



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

Препринты ИПМ • Препринт № 123 за 2018 г.



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

Бахвалов П.А.

Метод нестационарного
корректора для анализа
точности линейных
полудискретных схем

Рекомендуемая форма библиографической ссылки: Бахвалов П.А. Метод нестационарного корректора для анализа точности линейных полудискретных схем // Препринты ИПМ им. М.В.Келдыша. 2018. № 123. 38 с. doi:[10.20948/prepr-2018-123](https://doi.org/10.20948/prepr-2018-123)
URL: <http://library.keldysh.ru/preprint.asp?id=2018-123>

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

Метод нестационарного корректора
для анализа точности
линейных полудискретных схем

Москва — 2018

Бахвалов П. А.

Метод нестационарного корректора для анализа точности линейных полудискретных схем

Метод нестационарного корректора предназначен для исследования превышения порядка точности над порядком аппроксимации и специфики поведения ошибки решения при длительном счёте. В настоящей работе он излагается применительно к линейным полудискретным схемам для решения уравнения переноса, что позволяет дать более простое изложение. Предлагается обобщение метода нестационарного корректора на случай схем с матрицей перед временными производными. С использованием полученного метода получаются новые оценки точности 4-точечной схемы R3.

Ключевые слова: геометрический корректор, аппроксимация и точность, неравномерная сетка, суперсходимость

Pavel Alexeevich Bakhvalov

Unsteady corrector method for accuracy analysis of linear semidiscrete schemes

Unsteady corrector method is used to find the order of accuracy when it is greater than the order of truncation error and investigate the long-time evolution of solution error. In this paper this method is applied to linear semidiscrete schemes, which allows to simplify its description. We also present its generalization for schemes with matrices next to temporal derivatives. Using this method we obtain new accuracy estimates for 4-point difference scheme R3.

Key words: geometric corrector, consistency and accuracy, non-uniform mesh, superconvergence

Оглавление

Введение	3
Постановка задачи	4
Разностные схемы	5
Нестационарный корректор	10
Нестационарный корректор и точность схемы	15
Блочное измельчение	21
Исследование 4-точечной схемы	27
Заключение	37
Список литературы	38

Введение

Хорошо известно, что на неравномерных и неструктурированных сетках порядок точности разностных схем может превосходить порядок аппроксимации. Впервые этот эффект был обнаружен А. Н. Тихоновым и А. А. Самарским [1] на примере уравнения конвекции-диффузии. Для уравнения переноса и гиперболических систем этот эффект наблюдается в расчётах по схемам разных классов, см., например, [2] и список литературы в ней. Однако теоретически обосновать его обычно удаётся только для некоторых конечно-элементных схем, в том числе для простейшей конечно-объёмной схемы [3, 4] и др.

Со сверхсходимостью связан другой эффект: различие в точности схемы при малом и большом временах счёта. Для метода Галёркина с разрывными базисными функциями в одномерном случае ошибка решения имеет оценку $O(h_{\max}^{k+1} + h_{\max}^{2k+1}t)$, то есть при коротком счёте доминирующий член ошибки имеет $(k+1)$ -й, а при длительном счёте – $(2k+1)$ -й порядок малости по шагу сетки, несмотря на k -й порядок аппроксимации. Этот эффект был достаточно давно замечен в численных расчётах (см., например, [5]), но теоретически его удалось доказать только несколько лет назад [6].

Одним из способов доказательства повышенной точности является поиск такого способа проецирования решения на расчётную сетку, в смысле которого схема обладала бы повышенным порядком аппроксимации. Для разрывного метода Галёркина этот метод применялся, например, в [6, 7], причём в [6] этот оператор искался в виде ряда по производным от точного решения.

В [8] автором был предложен метод нестационарного корректора, в основе которого лежит этот же принцип. В отличие от большинства работ по конечно-элементным схемам, опирающихся на глубокое знание природы численного метода, метод нестационарного корректора основывается на методе неопределённых коэффициентов и поэтому применим для широкого класса разностных схем.

Настоящая работа продолжает разработку метода нестационарного корректора. В отличие от [8], интегрирование по времени предполагается точным, что существенно упрощает описание и использование метода. Даются несколько его технических обобщений: допускается наличие матрицы в разностной схеме перед производной по времени и рассматривается случай более сложного оператора проецирования, включающего в себя L_2 -проекцию.

Также в настоящей работе метод нестационарного корректора применяется к анализу точности схемы R3. Доказывается, что отличие численного решения от точного оценивается величиной порядка $O(h_{\max}^2 + (\Delta h)_{\max}^2 t + h_{\max}^3 t)$. Ранее аналогичная оценка была получена автором в [9] применительно к стационарному уравнению.

Постановка задачи

В настоящей работе будем рассматривать уравнение переноса с постоянной скоростью в конечной области, допуская задание граничных условий 1-го рода и периодических граничных условий. Задачу с периодическими условиями будем интерпретировать как задачу в бесконечной области с решением, инвариантным относительно пространственных трансляций.

Пусть $\mathbf{a} \in \mathbb{R}^d$ – некоторый вектор, постоянный во времени и пространстве, а $G \subseteq \mathbb{R}^d$ – область с кусочно-гладкой границей. Обозначим $(\partial G)^- = \{\mathbf{r} \in \partial G : \mathbf{a} \cdot \mathbf{n} < 0\}$, где \mathbf{n} – внешняя нормаль к границе. Рассмотрим уравнение переноса

$$\frac{\partial u}{\partial t} + \mathbf{a} \cdot \nabla u = 0, \quad t \in (0, t_{\max}), \quad \mathbf{r} \in G, \quad (1)$$

с начальными условиями $u(0, \mathbf{r}) = u_0(\mathbf{r})$, $\mathbf{r} \in G$, и граничными условиями на входной границе $u(t, \mathbf{r}) = u_b(t, \mathbf{r})$, $\mathbf{r} \in (\partial G)^-$. Будем считать, что начальные и граничные условия таковы, что решение (1) является непрерывной функцией.

Будем предполагать, что существует такой набор линейно независимых векторов $\mathbf{b}_i^s \in \mathbb{R}^d$, $i = 1, \dots, \bar{d}$, где $0 \leq \bar{d} \leq d$, что расчётная область G (и, следовательно, ∂G^-), начальные данные u_0 и граничные условия u_b инвариантны относительно пространственных трансляций на векторы \mathbf{b}_i^s .

Схема (дискретная или полудискретная) численного решения уравнения (1) обычно заключается в следующем.

1. Определяется пространство \mathbb{R}^M (M – конечное или счётное), которое будем называть сеточным.
2. Начальные данные $u_0(\mathbf{r})$ проецируются на сеточное пространство при помощи некоторого оператора $\Pi : C(\bar{G}) \rightarrow \mathbb{R}^M$ (мы ограничиваемся случаем непрерывных решений; в общем случае Π может действовать, например, на $L_1^{loc}(G)$).
3. С использованием граничных условий, заданных на множестве $M^b \subset M$ при $0 < t < t_{\max}$, и начальных данных некоторым образом находится численное решение на момент времени t_{\max} .

Таким образом, схема сочетает в себе как алгоритм эволюции решения в сеточном пространстве (также обычно называемый схемой), так и оператор Π . Анализ точности этой схемы заключается в добавлении к пунктам 1–3 ещё одного:

4. Точное решение на момент времени t_{\max} проецируется на сетку при помощи оператора Π (см. п. 2), после чего эта проекция вычитается из численного решения. Полученная величина является ошибкой решения и измеряется в некоторой норме.

Отметим, что обычно под оператором проецирования (или проектором) понимается оператор, действующий из пространства на его подпространство и удовлетворяющий свойству $\Pi^2 = \Pi$. В этом смысле оператор Π и другие операторы, которые в настоящей работе будут называться проекторами, вообще говоря, проекторами не являются.

Строго говоря, если M счётно, то на пространстве \mathbb{R}^M невозможно определить никакую норму. Если рассматривать множество ограниченных числовых последовательностей

$$V_\infty = \{f \in \mathbb{R}^M : \sup_j |f_j| < \infty\}, \quad (2)$$

то норму определить можно, например, как $\|f\| = \sup_j |f_j|$. Однако в этой норме может быть сложно доказать устойчивость (или получить оптимальную оценку на константу устойчивости, если схема обладает слабой неустойчивостью). Поэтому интерес могут представлять также и интегральные нормы, определить которые на V_∞ затруднительно. Чтобы избежать этих трудностей, в настоящей работе будем рассматривать конечномерные пространства \mathbb{R}^{M^s} , $M^s \subseteq M$, $|M^s| < \infty$, функций, периодических по \bar{d} направлениям. На \mathbb{R}^{M^s} уже можно определить, в том числе, интегральные и энергетические нормы. Все оценки, которые будут получены в настоящей работе, равномерны по периоду решения и $|M^s|$. Отметим, что при каждом фиксированном M^s различные нормы на \mathbb{R}^{M^s} , конечно, будут эквивалентными, но для констант, оценивающих нормы друг через друга, уже не будет равномерных оценок при $|M^s| \rightarrow \infty$.

Разностные схемы

Всюду в настоящей работе будет использоваться величина h , имеющая смысл максимального линейного размера сеточного элемента. При изложении теории для схемы общего вида она не будет определяться, как не будет определяться и понятие расчётной сетки. Поэтому формально её следует рассматривать как одну из констант, встречающихся в определениях. Физический смысл она будет приобретать только при анализе конкретных разностных схем.

Опишем класс схем, которые будут рассматриваться в настоящей работе. Начнём с операторов, проецирующих непрерывные функции на сеточное пространство.

Определение 1. Будем называть оператор $\Pi : C(\bar{G}) \rightarrow \mathbb{R}^M$ допустимым, если он представляется в виде

$$(\Pi(f))_j = \int_{\mathbb{R}^d} f(\mathbf{r}) d\mu_j, \quad (3)$$

где μ_j – некоторая, вообще говоря, знакопеременная мера, удовлетворяющая условиям нормировки,

$$\int_{\mathbb{R}^d} d\mu_j = 1, \quad j \in M, \quad (4)$$

конечности вариации

$$\sup_{j \in M} \int_{\mathbb{R}^d} |d\mu_j| = C_\mu < \infty. \quad (5)$$

Кроме того, её носитель которой должен лежать в некотором шаре радиуса $C_r h/2$:

$$\text{supp} \mu_j \subseteq \bar{B}_{C_r h/2}(\mathbf{c}_j) \cap \bar{G}, \quad j \in M, \quad (6)$$

а значения в $j \in M^b$ должны определяться только граничными условиями:

$$\text{supp} \mu_j \subseteq (\partial G)^-, \quad j \in M^b. \quad (7)$$

Далее, чтобы избежать лишних скобок в формулах, будем обозначать компоненты вектора $\Pi(f)$ через $\Pi_j(f)$, $j \in M$. Также введём обозначение

$$\mathbf{r}_j = \Pi_j(\mathbf{r}). \quad (8)$$

Из определения 1 напрямую вытекают следующие свойства допустимых операторов проецирования:

$$\begin{aligned} \Pi_j(1) &= 1, \quad j \in M; \\ \text{supp} \mu_j &\subseteq \bar{B}_{C_r h}(\mathbf{r}_j), \quad j \in M; \\ \Pi_j(f) &\leq C_\mu \sup |f(\mathbf{r})|, \quad j \in M. \end{aligned}$$

Примерами допустимых операторов Π на квазиравномерной сетке являются взятие точечного значения в сеточном узле или центре масс сеточной ячейки, взятие интегрального среднего по ячейке, взятие значений в точках коллокации L_2 -проекции на некоторое пространство кусочно-полиномиальных функций.

Теперь перейдём к описанию эволюции по времени. Пусть Π – некоторый допустимый оператор. Рассмотрим сеточное уравнение общего вида

$$\sum_{k \in M} Z_{jk} \frac{du_k}{dt} + \sum_{k \in M} L_{jk} u_k = 0, \quad j \in M \setminus M^b, \quad (9)$$

с действительными коэффициентами Z_{jk} и L_{jk} и граничными условиями

$$u_j(t) = \Pi_j(u_b), \quad j \in M^b. \quad (10)$$

Поскольку Π – допустимый, определение (10) корректно, так как в силу (7) значение $\Pi_j(f)$ определяется только значениями f на границе.

Определение 2. Будем называть $\mathfrak{S} = (M, M^b, \Pi, Z, L)$, где $M^b \subset M$, $\Pi : C(\tilde{G}) \rightarrow \mathbb{R}^M$, $Z = \{Z_{jk}\}$, $L = \{L_{jk}\}$, $j, k \in M$, допустимой схемой, если выполняются следующие условия:

- 1) Π – допустимый оператор;
- 2) $Z(V_\infty) = V_\infty$, где V_∞ определено (2);
- 3) Z обратим на V_∞ , то есть $\forall y \in V_\infty \exists! x \in V_\infty : Zx = y$;
- 4) $Z_{jk} = \delta_{jk}$ при $j \in M^b$; $L_{jk} = 0$ при $j \in M^b$. Здесь и далее δ – символ Кронекера.
- 5) имеет место точность на константе

$$\sum_{k \in M} L_{jk} = 0, \quad j \in M \setminus M^b; \quad (11)$$

- 6) существует константа C_s , такая что

$$\begin{aligned} \sum_{k \in M} |L_{jk}| |\mathbf{r}_k - \mathbf{r}_j|^m &\leq (C_s)^{m+1} h^{m-1}, \quad m \geq 0, \quad j \in M \setminus M^b; \\ \sum_{k \in M} |Z_{jk}| |\mathbf{r}_k - \mathbf{r}_j|^m &\leq (C_s)^{m+1} h^m, \quad m \geq 0, \quad j \in M \setminus M^b. \end{aligned} \quad (12)$$

Условие 6) фактически является ограничением на локальность шаблона. Отметим, что в условии 3) требуется обратимость $Z : V_\infty \rightarrow V_\infty$, что является более сильным условием, чем обратимость его ограничения на множество периодических функций с периодом \mathbb{R}^{M^s} при рассматриваемом периоде M^s .

Легко установить следующее.

Утверждение 1. Пусть $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ удовлетворяет условиям 1)–5) определения 2. Пусть коэффициенты L_{jk} и Z_{jk} , $j \in M \setminus M^b$, $k \in M$, удовлетворяют следующим условиям:

- 1) справедливы оценки $|L_{jk}| \leq C_1/h$ и $|Z_{jk}| \leq C_1$;
- 2) количество ненулевых коэффициентов в каждой строке матриц L и Z конечно: $\forall j \in M \setminus M^b \ |\{k : L_{jk} \neq 0 \text{ or } Z_{jk} \neq 0\}| \leq C_2$;
- 3) операторы L и Z локализованы в области диаметром порядка h :
 $L_{jk} = Z_{jk} = 0$ при $|\mathbf{r}_k - \mathbf{r}_j| > C_3 h$.

Тогда выполняется условие б) определения 2, причём $C_s = C_3 \max\{C_1 C_2, 1\}$, и поэтому \mathfrak{S} является допустимой.

Чтобы решение разностной задачи при задании периодических граничных условий действительно было периодическим, нужна периодичность схемы, которой оно будет вычисляться. Поэтому введём следующие определения.

Определение 3. Будем называть допустимую схему $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ периодической по \bar{d} направлениям, $\bar{d} \leq d$, с периодами \mathbf{b}_i , $i = 1, \dots, \bar{d}$, и разностным периодом M^0 , если верны следующие условия.

