2^2}{6(i\phi_1 + \phi_2^2)} - \frac{\phi_1\phi_2^2}{6(2\phi_2^2 + \sqrt{3\phi_1^2 + 4\phi_2^4})(i\phi_1 + \phi_2^2)}.$$

При любом линейном соотношении между ϕ_1 и ϕ_2 , исключая $\phi_1 = 0$, величина $(S^{-1}\mathbf{e})_1$ имеет первый порядок малости по $|\phi|$. В то же время, если положить $\phi_1 = 0$, то $x = y = 0$, поэтому все три собственных значения $A(\phi_1, \phi_2)$ обращаются в ноль и, следовательно, ошибка решения тождественно равна нулю. Таким образом, какое бы мы ни взяли направление, для волн, распространяющихся по этому направлению, имеет место оценка ошибки $O(h)$, и можно говорить, что порядок точности в длительном счёте равен бесконечности. Однако константа в этой оценке существенным образом зависит от направления. Действительно, если положить, например, $\phi_1 = \phi_2^2$, то $(S^{-1}\mathbf{e})_1 \rightarrow (i - 1/\sqrt{7})(1 - i)/12$ при $\phi = (\phi_2^2, \phi_2) \rightarrow 0$. Соответствующее собственное значение $\lambda_1(\phi_2^2, \phi_2) = ir \approx i\sqrt{7}\phi_2^4$. Таким образом, при $\phi = (\phi_2^2, \phi_2)$ в проколотой окрестности нуля выполняется

$$\sup_{\nu \geq 0} \|\hat{\varepsilon}(\phi, \nu, \mathring{\Pi}_h)\|_3 \geq C \sup_{\nu \geq 0} |e^{\nu\lambda_1(\phi)}(S^{-1}\mathbf{e})_1| = 2C|(S^{-1}\mathbf{e})_1| \geq \tilde{C} > 0.$$

Здесь $\|\cdot\|_{(3)}$ – полунорма на V_{per} , учитывающая только первые три компоненты решения в каждом сеточном блоке. Следовательно, оценки вида $\|\varepsilon_h(t, e^{i\alpha \cdot r}, \mathring{\Pi}_h)\| \leq Ch$, где C не зависит от t и h , для рассматриваемой схемы не существует.

Графики $\|\varepsilon_h(t, e^{i\alpha \cdot r}, \mathring{\Pi}_h)\|_{(3)}$ в зависимости от времени на сетках с $h = 1/60$, $h = 1/120$, $h = 1/240$ при $\alpha_x = 4\pi$, $\alpha_y = -2\pi$ приведены на рис. 2.

Рассмотрим теперь эту подсистему методом введения вспомогательного отображения. Понимая под W и E ограничения соответствующих операторов на трёхмерное векторное подпространство, имеем

$$\sum_{\eta} L_{\eta} = 0, \quad \sum_{\eta} L_{\eta}(e \cdot \eta) = 0, \quad (12.3)$$

