УДК 518:517.944/.947

выбор итерационных параметров в методе ричардсона

Е. С. НИКОЛАЕВ, А. А. САМАРСКИЙ

(Москва)

Предлагается упорядочение набора итерационных параметров метода Ричардсона, для которого он становится численно устойчивым. Число параметров произвольно.

В работе рассмотрен вопрос вычислительной устойчивости итерационного метода Ричардсона решения операторного уравнения I рода в гильбертовом пространстве. Этот метод обладает высокой скоростью сходимости, однако для задач с плохо обусловленным оператором на практике была выявлена его численная неустойчивость [1-3].

В [4] показано, что неустойчивость метода связана с порядком использования итерационных параметров и что предлагавшиеся ранее [$^{1-3}$] способы упорядочения набора параметров не устраняют численную неустойчивость, а только уменьшают ее.

Дальнейшие исследования показали, что существует такой порядок в наборе итерационных параметров, для которого метод становится численно устойчивым. Этот порядок был предложен в [4] и [5], гл. VIII, для случая, когда число параметров $n=2^p$.

В настоящей статье метод упорядочения параметров, изложенный в [5], обобщается на случай произвольного числа параметров n.

В § 1 дается описание метода Ричардсона и упорядочения набора параметров для случая $n=2^p$. Там же формулируются теоремы о вычислительной устойчивости. Подробное доказательство теорем будет дано отдельно. § 2 посвящен заданию порядка в наборе параметров для случая произвольного n. В § 3 приведены результаты экспериментального изучения вычислительной устойчивости метода с описанным набором параметров.

§ 1. Постановка задачи

1. Пусть в вещественном гильбертовом пространстве H дано операторное уравнение I рода с самосопряженным оператором ($A=A^*>0$)

$$(1.1) Au = f,$$

где f — заданный, u — искомый элемент из H.

Для приближенного решения задачи (1.1) рассмотрим неявную двухслойную итерационную схему

(1.2)
$$B \frac{y_{h+1} - y_h}{\tau_{h+1}} + Ay_h = f, \quad k = 0, 1, \dots, n-1,$$

с произвольным начальным приближением $y_0 \in H$. Семейство схем (1.2) выделим условием

$$(1.3) B = B^* \geqslant \beta E, \beta > 0,$$

и будем предполатать, что операторы A и B энергетически эквивалентны [5] с постоянными γ_1 и γ_2 :

$$(1.4) \gamma_1 B \leqslant A \leqslant \gamma_2 B, \gamma_2 > \gamma_1 > 0.$$

В предположениях (1.3), (1.4) решение задачи об оптимальном наборе итерационных параметров τ_h имеет вид [5]

(1.5)
$$\tau_{k} = \frac{\tau_{0}}{1 + \rho_{0}\mu_{k}},$$

$$\mu_{k} \in \mathfrak{M}_{n} = \left\{\cos\frac{2i - 1}{2n}\pi, \quad i = 1, 2, ..., n\right\}, \qquad k = 1, 2, ..., n,$$

где \mathfrak{M}_n — множество из n элементов, расположенных по возрастанию i,

$$\tau_{\scriptscriptstyle 0} = \frac{2}{\gamma_{\scriptscriptstyle 1} + \gamma_{\scriptscriptstyle 2}} \,, \qquad \rho_{\scriptscriptstyle 0} = \frac{1 - \xi}{1 + \xi} \,, \qquad \xi = \frac{\gamma_{\scriptscriptstyle 1}}{\gamma_{\scriptscriptstyle 2}} \,. \label{eq:tau_0}$$

При таком наборе параметров справедливы оценки

$$\|y_n - u\|_{D} \leqslant q_n \|y_0 - u\|_{D}, \qquad D = A$$
 или B .

Здесь норма в энергетическом пространстве $H_{\scriptscriptstyle D}$ определяется следующим образом:

$$\|x\|_D = V(Dx, x)$$
 для $D^* = D > 0$, $x \in H$, $q_n = \frac{2\rho_1^n}{1 + \rho_1^{2n}}$, $\rho_1 = \frac{1 - V\xi}{1 + V\xi}$.

Для уменьшения нормы начальной погрешности в пространстве $H_{\rm A}$ ($H_{\rm B}$) в $1/\epsilon$ раз достаточно n итераций:

$$n = n(\varepsilon, \xi) \approx \frac{\ln 0.5\varepsilon}{\ln \rho_1} \approx \frac{|\ln 0.5\varepsilon|}{2V\xi}.$$

Схему (1.2) с набором параметров (1.5) называют неявным итерационным методом Ричардсона.

2. При изучении сходимости метода (1.2) мы предполагали, что вычислительный процесс является идеальным, т. е. счет ведется с бесконечным числом знаков. Однако процесс округления результатов арифметических операций вносит в решение y_k некоторые погрешности на каждом этапе вычислений. Будем считать, что введение этих погрешностей эквивалентно возмущению входных данных задачи — начального приближения, правой части и операторов A и B.

Тогда реальное решение \widetilde{y}_k можно рассматривать, как точное решение задачи

$$(1.6) \widetilde{B} \frac{\widetilde{y}_{k+1} - \widetilde{y}_k}{\tau_{k+1}} + \widetilde{A}\widetilde{y}_k = \widetilde{f}_{k+1} + \frac{1}{\tau_{k+1}} \widetilde{w}_{k+1}, k = 0, 1, \ldots, n,$$

 $\widetilde{\boldsymbol{y}}_{0}$ задано.

При таком подходе задача о вычислительной погрешности метода сводится к исследованию устойчивости схемы (1.6) с возмущенными операторами \widetilde{A} и \widetilde{B} по начальным данным и правой части.

Вторым вопросом численной устойчивости метода (1.2) является попрос о росте промежуточных решений y_h при различных упорядочениях множества \mathfrak{M}_n . Исследование устойчивости схемы (1.6) позволяет получить оценку для величины промежуточного решения.