1. M представляется в виде $M = M^0 \times \mathbb{Z}^{\bar{d}}$. Далее для $j \in M$ будем использовать обозначение $j = [\xi, \boldsymbol{\eta}]$, где $\xi \in M^0$ – индекс неизвестной внутри блока, а $\boldsymbol{\eta} \in \mathbb{Z}^{\bar{d}}$ – индекс блока.
2. Если $M^b \neq \emptyset$, граничные условия равномерно распределены по блокам: $M^b = M^{b,0} \times \mathbb{Z}^{\bar{d}}$.
3. Операторы Z и L являются периодическими, то есть для любых $\xi_1, \xi_2 \in M^0$, $\boldsymbol{\eta}_1, \boldsymbol{\eta}_2, \boldsymbol{\vartheta} \in \mathbb{Z}^{\bar{d}}$ справедливо

$$Z_{[\xi_1, \boldsymbol{\eta}_1 + \boldsymbol{\vartheta}] [\xi_2, \boldsymbol{\eta}_2 + \boldsymbol{\vartheta}]} = Z_{[\xi_1, \boldsymbol{\eta}_1] [\xi_2, \boldsymbol{\eta}_2]},$$

$$L_{[\xi_1, \boldsymbol{\eta}_1 + \boldsymbol{\vartheta}] [\xi_2, \boldsymbol{\eta}_2 + \boldsymbol{\vartheta}]} = L_{[\xi_1, \boldsymbol{\eta}_1] [\xi_2, \boldsymbol{\eta}_2]}.$$

4. Π является периодическим, то есть мера $\mu_{[\xi, \boldsymbol{\eta} + \boldsymbol{\vartheta}]}$, входящая в определение допустимого оператора Π , является пространственной трансляцией меры $\mu_{[\xi, \boldsymbol{\eta}]}$ на вектор $\sum_{i=1}^{\bar{d}} \vartheta^i \mathbf{b}_i$, где ϑ^i – компонента вектора $\boldsymbol{\vartheta}$.

При этом будем называть допустимый оператор \mathcal{P} периодическим относительно \mathfrak{S} , если он периодический с теми же периодами \mathbf{b}_i и M^0 .

Очевидно, что периодичность Π равносильна тому, что для всех $f \in C(\bar{G})$ и всех $\xi \in M^0$, $\boldsymbol{\eta}_1 = \{\eta_1^i, i = 1, \dots, \bar{d}\}$, $\boldsymbol{\eta}_2 = \{\eta_2^i, i = 1, \dots, \bar{d}\}$ выполняется

$$\int f \left(\mathbf{r} - \sum_{i=1}^{\bar{d}} \eta_1^i \mathbf{b}_i \right) d\mu_{[\xi, \boldsymbol{\eta}_1]} = \int f \left(\mathbf{r} - \sum_{i=1}^{\bar{d}} \eta_2^i \mathbf{b}_i \right) d\mu_{[\xi, \boldsymbol{\eta}_2]}.$$

Пусть $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ – допустимая схема, периодическая по \bar{d} направлениям с периодами b_i и разностным периодом M^0 . По определению $j \in M$ представляются в виде $j = [\xi, \eta]$, $\xi \in M^0$, $\eta \in \mathbb{Z}^{\bar{d}}$. Пусть $N \in \mathbb{N}$. Очевидно, что \mathfrak{S} также является периодической с периодом $b_i^s = Nb_i$, $i = 1, \dots, \bar{d}$, и разностным периодом $M^s = M^0 \times \{0, \dots, N-1\}^{\bar{d}}$, причём $j \in M$ теперь представляются в виде $j = [(\xi, \zeta), \theta]$, где компоненты вектора θ являются результатом, а ζ – остатком от деления η на N .

Определим операторы \hat{Z} и \hat{L} , действующие из \mathbb{R}^{M^0} в \mathbb{R}^{M^0} , и операторы \check{Z} и \check{L} , действующие из \mathbb{R}^{M^s} в \mathbb{R}^{M^s} , равенствами

$$\hat{Z}_{\xi\xi'} = \sum_{\eta \in \mathbb{Z}^{\bar{d}}} Z_{[\xi, 0][\xi', \eta]}, \quad \hat{L}_{\xi\xi'} = \sum_{\eta \in \mathbb{Z}^{\bar{d}}} L_{[\xi, 0][\xi', \eta]}. \quad (13)$$

$$\check{Z}_{(\xi, \zeta)(\xi', \zeta')} = \sum_{\theta \in \mathbb{Z}^{\bar{d}}} Z_{[(\xi, \zeta), 0][(\xi', \zeta'), \theta]}, \quad \check{L}_{(\xi, \zeta)(\xi', \zeta')} = \sum_{\theta \in \mathbb{Z}^{\bar{d}}} L_{[(\xi, \zeta), 0][(\xi', \zeta'), \theta]}. \quad (14)$$

В частности, при $\xi \in M^b \cap M^0$ выполняется $\hat{Z}_{\xi\xi'} = \delta_{\xi\xi'}$, $\check{Z}_{(\xi, \zeta)(\xi', \zeta')} = \delta_{\xi\xi'} \delta_{\zeta\zeta'}$.

Будем говорить, что вектор $f \in \mathbb{R}^M$ периодичен с периодом M^0 , если $f_{[\xi, \eta_1]} = f_{[\xi, \eta_2]}$ для всех $\xi \in M^0$, $\eta_1, \eta_2 \in \mathbb{Z}^{\bar{d}}$. Далее будем отождествлять такие вектора с соответствующими элементами \mathbb{R}^{M^0} . Аналогичное определение даётся для M^s .

Утверждение 2. Пусть $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ – допустимая периодическая схема. Пусть V_∞^0 и V_∞^s – множества векторов из V_∞ , периодических с периодами M^0 и M^s соответственно. Тогда существуют операторы $\hat{Z}^{-1} : V_\infty^0 \rightarrow V_\infty^0$ и $\check{Z}^{-1} : V_\infty^s \rightarrow V_\infty^s$, такие что $\hat{Z}^{-1}\hat{Z} = \hat{Z}\hat{Z}^{-1} = I_0$, $\check{Z}^{-1}\check{Z} = \check{Z}\check{Z}^{-1} = I_s$, где I_0 и I_s – тождественные операторы на V_∞^0 и V_∞^s соответственно.

Докажем утверждение для операторов с крышками; для операторов с рогами оно доказывается аналогично. Рассмотрим систему $\hat{Z}x = y$, $x, y \in V_\infty^0$. Сопоставим ей систему $Zx = y$, $x \in V_\infty$, $y \in V_\infty^0 \subseteq V_\infty$. По условию допустимости \mathfrak{S} у этой системы будет существовать единственное решение $Z^{-1}y$. Покажем, что оно будет лежать в V_∞^0 . Действительно, пусть T_ϑ – оператор, сопоставляющий вектору $g \in \mathbb{R}^M$ вектор $(T_\vartheta g)_{[\xi, \eta]} = g_{[\xi, \eta + \vartheta]}$. Из определения периодической схемы легко убедиться, что для всех g верно $ZT_\vartheta g = T_\vartheta Zg$. Подставим в это равенство $g = Z^{-1}y$. Поскольку $y = Zg$ является периодическим, $T_\vartheta Zg = Zg$, следовательно, $ZT_\vartheta g = Zg$. Поскольку $g \in V_\infty$, решение системы $Zx = g$ существует и единственно. Следовательно, для всех T_ϑ выполняется $T_\vartheta g = g$. Это и означает, что $g = Z^{-1}y \in V_\infty^0$. А поскольку для векторов $g \in V_\infty^0$ справедливо $Zg = \hat{Z}g$, получаем, что $\hat{Z}g = y$.

Следующее утверждение показывает, что периодические (в смысле определения 3) допустимые схемы действительно позволяют решать задачи с периодическими условиями, то есть задача с множеством переменных $M = M^s \times \mathbb{Z}^{\bar{d}}$ сводится к задаче с множеством переменных M^s .

Утверждение 3. Пусть $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ – допустимая периодическая схема с периодами \mathbf{b}_i и разностным периодом M^0 . Пусть начальные данные $u_0(\mathbf{r})$ и граничные условия $u_b(t, \mathbf{r})$ согласованы друг с другом и периодичны с периодами $\mathbf{b}_i^s = N\mathbf{b}_i$. Тогда существует единственное решение $u \in V_\infty^s$ задачи

$$\sum_{k \in M^s} \check{Z}_{jk} \frac{du_k}{dt} + \sum_{k \in M^s} \check{L}_{jk} u_k = 0, \quad j = (\xi, \zeta), \quad \xi \in M^0 \setminus M^{b,0}, \quad (15)$$

$$u_j(t) = \Pi_j(u_b(t, \cdot)), \quad j = (\xi, \zeta), \quad \xi \in M^{b,0},$$

с начальными условиями $u_j(0) = \Pi_j(u_0)$. При этом функция $v \in V_\infty$, определённая $v_{[(\xi, \zeta), \eta]} = u_{(\xi, \zeta)}$, является решением задачи (9)–(10) при тех же начальных условиях.

Нестационарный корректор

В настоящем разделе будем рассматривать операторы $\tilde{\Pi}$ общего вида, проецирующие гладкое временное сечение функции на сеточное пространство. При этом способ проецирования может гладким образом меняться во времени, то есть $\tilde{\Pi} : [0, t_{\max}] \times C^q(\bar{G}) \rightarrow \mathbb{R}^M$.

Определение 4. Пусть $u \in C^{q+1}([0, t_{\max}] \times \bar{G})$ – точное решение дифференциальной задачи (1) при некоторых начальных и граничных условиях. Ошибкой аппроксимации схемы $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ в момент времени t на функции $u(t, \mathbf{r})$ в смысле оператора $\tilde{\Pi}$ будем называть вектор $\epsilon(t, u, \tilde{\Pi}) = \{\epsilon_j(t, u, \tilde{\Pi}), j \in M\}$ с компонентами

$$\epsilon_j(t, u, \tilde{\Pi}) = \sum_{k \in M} Z_{jk} \frac{d}{dt} \left(\tilde{\Pi}(t) u(t, \cdot) \right)_k + \sum_{k \in M} L_{jk} \left(\tilde{\Pi}(t) u(t, \cdot) \right)_k \quad (16)$$

при $j \in M \setminus M^b$ и $\epsilon_j(t, u, \tilde{\Pi}) = 0$ при $j \in M^b$. Будем называть схему \mathfrak{S} точной на функции $u(t, \mathbf{r})$, являющейся решением уравнения (1), если при всех $j \in M$ и $t \in (0, t_{\max})$ выполняется $\epsilon_j(t, u, \tilde{\Pi}) = 0$. Будем называть схему \mathfrak{S} точной на полиноме порядка p (например, константе, линейной, квадратичной функции) в смысле оператора $\tilde{\Pi}$, если она точна на всех соответствующих функциях u , являющихся решениями уравнения (1).

Пусть $m = (m_1, \dots, m_d)$ – мультииндекс: $m_i \geq 0$, $|m| = m_1 + \dots + m_d$, $m! = m_1! \dots m_d!$. Под обозначением $l \leq m$ будем понимать, что для $i = 1, \dots, d$ верно $l_i \leq m_i$. Под обозначением $l < m$ будем понимать, что $l \leq m$ и хотя бы для одного $i = 1, \dots, d$ верно $l_i < m_i$. Введём обозначения

$$\mathbf{r}^m = x_1^{m_1} \dots x_d^{m_d}, \quad D^m = \frac{\partial^{|m|}}{\partial x_1^{m_1} \dots \partial x_d^{m_d}}.$$

Будем говорить, что непустое множество мультииндексов \mathfrak{M} является *монотонным*, если $\forall m \in \mathfrak{M} \forall l < m \ l \in \mathfrak{M}$; положим $|\mathfrak{M}| = \max_{m \in \mathfrak{M}} |m|$. Введём обозначения $\mathfrak{M}_m = \{l : l \leq m\}$, $\mathfrak{M}_{< m} = \{l : l < m\}$. Символами P_m и $P_{< m}$ будем обозначать пространство многочленов вида, соответственно, $\sum_{l \leq m} \alpha_l (\mathbf{r} - \mathbf{a}t)^l$ и $\sum_{l < m} \alpha_l (\mathbf{r} - \mathbf{a}t)^l$, где $\alpha_l \in \mathbb{R}$. Очевидно, что при $l < m$ выполняется $P_l \subset P_m$, $P_{< l} \subset P_{< m}$.

Определение 5. Пусть $\Pi, \mathcal{P} : C(\bar{G}) \rightarrow \mathbb{R}^M$ – некоторые допустимые операторы. Пусть \mathfrak{M} – монотонное множество мультииндексов, и $\mathfrak{C}_j^m(t)$, $m \in \mathfrak{M} \setminus \{0\}$, $j \in M$, – некоторые функции времени. Тогда будем говорить, что оператор $\tilde{\Pi}_{\mathfrak{M}} : [0, t_{\max}] \times C^{|\mathfrak{M}|}(\bar{G}) \rightarrow \mathbb{R}^M$, определяемый равенством

$$\left(\tilde{\Pi}_{\mathfrak{M}}(t)f \right)_j = \Pi_j(f) + \sum_{m \in \mathfrak{M} \setminus \{0\}} \mathfrak{C}_j^m(t) \mathcal{P}_j(D^m f), \quad (17)$$

порождён операторами Π , \mathcal{P} и набором функций $\mathfrak{C}_j^m(t)$, $m \in \mathfrak{M} \setminus \{0\}$.

Из определения 5, в частности, следует $\tilde{\Pi}_0(t) \equiv \Pi$. Далее для краткости будем пользоваться обозначениями $\Pi_q = \Pi_{\mathfrak{M}_q}$ и $\Pi_{< q} = \Pi_{\mathfrak{M}_{< q}}$.

В большинстве случаев оператор \mathcal{P} выбирается совпадающим с Π , но формулируемые ниже утверждения справедливы и в том случае, если они различаются. В формулировках некоторых утверждений оператор \mathcal{P} будем опускать.

Определение 6. Пусть $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ – допустимая схема, \mathcal{P} – допустимый оператор, \mathfrak{M} – монотонное множество мультииндексов. Будем называть набор коэффициентов $\{\mathfrak{C}_j^m(t), j \in M, m \in \mathfrak{M} \setminus \{0\}\}$, цепочкой нестационарных корректоров для схемы \mathfrak{S} с использованием \mathcal{P} , если

- для всех $m \in \mathfrak{M}$ и $f \in P_m$ схема \mathfrak{S} точна на функции f в смысле $\tilde{\Pi}_m$ (17), порождённого Π , \mathcal{P} и $\{\mathfrak{C}^l\}$;
- $\mathfrak{C}_j^m(t) = 0$, $j \in M^b$, $m \in \mathfrak{M} \setminus \{0\}$.

Набор коэффициентов $\mathfrak{C}^m(t) = \{\mathfrak{C}_j^m(t), j \in M\}$ при фиксированном мультииндексе m будем называть нестационарным корректором. Будем называть его ограниченным, если при всех $t \in [0, t_{\max}]$ выполняется $\sup_j |\mathfrak{C}_j^m(t)| < \infty$.