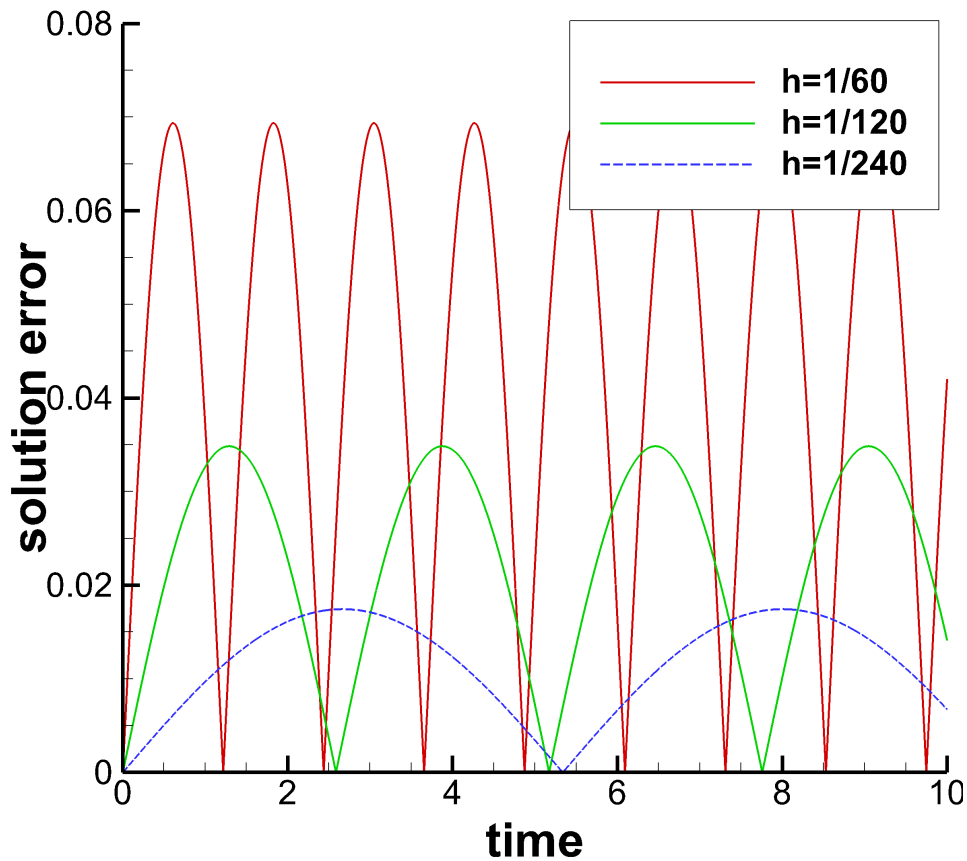


Рис. 2. Зависимость ошибки решения от времени

$$\sum_{\eta} L_{\eta} \frac{(\mathbf{e} \cdot \boldsymbol{\eta})^2}{2} = e_x^2 W, \quad \sum_{\eta} L_{\eta} \frac{(\mathbf{e} \cdot \boldsymbol{\eta})^3}{6} = \sum_{\eta} L_{\eta} \frac{(e_x \eta_1 + e_y \eta_2)^3}{6} = 2e_x e_y^2 E.$$

$$\sum_{\eta} L_{\eta} \frac{(\mathbf{e} \cdot \boldsymbol{\eta})^4}{24} = \sum_{\eta} L_{\eta} \frac{(e_x \eta_1 + e_y \eta_2)^4}{24} = \frac{e_x^4}{12} W. \quad (12.4)$$

Поскольку

$$\sum_{\eta} L_{\eta} \boldsymbol{\epsilon} = 0, \quad \sum_{\eta} L_{\eta} \boldsymbol{\eta} \boldsymbol{\epsilon} = 0, \quad \sum_{\eta} L_{\eta} \boldsymbol{\eta} \otimes \boldsymbol{\eta} \boldsymbol{\epsilon} = 0,$$

порядок аппроксимации в смысле $\mathring{\Pi}_h$ равен 2.

Выпишем уравнение на нахождение коэффициентов оператора $\mathring{\Pi}_{h,\mathbf{e}}$, в смысле которого может быть обеспечен 3-й порядок аппроксимации на плоских волнах, распространяющихся по направлению \mathbf{e} . Поскольку выполняется (12.3), в уравнение будут входить только коэффициенты 1-го порядка $\mathfrak{C}_{\mathbf{e}}$, но не

2-го и не 3-го. Имеем

$$\begin{aligned} \left(\epsilon_1((\mathbf{e} \cdot \mathbf{r})^3/6, \overset{\circ}{\Pi}_{1,e}) \right)_0 &= \sum_{\eta} L_{\eta} \left(\frac{(\mathbf{e} \cdot \boldsymbol{\eta})^3}{6} \boldsymbol{\epsilon} + \frac{(\mathbf{e} \cdot \boldsymbol{\eta})^2}{2} \mathfrak{C}_e \boldsymbol{\epsilon} \right) = \\ &= 2e_x e_y^2 E \boldsymbol{\epsilon} + e_x^2 W \mathfrak{C}_e \boldsymbol{\epsilon}. \end{aligned}$$

Приравнивая к нулю, получаем

$$e_x \left[e_x \begin{pmatrix} 0 & 1 & -1 \\ -1 & 0 & 1 \\ 1 & -1 & 0 \end{pmatrix} (\mathfrak{C}_e \boldsymbol{\epsilon}) + 2e_y^2 \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \right] = 0.$$