Некоторые примеры, показывающие влияние порядка использования параметров τ_k на рост решения и точность метода, приведены в [5] и § 3 этой работы.

Будем предполагать, что схема (1.6) принадлежит к исходному семейству схем, т. е. выполнены условия

(1.7)
$$\widetilde{A} = \widetilde{A}^* > 0, \quad \widetilde{B} = \widetilde{B}^* \geqslant \widetilde{\beta}E, \quad \widetilde{\beta} > 0.$$

Мерой возмущения операторов A и B будем считать относительное изменение их энергии (0 $\leqslant \alpha_1, \alpha_2 < 1$)

$$(1.8) \qquad |((A-\widetilde{A})x,x)| \leqslant \alpha_1(Ax,x), \quad |((B-\widetilde{B})x,x)| \leqslant \alpha_2(Bx,x).$$

Для исследования роста промежуточных решений \widetilde{y}_k перейдем от неявной схемы (1.6) к эквивалентной ей явной

$$x_{k+1} = \widetilde{S}_{k+1}x_k + au_{k+1}\phi_{k+1} + \psi_{k+1}, \qquad k = 0, 1, \dots, n-1,$$
 (1.9) $S_{k+1} = E - au_{k+1}\widetilde{C},$ где $x_k = \widetilde{B}^{-1/2}\widetilde{y}_k, \qquad \widetilde{C} = \widetilde{B}^{-1/2}\widetilde{A}\widetilde{B}^{-1/2}, \qquad \phi_k = \widetilde{B}^{-1/2}\widetilde{f}_k,$ $\psi_k = \widetilde{B}^{-1/2}\widetilde{w}_k.$

Для изучения вычислительной точности метода рассмотрим задачу для погрешности $z_k = \tilde{y}_k - u$:

$$\widetilde{B} \frac{z_{k+1} - z_k}{\tau_{k+1}} + \widetilde{A} z_k = \widetilde{f}_{k+1} - f + \frac{1}{\tau_{k+1}} \widetilde{w}_{k+1} + (A - \widetilde{A}) u,$$

$$z_0 = \widetilde{y}_0 - u.$$

Эквивалентная явная схема имеет вид

(1.10)
$$x_{k+1} - \varphi = \widetilde{S}_{k+1}(x_k - \varphi) + \tau_{k+1}\varphi_{k+1} + \psi_{k+1}, \quad k = 0, 1, \dots, n-1,$$

$$\widetilde{S}_{k+1} = E - \tau_{k+1}\widehat{C},$$

где
$$x_k=ar{B}^{\prime_k}z_k, \qquad ilde{C}=ar{B}^{-\prime_k}\widetilde{A}\widetilde{B}^{-\prime_k}, \qquad \phi_k=ar{B}^{-\prime_k}(ilde{f}_k-f), \ \psi_k=ar{B}^{-\prime_k}\widetilde{w}_k, \qquad \phi=ar{B}^{\prime_k}\widetilde{A}^{-1}(A-\widetilde{A})u.$$

Из (1.9), (1.10) следует, что для получения необходимых нам оценок достаточно исследовать устойчивость схемы (1.9) по начальным данным и правой части.

Пользуясь рекуррентной формулой (1.9), найдем

$$x_{k} = \widetilde{T}_{k,0}x_{0} + \sum_{j=1}^{k} \tau_{j}\widetilde{T}_{k,j}\psi_{j} + \sum_{j=1}^{k} \widetilde{T}_{k,j}\psi_{j}, \quad k = 1, 2, \dots, n$$

$$\widetilde{T}_{k,j} = \prod_{i=j+1}^{k} \widetilde{S}_{i}, \quad \widetilde{T}_{k,k} \equiv E.$$

Оператор $T_{k,j}$ называется разрешающим оператором со слоя j на слой k. О пределение. Назовем набор параметров $\{\tau_k\}$ устойчивым, если существуют положительные C_1 , C_2 , C_3 , зависящие, быть может, от γ_1 , γ_2 , α_1 , α_2 , но не зависящие от n, такие, что

$$\max_{1\leqslant k\leqslant n} \|\widetilde{T}_{k,0}\| \leqslant C_1, \qquad \max_{1\leqslant k\leqslant n} \sum_{j=1}^k \tau_j \|\widetilde{T}_{k,j}\| \leqslant C_2, \qquad \max_{1\leqslant k\leqslant n} \sum_{j=1}^k \|\widetilde{T}_{k,j}\| \leqslant C_3.$$

Для устойчивого набора параметров при любом k справедливы оценки

$$\|x_h\| \leqslant C_1 \|x_0\| + C_2 \max_{1 \leqslant j \leqslant h} \|\varphi_j\| + C_3 \max_{1 \leqslant j \leqslant h} \|\psi_j\|$$

для задачи (1.9) и

$$\|x_h\| \leqslant C_1 \|x_0\| + C_2 \max_{1 \leqslant j \leqslant h} \|\varphi_j\| + C_3 \max_{1 \leqslant j \leqslant h} \|\psi_j\| + (1 + C_1) \|\varphi\|$$

для задачи (1.10), из которых следует, что схемы (1.9), (1.10) устойчивы.

3. Перейдем к заданию порядка элементов в множестве \mathfrak{M}_n , который порождает устойчивый набор итерационных параметров $\{\tau_k\}$ метода Ричардсона. Рассмотрим здесь случай, когда $n=2^p$, p>0.

При ностроении набора параметров $\{\tau_k\}$ будем, следуя [5], исходить из минимального $\beta=\beta_1=\pi/2n$ и рекуррентным образом строить множества

$$M_1(eta) = \{-\coseta\},$$
 (1.12) $M_{2^k}(eta) = M_{2^{k-1}}(eta) \cup M_{2^{k-1}}(\pi/2^{k-1}-eta), \quad k=1,\ 2,\dots,p.$ Тогда (1.13) $\mathfrak{M}_{2^p} = M_{2^p}(eta_1).$

Следующая теорема утверждает, что множество \mathfrak{M}_{2^p} , упорядоченное таким образом, порождает устойчивый набор $\{\tau_k\}$ (μ_k в формуле (1.5) есть k-й элемент множества \mathfrak{M}_{2^p}).