Утверждение 4. Пусть для всех t зафиксированы начальные условия $\mathfrak{C}^m(0) \in V_\infty$, согласованные с граничными условиями: $\mathfrak{C}_j^m(0) = 0$, $j \in M^b$. Тогда ограниченные нестационарные корректоры \mathfrak{C}^m для схемы $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ находятся однозначно из последовательного решения системы уравнений

$$\sum_{k \in M} Z_{jk} \frac{d\mathfrak{C}_k^m}{dt} + \sum_{k \in M} L_{jk} \mathfrak{C}_k^m = \epsilon_j(t, f_m, \tilde{\Pi}_{< m}), \quad j \in M \setminus M^b, \quad (18)$$

$$\mathfrak{C}_j^m(t) = 0, \quad j \in M^b, \quad (19)$$

где $\epsilon(t, f_m, \tilde{\Pi}_{< m})$ – ошибка аппроксимации схемы \mathfrak{S} на функции

$$f_m(t, \mathbf{r}) = -\frac{1}{m!} (\mathbf{r} - \mathbf{r}_0 - \mathbf{a}(t - t_0))^m \quad (20)$$

в смысле $\tilde{\Pi}_{< m}$, причём она не зависит от выбора значений \mathbf{r}_0 и t_0 в f_m .

Доказательство. Предположим по индукции, что все корректоры $\mathfrak{C}^l \in [0, t_{\max}] \times V_\infty$, $l < m$, определены однозначно, и покажем это для корректора \mathfrak{C}^m .

Прежде всего, заметим, что $\epsilon_j(t, f_m, \tilde{\Pi}_{< m})$ не зависит от выбора \mathbf{r}_0 и t_0 в (20). Действительно, рассмотрим две функции $f_m^{(1)}(t, \mathbf{r})$ и $f_m^{(2)}(t, \mathbf{r})$, различающиеся выбором \mathbf{r}_0 и t_0 . Тогда $f_m^{(1)}(t, \mathbf{r}) - f_m^{(2)}(t, \mathbf{r})$ лежит в $P_{< m}$. Отсюда по определению 6 справедливо $\epsilon_j(t, f_m^{(1)} - f_m^{(2)}, \tilde{\Pi}_{< m}) = 0$.

Далее нужно отметить, что $\epsilon(t, f_m, \tilde{\Pi}_{< m}) \in V_\infty$. Для доказательства этого факта удобно зафиксировать $j \in M$, момент времени t и выбрать $\mathbf{r}_0 = \mathbf{r}_j$ и $t_0 = t$. Тогда все члены, входящие в формулу для $\epsilon_j(t, f_m, \tilde{\Pi}_{< m})$, будут содержать некоторые коэффициенты корректоров \mathfrak{C}^l , ограниченные по предположению индукции, и величины вида $\Pi_k D^{m-l}(\mathbf{r} - \mathbf{r}_j)^l$, сумма модулей которых ограничена по определению допустимости схемы.

Пусть f_m – произвольная функция вида (20). По определению 4

$$\epsilon_j(t, f_m, \tilde{\Pi}_m) = \sum_{k \in M} Z_{jk} \frac{d}{dt} \left(\tilde{\Pi}_m(t) f_m(t, \cdot) \right)_k + \sum_{k \in M} L_{jk} \left(\tilde{\Pi}_m(t) f_m(t, \cdot) \right)_k.$$

Комбинируя формулы (17) и (20), имеем

$$\left(\tilde{\Pi}_m(t) f_m(t, \cdot) \right)_k = \left(\tilde{\Pi}_{< m}(t) f_m(t, \cdot) \right)_k - \mathfrak{C}_k^m(t).$$

Отсюда

$$\epsilon_j(t, f_m, \tilde{\Pi}_m) = \epsilon_j(t, f_m, \tilde{\Pi}_{<m}) - \sum_{k \in M} Z_{jk} \frac{d}{dt} \mathfrak{E}_k^m(t) - \sum_{k \in M} L_{jk} \mathfrak{E}_k^m(t). \quad (21)$$

Предположим, что $\{\mathfrak{E}^l, l \leq m\}$ – цепочка нестационарных корректоров. Тогда по определению 6 справедливо $\epsilon_j(t, f_m, \tilde{\Pi}_m) = 0$, и поэтому в силу (21) выполняется (18). Обратно, пусть выполняется (18)–(19). Поскольку V_∞ нормируемо, на нём система ОДУ (18)–(19) имеет единственное решение. Покажем, что оно является нестационарным корректором. Рассмотрим произвольную функцию $f \in P_m$, и пусть f_m – некоторая функция вида (20). Заметим, что f представима в виде $f = \alpha f_m + g$, $g \in P_{<m}$. Поэтому

$$\epsilon_j(t, f, \tilde{\Pi}_m) = \alpha \epsilon_j(t, f_m, \tilde{\Pi}_m) + \epsilon_j(t, g, \tilde{\Pi}_{<m}).$$

Первое слагаемое равно нулю в силу (21) и (18), а второе – по предположению индукции. Это заканчивает доказательство.

Утверждение 5. Пусть m – мультииндекс. Пусть $\{\mathfrak{E}^l, 0 < l < m\}$ – цепочка корректоров для схемы \mathfrak{S} , не зависящих от времени. Тогда $\epsilon_j(t, f_m, \tilde{\Pi}_{<m})$ также не зависит от времени.

Обозначим через $f_m^{(\tau)}(t, \mathbf{r})$ функцию, определенную (20) с подстановкой $t_0 = \tau$. Заметим что для произвольного не зависящего от времени $\tilde{\Pi}$ справедливо $\epsilon_j(t, f_m^{(t_0)}, \tilde{\Pi}) = \epsilon_j(0, f_m^{(t_0-t)}, \tilde{\Pi})$. Тогда $\epsilon_j(t, f_m^{(t_0)}, \tilde{\Pi}_{<m}) = \epsilon_j(0, f_m^{(t_0-t)}, \tilde{\Pi}_{<m})$, что, в свою очередь, не зависит от выбора параметра u функции f_m в силу утверждения 4.

Легко установить следующее.

Утверждение 6. Пусть $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ – допустимая периодическая схема и $\{\mathfrak{E}^m(t), m \in \mathfrak{M} \setminus \{0\}\}$ – ограниченная цепочка нестационарных корректоров для \mathfrak{S} с использованием \mathcal{P} , периодичного относительно \mathfrak{S} . Предположим, что в начальный момент времени цепочка периодическая, то есть для всех m и всех ξ, η_1 и η_2 справедливо $\mathfrak{E}_{[\xi, \eta_1]}^m(0) = \mathfrak{E}_{[\xi, \eta_2]}^m(0)$. Тогда она в любой момент времени остаётся периодической:

$$\mathfrak{E}_{[\xi, \eta_1]}^m(t) = \mathfrak{E}_{[\xi, \eta_2]}^m(t). \quad (22)$$

Отметим, что в утверждении 6 не фигурирует период дифференциальной задачи. Таким образом, нестационарные корректоры периодичны с периодом, равным периоду схемы. Они находятся из системы (18)–(19), которую достаточно решать при $j \in M^0$ и дополнить алгебраическими вида соотношениями (22).

Таким образом, систему уравнений (18)–(19) нахождение нестационарного корректора можно записать в операторном виде

$$\hat{Z} \frac{d\mathfrak{C}^m}{dt} + \hat{L}\mathfrak{C}^m = \epsilon(t, f_m, \tilde{\Pi}_{<m}), \quad (23)$$

причём при $\xi \in M^{b,0}$ правая часть доопределяется как $\epsilon_\xi(t, f_m, \tilde{\Pi}_{<m}) = 0$.

Предположим, что схема \mathfrak{S} не точна на функции f_m (20) в смысле оператора проецирования Π , но точна на всех функциях f_l при $l < m$. Тогда нестационарные корректоры \mathfrak{C}^l , $l < m$, могут быть выбраны тождественно равными нулю, а нестационарный корректор \mathfrak{C}^m отличен от нуля.

Определение 7. *Тождественно не равный нулю набор коэффициентов \mathfrak{C}^m , такой что $\mathfrak{C}^l \equiv 0$ для всех $l < m$, будем называть главным корректором.*

В силу утверждения 4 коэффициенты главного корректора подчиняются системе (18)–(19), причём правая часть в (18) равна

$$\epsilon_j(t, f_m, \tilde{\Pi}_{<m}) = \epsilon_j(0, f_m, \Pi) = \sum_{k \in M} Z_{jk} \Pi_k \left(\frac{\partial f_m}{\partial t}(0, \cdot) \right) + \sum_{k \in M} L_{jk} \Pi_k f_m(0, \cdot)$$

и не зависит от времени.

При определении нестационарного корректора мы оставили свободу выбора начальных условий. Один из естественных способов задать их – положить $\mathfrak{C}^m(0) = 0$ для всех m . Однако в некоторых случаях этот выбор не является оптимальным.

Предположим, что система уравнений (18)–(19) для нахождения главного корректора \mathfrak{C}^m (без наложения начальных условий) имеет стационарное решение $\mathfrak{C}^{m,\infty}$. Тогда можно отказаться от зависимости от времени и положить $\mathfrak{C}^m(t) \equiv \mathfrak{C}^{m,\infty}$. Если случилось так, что для некоторого мультииндекса n все корректоры \mathfrak{C}^l , $l < n$, стационарны, то в силу утверждения 5 правая часть уравнений (18) не будет зависеть от времени. Если у неё будет стационарное решение $\mathfrak{C}^{n,\infty}$, то снова можно положить $\mathfrak{C}^n(t) \equiv \mathfrak{C}^{n,\infty}$, и так далее.

Нестационарный корректор и точность схемы

Далее будем предполагать, что период схемы M^0 конечен.

Определение 8. Пусть $\| \cdot \|$ – некоторая норма на множестве \mathbb{R}^{M^s} периодических сеточных функций. Будем называть неубывающую функцию $K(t)$ константой устойчивости для периодической схемы $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ в этой норме, если периодическое решение уравнения

$$\begin{aligned} \sum_{k \in M} Z_{jk} \frac{du_k}{dt} + \sum_{k \in M} L_{jk} u_k &= 0, \quad j \in M \setminus M^b, \\ u_j(t) &= 0, \quad j \in M^b, \end{aligned} \quad (24)$$

с периодическими начальными условиями $u_j(0) = u_j^0$, $j \in M$, такими что $u_j^0 = 0$ при $j \in M^b$, удовлетворяет оценке

$$\|u(t)\| \leq K(t) \|u^0\|.$$

Особо отметим, что если устойчивость понимать в смысле определения 8, то константа устойчивости не зависит от выбора проектора Π .

Хорошо известно, что устойчивость по правой части является следствием устойчивости по начальным данным.

Утверждение 7. Пусть $\| \cdot \|$ – некоторая норма в \mathbb{R}^{M^s} и $K(t)$ – константа устойчивости для схемы $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ в ней. Тогда периодическое решение

$$\begin{aligned} \sum_{k \in M} Z_{jk} \frac{du_k}{dt} + \sum_{k \in M} L_{jk} u_k &= f_j(t), \quad j \in M \setminus M^b, \\ u_j(t) &= 0, \quad j \in M^b, \end{aligned} \quad (25)$$

с начальными условиями $u_j(0) = u_j^0$, $j \in M$, при выполнении условий совместности $u_j^0 = 0$ и $f_j(t) = 0$ при $j \in M^b$, удовлетворяет оценке

$$\|u(t)\| \leq K(t) \left(\|u^0\| + \int_0^t \|\check{Z}^{-1} f(\tau)\| d\tau \right). \quad (26)$$

Действительно, пусть $u(t, u^0, 0)$ – периодическое решение однородной задачи (24) с начальными данными u^0 . Непосредственной подстановкой можно убедиться, что решение $u(t, u^0, f)$ неоднородной задачи (25) имеет вид

$$u(t, u^0, f) = u(t, u^0, 0) + \int_0^t u(t - \tau, \check{Z}^{-1} f(\tau), 0) d\tau.$$

Отсюда с учётом неубывания функции $K(t)$ следует искомая оценка (26).

Следующая теорема лежит в основе метода нестационарного корректора.

Теорема 1. Пусть выполнены следующие условия:

- 1) решение $u(\mathbf{r}, t)$ уравнения (1) q раз дифференцируемо и имеет липшицевы q -е пространственные производные с константой Липшица \mathbb{L} , и $\mathbb{L}_m(t) = \sup_{\mathbf{r} \in G} |D^m u(t, \mathbf{r})| < \infty$ для $m \leq q$;
- 2) $\mathfrak{S} = (M, M^b, \Pi, Z, L)$ – допустимая схема с константами C_s, C_r, C_μ, h , периодическая с периодами $\mathbf{b}_i^s, i = 1, \dots, \bar{d}$, и $M^0 = M^s$;
- 3) \mathcal{P} – допустимый оператор с константами C_r, C_μ, h , периодический относительно \mathfrak{S} ;
- 4) $\|\cdot\|$ и $\|\cdot\|_*$ – нормы в \mathbb{R}^{M^s} , такие что для $e = (1, \dots, 1)^T$ выполняется $\|e\| = \|e\|_* = 1$, для всех $f \in \mathbb{R}^{M^s}$ верно $\|f\| \leq C_n \|f\|_*$, и если $a, b \in \mathbb{R}^{M^s}$ такие, что $|a_j| \leq |b_j| \forall j \in M^s$, то $\|a\| \leq C_n \|b\|$;
- 5) если $a, b \in \mathbb{R}^{M^s}$ удовлетворяют условию

$$|a_j| \leq \max_{k: |\check{Z}_{jk}| + |\check{L}_{jk}| \neq 0} |b_k|, \quad j \in M^s,$$

то выполняется $\|a\|_* \leq C_w \|b\|_*$;

- 6) $\|\check{Z}^{-1}\|_* \leq C_z < \infty, \|\check{L}\|_* \leq C_l/h$;
- 7) $K(t)$ – константа устойчивости \mathfrak{S} в $\|\cdot\|$;
- 8) $\{\mathfrak{E}^m(t)\}, 0 < |m| \leq q$, – цепочка нестационарных корректоров для схемы \mathfrak{S} с использованием \mathcal{P} с некоторыми начальными данными $\mathfrak{E}^m(0)$, такими что $\mathfrak{E}_{[\xi, \eta]}^m(0) = \mathfrak{E}_{[\xi, 0]}^m(0)$ и $\mathfrak{E}_j^m(0) = 0$ при $j \in M^b$, и $\tilde{\Pi}_q$ – оператор, порождённый Π, \mathcal{P} и этой цепочкой корректоров;
- 9) $u_h(t)$ – решение по схеме \mathfrak{S} , а $\dot{u}_h(t) = \{u_j(t), j \in M\}$ – решение по схеме $\mathring{\mathfrak{S}} = (M, M^b, \mathring{\Pi}, Z, L)$, где $\mathring{\Pi} : C^q(\bar{G}) \rightarrow \mathbb{R}^M$ – некоторый проектор, периодический относительно \mathfrak{S} (например, Π или $\tilde{\Pi}_q$).