Видно, что при всех e система совместна. При $e_x \neq 0$ она имеет решение

$$\mathfrak{C}_e = -2 \frac{e_y^2}{e_x} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} + c_e I, \quad (12.5)$$

тогда как при $e_x = 0$ ей удовлетворяют любые значения \mathfrak{C}_e . Положим $c_e = 0$. Рассмотрим аппроксимационную ошибку на полиноме 4-го порядка в смысле полученного оператора $\overset{\circ}{\Pi}_{h,e}$:

$$\left(\epsilon_1 \left(\frac{(\mathbf{e} \cdot \mathbf{r})^4}{24}, \overset{\circ}{\Pi}_{1,e} \right) \right)_0 = \sum_{\eta} L_{\eta} \left(\frac{(\mathbf{e} \cdot \boldsymbol{\eta})^4}{24} \boldsymbol{\epsilon} + \frac{(\mathbf{e} \cdot \boldsymbol{\eta})^3}{6} \mathfrak{C}_e \boldsymbol{\epsilon} \right).$$

Используя (12.4) и заметив, что $W \boldsymbol{\epsilon} = 0$ и $E \mathfrak{C}_e \boldsymbol{\epsilon} = 0$, получаем, что $(\epsilon_1((\mathbf{e} \cdot \mathbf{r})^4, \overset{\circ}{\Pi}_{1,e}))_0 = 0$, то есть в смысле отображения $\overset{\circ}{\Pi}_{h,e}$ схема обладает 4-м порядком аппроксимации на волнах, распространяющихся по e .

Таким образом, какое бы ни было выбрано направление, для плоских волн, распространяющихся в этом направлении, можно записать оценку ошибки

$$\|\varepsilon_h(t, e^{i\boldsymbol{\alpha} \cdot \mathbf{r}}, \overset{\circ}{\Pi}_h)\|_{(3)} \leq C_1 \left(\frac{\boldsymbol{\alpha}}{|\boldsymbol{\alpha}|} \right) h |\boldsymbol{\alpha}| + C_2 h^4 (t + h) |\boldsymbol{\alpha}|^5.$$

Коэффициент $C_1(\cdot)$ не является ограниченной функцией на единичной окружности. Коэффициент C_2 же на единичной окружности ограничен, поскольку множитель $1/e_x = |\boldsymbol{\phi}|/\phi_1$ в (12.5) компенсируется тем, что $A(\phi_1, \phi_2) = O(\phi_1)$ при $\phi_1 \rightarrow 0$.

Проведённые рассуждения не дали ответ на вопрос, каково оптимальное значение порядка точности в длительном счёте при формальном порядке точности, равном 1. Получим этот ответ при помощи утверждения 10.4. Пусть матрица A определена (12.2). Выпишем критерий (10.2) с подстановкой $C_1 = |\boldsymbol{\phi}|$,

$C_2 = |\phi|^{Q+1}$, $v = \epsilon$. Имеем

$$\exists \tilde{\delta} : \quad \epsilon^*(A^*A + |\phi|^{2Q})^{-1}A^*A\epsilon \leq \tilde{\delta}|\phi|^2\|\epsilon\|^2.$$

Подставляя выражение для матрицы A и упрощая, получаем

$$\exists c : \quad 2y^2 \leq c|\phi|^2(3x^2 + y^2 + |\phi|^{2Q}).$$

Подставляя выражения для x и y , отбрасывая члены, заведомо мажорируемые соседними, и опуская несущественные коэффициенты, получаем

$$\exists c : \quad \phi_1^2\phi_2^4 \leq c(\phi_1^2 + \phi_2^2)(\phi_1^4 + \phi_1^2\phi_2^4 + (\phi_1^2 + \phi_2^2)^Q).$$

Выражение $(x + y)^Q$ для неотрицательных x и y эквивалентно $x^Q + y^Q$, поэтому можно переписать критерий в виде

$$\exists c : \quad \phi_1^2\phi_2^4 \leq c(\phi_1^2 + \phi_2^2)(\phi_1^4 + \phi_1^2\phi_2^4 + \phi_1^{2Q} + \phi_2^{2Q}). \quad (12.6)$$