Теорема 1. Если выполнены условия (1.3), (1.4) и множество \mathfrak{M}_{2^p} упорядочено согласно (1.13), то при условии $\mathfrak{a}_1=\mathfrak{a}_2=0$ справедливы сле-

дующие оценки, не зависящие от $n\ (m=j\cdot 2^k,j-$ нечетное число, $k\geqslant 0)$:

$$||T_{m,0}|| \leqslant 1/\xi, \qquad m = 1, 2, \dots, 2^{p},$$

$$\sum_{i=1}^{m} \tau_{i} ||T_{m,i}|| \leqslant \frac{1}{\gamma_{1}} \left[1 + \left(1 - \delta_{1,i} \right) \frac{1}{\xi^{2}} \right],$$

$$\sum_{i=1}^{m} ||T_{m,i}|| \leqslant \frac{1}{\xi} \left[1 + (1 - \delta_{1,j}) \frac{1}{\xi^{2}} \right],$$

где $\delta_{i,j}$ — символ Кронекера. Для $m=2^p$ справедливы более точные оценки

$$(1.14) ||T_{2^{p},0}|| \leqslant q_{2^{p}} < 1, \sum_{j=1}^{2^{p}} ||T_{2^{p},j}||_{\tau_{j}} \leqslant \frac{1 - q_{2^{p}}}{\gamma_{1}} < \frac{1}{\gamma_{1}},$$

(1.15)
$$\sum_{j=1}^{2^{p}} \|T_{2^{p},j}\| \leqslant \frac{4}{3^{\gamma} \xi}.$$

Теорема 1 выражает тот факт, что если в схеме (1.6) пренебречь возмущением операторов A и B, то промежуточные решения ограничены по норме $(m=j\cdot 2^k)$

$$\|\widetilde{y}_{m}\|_{B} \leqslant \frac{1}{\xi} \|\widetilde{y}_{0}\|_{B} + \left[1 + (1 - \delta_{1,j}) \frac{1}{\xi^{2}}\right] \left(\frac{1}{\gamma_{1}} \max_{1 \leqslant i \leqslant m} \|f_{i}\|_{B^{-1}} + \frac{1}{-\xi} \max_{1 \leqslant i \leqslant m} \|\widetilde{w}_{i}\|_{B^{-1}}\right)$$

и для погрешности решения после п итераций справедлива оценка

Теорема 2. Если выполнены условия (1.3), (1.4), (1.7), (1.8) и множество \mathfrak{M}_{2}^{r} упорядочено согласно (1.13), то при условии

$$\alpha = \frac{\alpha_1 + \alpha_2}{1 - \alpha_2} \leqslant \frac{\xi}{2}$$

для погрешности решения задачи (1.6) после п итераций справедлива оценка

$$\begin{split} (1.17) \qquad & \|\widetilde{y}_n - u\|_{\widetilde{B}} \leqslant \frac{q_n}{\widetilde{q}_n} \|\widetilde{y}_0 - u\|_{\widetilde{B}} + \frac{1}{\widetilde{\gamma}_1} \left(1 - \frac{q_n}{\widetilde{q}_n}\right) \max_{1 \leqslant i \leqslant n} \|\widetilde{f}_i - f\|_{\widetilde{B}^{-1}} + \\ & + \frac{4}{3} \frac{1 + \sqrt{\alpha}}{\sqrt{\xi - \sqrt{\alpha}}} \max_{1 \leqslant i \leqslant n} \|\widetilde{w}_i\|_{\widetilde{B}^{-1}} + \frac{\alpha_1}{\widetilde{\gamma}_1} \left(1 + \frac{q_n}{\widetilde{q}_n}\right) \|f\|_{\widetilde{B}^{-1}}. \end{split}$$

З∂есь

$$\begin{split} &\tilde{\gamma}_{\scriptscriptstyle 1} = \gamma_{\scriptscriptstyle 1} - \alpha \gamma_{\scriptscriptstyle 2} \geqslant 0.5 \; \gamma_{\scriptscriptstyle 1}, \\ &\tilde{q}_{\scriptscriptstyle n} = \frac{2 \tilde{\rho}_{\scriptscriptstyle 1}{}^{\scriptscriptstyle n}}{1 + \tilde{\rho}_{\scriptscriptstyle 1}{}^{\scriptscriptstyle 2n}} > q_{\scriptscriptstyle n}, \qquad \tilde{\rho}_{\scriptscriptstyle 1} = \frac{1 - \mathcal{V} \xi_{\scriptscriptstyle 0}}{1 + \mathcal{V} \xi_{\scriptscriptstyle 0}}, \qquad \xi_{\scriptscriptstyle 0} = \frac{\alpha}{1 + \alpha - \xi} \;. \end{split}$$

Теорема 2 выражает вычислительную устойчивость итерационного метода Ричардсона для случая, когда $n=2^p$. Оценка (1.17) переходит при

 α_1 , $\alpha_2 \to 0$ в полученную ранее (1.16). Оценки (1.14), даваемые теоремой 1, являются неулучшаемыми для любых упорядочений множества \mathfrak{M}_n , а (1.15) точна по порядку малости ξ при больших n (ср. с леммой 6 из [4]).

Замечание. Если предположить, что ψ_j в (1.11) имеет вид $\psi_j=T_{j,\;0}v_j$, то такое возмущение эквивалентно возмущению x_0 на величину $\sum\limits_{j=1}^{k}v_j$.

§ 2. Построение набора параметров для произвольного $n^{*)}$

1. Идея построения порядка в множестве \mathfrak{M}_n для произвольного n основана на двух соображениях: на переходе от упорядоченного множества \mathfrak{M}_{2k} к заданию порядка в множестве \mathfrak{M}_{2k+1} и на переходе от \mathfrak{M}_n к \mathfrak{M}_{2k} , где k — произвольное целое число.

