

УДК 519.62

ЧИСЛЕННОЕ РЕШЕНИЕ СИСТЕМЫ ЛИНЕЙНЫХ  
АЛГЕБРАИЧЕСКИХ УРАВНЕНИЙ С ТРЕХДИАГОНАЛЬНОЙ МАТРИЦЕЙ  
МЕТОДОМ ОРТОГОНАЛЬНЫХ ВСТРЕЧНЫХ ПРОГОНОВ

С.И. Фадеев

**1. Краткая характеристика метода.**  
Предлагаемый метод решения системы линейных алгебраических уравнений с трехдиагональной матрицей является вариантом метода встречных прогонков с использованием последовательности ортогональных преобразований Хаусхолдера. С тем же успехом мы могли бы обратиться к последовательности плоских вращений. Наш выбор связан с возможностью простого обобщения метода на случай, когда система имеет ленточную матрицу общего вида. В результате встречных прогонков формируются подсистемы относительно пар неизвестных. Так как при ортогональных преобразованиях сохраняются значения сингулярных чисел матрицы системы, то обусловленность матрицы каждой из подсистем как отношение наибольшего сингулярного числа к наименьшему по крайней мере не превосходит обусловленности матрицы исходной системы. В условиях машинной реализации алгоритма это утверждение нуждается в информации об оценках погрешности округлений. Благодаря использованию ортогональных преобразований, технология определения оценок оказывается достаточно простой. Трудоемкость метода характеризуется числом арифметических операций порядка  $N$ , где  $N$  - размерность системы.

По затрагиваемым в статье вопросам имеется обширная литература (см., например, [1]). В частности, сошлемся на работу [2], где приводятся алгоритмы встречных прогонков на основе исключения неизвестных по методу Гаусса, и работу [3], в которой подробно описан метод вычисления оценок погрешностей, возникающих при преобразованиях Хаусхолдера или, как принято называть, преобразованиях отражения.

2. **О п р е д е л е н и я.** Пусть в  $N$ -мерном евклидовом пространстве даны вектор  $a = (a_1, a_2, \dots, a_N)$ , гиперплоскость с нормалью  $v = (v_1, v_2, \dots, v_N)$  и пусть  $b = (b_1, b_2, \dots, b_N)$  - отраженный вектор:  $b = a - 2 \frac{(v, a)}{(v, v)} v = Pa$ . Здесь  $P$  - ортогональная симметрическая матрица отражения размерности  $N \times N$  с элементами:  $P_{i,j} = \delta_{i,j} - 2 \frac{v_i v_j}{(v, v)}$ , где  $\delta_{i,j}$  - индекс Кронекера. Обозначим через  $P[i, i+1]$  матрицу отражения, если отраженный вектор  $b$  таков, что

$$b_i = -\sigma \sqrt{a_i^2 + a_{i+1}^2}, \quad \sigma = \begin{cases} +1, & a_i \geq 0 \\ -1, & a_i < 0 \end{cases}, \quad |b_i| \neq 0, \quad b_{i+1} = 0, \quad (1)$$

а все остальные элементы совпадают с соответствующими элементами вектора  $a$ :  $b_k = a_k$ ,  $k \neq i$ ,  $k \neq i+1$ . Наряду с этим мы будем использовать обозначение  $P[i+1, i]$ , если

$$b_{i+1} = -\sigma \sqrt{a_i^2 + a_{i+1}^2}, \quad \sigma = \begin{cases} +1, & a_{i+1} \geq 0 \\ -1, & a_{i+1} < 0 \end{cases}, \quad |b_{i+1}| \neq 0, \quad b_i = 0, \quad (2)$$

а остальные элементы векторов  $a$  и  $b$  совпадают. Матрицы  $P[i, i+1]$  и  $P[i+1, i]$  отличаются от единичной матрицы лишь элементами с индексами  $(i, i)$ ,  $(i, i+1)$ ,  $(i+1, i)$  и  $(i+1, i+1)$ . Эти элементы образуют подматрицы  $Q[i, i+1]$  и  $Q[i+1, i]$  соответственно:

$$Q[i, i+1] = \begin{bmatrix} -p_i & -q_i \\ -q_i & p_i \end{bmatrix}, \quad p_i = -\frac{a_i}{b_i}, \quad q_i = -\frac{a_{i+1}}{b_i}, \quad (3)$$