Тогда существует \tilde{C} , зависящая только от $q, C_s, C_l, C_r, C_n, C_w, C_z$ и C_μ , такая что выполняются оценки

$$\begin{aligned} \|\dot{u}_h(t) - \tilde{\Pi}_q(t)u(t, \cdot)\| &\leq K(t) \|\tilde{\Pi}_q(0)u(0, \cdot) - \mathring{\Pi}u(0, \cdot)\| + \\ &+ tK(t)\tilde{C}\mathbb{L}h^q + K(t)\tilde{C}\mathbb{L} \sum_{0 < |m| \leq q} h^{q-|m|} \int_0^t \|\mathfrak{E}^m(\tau)\|_* d\tau, \end{aligned} \quad (27)$$

$$\|u_h(t) - \Pi u(t, \cdot)\| \leq \tilde{C}K(t)\mathbb{L} \sum_{0 < |m| \leq q} h^{q-|m|} \int_0^t \|\mathfrak{E}^m(\tau)\|_* d\tau + \quad (28)$$

$$+ C_n C_\mu \sum_{0 < |m| \leq q} (\|\mathfrak{E}^m(t)\| \mathbb{L}_m(t) + K(t) \|\mathfrak{E}^m(0)\| \mathbb{L}_m(0)) + \tilde{C}tK(t)\mathbb{L}h^q.$$

Отметим, что u_h и \dot{u}_h являются решением одной и той же системы уравнений (9)–(10) с разными начальными данными.

Доказательство. Прежде всего, установим существование такой константы \tilde{C} , зависящей только от $q, C_s, C_l, C_r, C_n, C_w, C_z$ и C_μ , такой что для всех решений $u(t, \mathbf{r})$ уравнения (1), имеющих q -е липшицевы пространственные производные с константой \mathbb{L}_{q+1} , выполняется оценка

$$\left\| \epsilon(t, u, \tilde{\Pi}_q) \right\|_* \leq \tilde{C} \mathbb{L}_{q+1} \left(h^q + \sum_{0 < |m| \leq q} h^{q-|m|} \|\mathfrak{e}^m(t)\|_* \right), \quad (29)$$

Доказательство этого утверждения проведём индукцией по q . Рассмотрим произвольную $j \in M \setminus M^b$ и произвольный момент времени τ . Представим функцию $u(t, \mathbf{r})$ в виде

$$u(t, \mathbf{r}) = p_{j,\tau}(t, \mathbf{r}) + g_{j,\tau}(t, \mathbf{r}),$$

где $p_{j,\tau}(t, \mathbf{r})$ – полином порядка q от t и \mathbf{r} , представляющий собой первые $q + 1$ членов разложения функции $u(t, \mathbf{r})$ в ряд Тейлора около точки (τ, \mathbf{r}_j) :

$$p_{j,\tau}(t, \mathbf{r}) = \sum_{0 \leq |m| \leq q} \frac{1}{m!} (\mathbf{r} - \mathbf{r}_j - \mathbf{a}(t - \tau))^m D^m u(\tau, \mathbf{r}_j).$$

Этот полином также является решением (1). Поскольку p – полином порядка q , то $g_{j,\tau} = u - p_{j,\tau}$ является q раз дифференцируемой функцией с липшицевыми q -ми производными с константой Липшица $2\mathbb{L}$. Кроме того, $D^m g_{j,\tau}(\tau, \mathbf{r}_j) = 0$ при $|m| \leq q$. Следовательно, при $|m| \leq q$ имеет место оценка

$$|D^m g_{j,\tau}(\tau, \mathbf{r})| \leq 2\mathbb{L} |\mathbf{r} - \mathbf{r}_j|^{q+1-|m|}, \quad (30)$$

а поскольку $g_{j,\tau}$ – решение (1), при $|m| < q$ также справедливо

$$\left| D^m \frac{\partial g_{j,\tau}}{\partial t}(\tau, \mathbf{r}) \right| \leq \sum_{|i|=1} |\mathbf{a}^i| |D^{m+i} g_{j,\tau}(\tau, \mathbf{r})| \leq 2\mathbb{L} d |\mathbf{a}| |\mathbf{r} - \mathbf{r}_j|^{q-|m|}. \quad (31)$$

Здесь d – размерность пространства. Отсюда в силу допустимости оператора \mathcal{P} и неравенства треугольника следует

$$|(\mathcal{P} D^m g_{j,\tau}(\tau, \cdot))_k| \leq 2\mathbb{L} C_\mu (|\mathbf{r}_k - \mathbf{r}_j| + C_r h)^{q+1-|m|}.$$

$$\left| \left(\mathcal{P} D^m \frac{\partial g_{j,\tau}}{\partial t}(\tau, \cdot) \right)_k \right| \leq 2\mathbb{L} C_\mu d |\mathbf{a}| (|\mathbf{r}_k - \mathbf{r}_j| + C_r h)^{q-|m|}.$$

Просуммируем эти выражения с весом $|L_{jk}|$ и $|Z_{jk}|$ соответственно:

$$\begin{aligned} \sum_{k \in M} |L_{jk}| |(\mathcal{P}D^m g_{j,\tau}(\tau, \cdot))_k| &\leq 2\mathbb{L}C_\mu \sum_{k \in M} |L_{jk}| (|\mathbf{r}_k - \mathbf{r}_j| + C_r h)^{q+1-|m|}. \\ \sum_{k \in M} |Z_{jk}| |(\mathcal{P}D^m g_{j,\tau}(\tau, \cdot))_k| &\leq 2\mathbb{L}C_\mu \sum_{k \in M} |Z_{jk}| (|\mathbf{r}_k - \mathbf{r}_j| + C_r h)^{q+1-|m|}. \\ \sum_{k \in M} |Z_{jk}| \left| \left(\mathcal{P}D^m \frac{\partial g_{j,\tau}}{\partial t}(\tau, \cdot) \right)_k \right| &\leq 2\mathbb{L}C_\mu d|\mathbf{a}| \sum_{k \in M} |Z_{jk}| (|\mathbf{r}_k - \mathbf{r}_j| + C_r h)^{q-|m|}. \end{aligned}$$

Представим выражения в правых частях в виде сумм одночленов и применим для каждого из них условие (12) допустимости разностной схемы. Таким образом, существует такая константа \tilde{C} , зависящая от q , C_s и C_r , что

$$\begin{aligned} \sum_{k \in M} |L_{jk}| |(\mathcal{P}D^m g_{j,\tau}(\tau, \cdot))_k| &\leq 2\mathbb{L}C_\mu \tilde{C} h^{q-|m|}. \\ \sum_{k \in M} |Z_{jk}| |(\mathcal{P}D^m g_{j,\tau}(\tau, \cdot))_k| &\leq 2\mathbb{L}C_\mu \tilde{C} h^{q-|m|+1}. \\ \sum_{k \in M} |Z_{jk}| \left| \left(\mathcal{P}D^m \frac{\partial g_{j,\tau}}{\partial t}(\tau, \cdot) \right)_k \right| &\leq 2\mathbb{L}C_\mu \tilde{C} h^{q-|m|}. \end{aligned} \tag{32}$$

Оценки (32) остаются справедливыми при подстановке Π вместо \mathcal{P} , поскольку Π является допустимым с теми же константами C_r , C_μ , h .

Поскольку по определению нестационарного корректора q -го порядка в смысле $\tilde{\Pi}_q$ схема \mathfrak{S} точна на полиномах порядка q , то $\epsilon(\tau, p_{j,\tau}, \tilde{\Pi}_q) = 0$ и $\epsilon(\tau, u, \tilde{\Pi}_q) = \epsilon(\tau, g_{j,\tau}, \tilde{\Pi}_q)$. По определению $\epsilon_j(\tau, g_{j,\tau}, \tilde{\Pi}_q)$

$$\begin{aligned} \epsilon_j(\tau, u, \tilde{\Pi}_q) &= \epsilon_j(\tau, g_{j,\tau}, \tilde{\Pi}_q) = \\ &= \sum_{k \in M} Z_{jk} \frac{d(\tilde{\Pi}_q(t) g_{j,\tau}(t, \cdot))_k}{dt}(\tau) + \sum_{k \in M} L_{jk} (\tilde{\Pi}_q(t) g_{j,\tau}(\tau, \cdot))_k. \end{aligned}$$

Подставим определение $\tilde{\Pi}_q$ (17):

$$\begin{aligned} \epsilon_j(\tau, u, \tilde{\Pi}_q) &= \sum_{k \in M} Z_{jk} \left(\Pi \frac{\partial g_{j,\tau}}{\partial t}(\tau, \cdot) \right)_k + \sum_{k \in M} L_{jk} (\Pi g_{j,\tau}(\tau, \cdot))_k + \\ &+ \sum_{0 < |m| \leq q} \left[\sum_{k \in M} L_{jk} \mathfrak{e}_k^m(\tau) (\mathcal{P}D^m g_{j,\tau}(\tau, \cdot))_k \right. \\ &+ \left. \sum_{k \in M} Z_{jk} \left(\frac{d\mathfrak{e}_k^m}{dt} (\mathcal{P}D^m g_{j,\tau}(\tau, \cdot))_k + \mathfrak{e}_k^m(\tau) \left(\mathcal{P}D^m \frac{\partial g_{j,\tau}}{\partial t}(\tau, \cdot) \right)_k \right) \right]. \end{aligned} \tag{33}$$

Подставляя теперь оценки (32) в (33), получаем

$$\begin{aligned} & \left| \epsilon_j(\tau, u, \tilde{\Pi}_q) \right| \leq 4\mathbb{L}C_\mu \tilde{C} h^q + \\ & + 4\mathbb{L}C_\mu \tilde{C} \sum_{0 < |m| \leq q} h^{q-|m|} \max_{k: |\check{Z}_{jk}| + |\check{L}_{jk}| \neq 0} \left(|\mathfrak{e}_k^m(\tau)| + h \left| \frac{d\mathfrak{e}_k^m}{dt}(\tau) \right| \right). \end{aligned} \quad (34)$$

Оценка (34) получена при $j \in M \setminus M^b$; при $j \in M^b$ по определению $\epsilon_j(\tau, u, \tilde{\Pi}_q) = 0$, и поэтому (34) также выполняется. Согласно утверждению 6, $\mathfrak{e}^m(\tau)$ является периодическим, поскольку по условию он является таковым в начальный момент времени. Следовательно, правая часть (34) лежит в \mathbb{R}^{M^s} . Применим теперь к этому неравенству условие 5) настоящей теоремы. Учитывая, что норма вектора из единиц по условию 4) равна 1, получаем

$$\begin{aligned} & \left\| \epsilon(\tau, u, \tilde{\Pi}_q) \right\|_* \leq 4\mathbb{L}C_\mu \tilde{C} h^q + \\ & + 4\mathbb{L}C_\mu \tilde{C} C_w \sum_{0 < |m| \leq q} h^{q-|m|} \left(\left\| \mathfrak{e}^m(\tau) \right\|_* + h \left\| \frac{d\mathfrak{e}^m}{dt}(\tau) \right\|_* \right). \end{aligned} \quad (35)$$

Согласно утверждению (4), нестационарные корректоры подчиняются системе (18)–(19). По условию $M^0 = M^s$, поэтому $\hat{L} = \check{L}$. Следовательно,

$$\left\| \frac{d\mathfrak{e}^m}{dt}(\tau) \right\|_* \leq \left\| \check{Z}^{-1} \right\|_* \left(\left\| \check{L} \right\|_* \left\| \mathfrak{e}^m(\tau) \right\|_* + \left\| \epsilon(\tau, f_m, \tilde{\Pi}_{< m}) \right\|_* \right).$$

Далее, функция f_m является решением (1), причём константы Липшица её пространственных производных порядка $|m| - 1$ не превосходят 1. По предположению индукции для $q = |m| - 1$ имеем

$$\left\| \epsilon(\tau, f_m, \tilde{\Pi}_{< m}) \right\|_* = \left\| \epsilon(\tau, f_m, \tilde{\Pi}_{|m|-1}) \right\|_* \leq \tilde{C} h^{|m|-1} + \tilde{C} \sum_{0 < |r| \leq |m|-1} h^{|m|-1-|r|} \left\| \mathfrak{e}^r(\tau) \right\|_*.$$

Поскольку по условию $\left\| \check{Z}^{-1} \right\|_* \leq C_z$ и $\left\| \check{L} \right\|_* \leq C_l/h$, получаем

$$\left\| \frac{d\mathfrak{e}^m}{dt}(\tau) \right\|_* \leq \frac{C_z}{h} \left(C_l \left\| \mathfrak{e}^m(\tau) \right\|_* + \tilde{C} h^{|m|} + \tilde{C} \sum_{0 < |r| \leq |m|-1} h^{|m|-|r|} \left\| \mathfrak{e}^r(\tau) \right\|_* \right).$$

С учётом этой оценки из (35) следует утверждение индукции.

Теперь представим численное решение в виде

$$u_j(t) = \left(\tilde{\Pi}_q(t)u(t, \cdot) \right)_j + \varepsilon_j(t). \quad (36)$$

Подставляя выражение (36) в схему (9), получим уравнение для ε :

$$\sum_k Z_{jk} \frac{d\varepsilon_k}{dt} + \sum_k L_{jk} \varepsilon_k = \varepsilon_j(\tau, u, \tilde{\Pi}_q). \quad (37)$$

Поскольку начальное значение численного решения было получено оператором $\tilde{\Pi}$, имеем

$$\varepsilon_j(0) = \left(\dot{\Pi}u(0, \cdot) \right)_j - \left(\tilde{\Pi}_q(0)u(0, \cdot) \right)_j. \quad (38)$$

В силу утверждения 7 справедлива оценка

$$\|u_h - \tilde{\Pi}_q(t)u(t, \cdot)\| = \|\varepsilon(t)\| \leq K(t)\|\varepsilon(0)\| + K(t) \int_0^t \|\check{Z}^{-1}\varepsilon(\tau, u, \tilde{\Pi}_q)\| d\tau.$$

Поскольку по условию 4) для всех f верно $\|f\| \leq C_n \|f\|_*$, получаем

$$\|u_h - \tilde{\Pi}_q(t)u(t, \cdot)\| \leq K(t)\|\varepsilon(0)\| + K(t)C_n C_z \int_0^t \|\varepsilon(\tau, u, \tilde{\Pi}_q)\|_* d\tau.$$

Подставляя оценку (29) для $\|\varepsilon(\tau, u, \tilde{\Pi}_q)\|_*$, получаем искомую оценку (27).

Далее, по неравенству треугольника

$$\|u_h(t) - \Pi u(t, \cdot)\| \leq \|u_h(t) - \tilde{\Pi}_q(t)u(t, \cdot)\| + \|\tilde{\Pi}_q(t)u(t, \cdot) - \Pi u(t, \cdot)\|. \quad (39)$$

Первое слагаемое в правой части оценивается по (27). Оценим второе. По определению $\tilde{\Pi}_q$ имеем

$$\begin{aligned} & \left| \left(\tilde{\Pi}_q(t)u(t, \cdot) \right)_k - (\Pi u(t, \cdot))_k \right| \leq \\ & \leq \sum_{0 < |m| \leq |q|} |\mathfrak{e}_k^m(t)| |(\mathcal{P}D^m u(t, \cdot))_k| \leq \sum_{0 < |m| \leq |q|} |\mathfrak{e}_k^m(t)| C_\mu \mathbb{L}_m. \end{aligned}$$

По условию 4) отсюда следует

$$\left\| \tilde{\Pi}_q(t)u(t, \cdot) - \Pi u(t, \cdot) \right\| \leq C_n C_\mu \sum_{0 < |m| \leq |q|} \mathbb{L}_m \|\mathfrak{e}^m(t)\|.$$

Подставляя это неравенство в (39) и пользуясь им же при $t = 0$ для оценки первого слагаемого, даваемого (27), получаем искомую оценку (28).