Предположим, что множество \mathfrak{M}_{2k} уже упорядочено необходимым нам образом. Тогда представим его в виде следующей суммы:

$$\mathfrak{M}_{2k} = \bigcup_{i=1}^{2k} M_1(\beta_{2k,i}) = M_1(\beta_{2k,1}) \cup M_1(\beta_{2k,2}) \cup \ldots \cup M_1(\beta_{2k,2k}),$$

где $\beta_{2h,i}$ при любом i принадлежит множеству

$$\left\{ \frac{2j-1}{4k} \pi, \ j=1,2,\ldots,2k \right\}.$$

Тогда множество \mathfrak{M}_{2k+1} упорядочим следующим образом:

$$\mathfrak{M}_{2k+1} = \bigcup_{i=1}^{2k} M_1(\beta_{2k+1,i}) \cup M_1\left(\frac{\pi}{2}\right),$$

где $\beta_{2k+1,\ i}$ — ближайшее к $\beta_{2k,\ i}$ число из множества

$$\left\{\frac{2j-1}{2(2k+1)}\pi, \ j=1,2,\ldots,2k+1, \ j\neq k+1\right\}.$$

Переход от \mathfrak{M}_{2k} к \mathfrak{M}_{2k+1} осуществлен.

Пусть, далее, упорядочено множество \mathfrak{M}_h :

$$\mathfrak{M}_{k} = \bigcup_{i=1}^{k} M_{1}(\beta_{k,i}) = M_{1}(\beta_{k,1}) \cup \ldots \cup M_{1}(\beta_{k,k}),$$
$$\beta_{k,i} \in \left\{ \frac{2j-1}{2k} \pi, \ j=1,2,\ldots,k \right\}.$$

Тогда, используя формулы (1.12), зададим порядок в множестве \mathfrak{M}_{2k} следующим образом:

$$\mathfrak{M}_{2k} = \bigcup_{i=1}^{k} M_{2}(\beta_{2k,i}) = \bigcup_{i=1}^{k} (M_{1}(\beta_{2k,i}) \cup M_{1}(\pi - \beta_{2k,i})),$$

$$\beta_{2k,i} = 0.5\beta_{k,i}, \quad i = 1, 2, \dots, k.$$

^{*)} Как сообщил нам В. И. Лебедев, в $[^{6,7}]$ предложены другие наборы параметров для любого n.

Эти соображения позволяют перейти от упорядоченного множества $\mathfrak{M}_1 = M_1(\pi/2)$, состоящего из одного элемента, к множеству \mathfrak{M}_n с произвольным n, чередуя должным образом переходы от множества с четным числом элементов к множеству с печетным числом и от множества из k элементов к множеству из 2k элементов.

Эта процедура упорядочения множества \mathfrak{M}_n для произвольного n может быть формализована следующим образом.

Представим n в виде разложения в сумму по степеням 2 с целыми показателями k_i :

$$n = 2^{h_1} + 2^{h_2} + \ldots + 2^{h_t}, \quad k_j \leqslant k_{j-1} - 1, \quad k_t \geqslant 0.$$

Здесь t — целый индекс. Введем величины

(2.1)
$$n_j = \sum_{i=1}^j 2^{k_i - k_j}, \quad \delta_j = \frac{n_j}{n} 2^{k_j}, \quad j = 1, 2, \dots, t,$$

и положим $k_{t+1} = -1$. Отметим, что все n_j — нечетные числа.

При построении набора параметров будем исходить из минимального

$$\beta = \beta_1 = \pi / 2n$$
.

Образуем упорядоченную сумму множеств

$$(2.2) M_n(\beta) = \bigcup_{i=1}^t M_{2^{k_i}}(n_i\beta) = M_{2^{k_i}}(n_1\beta) \cup \ldots \cup M_{2^{k_t}}(n_t\beta),$$

где $M_{2^k}(\beta)$ определяется рекуррентным образом $(j=1,2,\ldots,t)$:

$$M_1(\beta) = \{-\cos \beta\},\$$

(2.3)

$$M_{2^k}(\beta) = M_{2^{k-1}}(\beta) \cup M_{2^{k-1}}\left(\frac{\pi}{2^{k-1}}\delta_j - \beta\right),$$

если $k_{i+1} + 2 \leqslant k \leqslant k_i + 1$.

Тогда

$$\mathfrak{M}_n = M_n(\beta_1).$$

Заметим, что если $n=2^p$, то t=1, $n_1=1$, $\delta_1=1$, $k_1=p$. При этом формула (2.2) переходит в

$$M_n(\beta) = M_{2p}(\beta)$$

и рекуррентные соотношения (2.3), определяющие $M_{zp}(\beta)$, переходят в приведенные выше для $n=2^p$ формулы (1.12). Следовательно, упорядочение множества \mathfrak{M}_n согласно (2.4) является обобщением построений для случая $n=2^p$.

Мы изложим сейчас алгоритм, позволяющий упорядочить множество \mathfrak{M}_n в соответствии с (2.4). Построение множества \mathfrak{M}_n непосредственно по формулам (2.2), (2.3) затруднительно; мы используем их для описания упорядочения \mathfrak{M}_n в теоремах 1, 2.

Пусть θ_m — множество из m целочисленных элементов

$$\theta_m = \{\theta_m(1), \theta_m(2), \dots, \theta_m(m)\}.$$

Положим $n_{t+1} = 2n + 1$. Пусть j = 1. Строим множества

(2.5)
$$\theta_{n_i} = \{\theta_{n_i}(i) = \theta_{n_i-1}(i), \quad i = 1, 2, \dots, n_j - 1; \ \theta_{n_i}(n_j) = n_j\},$$

(2.6)
$$\theta_{2m} = \{\theta_{2m}(2i) = 4m - \theta_m(i), \ \theta_{2m}(2i-1) = \theta_m(i), \ i=1, 2, \dots, m\}, \\ m = n_i, 2n_i, 4n_i, \dots, 0.25(n_{i+1}-1).$$

Если j=t, то необходимое множество θ_n уже построено, иначе строится множество

(2.7)
$$\theta_{n_{j+1}-1} = \{\theta_{n_{j+1}-1}(2i) = 2n_{j+1} - \theta_{0.5(n_{j+1}-1)}(i), \quad \theta_{n_{j+1}-1}(2i-1) = \theta_{0.5(n_{j+1}-1)}(i), \quad i = 1, 2, \dots, 0.5(n_{j+1}-1)\}.$$

Затем j увеличивается на 1 и процесс повторяется, начиная с (2.5). В результате будет построено множество θ_n .

