

## РЕШЕНИЕ НАЧАЛЬНЫХ ЗАДАЧ НА ОСНОВЕ ПРИНЦИПА ДИФФЕРЕНЦИАЛЬНЫХ НЕВЯЗОК

We describe the way of constructing the methods of solution of the IVP for systems of ODEs based on the differential defect principle. Different ways of the error estimation are provided. The scheme of the machine implementation and the results of the numerical experiments are also discussed.

Рассмотрим задачу Коши для системы ОДУ

$$u' = f(x, u), \quad u(x_0) = u_0. \quad (1)$$

Наша цель - построить метод, который (в идеале) позволял бы приближенно решать эту задачу со сколь угодно высокой точностью на фиксированном отрезке  $[x_0, x_0+h]$ , не прибегая при этом к неограниченному уменьшению шага  $h$ .

Пусть нам известно некоторое приближение  $y \in C^1[x_0, x_0+h]$  к решению задачи (1), удовлетворяющее начальному условию  $y(x_0)=u_0$ . На отрезке  $[x_0, x_0+h]$  определим также погрешность  $\varepsilon(x)=u(x) - y(x)$  и дифференциальную невязку

$$\rho(x) = f(x, y(x)) - y'(x). \quad (2)$$

Функция  $\rho$  характеризует близость  $y$  к точному решению (1), и, так как ее значения могут быть вычислены на практике, естественно попытаться с их помощью, во-первых, оценить погрешность  $\varepsilon$  и, во-вторых, улучшить уже имеющееся приближение.

**Утверждение.** Если функция  $(x, u)$  удовлетворяет условию Липшица по второму аргументу, то справедливы оценки

$$1 - L(x - x_0) \leq \frac{\mu_\rho(x)}{\mu_\varepsilon(x)} \leq 1 + L(x - x_0), \quad (3.1)$$

$$1 - L(x - x_0) \leq \frac{\mu_\delta(x)}{\mu_\varepsilon(x)} \leq 1 + L(x - x_0), \quad (3.2)$$

где  $\mu_g(x) = \max_{\xi \in [x_0, x]} m_g(\xi)$ ,  $m_g(x) = \|g(x)\|$ ,  $\delta(x) = \int_{x_0}^x \rho(\xi) d\xi$ ,  $L$  - константа Липшица.

Для доказательства рассмотрим дифференциальное уравнение, которому удовлетворяет погрешность:

$$\varepsilon' = \rho + \Delta f. \quad (4)$$

Здесь  $\Delta f = f(x, y + \varepsilon) - f(x, y)$ . Применяя неравенство треугольника к (4) и используя условие Липшица, можно записать:

$$\mu_\varepsilon(x) < \mu_\rho(x) + \mu_{\Delta f}(x) < \mu_\rho(x) + L\mu_\varepsilon(x) \leq \mu_\rho(x) + L(x - x_0)\mu_\varepsilon(x).$$

Отсюда находим, что

$$\mu_\varepsilon(x) (1 - L(x - x_0)) \leq \mu_\rho(x). \quad (5.1)$$

Выразив из (4)  $\rho$ , по аналогии имеем:

$$\mu_\rho(x) \leq \mu_\varepsilon(x) + \mu_{\Delta f}(x) \leq \mu_\varepsilon(x) (1 + L(x - x_0)). \quad (5.2)$$

Оценки (3.1) получаются объединением неравенств (5.1) и (5.2). После интегрирования (4) по такой же схеме доказывается (3.2).

*Следствие.* Функции  $\rho$  и  $\delta$  представляют собой главные части погрешностей  $\varepsilon'$  и  $\varepsilon$  соответственно:

$$\lim_{h \rightarrow 0} \frac{\mu_\rho(x_0 + h)}{\mu_{\varepsilon'}(x_0 + h)} = \lim_{h \rightarrow 0} \frac{\mu_\delta(x_0 + h)}{\mu_\varepsilon(x_0 + h)} = 1. \quad (6)$$

Эти результаты позволяют сделать следующие выводы. Во-первых, из (4), (6) вытекает, что приближение  $y$  можно уточнять, прибавляя к нему  $\delta$  - главную часть погрешности. Во-вторых, неравенства (3.2) позволяют получить оценки

$$\tilde{\mu}(x) = \frac{\mu_\delta(x)}{1 + L(x - x_0)} \leq \mu_\varepsilon(x) \leq \frac{\mu_\delta(x)}{1 - L(x - x_0)}. \quad (7)$$

К сожалению, оценка сверху в (7) имеет смысл лишь при  $0 < hL < 1$ . Отметим также, что здесь фактически оценивается не норма погрешности  $m_\varepsilon(x)$ , а ее максимум на отрезке  $[x_0, x]$ .