Нас интересует наибольшее значение Q , при котором это условие выполняется. При $Q = 2$ это верно: раскрыв скобки в правой части, можно увидеть в ней слагаемое $\phi_1^2\phi_2^4$, совпадающее с левой частью. Покажем при помощи утверждения 11.5, что неравенство (12.6) не выполняется при $Q > 2$. Действительно, в левой части неравенства находится ϕ^m , где $m = (2, 4)$, а в его правой части находятся слагаемые

$$\mathbf{n}_1 = (6, 0), \quad \mathbf{n}_2 = (4, 2), \quad \mathbf{n}_3 = (2, 2Q), \quad \mathbf{n}_4 = (2, 6), \quad \mathbf{n}_5 = (0, 2 + 2Q).$$

Остальные слагаемые мажорируются перечисленными. Изобразив точку m и точки \mathbf{n}_j на координатной плоскости, можно заметить, что при всех $Q > 2$ точка m лежит ниже выпуклой оболочки точек \mathbf{n} . Следовательно, неравенство (12.6) не выполняется.

Таким образом, оптимальное значение порядка точности в длительном счёте при 1-м формальном порядке равно 2. Можно записать оценку

$$\|\varepsilon_h(t, e^{i\alpha \cdot \mathbf{r}}, \overset{\circ}{\Pi}_h)\|_{(3)} \leq \min \left\{ C_1 \left(\frac{\alpha}{|\alpha|} \right) h|\alpha| + C_2 h^4(t + h)|\alpha|^5, Ch^2(t + h)|\alpha|^2 \right\}.$$

Теперь перейдём к рассмотрению системы для последних двух компонент сеточной функции. Блок матрицы $A(\phi)$, определяющий поведение последних двух компонент решения, зависит только от ϕ_1 . Этот блок имеет вид (11.6), где $a = 2 \sin(\phi_1/2)$, $b = -16 \sin^4(\phi_1/2)$. Выпишем выражение (11.8) для функции численной ошибки. Имеем

$$\begin{aligned} \hat{\varepsilon}(\phi, \nu, \overset{\circ}{\Pi}_h) &= S(e^{\nu\Lambda} - I)S^{-1}(\overset{\circ}{\Pi}_1 \exp(i\phi \cdot \mathbf{r}))_0 = \\ &= S \begin{pmatrix} e^{\nu(i\alpha + im)} - 1 & 0 \\ 0 & e^{\nu(i\alpha - im)} - 1 \end{pmatrix} \frac{1}{2} \begin{pmatrix} -1 & m/\nu \\ 1 & m/\nu \end{pmatrix} \begin{pmatrix} 1 \\ \exp(i\alpha_1 h) \end{pmatrix}, \end{aligned}$$

где $\nu = t/h$, $v = a + ib$, $\phi = \alpha h$, $m = |v|$. Поскольку S унитарная, домножение на S не оказывает влияния на норму ошибки. Вычислим произведение последних двух сомножителей. Имеем

$$2S^{-1}(\mathring{\Pi}_h \exp(i\alpha \cdot \mathbf{r}))_0 = \begin{pmatrix} -1 + e^{i\alpha_1 h m/v} \\ 1 + e^{i\alpha_1 h m/v} \end{pmatrix}. \quad (12.7)$$

Далее без ограничения общности будем считать, что $\phi_1 > 0$; если $\phi_1 < 0$, то все рассуждения проводятся аналогично с тем отличием, что две компоненты (12.7) обмениваются свойствами. При $\phi_1 > 0$ имеем $a/m = v/m = 1 + O(h^2)$, поэтому в некоторой окрестности $h = 0$ первая компонента (12.7) по модулю не превосходит $2\alpha_1 h$. При этом $|\exp(\nu(ia + im)) - 1| \leq 2$, следовательно, произведение не превосходит $4\alpha_1 h$. Вторая компонента (12.7) по модулю не превосходит 2, а

$$\begin{aligned} |\exp(\nu(ia - im)) - 1| &= |\exp(it(a - m)/h) - 1| \leq t|a - m|/h = \\ &= \frac{t}{h} \left| a - \sqrt{a^2 + b^2} \right| = \frac{t}{h} \frac{b^2}{a + \sqrt{a^2 + b^2}} \leq \frac{t}{h} \frac{b^2}{a} = \\ &= \frac{t}{h} \frac{256 \sin^8(\phi_1/2)}{2 \sin(\phi_1/2)} = 128 \frac{t}{h} \sin^7(\alpha_1 h/2) \leq th^6. \end{aligned}$$