Тогда

(2.8)
$$\mathfrak{M}_n = \left\{ -\cos \beta_i, \ i = 1, 2, \dots, n, \ \beta_i = \theta_n(i) \frac{\pi}{2n} \right\}$$

и μ_k в формуле (1.5) есть k-й элемент множества \mathfrak{M}_n .

Для случая $n=2^p$ алгоритм (2.5)-(2.7) упрощается:

$$\theta_{1} = \{\theta_{1}(1) = 1\},\$$

$$\theta_{2m} = \{\theta_{2m}(2i) = 4m - \theta_{m}(i), \ \theta_{2m}(2i - 1) = \theta_{m}(i), \ i = 1, 2, \ldots, m\},\$$

$$m = 1, 2, 4, \ldots, 2^{p-1},$$

и после нахождения θ_{2^p} множество \mathfrak{M}_{2^p} строится согласно (2.8). Приведем некоторые примеры:

$$\begin{split} \theta_8 &= \{1,\, 15,\, 7,\, 9,\, 3,\, 13,\, 5,\, 11\},\\ \theta_9 &= \{1,\, 17,\, 7,\, 11,\, 3,\, 15,\, 5,\, 13,\, 9\},\\ \theta_{12} &= \{1,\, 23,\, 11,\, 13,\, 5,\, 19,\, 7,\, 17,\, 3,\, 21,\, 9,\, 15\},\\ \theta_{16} &= \{1,\, 31,\, 15,\, 17,\, 7,\, 25,\, 9,\, 23,\, 3,\, 29,\, 13,\, 19,\, 5,\, 27,\, 11,\, 21\},\\ \theta_{18} &= \{1,\, 35,\, 17,\, 19,\, 7,\, 29,\, 11,\, 25,\, 3,\, 33,\, 15,\, 21,\, 5,\, 31,\, 13,\, 23,\, 9,\, 27\}, \end{split}$$

§ 3. Вычислительная устойчивость метода

1. Для изучения устойчивости набора итерационных параметров $\{\tau_k\}$, построенного в соответствии с упорядочением множества \mathfrak{M}_n по формуле (2.4), используется численный эксперимент.

Исследование характера итерационного метода, очевидно, может быть проведено на простейшей модели, так как характер процесса определяется не конкретным видом оператора A, а его основными функциональными свойствами как оператора в гильбертовом пространстве.

Эксперимент на модельной задаче позволяет сравнить теоретические оценки, полученные для случая $n=2^p$ в теоремах 1, 2, с численными результатами.

В качестве модельной задачи была выбрана разностная аппроксимация краевой задачи

(3.1)
$$\frac{d^{4}v}{dx^{4}} = f(x), \qquad 0 < x < 1, \qquad v(0) = v_{1}, \\ \frac{d^{2}v}{dx^{2}}(0) = v_{2}, \qquad v(1) = v_{3}, \qquad \frac{d^{2}v}{dx^{2}}(1) = v_{4}.$$

На равномерной сетке с шагом h = 1/N строится разностная схема, аппроксимирующая задачу (3.1) с погрешностью $O(h^2)$:

(3.2)
$$\begin{array}{ll} u_{\overline{x}x\overline{x}x} = f, & 2h \leqslant x \leqslant 1 - 2h, \\ u(0) = v_1, & u_{\overline{x}x} - hu_{\overline{x}xx} = v_2, & x = h, \\ u(1) = v_3, & u_{\overline{x}x} + hu_{\overline{x}x\overline{x}} = v_4, & x = 1 - h. \end{array}$$

Оператор A, соответствующий задаче (3.2), является самосопряженным в пространстве H сеточных функций, определенных во внутренних узлах сетки. Скалярное произведение в H задается обычным образом:

$$(u,v) = \sum_{n=1}^{1-h} u(x)v(x)h.$$

При этом собственными функциями оператора A будут $\mu_k(x) = \sin k\pi x$ и соответствующими собственными значениями будут

$$\lambda_k = \frac{16}{h^4} \sin^4 \frac{k\pi h}{2}, \qquad k = 1, 2, \dots, N-1.$$

Заметим, что $\|A\|=\lambda_{N-1}\approx 16\,/\,h^4=1.6\cdot 10^5$ при N=10 и $\|A\|\approx 1.6\cdot 10^9$ при N=100.

Выбор задачи (3.2) в качестве объекта эксперимента позволяет моделировать на грубой сетке плохо обусловленные операторы A.

Рассматривался явный итерационный процесс (1.2) ($B \equiv E$). Тогда в условиях (1.4) постоянные энергетической эквивалентности γ_1 и γ_2 есть $\gamma_1 = \lambda_1$, $\gamma_2 = \lambda_{N-1}$. При этом

$$\xi = \operatorname{tg}^4 \frac{\pi h}{2}.$$

2. Проводилось несколько серий экспериментов. В первой серии изучалось влияние упорядочения множества \mathfrak{M}_n на рост промежуточных решений и достигаемую после n итераций точность.

В модельной задаче (3.2) бралось $v_1=1, v_2=v_3=v_4=0, f\equiv 0$. При таких входных данных точное решение задачи есть u(x)=1-x.