Более строгая оценка нормы погрешности имеет вид

$$m_\varepsilon(x) \leq e^{L(x-x_0)} \int_{x_0}^x e^{-L(\xi-x_0)} m_\rho(\xi) d\xi = M_L(x). \quad (8)$$

Она является тривиальным следствием теоремы 10.2 [1, с. 63]. Кроме того, если функция  $f$  удовлетворяет одностороннему условию Липшица

$$\langle f(x, u) - f(x, z), u - z \rangle \leq v \|u - z\|^2,$$

то (8) можно усилить, заменив  $L$  на  $v$  - одностороннюю константу Липшица (см. [1, с. 65, 68], [2, с. 207-208]). С учетом (7), (8) и последнего замечания получим оценки для  $\mu_\varepsilon(x_0 + h)$ :

$$\frac{\mu_\delta(x_0 + h)}{1 + hL} \leq \mu_\varepsilon(x_0 + h) \leq e^{hv} \int_{x_0}^{x_0+h} e^{-v(\xi-x_0)} \mu_\rho(\xi) d\xi. \quad (9)$$

Таким образом, неравенства (8) и (9) дают нам возможность оценивать погрешность  $\varepsilon$  по значениям дифференциальной невязки (2). Приступим теперь непосредственно к обсуждению метода решения задачи (1).

Как уже упоминалось, основной идеей метода является последовательное уточнение приближенного решения путем добавления к нему главной части погрешности. При машинной реализации такого процесса необходим вычислительно-надежный способ нахождения значений  $p(x)$ . Однако непосредственное использование формулы (2) затрудняется необходимостью дифференцирования

приближенного решения. Эта операция накладывает ограничение на гладкость  $u$  и, кроме того, известна своей вычислительной неустойчивостью. В случае,

когда интеграл  $\int_{x_0}^x f(\xi, y(\xi))d\xi$  может быть вычислен аналитически, подобные

трудности отпадают, и мы приходим к известному процессу Пикара. В общем же случае разумнее строить сначала приближение к производной  $v=y' \approx u'$  и уточнять уже его аналогичным образом. При этом способ построения  $v$  выбирается так, чтобы можно было легко вычислить точное значение

$$y(x) = u_0 + \int_{x_0}^x v(\xi) d\xi \text{ в любой точке отрезка } [x_0, x_0+h].$$

В этом случае дифференциальная невязка (2) принимает вид

$$\rho(x) = f(x, u_0 + \int_{x_0}^x v(\xi) d\xi) - v(x),$$

и процесс уточнений производной можно сформулировать, например, так:

$$v_{i+1} = v_i + \omega_i \rho_i \tag{10}$$

Здесь  $\omega_i = \omega_i(x)$  - некоторая скалярная весовая функция, с помощью которой можно регулировать свойства сходимости процесса. В простейшем случае можно положить  $\omega_i(x) = 1$ . Неравенства (3.1) позволяют получать другие, иногда

более рациональные, варианты. Так, при  $\omega_i(x) = \frac{1}{1 + \mu_i(x - x_0)}$  будет выполнено

условие  $\mu_{\rho_i}(x) \geq \omega_i(x) \mu_{\rho}(x)$ , что гарантирует малость добавки в (10) и, как следствие, повышает вычислительную устойчивость алгоритма. Следует, однако, заметить, что и этот выбор не является оптимальным, так как в некоторых случаях может ощутимо замедлять сходимость.