Таким образом, 4-я и 5-я компоненты ошибки удовлетворяют оценке $O(h + th^6)$.

Складывая оценки ошибки для первых трёх и последних двух компонент решения, получаем оценку ошибки вида $O(h + th^4)$, зависящую от направления. Одновременно с этим, оценивая ошибку по первым трём компонентам напрямую исходя из свойств аппроксимации и устойчивости, получаем суммарную оценку $O(h + th^2)$, от направления не зависящую.

Разобранный пример показывает возможность такого случая, что на всех решениях в виде плоских волн имеется порядок точности P и порядок точности в длительном счёте Q , но при этом на $H_{per}^{Q+1}(\mathbb{R}^d)$ порядок точности в длительном счёте Q не имеет места.

13. Зависимость от скорости переноса

Выше мы полагали скорость переноса ω в уравнении (2.3) и схему (2.5) для аппроксимации этого уравнения заданными. Теперь зададимся вопросом, как свойства схемы зависят от ω (по-прежнему полагаем ω постоянным во времени и пространстве). При этом будем считать, что коэффициенты схемы линейно зависят от вектора ω , и, следовательно, рассматривать схемы вида

$$\sum_{\zeta \in \mathcal{S}} Z_\zeta \frac{du_{\eta+\zeta}}{dt}(t) + \frac{1}{h} \sum_{l=1}^d \omega_l \sum_{\zeta \in \mathcal{S}} L_\zeta^{(l)} u_{\eta+\zeta}(t) = 0. \quad (13.1)$$

Можно показать, что схема вида (13.1), устойчивая для всех $\omega \in \mathbb{R}^d$, удовлетворяет условию $\sum_{\eta} L_{\zeta}^{(l)} = 0$. В силу следствия 5.2, для такой схемы порядок точности не может превосходить порядок аппроксимации. В большинстве же используемых на практике схем используется информация о направлении характеристики, что делает зависимость схемы от ω кусочно-линейной. Коэффициенты L_{ζ} в (2.5) линейно зависят не только от компонент вектора ω , но и от величин вида $|\omega \cdot \mathbf{n}_m|$, где \mathbf{n}_m – некоторые векторы. При каждом фиксированном направлении $\omega/|\omega|$ зависимость коэффициентов от $|\omega|$ остаётся линейной.

Легко построить схему вида (13.1), такую что для некоторых направлений переноса имеют место более высокие порядки аппроксимации и точности, чем для остальных направлений. Для этого достаточно взять равномерную сетку, размер блока $|M^0| = 1$ и записать конечно-разностную схему, аппроксимировав производную по x с порядком Q_x , а производную по y – с порядком $Q_y > Q_x$. Тогда для всех направлений, кроме совпадающего с осью OY и противоположному ему, схема будет обладать порядком аппроксимации и порядком точности Q_x , а для двух выделенных направлений – порядком Q_y .

Более интересно, что возможна и обратная ситуация: для некоторых выделенных направлений порядок может быть меньшим, чем для остальных. Такой эффект ранее упоминался для метода Галёркина с разрывными базисными функциями, но не в точности схемы, а лишь в её оценках (см [8]), не являющихся оптимальными. В то же время этот эффект нами наблюдался на практике при оценке точности метода спектральных разностей [9]. Приведём простой пример, иллюстрирующий такую возможность.