Пусть n кратно 8. Рассматривались следующие упорядочения множества \mathfrak{M}_n :

1)
$$\mathfrak{M}_{n} = \bigcup_{k=1}^{n} M_{1} \left(\frac{2k-1}{2n} \pi \right) =$$

$$= M_{1} \left(\frac{\pi}{2n} \right) \cup M_{1} \left(\frac{3\pi}{2n} \right) \cup \ldots \cup M_{1} \left(\frac{2n-1}{2n} \pi \right);$$

Этому упорядочению соответствует обычный «обратный» набор {ть}:

$$au_{h} = rac{ au_{0}}{1 +
ho_{0}\mu_{h}}, \qquad \mu_{h} = -\cosrac{2k-1}{2n}\pi, \qquad k = 1, 2, \dots, n;$$
2) $extit{$\mathfrak{M}_{n} = igcup_{h=1}^{n} M_{1}\left(\pi - rac{2k-1}{2n}\pi\right)$, здесь $\mu_{h} = \cosrac{2k-1}{2n}\pi$,}$

что соответствует обычному «прямому» набору параметров $\{\tau_k\}$;

3)
$$\mathfrak{M}_{n} = \bigcup_{k=1}^{n/2} M_{2} \left(\frac{2k-1}{2n} \pi \right) =$$

$$= \bigcup_{k=1}^{n/2} \left(M_{1} \left(\frac{2k-1}{2n} \pi \right) \cup M_{1} \left(\pi - \frac{2k-1}{2n} \pi \right) \right) ;$$

такому упорядочению соответствует разбиение набора $\{\tau_k\}$ на блоки из двух элементов $(\mu_{2k-1} = -\cos [(2k-1)/2n]\pi, \mu_{2k} = \cos [(2k-1)/2n]\pi)$; здесь и ниже используются рекуррентные формулы (1.12);

4)
$$\mathfrak{M}_{n} = \bigcup_{k=1}^{n/2} M_{2} \left(\frac{\pi}{2} - \frac{2k-1}{2n} \pi \right) =$$

$$= \bigcup_{k=1}^{n/2} \left(M_{1} \left(\frac{\pi}{2} - \frac{2k-1}{2n} \pi \right) \cup M_{1} \left(\frac{\pi}{2} + \frac{2k-1}{2n} \pi \right) \right).$$

Упорядочение, рекомендованное в [1], в наших обозначениях принимает вид

$$\mathfrak{M}_{n} = \bigcup_{k=1}^{n/2} M_{2} \left(\frac{\pi}{2} + \frac{2k-1}{2n} \pi \right);$$

$$\mathfrak{M}_{n} = \bigcup_{k=1}^{n/4} M_{4} \left(\frac{2k-1}{2n} \pi \right) =$$

$$= \bigcup_{k=1}^{n/4} \left(M_{4} \left(\frac{2k-1}{2n} \pi \right) \cup M_{4} \left(\pi - \frac{2k-1}{2n} \pi \right) \cup M_{4} \left(\frac{\pi}{2} + \frac{2k-1}{2n} \pi \right) \right).$$

Здесь производится организация блоков из 4 элементов:

$$\left(-\cos\beta_{k}, \cos\beta_{k}, -\sin\beta_{k}, \sin\beta_{k}; \quad \beta_{k} = \frac{2k-1}{2n}\pi\right);$$

$$\mathfrak{M}_{n} = \bigcup_{k=1}^{n/8} M_{8} \left(\frac{2k-1}{2n}\pi\right) =$$

$$= \bigcup_{k=1}^{n/8} \left(M_{4} \left(\frac{2k-1}{2n}\pi\right) \cup M_{4} \left(\frac{\pi}{4} - \frac{2k-1}{2n}\pi\right)\right);$$

здесь блок состоит из 8 элементов;

- 7) упорядочение множества \mathfrak{M}_n по формуле (2.4).
- 9 ЖВМ и МФ, № 4

Чтобы выявить влияние обусловленности оператора A на вычислительную устойчивость метода, расчеты проводились на последовательности сеток $h = \frac{1}{10}, \frac{1}{12}, \frac{1}{14}$.

Задавалось n с шагом 8 от 8 до 512. Для каждого n и каждого из укаванных упорядочений \mathfrak{M}_n проводились итерации по схеме (1.2). Определялось критическое значение $n_{\text{точ}}$, для которого реальная относительная точность еще не превышала теоретическую точность, т. е. выполнялось неравенство

$$\frac{\|y_n-u\|}{\|y_0-u\|}=\varepsilon_{\text{peam}}\leqslant q_n.$$

Так же определялось $n_{\text{авост}}$, при котором на какой-то промежуточной итерации решение выходило за границу максимально возможного для представления в машине числа (10^{19}) .

Расчеты показали, что для упорядочений 1)-6) при увеличении n сначала происходит потеря точности, когда n становится больше $n_{\text{точ}}$, а затем аварийный останов, вызываемый ростом промежуточных решений.

Использование блоков более сложной структуры приводит к уменьшению численной неустойчивости метода, и чем больше размер блока, тем больше $n_{\rm авост}$.

При использовании упорядочения 7) реальная относительная точность $\varepsilon_{\rm pean}$ для всех n не превышала теоретическую, которая для n=512 равна $q_n=1.4\cdot 10^{-11},\ 3.9\cdot 10^{-8},\ 4.5\cdot 10^{-6}$ для $h={}^4/_{10},\ {}^4/_{12},\ {}^4/_{14}$. Промежуточные решения были ограничены и величина

$$R_n = \max_{\substack{0 \le x \le 1 \\ 1 \le k \le n}} |y_k(x)|$$

была монотонно возрастающей функцией n и для больших n не менялась с ростом n, т. е. выходила на асимптотическое значение.

Расчеты показали, что упорядочение \mathfrak{M}_n согласно (2.4) для n, отличных от степени 2, сохраняет те же характеристики вычислительной устойчивости метода, как и в случае $n=2^p$.