Блочное измельчение

Теорема 1 позволяет оценить ошибку решения через величины нестационарных корректоров на фиксированной сетке. Однако с теоретической точки зрения конкретное значение ошибки неинформативно; интерес представляет поведение ошибки на последовательности измельчающихся сеток. При произвольном измельчении провести такой анализ для схемы общего вида принципиально невозможно, поскольку на каждой новой сетке операторы Z и L будут мало похожи на предыдущие. Однако если сетка задана во всём пространстве \mathbb{R}^d , то её измельчение можно осуществлять преобразованием гомотетии относительно начала координат: радиус-векторы узлов будут преобразовываться по правилу $\mathbf{r}_j \rightarrow (h/h_0)\mathbf{r}_j$. Если сетка периодическая, а число $h_0/h = K$ целое, то на эту процедуру можно взглянуть иначе: возьмём один сеточный блок (период), уменьшим его в K раз по линейному размеру вдоль каждой оси и заполним всё пространство полученными блоками. При необходимости удалим дублируемые сеточные примитивы, лежащие на склеиваемых гранях кубических блоков. Будем называть такое измельчение блочным.

Пример блочного измельчения приведён на рис. 1.

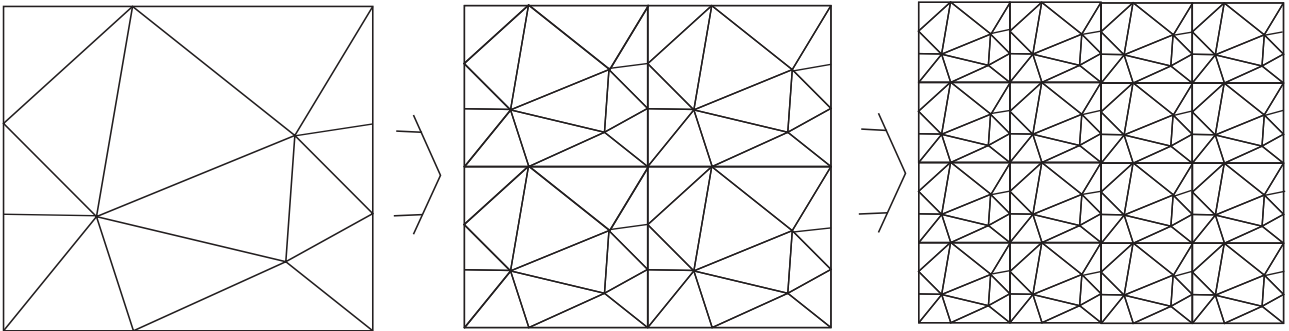


Рис. 1. Блочное измельчение

Дадим блочному измельчению формальное определение.

Определение 9. Рассмотрим допустимую схему $\mathfrak{S} = (M, M^b, \Pi, Z, L)$, где $M^b = \emptyset$. Пусть h_0 – произвольное положительное число. Будем называть её блочным измельчением семейство схем \mathfrak{S}_h , характеризуемых параметром $h > 0$, таких что $\mathfrak{S}_h = (M, \emptyset, \Pi_h, Z_h, L_h)$,

$$(Z_h)_{jk} = Z_{jk}, \quad (L_h)_{jk} = \frac{h_0}{h} L_{jk}, \quad j, k \in M, \quad (40)$$

$$(\Pi_h(f))_j = \int_{\mathbb{R}^d} f\left(\frac{h}{h_0} \mathbf{r}\right) d\mu_j, \quad (41)$$

где μ_j – мера, входящая в оператор проецирования Π .

Легко убедиться, что если $\mathfrak{S} = (M, \emptyset, \Pi, Z, L)$ является допустимой схемой с константами C_μ, C_r, C_s, h_0 , то при любом h схема \mathfrak{S}_h является допустимой с константами C_μ, C_r, C_s, h . Если \mathfrak{S} является периодической схемой с периодами \mathbf{b}_i и M^0 , то при любом h схема \mathfrak{S}_h является периодической с периодами $\mathbf{b}_i h/h_0$ и M^0 . Если $h = h_0/N$, то она также является периодической с кратными периодами \mathbf{b}_i и $M^0 \times (0, \dots, N-1)^{\bar{d}}$. В противном случае она позволяет решать задачи с другим периодом в решении.

Всюду далее будем предполагать, что \mathcal{P} периодичен с тем же набором векторов \mathbf{b}_i , что и Π , а при блочном измельчении он удовлетворяет соотношению, аналогичному (41).

Поскольку при блочном измельчении множество M^0 не меняется, то \hat{Z}_h и \hat{L}_h принадлежат тому же пространству, что и \hat{Z} и \hat{L} . Более того, справедливо $\hat{Z}_h = \hat{Z}$, $\hat{L}_h = (h_0/h)\hat{L}$, и уравнение на нахождение нестационарного корректора (23) приобретает вид

$$\hat{Z} \frac{d\mathfrak{C}_h^m}{dt} + \frac{h_0}{h} \hat{L} \mathfrak{C}_h^m = \epsilon_h(t, f_m, \tilde{\Pi}_{h, < m}), \quad (42)$$

где ϵ_h задаётся определением 4 для схемы \mathfrak{S}_h . Исследование уравнения (42) по сравнению с (23) является более интересным, поскольку позволяет исследовать поведение нестационарных корректоров \mathfrak{C}_h^m при $h \rightarrow 0$ и, следовательно, делать выводы о свойствах ошибки решения при измельчении сетки.

Далее всюду будем считать, что схема \mathfrak{S} периодическая по \bar{d} направлениям; случай непериодической схемы соответствует $\bar{d} = 0$.

Основным результатом, который удаётся получить для схемы общего вида при блочном измельчении, является автомодельность корректоров.

Докажем вначале вспомогательное утверждение.

Утверждение 8. Пусть $\mathfrak{S}_h = (M, \emptyset, \Pi_h, Z_h, L_h)$ – блочное измельчение некоторой допустимой схемы \mathfrak{S} . Пусть \mathfrak{M} – монотонное множество мультииндексов, и $\tilde{\Pi}_h : C^{|\mathfrak{M}|}(\mathbb{R}^d) \rightarrow \mathbb{R}^M$ – семейство проекторов вида

$$\left(\tilde{\Pi}_h f\right)_j = (\Pi_h f)_j + \sum_{l \in \mathfrak{M} \setminus \{0\}} (\mathfrak{C}_h^l(t))_j (\mathcal{P}_h(D^l f))_j, \quad (43)$$

где $\mathfrak{C}_h^l(t)$ – цепочка функций (не обязательно корректоров), таких что

$$\left(\mathfrak{C}_h^l(t)\right)_j = \left(\frac{h}{h_0}\right)^{|l|} \left(\mathfrak{C}_{h_0}^l\left(t \frac{h_0}{h}\right)\right)_j,$$

а \mathcal{P}_h связаны с \mathcal{P} соотношением, аналогичным (41). Пусть m – произвольный мультииндекс, и $f_m(t, \mathbf{r})$ определена (20) с параметрами t_0 и \mathbf{r}_0 . Тогда

$$\epsilon_h(t, f_m, \tilde{\Pi}_h) = \left(\frac{h}{h_0}\right)^{|m|-1} \epsilon_{h_0}\left(t \frac{h_0}{h}, f_m^0, \tilde{\Pi}_{h_0}\right), \quad (44)$$

где ϵ_h даётся определением 4 для \mathfrak{S}_h , а $f_m^0(t, \mathbf{r})$ определена (20) с параметрами $t_0 h_0/h$ и $\mathbf{r}_0 h_0/h$.

Отметим, что если \mathfrak{C}^l не являются корректорами, то значения ϵ_h вообще говоря, зависят как от t , как и от выбора констант t_0 и \mathbf{r}_0 в функции f_m .

Доказательство этого утверждения примитивно и заключается в последовательном применении всех его условий. По условию функция $f_m(t, \mathbf{r})$ определяется (20), поэтому для всех $l \leq m$

$$D^l f_m(t, \mathbf{r}) = -\frac{1}{(m-l)!} (\mathbf{r} - \mathbf{r}_0 - \mathbf{a}(t - t_0))^{m-l},$$

$$\frac{\partial D^l f_m}{\partial t}(t, \mathbf{r}) = \sum_{|i|=1, i \leq m-l} \mathbf{a}^i \frac{1}{(m-l-i)!} (\mathbf{r} - \mathbf{r}_0 - \mathbf{a}(t - t_0))^{m-l-i}.$$

Для функции f_m^0 , аналогично,

$$\begin{aligned} D^l f_m^0\left(\frac{h_0}{h}t, \frac{h_0}{h}\mathbf{r}\right) &= -\frac{1}{(m-l)!} (\mathbf{r} - \mathbf{r}_0 - \mathbf{a}(t - t_0))^{m-l} \left(\frac{h_0}{h}\right)^{|m-l|}, \\ \frac{\partial D^l f_m^0}{\partial t}\left(\frac{h_0}{h}t, \frac{h_0}{h}\mathbf{r}\right) &= \frac{h}{h_0} \frac{d}{dt} \left[D^l f_m^0\left(\frac{h_0}{h}t, \frac{h_0}{h}\mathbf{r}\right) \right] = \\ &= \sum_{|i|=1, i \leq m-l} \mathbf{a}^i \frac{1}{(m-l-i)!} (\mathbf{r} - \mathbf{r}_0 - \mathbf{a}(t - t_0))^{m-l-i} \left(\frac{h_0}{h}\right)^{|m-l|-1}. \end{aligned}$$

По определению блочного измельчения проекторы Π_h и Π_{h_0} связаны соотношением (41), откуда получаем

$$\begin{aligned} (\Pi_h D^l f_m(t, \cdot))_k &= \left(\frac{h}{h_0}\right)^{|m-l|} \left(\Pi_{h_0} D^l f_m^0\left(\frac{h_0}{h}t, \cdot\right)\right)_k, \\ \left(\Pi_h \frac{\partial D^l f_m}{\partial t}(t, \cdot)\right)_k &= \left(\frac{h}{h_0}\right)^{|m-l|-1} \left(\Pi_{h_0} \frac{\partial D^l f_m^0}{\partial t}\left(\frac{h_0}{h}t, \cdot\right)\right)_k. \end{aligned} \quad (45)$$

Оценки вида (45) справедливы и при подстановке \mathcal{P} вместо Π . Теперь воспользуемся определением 4, в которое подставим проектор (43).

$$\begin{aligned} &(\epsilon_h(t, f_m, \tilde{\Pi}_h))_j = \\ &= \sum_{k \in M} (Z_h)_{jk} \left[\left(\Pi_h \left(\frac{\partial f_m}{\partial t}(t, \cdot)\right)\right)_k + \sum_{0 < |l| \leq q} (\mathfrak{E}_h^l)_k(t) \left(\mathcal{P}_h \left(\frac{\partial D^l f_m}{\partial t}(t, \cdot)\right)\right)_k \right] + \\ &+ \sum_{k \in M} (L_h)_{jk} \left[(\Pi_h f_m(t, \cdot))_k + \sum_{0 < |l| \leq q} (\mathfrak{E}_h^l)_k(t) (\mathcal{P}_h D^l f_m(t, \cdot))_k \right]. \end{aligned}$$

Остаётся выразить все значения, зависящие от h , через аналогичные значения при $h = h_0$, для чего воспользоваться условием $\mathfrak{E}_h^l(t) = (h/h_0)^{|l|} \mathfrak{E}_{h_0}^l(th_0/h)$, свойствами проектора (45) и свойствами операторов схемы (40). Имеем

$$\begin{aligned} &(\epsilon_h(t, f_m, \tilde{\Pi}_h))_j = \\ &= \left(\frac{h}{h_0}\right)^{|m|-1} \sum_{k \in M} Z_{jk} \left[\left(\Pi_{h_0} \left(\frac{\partial f_m^0}{\partial t}\right)\right)_k + \sum_{0 < |l| \leq q} (\mathfrak{E}_{h_0}^l)_k \left(\mathcal{P}_{h_0} \left(\frac{\partial D^l f_m^0}{\partial t}\right)\right)_k \right] + \\ &+ \left(\frac{h}{h_0}\right)^{|m|-1} \sum_{k \in M} L_{jk} \left[(\Pi_{h_0} f_m^0)_k + \sum_{0 < |l| \leq q} (\mathfrak{E}_{h_0}^l)_k (\mathcal{P}_{h_0} D^l f_m^0)_k \right], \end{aligned}$$

где все функции $\mathfrak{E}_{h_0}^l$ и временные сечения функции f_m^0 берутся на момент времени th_0/h . Снова пользуясь определением 4 для ошибки в смысле проектора $\tilde{\Pi}_{h_0}$ (43), получаем искомое равенство (44).

Теорема 2 (об автомодельности корректоров при блочном измельчении). Пусть $\mathfrak{S}_h = (M, \emptyset, \Pi_h, Z_h, L_h)$ – блочное измельчение некоторой допустимой схемы \mathfrak{S} . Пусть $\{\mathfrak{C}_{h_0}^m(t), m \in \mathfrak{M} \setminus \{0\}\}$, – цепочка нестационарных корректоров для схемы \mathfrak{S} . Тогда

$$\mathfrak{C}_h^m(t) = \left(\frac{h}{h_0}\right)^{|m|} \mathfrak{C}_{h_0}^m\left(t\frac{h_0}{h}\right), \quad m \in \mathfrak{M} \setminus \{0\}, \quad (46)$$

будет цепочкой корректоров для \mathfrak{S}_h .

Доказательство теоремы 2 проведём по индукции. Пусть $m \in \mathfrak{M} \setminus \{0\}$. Предположим, что $\{\mathfrak{C}_h^l(t), l < m\}$, определённая формулой (46), является цепочкой корректоров. Покажем, что $\{\mathfrak{C}_h^l(t), l \leq m\}$, также будет цепочкой.

Действительно, чтобы \mathfrak{C}_h^m был корректором, необходимо и достаточно, чтобы он подчинялся уравнению (42). В силу утверждения 4 правая часть не зависит от выбора параметров t_0 и r_0 у функции f_m . С использованием утверждения 8 уравнение (42) на нестационарный корректор переписывается в виде

$$\hat{Z} \frac{d\mathfrak{C}_h^m}{dt} + \frac{h_0}{h} \hat{L} \mathfrak{C}_h^m = \epsilon_h(t, f_m, \tilde{\Pi}_{h, < m}) = \left(\frac{h}{h_0}\right)^{|m|-1} \epsilon_{h_0}\left(t\frac{h_0}{h}, f_m, \tilde{\Pi}_{h_0, < m}\right). \quad (47)$$

Подставляя $\mathfrak{C}_h^m(t)$, заданный (46), в это уравнение, и сокращая общий множитель $(h/h_0)^{|m|-1}$ и вводя $\tau = th_0/h$, получаем

$$\hat{Z} \frac{d\mathfrak{C}_{h_0}^m}{d\tau} + \hat{L} \mathfrak{C}_{h_0}^m = \epsilon_{h_0}\left(\tau, f_m, \tilde{\Pi}_{h_0, < m}\right), \quad (48)$$

что верно, поскольку по условию $\mathfrak{C}_{h_0}^m$ является нестационарным корректором для $\mathfrak{S} = \mathfrak{S}_{h_0}$.

Отметим, что в связи с произвольностью начальных данных $\mathfrak{C}^m(0)$ для каждого мультииндекса $m \neq 0$ цепочка корректоров, определённых формулой (46), не является единственно возможной.