Рассмотрим уравнение переноса (2.3) при $d = 2$ со скоростью переноса $(\omega_x, \omega_y)^T$, где $\omega_x, \omega_y \geq 0$. Пусть векторы \mathbf{a}_j совпадают с единичными ортами. Положим $M^0 = \{L, R\}$ и определим отображение

$$(\Pi_h f)_{\eta,L} = f(h\eta_x, h(\eta_y + 1)), \quad (\Pi_h f)_{\eta,R} = f(h(\eta_x + 1), h\eta_y).$$

Далее блоки будем нумеровать привычными индексами $i \equiv \eta_x, j \equiv \eta_y$. Производную для степени свободы R будем определять центральной разностной производной:

$$\frac{du_{i,j}^R}{dt} + \omega_x \frac{u_{i+1,j}^R - u_{i-1,j}^R}{2h} + \omega_y \frac{u_{i,j+1}^R - u_{i,j-1}^R}{h} = 0.$$

Производную для степени свободы L будем определять при помощи направленных разностей. Воспользуемся тем, что точка, где определяется значение L на ячейке $i, j - 1$, совпадает с точкой определения значения R на ячейке $i - 1, j$, и “завяжем” эти компоненты друг на друга:

$$\frac{du_{i,j}^L}{dt} + \omega_x \frac{u_{i,j}^L - u_{i-1,j}^L}{h} + \omega_y \frac{u_{i,j}^L - u_{i-1,j}^R}{2h} = 0.$$

В матричном виде эта схема примет вид (2.5), где

$$L_{-1,0} = \begin{pmatrix} -\omega_x & -\omega_y \\ 0 & -\omega_x/2 \end{pmatrix}, \quad L_{0,0} = \begin{pmatrix} \omega_x + \omega_y & 0 \\ 0 & 0 \end{pmatrix},$$

$$L_{1,0} = \begin{pmatrix} 0 & 0 \\ 0 & \omega_x/2 \end{pmatrix}, \quad L_{0,1} = \begin{pmatrix} 0 & 0 \\ 0 & \omega_y/2 \end{pmatrix}, \quad L_{0,-1} = \begin{pmatrix} 0 & 0 \\ 0 & -\omega_y/2 \end{pmatrix}.$$

Если обозначить $\phi = (\phi, \psi)$, то матрица $L(\phi)$ примет вид

$$L(\phi) = \begin{pmatrix} \omega_x(1 - e^{-i\phi}) + \omega_y & -\omega_y e^{-i\phi} \\ 0 & i(\omega_x \sin \phi + \omega_y \sin \psi) \end{pmatrix},$$

и $A(\phi) = i(\omega_x \phi + \omega_y \psi) - L(\phi)$.

Покажем, что эта схема устойчива при $\omega_x, \omega_y \geq 0$, то есть что $\|\exp(\nu A(\phi))\|$ ограничено по $\nu \geq 0$ равномерно по ϕ . Матрица $A(\phi)$ имеет вид

$$A(\phi) = \begin{pmatrix} \lambda_1(\phi) & a(\phi) \\ 0 & \lambda_2(\phi) \end{pmatrix},$$

при

$$\lambda_1(\phi) = i(\omega_x \phi + \omega_y \psi) - (\omega_x(1 - e^{-i\phi}) + \omega_y),$$

$$\lambda_2(\phi) = i(\omega_x \phi + \omega_y \psi) - i(\omega_x \sin \phi + \omega_y \sin \psi),$$

$a(\phi) = -\omega_y e^{-i\phi}$. Можно показать, что

$$e^{\nu A(\phi)} = \exp \left[\nu \begin{pmatrix} \lambda_1(\phi) & a(\phi) \\ 0 & \lambda_2(\phi) \end{pmatrix} \right] = \begin{pmatrix} e^{\nu \lambda_1(\phi)} & \frac{a(\phi)}{\lambda_2(\phi) - \lambda_1(\phi)} [e^{\nu \lambda_2(\phi)} - e^{\nu \lambda_1(\phi)}] \\ 0 & e^{\nu \lambda_2(\phi)} \end{pmatrix}.$$