Для краткой записи конструируемого метода перепишем (1) в операторном виде:

$$Du = F(u), u(x_0) = u_0.$$

Здесь

$$Du = u', F(u)(x) = f(x, u(x)).$$

Вместо (10) рассмотрим процесс

$$v_{i+1} = v_i + A_i (\omega_i(F(D^{-1}v_i) - v_i)), \tag{11}$$

где для соблюдения начальных условий исходной задачи полагается

$$v_i(x_0) = f(x_0, u_0), \text{ и первообразная берется с учетом условия } D^{-1}v_i(x_0) = u_0.$$

В формуле (11) введены новые объекты - операторы аппроксимации  $A_i$ . Предполагается, что они обладают тем свойством, что полученные с их помощью приближения легко интегрируются аналитически. За счет последовательного повышения уровня такой аппроксимации и предполагается достижение поставленной цели - нахождения решения с любой заданной точностью на фиксированном отрезке. Поэтому выбор операторов  $A_i$  играет существенную роль при построении методов, основанных на формулах типа (11). Необходимо подчеркнуть, что на данный момент вопрос об оптимальности подобного выбора в общем случае остается открытым.

Рассмотрим подробнее один из возможных вариантов реализации алгоритма (11) на примере использования кусочно-линейной аппроксимации.

Итак, на каждом шаге итерационного процесса нам известно кусочно-линейное приближение  $v_i$  к производной  $u'$  точного решения, определенное на всем отрезке  $[x_0, x_0+h]$  набором узлов  $X_i = \{x_{ij}\}, j = \overline{0, n_i}$ , и значениями  $V_i = \{v_i(x_{ij})\}$ . Переход к следующему приближению  $v_{i+1}$  осуществляется таким образом.

1. Находим  $y_i = D^{-1}v_i$ . Для этого достаточно вычислить  $Y_i = \{y_i(x_{ij})\}$  по квадратурной формуле трапеций. Значения  $y_i(x)$  в точках  $x \in X_i$  легко вычисляются с помощью  $Y_i$  и  $V_i$ . Очевидно, что найденное таким образом  $y_i$  представляет собой квадратичный сплайн, а  $v_i$  - его точная производная.

2. Аппроксимируем добавку к производной, т. е. находим  $\Delta v_i = A_i((\omega_i(F(v_i) - v_i))$  - кусочно-линейную вектор-функцию, которая будет задана набором ее значений  $\{\Delta v_{ij}\}$  в узлах сетки  $Z_i = \{z_{ij}\}, j = \overline{0, m_i}$ . Величина  $\Delta_i = \max_j \|\Delta v_{ij}\|$  служит индикатором сходимости итерационного процесса. Верхнюю границу погрешности текущего приближения в любой точке  $x \in [x_0, x_0 + h]$  можно вычислить с помощью (8), а оценить с обеих сторон ее максимум на всем отрезке позволяет формула (9). Заметим, что вычисление этих оценок является достаточно трудоемким, и, кроме этого, при возрастании значений параметров  $L, v, h$  они становятся все более грубыми. Поэтому о погрешности можно судить также эмпирически по величинам  $\Delta_i$  и  $\mu_\delta(h)$ . Экспериментальные данные для сравнения этих способов оценки погрешности будут приведены ниже.

3. Находим  $v_{i+1} = v_i + \Delta v_i$ , или, что то же самое,  $X_{i+1}$  и  $V_{i+1}$ . В точном варианте  $X_{i+1} = X_i \cup Z_i, V_{i+1} = \{v_i(x_{i+1,j}) + \Delta v_i(x_{i+1,j})\}, j = \overline{0, n_{i+1}}$ . Однако в общем случае произвольного выбора сеток такая процедура может привести к очень быстрому росту числа узлов  $X_{i+1}$  и, как следствие, к большой вычислительной трудоемкости. Поэтому можно выбирать  $X_{i+1}$  из какого-либо параметрического семейства сеток.

Приведем теперь некоторые результаты практического применения описанного метода. В качестве тестовой возьмем двумерную линейную задачу

$$u' = Qu, u(0) = (2, -2)^T, x \in [0, h],$$

где  $Q = \begin{pmatrix} -3 & -4 \\ 2 & 3 \end{pmatrix}, h=1$ . Задача решалась с использованием в качестве  $A_i$  оператора кусочно-линейной аппроксимации по  $N$  равноотстоящим узлам. На  $i$  й итерации вычислялись следующие величины:  $\Delta_i$  - норма добавки в (11);  $\mu_\delta(h)$  - норма главной части погрешности;  $\check{\mu}(h)$  - оценка снизу для нормы погрешности (см. (7));  $M_L(h)$  и  $M_r(h)$  - оценки сверху (8) с использованием соответственно классической и односторонней констант Липшица;  $m_\epsilon(h)$  - норма погрешности, вычисленная на точном решении. Во всех случаях при оценке погрешностей использовалась евклидова норма. Результаты эксперимента приведены в таблицах 1-3.