Легко показать, что если $K_{h_0}(t)$ – наименьшая функция, являющаяся константой устойчивости $\mathfrak{S} \equiv \mathfrak{S}_{h_0}$, то константа устойчивости $K_h(t)$ для схемы \mathfrak{S}_h при $h = h_0/N$ удовлетворяет оценке

$$K_h(t) \geq K_{h_0}\left(t\frac{h_0}{h}\right).$$

Рассмотрим теперь, что даёт теорема 1 для последовательности схем \mathfrak{S}_h , являющихся блочным измельчением.

Утверждение 9. Пусть выполнены следующие условия:

- 1) решение $u(r, t)$ уравнения (1) q раз дифференцируемо и имеет липшицевы q -е пространственные производные с константой Липшица \mathbb{L} .
- 2) $\mathfrak{S} = (M, \emptyset, \Pi, Z, L)$ – допустимая схема с константами C_s, C_r, C_μ, h_0 , периодическая с периодами $\mathbf{b}_i = \mathbf{b}_i^s, i = 1, \dots, \bar{d}$, и $M^0 = M^s$;
- 3) $\mathfrak{S}_h = (M, \emptyset, \Pi_h, Z_h, L_h)$ – блочное измельчение \mathfrak{S} ;
- 4) \mathcal{P} – допустимый оператор с константами C_r, C_μ, h , периодический относительно \mathfrak{S} , а \mathcal{P}_h связаны с \mathcal{P} соотношением, аналогичным (41);
- 5) $\|\cdot\|$ – некоторая норма в \mathbb{R}^{M^s} , такая что для $e = (1, \dots, 1)^T$ выполняется $\|e\| = 1$ и если $a, b \in \mathbb{R}^{M^s}$ удовлетворяют условию $|a_j| \leq |b_j| \forall j \in M^s$, то $\|a\| \leq C_n \|b\|$;
- 6) если $a, b \in \mathbb{R}^{M^s}$ удовлетворяют условию $|a_j| \leq \max_{k: |\check{Z}_{jk}| + |\check{L}_{jk}| \neq 0} |b_k|$ при $j \in M^s$, то выполняется $\|a\| \leq C_w \|b\|$;
- 7) $\|\check{Z}^{-1}\| \leq C_z < \infty, \|\check{L}\| \leq C_l/h$;
- 8) константа устойчивости \mathfrak{S}_h удовлетворяет оценке $K_h(t) \leq K$, где K не зависит от t и h ;
- 9) $\{\mathfrak{E}^m(t)\}, 0 < |m| \leq q$, – цепочка нестационарных корректоров для схемы \mathfrak{S} с использованием \mathcal{P} с некоторыми начальными данными $\mathfrak{E}^m(0)$, такими что $\mathfrak{E}_{[\xi, \eta]}^m(0) = \mathfrak{E}_{[\xi, 0]}^m(0), \mathfrak{E}_h^m(t)$ определены формулой (46), и $\tilde{\Pi}_{h, q}$ – оператор, порождённый Π_h, \mathcal{P}_h и этой цепочкой корректоров;
- 10) $u_h(t) = \{u_j(t), j \in M\}$ – решение по схеме $\mathfrak{S}_h = (M, 0, \tilde{\Pi}_h, Z_h, L_h)$, где $\tilde{\Pi}_h$ – либо Π_h , либо $\tilde{\Pi}_{h, q}$.

Тогда выполняется оценка

$$\begin{aligned} \|u_h(t) - \tilde{\Pi}_{h, q}(t)u(t, \cdot)\| &\leq K_h(t) \|\tilde{\Pi}_{h, q}(0)u(0, \cdot) - \mathring{\Pi}_h u(0, \cdot)\| + \\ &+ tK_h(t)\tilde{C}\mathbb{L}h^q + K_h(t)\tilde{C}\mathbb{L} \sum_{0 < |m| \leq q} h^{q-|m|} \int_0^t \|\mathfrak{E}_h^m(\tau)\| d\tau, \end{aligned} \quad (49)$$

где константа \tilde{C} зависит только от $q, C_s, C_r, C_n, C_l, C_w, C_z$ и C_μ .

Следствие 10. Пусть выполняются все условия утверждения 9. Пусть величины $\mathfrak{E}^m(t)$ являются корректорами с нулевыми начальными условиями, и пусть скорость их роста меньше, чем $t^{|m|}$: существует такое $\delta, 0 \leq \delta < 1$, что для всех t справедливо $\|\mathfrak{E}^m(t)\| \leq C_m t^{\delta|m|}$. Обозначим символом \times почленное умножение, то есть $(a \times b)_j = a_j b_j, j \in M^s$. Тогда для решения u_h по схеме \mathfrak{S}_h при всех $t \in [0, t_{\max}]$ верно

$$\lim_{h \rightarrow 0} \frac{1}{h^{(1-\delta)(q-1)}} \left\| u_h(t) - \Pi_h u(t, \cdot) - \sum_{0 < |m| \leq q} \mathfrak{E}_h^m(t) \times \mathcal{P}_h D^m u(t, \cdot) \right\| = 0.$$

Исследование 4-точечной схемы

В [9] была рассмотрена 4-точечная схема R3 для аппроксимации производной для решения задачи Коши для обыкновенного дифференциального уравнения $u' + \lambda u = 0$, $u(0) = 1$. Такая постановка является противоестественной, поскольку для решения задачи Коши используется схема не маршевого счёта. Дополнительную неясность вносит отсутствие доказательства устойчивости этой разностной схемы.

В настоящем разделе мы рассмотрим схему R3, применённую к уравнению переноса (1) в одномерном ($d = 1$) случае с периодическими условиями $u(t, x + nx_{\max}) = u(t, x)$, $n \in \mathbb{Z}$. Пусть для простоты скорость переноса $a = 1$.

Пусть на отрезке $[0, x_{\max}]$ задана расчётная сетка, состоящая из $N + 1$ узлов с координатами $0 = x_0 < \dots < x_N = x_{\max}$. Положим $M^0 = M^s = \{0, \dots, N - 1\}$, $M = \mathbb{Z}$, и координаты доопределим по периодике: $x_{i \pm N} = x_i \pm x_{\max}$. Обозначим

$$h_{j+1/2} = x_{j+1} - x_j, \quad x_{j+1/2} = \frac{x_j + x_{j+1}}{2}, \quad \bar{h}_j = \frac{x_{j+1} - x_{j-1}}{2}.$$

Введём обозначения $h_{\max} = \max_j h_{j-1/2}$, $h_{\min} = \min_j h_{j-1/2}$, $(\Delta h)_{\max} = \max_j |h_{j+1/2} - h_{j-1/2}|$.

Схема R3 записывается в консервативном виде

$$\frac{du_j}{dt} + \frac{1}{\bar{h}_j} (F_{j+1/2}|_u - F_{j-1/2}|_u) = 0, \quad (50)$$

где потоки $F_{j+1/2}$ определяются формулой

$$F_{j+1/2}|_u = u_j + \frac{x_{j+1} - x_j}{2} \left(\frac{2u_{j+1} - u_j}{3x_{j+1} - x_j} + \frac{1u_j - u_{j-1}}{3x_j - x_{j-1}} \right), \quad (51)$$

а определение $F_{j-1/2}$ получается из определения $F_{j+1/2}$ заменой j на $j - 1$. Используется точечный проектор Π , определяемый равенством $(\Pi f)_j = f(x_j)$. На равномерной сетке для периодической задачи эта схема обладает 3-м, а на неравномерной сетке – 1-м порядком аппроксимации.

Всюду в настоящем разделе мы будем рассматривать как операторы, действующие из \mathbb{R}^M в \mathbb{R}^M , так и их ограничения на множества функций, периодических с периодом $M^0 = M^s$, действующие из \mathbb{R}^{M^s} в \mathbb{R}^{M^s} . Чтобы их различать, последние будем помечать крышкой над переменными по аналогии с (13).

Сравнивая (50)–(51) с общим видом (9), получаем, что Z – тождественный оператор, а L представим в виде $L = H^{-1}DF$, где H – оператор, домножающий j -ю компоненту вектора на \bar{h}_j , D – оператор, заданный $(Df)_j = f_j - f_{j-1}$, и F – оператор, заданный $(Fu)_j = F_{j+1/2}|_u$. Для оператора \hat{L} , действующего из

\mathbb{R}^{M^s} в \mathbb{R}^{M^s} , справедливо аналогичное представление $\hat{L} = \hat{H}^{-1} \hat{D} \hat{F}$, где \hat{H} – диагональная матрица с компонентами \hat{h}_j на диагонали, а $\hat{D}_{jk} = \hat{\delta}_{j,k} - \hat{\delta}_{j-1,k}$, где $\hat{\delta}_{j,k} = 1$, если $j = k$ по модулю N , и $\hat{\delta}_{j,k} = 0$ иначе. \hat{F} определяется через F по формуле, аналогичной (13).

Вопрос устойчивости схем (50)–(51) является открытым. На равномерной сетке она устойчива в норме L_2 , а на сетке с чередующимися шагами h_{\min} и h_{\max} она является устойчивой в L_2 , если $h_{\max}/h_{\min} \leq 3$. Этот результат устанавливается при помощи спектрального анализа.

Спектральный анализ на чередующейся сетке. Рассмотрим неравномерную сетку с шагами $h_{j+1/2}$, где для чётных j выполняется $h_{j+1/2} = h_{\min} = \hbar(1 - \Delta)$, а для нечётных j выполняется $h_{j+1/2} = h_{\max} = \hbar(1 + \Delta)$. Здесь \hbar не зависит от j и играет роль среднего шага сетки. Обозначим $\mathcal{H} = h_{\max}/h_{\min}$. Оператор \hat{Z} является тождественным оператором, а оператор \hat{L} , – блочно-циклической матрицей с блоками размера 2×2 :

$$\hat{L} = \begin{pmatrix} A_0 & A_1 & 0 & \dots & A_{N-1} \\ A_{N-1} & A_0 & A_1 & \dots & 0 \\ 0 & A_{N-1} & A_0 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ A_1 & 0 & 0 & \dots & A_0 \end{pmatrix},$$

$$A_0 = \frac{1}{6\hbar} \begin{pmatrix} 2 + \mathcal{H} & 2 \\ -4 - \mathcal{H} - \mathcal{H}^{-1} & 2 + \mathcal{H}^{-1} \end{pmatrix}, \quad A_1 = \frac{1}{6\hbar} \begin{pmatrix} 0 & 0 \\ 2 & 0 \end{pmatrix},$$

$$A_{N-1} = \frac{1}{6\hbar} \begin{pmatrix} \mathcal{H}^{-1} & -4 - \mathcal{H} - \mathcal{H}^{-1} \\ 0 & \mathcal{H} \end{pmatrix}.$$

Собственные вектора v_{kl} матрицы L и соответствующие им собственные значения λ_{kl} , где $k = 0, \dots, N - 1, l = 0, 1$, ищутся в виде

$$v_{kl} = \begin{pmatrix} w_{kl} \\ w_{kl} \exp(2\pi ik/N) \\ \vdots \\ w_{kl} \exp(2\pi ik(N-1)/N) \end{pmatrix},$$

где $w_{kl} \in \mathbb{R}^2$ являются собственными векторами матрицы $A(k/N)$, где

$$A(\phi) = A_{N-1} \exp(-2\pi i\phi) + A_0 + A_1 \exp(2\pi i\phi).$$

Для схемы R3 матрица $A(\phi)$ имеет вид

$$A(\phi) = \frac{1}{6\hbar} \begin{pmatrix} 2 + \mathcal{H} + e^{-i\phi} \mathcal{H}^{-1} & 2 - (4 + \mathcal{H} + \mathcal{H}^{-1})e^{-i\phi} \\ -4 - \mathcal{H} - \mathcal{H}^{-1} + 2e^{i\phi} & 2 + \mathcal{H}^{-1} + \mathcal{H}e^{-i\phi} \end{pmatrix}.$$

Можно показать, что при всех ϕ она имеет два линейно независимых собственных вектора, а, значит, базисом из собственных векторов обладает и матрица \hat{L} . Матрица $A(\phi)$ имеет два семейства собственных значений. Одно из них при $\phi \rightarrow 0$ стремится к $(\mathcal{H}^{-1} + 2 + \mathcal{H})/(3\hbar)$. Второе собственное значение имеет разложение

$$\lambda(\phi) = \frac{1}{2\hbar}i\phi + \frac{1}{24\hbar}i\frac{(\mathcal{H}-1)^2}{(\mathcal{H}+1)^2}\phi^3 - \frac{1}{48\hbar}\frac{\mathcal{H}}{(\mathcal{H}+1)^2}(3\mathcal{H}^2 - 10\mathcal{H} + 3)\phi^4 + O(\phi^5).$$

Для точного оператора дифференцирования собственное значение было бы равно $i\phi/2\hbar$ (здесь $2\hbar$ – геометрический размер блока). Видно, что при $\mathcal{H} < 3$ собственное значение вблизи $\phi = 0$ имеет положительную действительную часть, а при $\mathcal{H} > 3$ – отрицательную, что делает схему неустойчивой. Рассмотрение следующих членов разложения показывает, что при $\mathcal{H} = 3$ сохраняется $\text{Re}\lambda(\phi) > 0$ при $\phi \neq 0$. Численный анализ показывает, что в случае $\mathcal{H} < 3$ условие $\text{Re}\lambda(\phi) > 0$ сохраняется для всех $\phi \neq 0$.

Введём обозначение

$$\Lambda_j = h_{j+1/2}/h_{j-1/2}.$$

Проведённые автором численные эксперименты показывают, что при выполнении ограничения

$$\Lambda_j < \Lambda_{\max} < 3 \quad (52)$$

схема R3 остаётся устойчивой и на произвольной неравномерной сетке. Однако поскольку неустойчивость проявляется на длинноволновых возмущениях, её влияние может быть достаточно незаметным, и поэтому выявить её затруднительно.

Ниже будем рассматривать расчётные сетки, удовлетворяющие условию (52). Отметим, что условие ограниченности отношения h_{\max}/h_{\min} при измельчении сетки не накладывается.

Анализ разностного оператора. Докажем вспомогательное утверждение.

Утверждение 11. Пусть $A = \{a_{ij}\}$ – матрица со строчным диагональным преобладанием: $a_{ii} > \sum_{j \neq i} |a_{ij}|$. Тогда справедлива оценка

$$\|A^{-1}\|_{\infty} \leq \max_i \left\{ \frac{1}{a_{ii}} \right\} \left(1 - \max_i \left\{ \frac{1}{a_{ii}} \sum_{j \neq i} |a_{ij}| \right\} \right)^{-1}. \quad (53)$$

Действительно, разобьём A на диагональную и внедиагональную часть: $A = \mathcal{D} + \mathcal{M}$. Тогда $(\mathcal{D}^{-1}A) = I + \mathcal{D}^{-1}\mathcal{M}$. В силу диагонального преобладания

у матрицы A справедливо $\|\mathcal{D}^{-1}\mathcal{M}\|_\infty = \max_i \sum_j |(\mathcal{D}^{-1}\mathcal{M})_{ij}| < 1$, и поэтому выполняется

$$(\mathcal{D}^{-1}A)^{-1} = \sum_{k=0}^{\infty} (-\mathcal{D}^{-1}\mathcal{M})^k.$$

Отсюда

$$A^{-1} = \sum_{k=0}^{\infty} (-\mathcal{D}^{-1}\mathcal{M})^k \mathcal{D}^{-1}.$$

Простая оценка нормы даёт

$$\|A^{-1}\|_\infty \leq \frac{\|\mathcal{D}^{-1}\|_\infty}{1 - \|\mathcal{D}^{-1}\mathcal{M}\|_\infty}.$$

Вычисляя нормы матриц \mathcal{D}^{-1} и $\mathcal{D}^{-1}\mathcal{M}$, получаем искомую оценку (53).