При $\omega_x \geq 0, \omega_y \geq 0$ верно $\operatorname{Re} \lambda_2(\phi) = 0$ и $\operatorname{Re} \lambda_1(\phi) \leq 0$, поэтому диагональные компоненты $e^{\nu A(\phi)}$ ограничены. Величина в квадратных скобках при больших ν стремится к единице и, следовательно, тоже ограничена. Осталось показать, что ограниченной является величина

$$Y(\phi) = a(\phi) / (\lambda_2(\phi) - \lambda_1(\phi)).$$

Для проверки этого условия запишем

$$|Y(\phi)|^2 = \left| \frac{\omega_y}{i(\omega_x \sin \phi + \omega_y \sin \psi) - (\omega_x(1 - e^{-i\phi}) + \omega_y)} \right|^2 =$$

$$= \frac{\omega_y^2}{(\omega_x(1 - \cos \phi) + \omega_y)^2 + \omega_y^2 \sin^2 \psi} \leq 1,$$

поскольку в силу неотрицательности ω_x и ω_y минимум знаменателя достигается при $\phi = \psi = 0$. Таким образом, записанная схема при $\omega_x \geq 0, \omega_y \geq 0$ является устойчивой.

Выпишем теперь ошибку аппроксимации. Подставим в схему (2.5) функции $x^2/2, xy$ и $y^2/2$. Поскольку для компонент R используется схема 2-го порядка, производная на них точна и ошибка аппроксимации равна 0. Для компоненты L вычислим её явно. Для $f^{xx} \equiv -(\epsilon_1(x^2/2, \Pi_1))_0, f^{xy} \equiv -(\epsilon_1(xy, \Pi_1))_0, f^{yy} \equiv -(\epsilon_1(y^2/2, \Pi_1))_0$, получим

$$f^{xx} = h \begin{pmatrix} \omega_x/2 \\ 0 \end{pmatrix}, \quad f^{xy} = h \begin{pmatrix} -\omega_x/2 \\ 0 \end{pmatrix}, \quad f^{yy} = h \begin{pmatrix} -\omega_y/2 \\ 0 \end{pmatrix}.$$

Получим теперь уравнение на корректор. Заметим, что

$$L(0) = \begin{pmatrix} \omega_y & -\omega_y \\ 0 & 0 \end{pmatrix}.$$

При $\omega_y \neq 0$ матрица $L(0)$ имеет два собственных значения: 0 и $\omega_y \neq 0$. Поскольку нулевое собственное значение является простым, диагональные матрицы $\mathfrak{E}^{(m)}$ отображения $\tilde{\Pi}_h$ определяются однозначно с точностью до cI , где $c \in \mathbb{R}$, а I – единичная матрица. Отсюда следует, что $\mathfrak{E}^{(m)}, |\mathbf{m}| = 1$, можно полагать равными нулю. Коэффициенты второго порядка находятся из соотношений

$$L(0)\mathfrak{E}^{xx} = f^{xx}, \quad L(0)\mathfrak{E}^{xy} = f^{xy}, \quad L(0)\mathfrak{E}^{yy} = f^{yy}.$$

Все эти системы, очевидно, совместны. Таким образом, отображение $\tilde{\Pi}_h$ построено, что доказывает 2-й порядок точности рассматриваемой схемы (при $\omega_x \geq 0, \omega_y > 0$).

Если же $\omega_y = 0$, то $L(0) = 0$. Применяя следствие 5.2, получаем, что схема не может обладать порядком точности, превосходящим порядок аппроксимации, то есть её порядок точности равен 1.

Таким образом, построен пример схемы для уравнения (2.3) при $\omega_x \geq 0, \omega_y \geq 0$, такой что при $\omega_y \neq 0$ она обладает 1-м порядком аппроксимации и 2-м порядком точности, а при $\omega_y = 0$ – 1-м порядком аппроксимации и 1-м порядком точности.

14. Нерешенные задачи

Мы привели ряд примеров, показывающих, что в многомерном случае схемы могут обладать неожиданным поведением ошибки или странным образом откликаться на попытки исследовать это поведение методом введения вспомогательного отображения. Тем не менее, несколько вопросов пока остались неотвеченными. Приведём некоторые из них.