N=10

Таблица 1

$i$	$\Delta_i$	$\mu_\delta(h)$	$\check{\mu}(h)$	$M_L(h)$	$M_r(h)$	$m_\epsilon(h)$
10	0,209441	0,40122	0,0560185	2,61745	0,879389	0,495998
20	0,0717992	0,0926385	0,0129342	0,520613	0,175141	0,106122
50	0,00166674	0,0035819	0,00050011	0,247205	0,0255775	0,006499
100	$1,889 \cdot 10^{-6}$	0,0046942	0,00065540	0,248534	0,026423	0,007928

N=100

Таблица 2

$i$	$\Delta_i$	$\mu_\delta(h)$	$\check{\mu}(h)$	$M_L(h)$	$M_r(h)$	$m_\epsilon(h)$
10	0,209394	0,399464	0,05577	2,47761	0,865879	0,494225
20	0,071415	0,093682	0,01308	0,33354	0,162197	0,109318
50	0,001606	0,001091	0,000152	0,004030	0,0016535	0,001169
100	$1,65 \cdot 10^{-6}$	0,000038	$5,32 \cdot 10^{-6}$	0,002060	0,0002176	0,000065

$i$	$\Delta_i$	$\mu_\delta(h)$	$\tilde{\mu}(h)$	$M_L(h)$	$M_\varepsilon(h)$	$m_\varepsilon(h)$
10	0,209394	0,399449	0,0557712	2,47712	0,865814	0,494211
20	0,0714118	0,0936906	0,0130811	0,33272	0,162172	0,109345
50	0,0016054	0,0011282	0,0001575	0,002207	0,001528	0,001232
100	$1,64 \cdot 10^{-6}$	$2,86 \cdot 10^{-7}$	$3,99 \cdot 10^{-8}$	0,000021	$2,634 \cdot 10^{-6}$	$6,06 \cdot 10^{-8}$

Полученные результаты позволяют сделать следующие выводы.

Во-первых, как и следовало ожидать, точность приближенного решения повышается при увеличении числа уточнений  $i$  и числа узлов  $N$ . Это дает повод надеяться на то, что, управляя этими двумя рычагами, на практике можно достигать любой разумной точности на заданном шаге.

Во-вторых, сравнение значений  $M_L(h)$  и  $M_\varepsilon(h)$  показывает, что использование односторонней константы Липшица вместо классической делает верхнюю оценку погрешности существенно точнее. Несмотря на это, такая оценка обычно является достаточно завышенной. Что касается оценки снизу  $\tilde{\mu}(h)$ , то на данной задаче она дает неплохие результаты.

В-третьих, достаточно ясно прослеживается связь между величиной  $\mu_\delta(h)$ , которая фактически представляет собой норму главной части погрешности, и нормой погрешности  $m_\varepsilon(h)$ . Принимая во внимание также то, что значение  $\mu_\delta(h)$  может быть эффективно вычислено, представляется разумным на практике использовать его вместо (7) - (9) для оценки  $m_\varepsilon(h)$ .

В заключение сделаем еще одно наблюдение, основанное на сравнении погрешности и ее оценок для разных  $N$  при  $i=50$ . Мы видим, что при большой разнице в уровне аппроксимации (а следовательно, и в вычислительной трудоемкости) точность приближения практически одинакова. Следовательно, при высоких требованиях к точности можно избежать чрезмерных вычислительных затрат, если начинать процесс приближений с низкого уровня аппроксимации и затем повышать его по мере необходимости.

1. Хайрер Э., Нёрсетт С, Ваннер Г. Решение обыкновенных дифференциальных уравнений. Нежесткие задачи. М, 1990.
2. Хайрер Э., Ваннер Г. Решение обыкновенных дифференциальных уравнений. Жесткие и дифференциально-алгебраические задачи. М., 1999.

Поступила в редакцию 08.11.04.

**Владимир Васильевич Бобков** - доктор физико-математических наук, профессор кафедры вычислительной математики.

**Борис Викторович Фалейчик** - аспирант кафедры вычислительной математики. Научный руководитель - В.В. Бобков.