Для периодических решений запишем схему (50)–(51) в операторном виде: $du/dt + \hat{L}u = 0$. Оператор \hat{L} в силу консервативной формы (50) представляется в виде $\hat{L} = \hat{H}^{-1}\hat{D}\hat{F}$. Рассмотрим оператор \hat{F} . Покажем, что он обратим, и оценим $\|\hat{F}^{-1}\|_\infty$. По определению (51) имеем

$$\hat{F} = \left\| \begin{array}{ccccccc} \frac{2}{3} + \frac{\Lambda_0}{6} & \frac{1}{3} & 0 & 0 & 0 & 0 & -\frac{\Lambda_0}{6} \\ -\frac{\Lambda_1}{6} & \frac{2}{3} + \frac{\Lambda_1}{6} & \frac{1}{3} & 0 & 0 & 0 & 0 \\ 0 & -\frac{\Lambda_2}{6} & \frac{2}{3} + \frac{\Lambda_2}{6} & \frac{1}{3} & 0 & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \frac{1}{3} & 0 & 0 & 0 & 0 & -\frac{\Lambda_{N-1}}{6} & \frac{2}{3} + \frac{\Lambda_{N-1}}{6} \end{array} \right\|.$$

Применяя формулу (53), получаем оценку

$$\|\hat{F}^{-1}\|_\infty \leq \frac{3/2}{1 - \max_i \left\{ \frac{1/3 + \Lambda_i/6}{2/3 + \Lambda_i/6} \right\}} = 3 + \frac{3}{4} \max \Lambda_i.$$

Важно, что эта оценка не зависит ни от h , ни от размера матрицы.

Введём диагональную матрицу \tilde{H} размера $|M^s|$ с компонентами $h_{j-1/2}$ на диагонали. Представим \hat{L} в виде $\hat{L} = \hat{G}\tilde{H}^{-1}\hat{D}$. Матрица \hat{G} имеет вид

$$\hat{G} = \hat{H}^{-1} \left\| \begin{array}{cccccc} \left(\frac{2}{3} + \frac{\Lambda_0}{6}\right) h_{-1/2} & \frac{1}{3} h_{1/2} & 0 & 0 & -\frac{1}{6} h_{-1/2} \\ -\frac{1}{6} h_{1/2} & \left(\frac{2}{3} + \frac{\Lambda_1}{6}\right) h_{1/2} & \frac{1}{3} h_{3/2} & 0 & 0 \\ 0 & -\frac{1}{6} h_{3/2} & \left(\frac{2}{3} + \frac{\Lambda_2}{6}\right) h_{3/2} & \frac{1}{3} h_{5/2} & 0 \\ \dots & \dots & \dots & \dots & \dots \\ \frac{1}{3} h_{-1/2} & 0 & 0 & -\frac{1}{6} h_{N-3/2} & \left(\frac{2}{3} + \frac{\Lambda_{N-1}}{6}\right) h_{N-3/2} \end{array} \right\|.$$

Рассмотрим вопрос, при каких условиях матрица \hat{G} имеет сильное диагональное преобладание. Имеем

$$\begin{aligned} 1 - \frac{1}{a_{kk}} \sum_{m \neq k} |a_{km}| &= 1 - \left(\left(\frac{2}{3} + \frac{\Lambda_k}{6} \right) h_{k-1/2} \right)^{-1} \left(\frac{1}{6} h_{k-1/2} + \frac{1}{3} h_{k+1/2} \right) = \\ &= 1 - \frac{2\Lambda_k + 1}{\Lambda_k + 4} = \frac{3 - \Lambda_k}{\Lambda_k + 4}. \end{aligned} \quad (54)$$

Таким образом, для сильного диагонального преобладания необходимо и достаточно выполнения условия $\Lambda_k < 3$, которое выполняется по сделанному выше предположению (52). Есть ли связь условия диагонального преобладания матрицы \hat{G} с условием устойчивости схемы R3, неизвестно.

Оценим $\|\hat{G}^{-1}\|_\infty$ по формуле (53). Пользуясь оценкой (54) и

$$(G_{ii})^{-1} = \hbar_i \left(\left(\frac{2}{3} + \frac{\Lambda_i}{6} \right) h_{i-1/2} \right)^{-1} = \frac{1/2 + \Lambda_i/2}{2/3 + \Lambda_i/6} \leq 3,$$

получаем

$$\|\hat{G}^{-1}\|_\infty \leq 3 \frac{4 + \Lambda_{\max}}{3 - \Lambda_{\max}} := C_G. \quad (55)$$

Анализ ошибки аппроксимации. Рассмотрим аппроксимационную ошибку $\epsilon(0, f, \Pi)$, где f – некоторое решение (1). В (51) было введено обозначение $F_{j+1/2}|_u$ для $u \in \mathbb{R}^M$; для непрерывных функций f будем использовать обозначение $F_{j+1/2}|_f = F_{j+1/2}|_{\Pi f}$. Тогда аппроксимационная ошибка равна

$$\epsilon_j(0, f, \Pi) = \frac{1}{\hbar_j} (F_{j+1/2}|_f - F_{j-1/2}|_f) - f'(x_j).$$

Преобразуем это выражение следующим образом:

$$\begin{aligned} \epsilon_j(0, f, \Pi) &= \frac{1}{\hbar_j} (f''(x_{j+1/2}) F_{j+1/2}|_{(x-x_j)^2/2} - f''(x_{j-1/2}) F_{j-1/2}|_{(x-x_j)^2/2}) + \\ &+ \frac{1}{\hbar_j} \left(F_{j+1/2}|_{f(x) - f''(x_{j+1/2})(x-x_j)^2/2} - F_{j-1/2}|_{f(x) - f''(x_{j-1/2})(x-x_j)^2/2} \right) - f'(x_j). \end{aligned}$$

В первом слагаемом преобразуем выражение вида $F_{j\pm 1/2}|_{(x-x_j)^2/2}$:

$$F_{j\pm 1/2}|_{(x-x_j)^2/2} = F_{j\pm 1/2}|_{(x-x_{j\pm 1/2})^2/2} + F_{j\pm 1/2}|_{(x-x_j)^2/2 - (x-x_{j\pm 1/2})^2/2}.$$

Функция $(x-x_j)^2/2 - (x-x_{j\pm 1/2})^2/2$ линейна по x , а потоки по построению точны на линейной функции. Поэтому последнее слагаемое равно значению

этой линейной функции в $x = x_{j\pm 1/2}$, то есть равно $(x_j - x_{j\pm 1/2})^2/2 = h_{j\pm 1/2}^2/8$. Таким образом,

$$\begin{aligned} \epsilon_j(0, f, \Pi) = & \frac{1}{\bar{h}_j} \left(f''(x_{j+1/2}) \left(F_{j+1/2}|_{(x-x_{j+1/2})^2/2} + \frac{h_{j+1/2}^2}{8} \right) - \right. \\ & \left. - f''(x_{j-1/2}) \left(F_{j-1/2}|_{(x-x_{j-1/2})^2/2} + \frac{h_{j-1/2}^2}{8} \right) \right) + \\ & + \frac{1}{\bar{h}_j} \left(F_{j+1/2}|_{f(x)-f(x_j)-(x-x_j)f'(x_j)-f''(x_{j+1/2})(x-x_j)^2/2} - \right. \\ & \left. - F_{j-1/2}|_{f(x)-f(x_j)-(x-x_j)f'(x_j)-f''(x_{j-1/2})(x-x_j)^2/2} \right). \end{aligned} \quad (56)$$

Последнее слагаемое преобразовано с учётом точности на константе, линейной функции, а также определения $\bar{h}_j = (h_{j+1/2} + h_{j-1/2})/2$.

Теперь, пользуясь проведённым выше анализом, получим оценки на величины корректоров \mathfrak{E}^m , попутно выбирая для них подходящие начальные данные. Корректоры будем определять с использованием $\mathcal{P} \equiv \Pi$. Поскольку схема точна на линейной функции, можно положить $\mathfrak{E}^1 \equiv 0$.

Анализ корректора \mathfrak{E}^2 . Если подставить в формулу (56) функцию $f_{2,j} = -(x - x_j)^2/2$, последние два слагаемых зануляются, и получаем

$$\begin{aligned} \epsilon_j(0, f_{2,j}, \Pi) &= \frac{1}{\bar{h}_j} (\phi_{j+1/2} - \phi_{j-1/2}), \\ \phi_{j+1/2} &= -F_{j+1/2}|_{(x-x_{j+1/2})^2/2} - \frac{1}{8}h_{j+1/2}^2 = \\ &= -\frac{1}{8}h_{j+1/2}^2 - \frac{h_{j+1/2}}{2} \left(0 + \frac{1}{3h_{j-1/2}} \left(\frac{h_{j+1/2}^2}{8} - \frac{(h_{j+1/2} + 2h_{j-1/2})^2}{8} \right) \right) - \frac{1}{8}h_{j+1/2}^2 = \\ &= -\frac{1}{6}h_{j+1/2}^2 + \frac{1}{12}h_{j+1/2}h_{j-1/2}. \end{aligned}$$

Корректор \mathfrak{E}^2 находится из системы уравнений

$$\frac{d\mathfrak{E}^2}{dt} + \hat{H}^{-1} \hat{D} \hat{F} \mathfrak{E}^2 = \epsilon(0, f_2, \Pi) = \hat{H}^{-1} \hat{D} \phi.$$

Удобно выбрать его не зависящим от t и равным $\mathfrak{E}^2 = \hat{F}^{-1}(\phi + h_{\min}^2 e/12) = \hat{F}^{-1} \phi + h_{\min}^2 e/12$, где $e = (1, \dots, 1)^T$. Тогда

$$\begin{aligned} \left| \phi_j + \frac{1}{12} h_{\min}^2 \right| &= \frac{1}{12} |h_{j+1/2}(h_{j-1/2} - h_{j+1/2}) + (h_{\min} - h_{j+1/2})(h_{\min} + h_{j+1/2})| \leq \\ &\leq \frac{1}{4} h_{\max}(h_{\max} - h_{\min}). \end{aligned}$$

Отсюда

$$\|\mathfrak{E}^2\|_{\infty} \leq \left\| \hat{F}^{-1} \right\|_{\infty} \|\phi + h_{\min}^2 e/12\|_{\infty} \leq \left(3 + \frac{3}{4} \Lambda_{\max} \right) h_{\max}(h_{\max} - h_{\min}). \quad (57)$$

В то же время имеем

$$\begin{aligned} |\epsilon_j(0, f_{2,j}, \Pi)| &= \frac{1}{\hbar_j} |\phi_{j+1/2} - \phi_{j-1/2}| = \\ &= \frac{1}{\hbar_j} \left| -\frac{1}{6}(h_{j+1/2}^2 - h_{j-1/2}^2) + \frac{1}{12} h_{j-1/2}(h_{j+1/2} - h_{j-3/2}) \right| \leq (\Delta h)_{\max}, \end{aligned}$$

и поэтому в силу (55)

$$\left\| \tilde{H}^{-1} \hat{D} \mathfrak{E}^2 \right\|_{\infty} = \left\| \hat{G}^{-1} \epsilon(0, f_2, \Pi) \right\|_{\infty} \leq \left\| \hat{G}^{-1} \right\|_{\infty} \|\epsilon(0, f_2, \Pi)\|_{\infty} \leq C_G (\Delta h)_{\max}. \quad (58)$$

Анализ корректора \mathfrak{E}^3 . Нестационарный корректор \mathfrak{E}^3 , образующий цепочку вместе с выбранным корректором \mathfrak{E}^2 , находится из системы уравнений

$$\frac{d\mathfrak{E}^3}{dt} + \hat{H}^{-1} \hat{D} \hat{F} \mathfrak{E}^3 = \epsilon(t, f_3, \tilde{\Pi}_2).$$

Правая часть не зависит от времени, поскольку $\tilde{\Pi}_2$ стационарный. Зафиксируем узел k , в котором будем вычислять правую часть. Поскольку правая часть не зависит от выбора констант в f_3 , положим $\mathbf{r}_0 = \mathbf{r}_k$ и $t_0 = 0$. Имеем

$$\epsilon_k(0, f_3, \tilde{\Pi}_2) = -\mathfrak{E}_k^2 + \sum_m L_{km} \left(\frac{(x_m - x_k)^3}{6} + \mathfrak{E}_m^2(x_m - x_k) \right). \quad (59)$$

Представим $\epsilon_k(0, f_3, \tilde{\Pi}_2) = p_k + q_k$, причём в p_k отнесём члены, содержащие коэффициенты корректора \mathfrak{E}^2 , а в q_k – не содержащие. Для краткости обозначим $\alpha \equiv \mathfrak{E}^2$.

Рассмотрим вначале слагаемое p_k :

$$\begin{aligned} p_k &= -\alpha_k + \sum_m L_{km} \alpha_m (x_m - x_k) = \\ &= \frac{1}{\hbar_k} \left[-\alpha_k \hbar_k + \sum_m (F_{km} - F_{k-1,m}) \alpha_m (x_m - x_k) \right] = \frac{\eta_{k+1/2}^{(0)} - \eta_{k-1/2}^{(0)}}{\hbar_k} + \frac{1}{\hbar_k} r_k, \end{aligned}$$

где

$$\eta_{k+1/2}^{(0)} = \sum_m F_{km} \alpha_m (x_m - x_{k+1}), \quad r_k = -\alpha_k \bar{h}_k + (x_{k+1} - x_k) \sum_m F_{km} \alpha_m.$$

Теперь распишем r_k , зная коэффициенты F_{km} для схемы R3:

$$\begin{aligned} r_k &= -\alpha_k \bar{h}_k + h_{k+1/2} \left(-\frac{\Lambda_k}{6} \alpha_{k-1} + \left(\frac{2}{3} + \frac{\Lambda_k}{6} \right) \alpha_k + \frac{1}{3} \alpha_{k+1} \right) = \\ &= -\frac{1}{2} \alpha_k h_{k+1/2} - \frac{1}{2} \alpha_k h_{k-1/2} + \frac{1}{6} \frac{h_{k+1/2}^2}{h_{k-1/2}} (\alpha_k - \alpha_{k-1}) + \frac{2}{3} h_{k+1/2} \alpha_k + \frac{1}{3} h_{k+1/2} \alpha_{k+1} = \\ &= -\frac{1}{6} \frac{h_{k-1/2}^2 - h_{k+1/2}^2}{h_{k-1/2}} (\alpha_k - \alpha_{k-1}) + \frac{1}{6} (\alpha_k h_{k+1/2} - \alpha_{k-1} h_{k-1/2}) + \\ &\quad + \frac{1}{3} (\alpha_{k+1} h_{k+1/2} - \alpha_k h_{k-1/2}). \end{aligned}$$

Таким образом,

$$p_k = \frac{\eta_{k+1/2}^{(p)} - \eta_{k-1/2}^{(p)}}{\bar{h}_k} - \frac{1}{3} (h_{k-1/2} - h_{k+1/2}) \frac{\alpha_k - \alpha_{k-1}}{h_{k-1/2}}, \quad (60)$$

$$\eta_{k+1/2}^{(p)} = \sum_m F_{km} \alpha_m (x_m - x_{k+1}) + \frac{1}{6} \alpha_k h_{k+1/2} + \frac{1}{3} \alpha_{k+1} h_{k+1/2}. \quad (61)$$

В силу (58) последнее слагаемое в правой части (60) по модулю не превосходит $((\Delta h)_{\max})^2 C_G / 3$. В силу (57) имеем $\eta_{k+1/2}^{(p)} = O(h_{\max}^3)$.