1. Может ли оптимальное значение порядка точности в длительном счёте для схемы вида (2.5) в смысле \dot{P}_h вида (2.9) не быть целым числом?
2. Пусть схема обладает порядком точности P и порядком точности в длительном счёте Q . Может ли существовать такое направление e_∞ , что при $e \rightarrow e_\infty$ коэффициенты $\mathfrak{C}_e^{(Q)}$, вычисленные исходя из условия Q -го порядка аппроксимации на волнах вида $f(\mathbf{r} \cdot \mathbf{e})$, неограниченно растут?
3. Множество направлений, по которым могут распространяться волны, не покрывает всю единичную сферу. Например, при смещениях $\mathbf{a}_1 = (1,0)^T$, $\mathbf{a}_2 = (0,1)^T$ справедливо $|e_y|/|e_x| \leq N_0/h$, где N_0 – период решения. Каковы могут быть оценки ошибки, учитывающие зависимость от N_0 ?
4. Рассмотрим разрывный метод Галёркина на тетраэдральной трансляционно-инвариантной сетке и предположим, что он обладает оценкой $O(h^P + th^Q)$. Вопрос: найдётся ли оператор вида (2.10), доставляющий Q -й порядок аппроксимации? Если вектор переноса совпадает с вектором одного из рёбер сетки, имеем квазиодномерный случай, и ответ положительный. Если он не лежит ни в одной из плоскостей граней сетки, имеем простой случай, и ответ также положительный. А если вектор переноса лежит вдоль грани, но не вдоль ребра?
5. И, наконец, главный вопрос, мотивировавший авторов на проведение настоящего исследования: существует ли метод, позволяющий установить порядок точности в длительном счёте для заданной схемы?

15. Заключение

Было показано, что для схемы, обладающей p -м порядком точности и q -м порядком точности в длительном счёте, найдётся отображение, отличающееся от исходного на величину $O(h^p)$ и доставляющее q -й порядок аппроксимации. Однако это отображение, вообще говоря, нелокально; были доказаны два достаточных условия, при которых оно является локальным. Был приведён ряд искусственных примеров схем, демонстрирующих неожиданное поведение в длительном счёте.

Список литературы

1. Бахвалов П. А., Сурначёв М. Д. О спектральном анализе схем для линейного уравнения переноса // Препринты ИПМ им. М.В.Келдыша. 2019. № 70. 28 с.
2. Бахвалов П. А., Сурначёв М. Д. Об аналитических семействах матриц, порождающих ограниченные полугруппы // Сибирский журнал вычислительной математики. (представлено в редакцию).

3. Бахвалов П. А., Сурначёв М. Д. Линейные схемы с несколькими степенями свободы для одномерного уравнения переноса // Препринты ИПМ им. М.В.Келдыша. 2019. № 73. 40 с.
4. Бахвалов П. А., Сурначёв М. Д. О приведении устойчивых матриц к блочно-диагональному виду // Препринты ИПМ им. М.В.Келдыша. 2019. № 71. 15 с.
5. Cao W., Zhang Z., Zou Q. Superconvergence of discontinuous Galerkin methods for linear hyperbolic equations // SIAM Journal on Numerical Analysis. 2014. Vol. 52, no. 5. P. 2555–2573.
6. Superconvergence of discontinuous Galerkin methods for two-dimensional hyperbolic equations / Cao W., Shu C.-W., Yang Y. et al. // SIAM Journal on Numerical Analysis. 2015. Vol. 53, no. 4. P. 1651–1671.
7. Kato T. Perturbation theory for linear operators. Grund. math. Wiss., B. 132, Springer, 1966. P. XIX, 592.
8. Richter G. An optimal-order error estimate for the discontinuous Galerkin method // Mathematics of Computation. 1988. Vol. 50. P. 75–88.
9. Balan A., May G., Schöberl J. A stable high-order Spectral Difference method for hyperbolic conservation laws on triangular elements // J. Comput. Phys. 2012. Vol. 231. P. 2359–2375.

Оглавление

1	Введение	3
2	Постановка задачи	4
3	Спектральное представление схемы	7
4	Общие замечания	10
5	Определение формального порядка	15
6	Существование отображения	16
7	Простой случай	18
8	Квазиодномерный случай	20
9	Контрпример	22
10	Особые направления волнового вектора	27
11	Минус-особое направление по точности в длительном счёте	30
12	Странное особое направление	33
13	Зависимость от скорости переноса	39
14	Нерешенные задачи	42
15	Заключение	43
	Список литературы	43