Результаты приведены в табл. 1. Первые строки соответствуют начальному приближению

$$y_0(x) = \begin{cases} 1, & x = 0, \\ 0, & x \neq 0, \end{cases}$$

вторые соответствуют $y_{\scriptscriptstyle 0}(x) = \cos{(\pi x \, / \, 2)}$. Для этих приближений при упорядочении 7)

$$\max_{1 \leqslant n \leqslant 512} R_n = 208, 427, 784 \text{ m} \max_{1 \leqslant n \leqslant 512} R_n = 1.63, 2.73, 4.00$$

для $h = \frac{1}{10}, \frac{1}{12}, \frac{1}{14}$ соответственно.

• 3. Во второй серии находились оценки для

$$\|T_{n,0}\|, \qquad \sum_{j=1}^n \tau_j \|T_{n,j}\|, \qquad \sum_{j=1}^n \|T_{n,j}\|,$$

Таблица 1

	h=1/10,		$\xi = 6.29 \cdot 10^{-4}$		h=1/12,		$\xi = 3.004 \cdot 10^{-4}$		h = 1/14,		$\xi = 1.61 \cdot 10^{-6}$	
	$n_{ m TOH}$	$q_{m{n}}$	^ε реал	naboct	n_{TOH}	$q_{m{n}}$	^є реал	nabocт	$n_{\mathbf{T}04}$	$q_{m{n}}$	εpeaл	павост
1	24 24	0.55 0.55	0.495 0.548	48 48		0.732 0.732	0.729 0.731	48 48			0.773 0.837	40 48
2	24 16		0.481 0.746	72 64			0.713 0.86	64 64			0.765 0.918	64 64
3	48 48	$0.178 \\ 0.178$	0.17 0.177	88 96	48 48		0.364 0.366	80 88			0.376 0.542	80 88
4	48 40		0.171 0.2 6 3	136 136			0.42 0.469	136 128			0.482 0.639	128 128
5	88 104		$2.2 \cdot 10^{-2} \\ 8.68 \cdot 10^{-3}$			$7.2 \cdot 10^{-2} \ 5.4 \cdot 10^{-2}$					$0.141 \\ 0.105$	160 1 76
6	184 184	1.95·10 ⁻⁴ 1.95·10 ⁻⁴	$\begin{bmatrix} 1.24 \cdot 10^{-4} \\ 1.9 \cdot 10^{-4} \end{bmatrix}$	416 456	176 176	4.48·10 ⁻³ 4.48·10 ⁻³	3.74·10 ⁻³ 4.44·10 ⁻³	408 448	176 200	$2.29 \cdot 10^{-2} $ $1.25 \cdot 10^{-2}$	$1.98 \cdot 10^{-2} \\ 9 \cdot 10^{-3}$	

когда \mathfrak{M}_n упорядочено согласно (2.4). Для этого вычислялись величины

$$(3.3) I_1 = ||T_{n, 0}y|| / ||y|| \leqslant ||T_{n, 0}||,$$

(3.4)
$$I_2 = \sum_{j=1}^n \tau_j \frac{\|T_{n,j}y\|}{\|y\|} \leqslant \sum_{j=1}^n \tau_j \|T_{n,j}\|,$$

(3.5)
$$I_3 = \sum_{j=1}^n \frac{\|T_{n,j}y\|}{\|y\|} \leqslant \sum_{j=1}^n \|T_{n,j}\|.$$

Для случая $n=2^p$ в теореме 1 доказано, что равенства в (3.3)—(3.5) достигаются, если y является собственной функцией, соответствующей минимальному собственному значению задачи

$$Ay - \lambda By = 0.$$

При этом справедливы равенства

(3.6)
$$I_1 = ||T_{2^p,0}|| = q_{2^p}, \qquad I_2 = \sum_{j=1}^{2^p} \tau_j ||T_{2^p,j}|| = \frac{1 - q_{2^p}}{\gamma_1}$$

и оценка

(3.7)
$$I_3 = \sum_{j=1}^{2^p} ||T_{2^p,j}|| \leqslant \frac{4}{3^{\gamma} \xi}.$$

Для модельной задачи (3.2) имеем $y = \sin \pi x$.

Кроме того в теореме 1 доказано, что для произвольного n при любом упорядочении множества \mathfrak{M}_n справедлива оценка

$$\sum_{i=1}^n \tau_i \|T_{n,i}\| \geqslant \frac{1-q_n}{\gamma_1}.$$

Таблица 2

$$h = 1/10, \quad \xi = 6.29 \cdot 10^{-4}$$

n	I ₁	q_n	I_2	$(1-q_n)/\gamma_1$	I_3	<i>I₃</i> V ξ	
64	8.0451.10-2	$8.0451 \cdot 10^{-2}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	9.5968 · 10-3	42.726 27.171	1.072 0.6816	
96	1.6174.10-2	$1.6174 \cdot 10^{-2}$	$\begin{array}{ c c c c c c }\hline 1.0268 \cdot 10^{-2} \\ 3.6973 \cdot 10^{-4} \\ \hline\end{array}$	1.0268 · 10-2	45.034 28.641	1.1297 0.718	
128	$3.2467 \cdot 10^{-3}$	$3.2467 \cdot 10^{-3}$	$\begin{array}{ c c c c c }\hline 1.0403 \cdot 10^{-2} \\ 3.8662 \cdot 10^{-4} \\ \hline \end{array}$	1.0403.10-2	47.072 29.933	1.181 0.7509	
192	1.308.10-4	1.308.10-4	$\begin{array}{ c c c c c c }\hline 1.0435 \cdot 10^{-2} \\ 3.8184 \cdot 10^{-4} \\ \hline\end{array}$	1.0435.10-2	$46.5 \\ 29.57$	1.167 0.7418	
256	5.27.10-6	5.27.10-6	$\begin{array}{c c} 1.044 \cdot 10^{-2} \\ 3.868 \cdot 10^{-4} \end{array}$	1.044 · 10-2	47.098 29.95	1.182 0.7513	
344	6.37.10-8	6.37.10-8	$\begin{array}{ c c c c c }\hline 1.044 \cdot 10^{-2} \\ 4.3697 \cdot 10^{-4} \\ \hline\end{array}$	1.044.10-2	53.143 33.768	1.333 0.8471	
384	8.55.10-9	8.55.10-9	$\begin{array}{ c c c c c }\hline 1.044 \cdot 10^{-2} \\ 3.8787 \cdot 10^{-4} \\ \hline \end{array}$	1.044.10-2	47.225 30.03	1.185 0.753	