Рассмотрим теперь слагаемое $q_k = \epsilon_k(0, f_3, \Pi)$. Схема имеет коэффициенты $L_{k,k-2} = \Lambda_{k-1}/6$, $L_{k,k-1} = -(4 + \Lambda_{k-1} + \Lambda_k)/6$, $L_{k,k} = 1/3 + \Lambda_k/6$, $L_{k,k+1} = 1/3$, и остальные коэффициенты нулевые. Поэтому

$$q_k = \frac{1}{\bar{h}_k} \left(\frac{\Lambda_{k-1} (x_{k-2} - x_k)^3}{6} - \left(\frac{2}{3} + \frac{\Lambda_{k-1} + \Lambda_k}{6} \right) \frac{(x_{k-1} - x_k)^3}{6} + \frac{1}{3} \frac{(x_{k+1} - x_k)^3}{6} \right).$$

Рассмотрим слагаемые, содержащие множитель Λ_{k-1} :

$$\begin{aligned} &\frac{1}{36\bar{h}_k} \Lambda_{k-1} ((x_{k-2} - x_k)^3 - (x_{k-1} - x_k)^3) = \\ &= \frac{1}{36\bar{h}_k} \Lambda_{k-1} (x_{k-2} - x_{k-1}) ((x_{k-2} - x_k)^2 + (x_{k-2} - x_k)(x_{k-1} - x_k) + (x_{k-1} - x_k)^2) = \\ &= -\frac{1}{36\bar{h}_k} h_{k-1/2} \left((h_{k-1/2} + h_{k-3/2})^2 + (h_{k-1/2} + h_{k-3/2}) h_{k-1/2} + h_{k-1/2}^2 \right) = \\ &= -\frac{1}{36\bar{h}_k} h_{k-1/2} \left(h_{k-3/2}^2 + 3h_{k-1/2} h_{k-3/2} + 3h_{k-1/2}^2 \right). \end{aligned}$$

Подставляя этот результат в выражение для q_k и раскрывая Λ_k , получаем

$$q_k = -\frac{1}{36\hbar_k} h_{k-1/2} \left(h_{k-3/2}^2 + 3h_{k-1/2}h_{k-3/2} + 3h_{k-1/2}^2 \right) + \frac{1}{36\hbar_k} \left(4h_{k-1/2}^3 + h_{k-1/2}^2 h_{k+1/2} + 2h_{k+1/2}^3 \right).$$

Теперь нетрудно убедиться, что q_k представимо в виде

$$q_k = \frac{\eta_{k+1/2}^{(q)} - \eta_{k-1/2}^{(q)}}{\hbar_k} + \frac{1}{36} \frac{(2h_{k+1/2} + h_{k-1/2})(h_{k+1/2} - h_{k-1/2})^2}{\hbar_k}, \quad (62)$$

где $\eta_{k+1/2}^{(q)} = (h_{k-1/2}^2 h_{k+1/2} + 3h_{k-1/2} h_{k+1/2}^2)/36 = O(h_{\max}^3)$. Легко видеть, что второе слагаемое в (62) не превосходит $(\Delta h)_{\max}^2/9$.

Подставляя (60)–(62) в (59), получаем, что правая часть уравнения на \mathfrak{E}^3 приводится к виду

$$\epsilon_k(t, f_3, \tilde{\Pi}_2) = q_k + p_k = \frac{\eta_{k+1/2} - \eta_{k-1/2}}{\hbar_k} + g_k, \quad (63)$$

где $\eta_{k+1/2} = \eta_{k+1/2}^{(p)} + \eta_{k+1/2}^{(q)} = O(h_{\max}^3)$, а $g_k = O((\Delta h)_{\max}^2)$. Величины $\eta_{k+1/2}$ и g_k не зависят от времени.

Выберем в качестве начальных данных $\mathfrak{E}^3(0) = \hat{F}^{-1}\eta$. Тогда величина $w(t) = \mathfrak{E}^3(t) - \mathfrak{E}^3(0)$ будет удовлетворять уравнению

$$\frac{dw}{dt} + \hat{H}^{-1} \hat{D} \hat{F} w = g, \quad w(0) = 0.$$

Отсюда в любой норме $\|w(t)\| \leq t K_h(t) \|g\|$, где $K_h(t)$ – константа устойчивости в этой норме. Если дополнительно наложить условие, что для произвольного вектора $\|y\| \leq \|y\|_{\infty}$, получаем

$$\|\mathfrak{E}^3(t)\| = O((\Delta h)_{\max}^2 t + h_{\max}^3). \quad (64)$$

Теорема 3. Пусть $u(t, x) \in C^3(\mathbb{R})$ – периодическая функция с периодом x_{\max} , \mathbb{L} – константа Липшица её 3-й производной, $\mathbb{L}_q = \max D^q u$. Рассмотрим расчётную сетку, удовлетворяющую условию (52). Пусть $\|\cdot\|$ – некоторая норма на \mathbb{R}^{M^s} , удовлетворяющая условиям 4) и 5) теоремы 1 с константами C_n и C_w , и $K_h(t)$ – константа устойчивости схемы R3 на этой сетке в этой норме. Тогда существует константа C , зависящая только от Λ_{\max} , C_n и C_w , такая что

$$\|u_h(t) - \Pi u(t, \cdot)\| \leq C K_h(t) (\mathbb{L}_2 h (h_{\max} - h_{\min}) + \mathbb{L}_3 ((\Delta h)^2 t + h^3)) + C K_h(t) \mathbb{L} h^3 (h + t) + C K_h(t) \mathbb{L} (\Delta h)^2 t^2. \quad (65)$$

Заметим, что схема R3 является допустимой в смысле определения 2 с константами $h = h_{\max}$, $C_\mu = C_r = 1$, C_s оценивается по утверждению 1 при $C_1 = (2 + \Lambda_{\max})/3$, $C_2 = 4$, $C_3 = 2$, что даёт $C_s = 8(2 + \Lambda_{\max})/3$. Константа $C_z = 1$. Для доказательства искомого утверждения достаточно применить теорему 1 с учётом полученных оценок (57) и (64), заметив, что константа C_l не участвует в оценке, так как $Z = I$. Отметим, что оценка (65) содержит константу устойчивости K_h , ограниченность которой при $h \rightarrow 0$ не доказана.

Анализ корректора \mathfrak{E}^4 . Для получения искомой оценки на величину ошибки решения нам не потребовалось рассмотрения корректоров следующих порядков. Однако можно показать, что вся ошибка 2-го порядка малости по h содержится в корректорах \mathfrak{E}^2 и \mathfrak{E}^3 , то есть $\mathfrak{E}^4 = \bar{o}(h_{\max}^2)$. Мы покажем это при дополнительном предположении, что константа устойчивости схемы R3 в рассматриваемой норме $K_h(t) \leq K$, где K не зависит от t и h , а константа устойчивости $K_\infty(t)$ в норме L_∞ удовлетворяет оценке $K_\infty(t) \leq (t/h)^\delta K$, $\delta < 1$.

Продифференцируем уравнение на нахождение \mathfrak{E}^3 . Имеем

$$\ddot{\mathfrak{E}}^3 + \hat{L}\dot{\mathfrak{E}}^3 = 0,$$

$$\dot{\mathfrak{E}}^3(0) = \epsilon(0, f_3, \tilde{\Pi}_2) - \hat{L}\mathfrak{E}^3(0) = (\hat{H}^{-1}\hat{D}\eta + g) - (\hat{H}^{-1}\hat{D}\hat{F})(\hat{F}^{-1}\eta) = g.$$

По определению константы устойчивости

$$\|\dot{\mathfrak{E}}^3(t)\|_\infty \leq K_\infty(t)\|\dot{\mathfrak{E}}^3(0)\|_\infty = K_\infty(t)\|g\|_\infty.$$

Отсюда из уравнения $\dot{\mathfrak{E}}^3 + \hat{L}\mathfrak{E}^3 = \epsilon(0, f_3, \tilde{\Pi}_2)$ по неравенству треугольника

$$\|\hat{L}\mathfrak{E}^3(t)\|_\infty \leq \|\epsilon(0, f_3, \tilde{\Pi}_2)\|_\infty + K_\infty(t)\|g\|_\infty.$$

Подставляя разложение $\hat{L} = \hat{G}\tilde{H}^{-1}\hat{D}$, получаем

$$\|\hat{D}\mathfrak{E}^3(t)\|_\infty \leq \|\tilde{H}\hat{G}^{-1}\|_\infty(\|\epsilon(0, f_3, \tilde{\Pi}_2)\|_\infty + K_\infty(t)\|g\|_\infty) \leq Ct^\delta h_{\max}^{3-\delta}.$$

Рассмотрим теперь правую часть уравнения на нахождение \mathfrak{E}^4 . Поскольку корректор \mathfrak{E}^3 зависит от времени, она также будет зависеть от времени, но не будет зависеть от выбора параметров t_0 и \mathbf{r}_0 у функции f_4 . Зафиксируем момент времени t и узел j , в котором будем вычислять правую часть. Выберем у функции $t_0 = t$, $\mathbf{r} = \mathbf{r}_j$. Тогда

$$\epsilon_j(t, f_4, \tilde{\Pi}_3) = -\mathfrak{E}_j^3(t) + \sum_m L_{jm} \left(\frac{(x_m - x_j)^4}{24} + \mathfrak{E}_m^2 \frac{(x_m - x_j)^2}{2} + \mathfrak{E}_m^3(t)(x_m - x_j) \right).$$

В силу точности на линейной функции имеем

$$\epsilon_j(t, f_4, \tilde{\Pi}_3) = \sum_m L_{jm} \left(\frac{(x_m - x_j)^4}{24} + \mathfrak{E}_m^2 \frac{(x_m - x_j)^2}{2} + (\mathfrak{E}_m^3(t) - \mathfrak{E}_j^3(t))(x_m - x_j) \right).$$

Слагаемые без коэффициентов корректора и слагаемые, содержащие коэффициенты \mathfrak{E}^2 , имеют порядок $O(h^4)$, что после домножения на $L_{jm} \sim 1/h$ даёт $O(h^3)$. Таким образом,

$$|\epsilon_j(t, f_4, \tilde{\Pi}_3)| = O(h^3) + O\left(\|\hat{D}\mathfrak{E}^3(t)\|_\infty\right) = O(h^{3-\delta}t^\delta).$$

Следовательно, если положить $\mathfrak{E}^4(0) = 0$, получаем

$$\|\mathfrak{E}^4(t)\| \leq CK_h(t)h^{3-\delta}t^{1+\delta}. \quad (66)$$

Применяя теперь оценку (27) теоремы 1, для решения $u_h(t)$ получаем

$$\begin{aligned} & \|u_h(t) - \tilde{\Pi}_4(t)u(t, \cdot)\| \leq K(\mathbb{L}_2\|\mathfrak{E}^2\| + \mathbb{L}_3\|\mathfrak{E}^3(0)\|) + \\ & + tK\tilde{\mathbb{L}}h^4 + tK\tilde{\mathbb{L}}\max_{\tau < t}(h^2\|\mathfrak{E}^2\| + h\|\mathfrak{E}^3(\tau)\| + \|\mathfrak{E}^4(\tau)\|) \leq \\ & \leq \tilde{C}' [K(\mathbb{L}_2h_{\max}(h_{\max} - h_{\min}) + \mathbb{L}_3h_{\max}^3) + tK^2\mathbb{L}(h_{\max}^4 + h_{\max}^3t + h_{\max}^{3-\delta}t^{1+\delta})]. \end{aligned} \quad (67)$$

В последнем неравенстве использованы оценки (57), (64) и (66). Таким образом, растущий со временем член ошибки решения наименьшего порядка по h , имеющий величину $O((\Delta h)_{\max}^2 t)$, содержится в корректоре $\mathfrak{E}^3(t)$.

Также с использованием (67) оценка (65) уточняется до

$$\begin{aligned} \|u_h(t) - \Pi u(t, \cdot)\| & \leq \tilde{C}'K\mathbb{L}_2h(h_{\max} - h_{\min}) + \tilde{C}'K\mathbb{L}_3((\Delta h)^2t + h^3) + \\ & + \tilde{C}'K\mathbb{L}_4h^{3-\delta}t^{1+\delta} + \tilde{C}'tK^2\mathbb{L}(h_{\max}^4 + h_{\max}^3t + h_{\max}^{3-\delta}t^{1+\delta}). \end{aligned} \quad (68)$$

Отметим, что если при использовании метода нестационарного корректора ограничиться рассмотрением главного корректора (в нашем случае – \mathfrak{E}^2), то это было бы эквивалентно оценке точности схемы с использованием негативной нормы. При этом получается оценка $\|u_h(t) - \Pi u(t, \cdot)\| = O(h_{\max}^2 t + h_{\max}^2)$.

Заключение

В работе был изложен метод нестационарного корректора применительно к исследованию точности линейных разностных схем для уравнения переноса. Интегрирование по времени предполагалось точным. Это позволило упростить изложение по сравнению с предшествующей работой [8]. Метод был обобщён на случай схем с матрицей перед временными производными. Был проведён анализ схемы R3 на неравномерной сетке с периодическими условиями.

В следующей работе мы подробнее остановимся на случае периодических разностных схем с конечным периодом, не зависящим от периода решения.

Автор выражает благодарность М. Д. Сурначёву за внимательное прочтение работы и содержательные замечания к ней.

Работа выполнена при поддержке Российского фонда фундаментальных исследований, проект 16-31-60072 мол-а-дк.

Список литературы

1. Тихонов А. Н., Самарский А. А. Однородные разностные схемы на неравномерных сетках // Журнал вычислительной математики и математической физики. 1962. Т. 2. С. 812–832.
2. Diskin B., Thomas J.-L. Notes on accuracy of finite-volume discretization schemes on irregular grids // Applied Numerical Mathematics. 2010. Vol. 60, no. 3. P. 224–226.
3. Despres B. Lax theorem and finite volume schemes // Mathematics of Computation. 2003. Vol. 73. P. 1203–1234.
4. Bouche D., Ghidaglia J.-M., Pascal F. Error Estimate and the Geometric Corrector for the Upwind Finite Volume Method Applied to the Linear Advection Equation // SIAM Journal on Numerical Analysis. 2006. Vol. 43. P. 557–603.
5. Lowrie R. Compact higher-order numerical methods for hyperbolic conservation laws. Ph.D. thesis: The University of Michigan. 1996.
6. 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.
7. Cheng Y., Shu C.-W. Superconvergence and time evolution of discontinuous Galerkin finite element solutions // Journal of Computational Physics. 2008. Vol. 227, no. 22. P. 9612–9627.
8. Бахвалов П. А. Метод нестационарного корректора для анализа точности линейных разностных схем для уравнения переноса // Препринты ИПМ им. М.В.Келдыша. 2016. № 140. С. 1–32.
9. Бахвалов П. А., Козубская Т. К. Структура ошибки консервативного 4-точечного конечно-разностного оператора дифференцирования на неравномерных сетках // Препринты ИПМ им. М.В.Келдыша. 2014. № 74. С. 1–32.