Таблица 3

$$h = 1/20, \ \xi = 3.84 \cdot 10^{-5}$$

ņ	11	q_n	I_2	$(1-q_n)/\gamma_1$	I_3	<i>I</i> ₃ √ ξ		
64	0.75125	0.75 1 25	$2.5641 \cdot 10^{-3}$ $3.115 \cdot 10^{-5}$	$2.5641 \cdot 10^{-3}$	62.066 39.506	$0.384 \\ 0.245$		
96	0.55725	0.55725	$4.564 \cdot 10^{-3} \\ 4.48 \cdot 10^{-5}$	4.564.10-3	89.331 56.863	$\substack{0.553\\0.352}$		
128	0.39313	0.39313	$6.2558 \cdot 10^{-3} \\ 5.708 \cdot 10^{-5}$	6.2558 • 10-3	113.86 72.474	0.705 0.449		
192	0.1838	0.1838	$\begin{array}{c} 8.4137 \cdot 10^{-3} \\ 7.42 \cdot 10^{-5} \end{array}$	8.4137.10-3	$148.04 \\ 94.234$	$\begin{array}{c} \textbf{0.917} \\ \textbf{0.584} \end{array}$		
25 6	8.3747.10-2	8.3747.10-2	$9.445 \cdot 10^{-3} \\ 8.64 \cdot 10^{-5}$	9.445.10-3	$172.26 \\ 109.65$	1.067 0.679		
344	2.8197 · 10-2	2.8197.10-2	$\substack{1.002 \cdot 10^{-2} \\ 9.88 \cdot 10^{-5}}$	1.002 · 10-2	$\substack{197.03\\125.4}$	$\substack{1.22\\0.777}$		
3 84	1.7181.10-2	1.7181 · 10-2	$1.013 \cdot 10^{-2} \ 9.136 \cdot 10^{-5}$	1.013.10-2	182.23 116.0	1.129 0.718		
51 2	3.5191.10-3	3.5191 10-3	$\begin{array}{c} 1.027 \cdot 10^{-2} \\ 9.56 \cdot 10^{-5} \end{array}$	1.027.10-2	190.66 121.37	1.181 0.752		
768	1.4762.10-4	1.4762-10-4	$\begin{array}{c c} 1.0307 \cdot 10^{-2} \\ 9.43 \cdot 10^{-5} \end{array}$	1.0307.10-2	188.18 119.78	1.166 0.742		
1024	$6.192 \cdot 10^{-6}$	6.192.10-6	$\begin{array}{c c} 1.0308 \cdot 10^{-2} \\ 9.56 \cdot 10^{-5} \end{array}$	1.0308.10-2	$\begin{array}{c} 190.72 \\ 121.4 \end{array}$	1.181 0.752		

Чтобы исследовать, как меняются все три нормы I_1 , I_2 , I_3 в окрестности $n=2^p$, используется численный эксперимент. Для вычисления указанных выше норм достаточно найти

$$\max_{v} I_1$$
, $\max_{v} I_2$, $\max_{v} I_3$,

когда y пробегает множество собственных функций оператора A.

Расчеты проводились на последовательности сеток $h={}^1/_{10}, {}^1/_{12}, {}^1/_{14}, {}^1/_{20}$. Число итераций n задавалось, как и в первой серии. При этом оказалось, что в (3.3)-(3.5) равенство достигалось на первой собственной функции оператора A, как и для случая $n=2^p$, и были справедливы равенства (3.6).

Попутно мы убедились в том, что теоретические оценки (3.6), (3.7) при $n=2^p$ достигаются.

Некоторые результаты приведены в табл. 2, 3. В первых строках значения I_1 , I_2 , I_3 даны для $y = \sin \pi x$, во вторых строках I_2 , I_3 даны для сравнения при $y = \sin (N-1)\pi x$.

Рассмотренные выше эксперименты на модельной задаче показывают, что для n, отличного от степени 2, при упорядочении множества \mathfrak{M}_n согласно (2.4) сохраняются те же характеристики вычислительной устойчивости метода, как и для случая $n=2^p$, для которого в теоремах 1, 2 получены теоретические оценки.

Использование упорядочения \mathfrak{M}_n согласно (2.4) необходимо при решении задач с плохо обусловленным оператором, например разнестных задач, возникающих из аппроксимации эллиптических уравнений высокого порядка. В то же время, как показывает пример 1, если n не слишком велико и задача не очень плохо обусловлена, то можно пользоваться более простым упорядочением с помощью блоков малой длины.

Поступила в редакцию 24.03.1972

Цитированная литература

- 1. D. Young. On Richardson's method for solving linear system with positive definite matrices. J. Math. and Phys., 1954, 32, № 4, 243—255.
- 2. В. Вазов, Дж. Форсайт. Разностные методы решения дифференциальных уравнений в частных производных. М., Изд-во ин. лит., 1963.
- 3. Д. К. Фаддеев, В. Н. Фаддеева, Вычислительные методы линейной алгебры. М., Физматгиз, 1963.
- 4. В. Н. Лебедев, С. А. Финогенов. О порядке выбора итерационных параметров в чебышевском циклическом итерационном методе. Ж. вычисл. матем. и матем. физ. 1971, 11, \mathbb{N}_2 2, 425—438.
- 5. А. А. Самарский. Введение в теорию разностных схем. М., «Наука», 1971.
- 6. Г. И. Марчук, В. И. Лебедев. Численные методы в теории переноса нейтронов. М., Атомиздат, 1971.
- 7. В. И. Лебедев, С. А. Финогенов. Об одном алгоритме выбора параметров в чебышевских циклических методах. В сб. «Вычисл. методы линейной алгебры». Новосибирск, ВЦ СО АН СССР, 1972, 21—27.