согласно (1), и

$$Q[i+1, i] = \begin{bmatrix} p_i & -q_i \\ -q_i & -p_i \end{bmatrix}, \quad p_i = -\frac{a_{i+1}}{b_{i+1}}, \quad q_i = -\frac{a_i}{b_{i+1}}, \quad (4)$$

согласно (2). Подматрицы  $Q[i, i+1]$  и  $Q[i+1, i]$  являются матрицами отражения вектора  $(a_i, a_{i+1})$  в пространстве двух измерений, причем  $(b_i, b_{i+1})$  - отраженный вектор:

$$\begin{bmatrix} b_i \\ 0 \end{bmatrix} = Q[i, i+1] \begin{bmatrix} a_i \\ a_{i+1} \end{bmatrix}, \quad \begin{bmatrix} 0 \\ b_{i+1} \end{bmatrix} = Q[i+1, i] \begin{bmatrix} a_i \\ a_{i+1} \end{bmatrix}.$$

3. Оценка погрешности расчета отраженного вектора. Следуя [3], воспользуемся положительными параметрами  $\epsilon_1$  и  $\epsilon_2$  - характеристиками разрядности ЭВМ, использующей представление числа с плавающей запятой - при оценке погрешности вычислений. Здесь  $\epsilon_1$  таково, что  $1+\epsilon_1$  - наименьшее по модулю машинное число, большее единицы;  $\epsilon_2/2$  и  $1/\epsilon_2$  - соответственно нижняя и верхняя границы модулей машинных чисел. Предполагается, что все вычисления производятся с одинарной точностью.

Пусть, например, известна матрица  $P[i, i+1]$  и требуется найти вектор  $y = (y_1, y_2, \dots, y_N)$  как отражение вектора  $x = (x_1, x_2, \dots, x_N)$ :

$$y = P[i, i+1]x, \quad \|y\| = \|x\| = \sqrt{(x, x)}. \quad (5)$$

В силу структуры  $P[i, i+1]$  имеем:  $y_k = x_k, \quad k \neq i, k \neq i+1$ ,

$$\begin{bmatrix} y_1 \\ \vdots \\ y_{i+1} \end{bmatrix} = Q[i, i+1] \begin{bmatrix} x_1 \\ \vdots \\ x_{i+1} \end{bmatrix} = \begin{bmatrix} -p_i x_i - q_i x_{i+1} \\ \vdots \\ -q_i x_i + p_i x_{i+1} \end{bmatrix}, \quad (6)$$

где  $p_i$  и  $q_i$  определены по формулам (3). Из-за ошибок округлений, возникающих а) при вычислении  $p_i$  и  $q_i$  и б) при вычислении  $y_i$  и  $y_{i+1}$  по формулам (6), мы вместо  $y$  получим вектор  $\tilde{y} = (\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_N)$ :  $\tilde{y} = (P[i, i+1] + \delta P)(x + \delta x) = y + h$ , где вектор  $h = (h_1, h_2, \dots, h_N)$  моделирует влияние ошибок округлений, матрица  $\delta P$  и вектор  $\delta x$  имеют нормы порядка  $\epsilon_1$ , причем  $\tilde{y}_k = y_k, \quad h_k = 0, \quad k \neq i, k \neq i+1$ . С точностью до малых более высокого порядка

$$\|h\| \leq \|\delta P\| \|x\| + \|\delta x\|. \quad (7)$$

Рассмотрим оценки погрешностей в случае "а". Имеем:  $(p_i)_{\text{выч}} = p_i + \delta p_i + u_i, \quad (q_i)_{\text{выч}} = q_i + \delta q_i + v_i$ , где

$$|\delta p_i| \leq 4\epsilon_1 |p_i|, \quad |\delta q_i| \leq 4\epsilon_1 |q_i|, \quad |u_i| \leq \frac{\epsilon_2}{2}, \quad |v_i| \leq \frac{\epsilon_2}{2}.$$

Здесь использовано предположение о том, что, привлекая при необходимости масштабирование элементов  $a_i$  и  $a_{i+1}$  (значения  $p_i$  и  $q_i$  от этого не изменяются), можно считать ограниченным снизу модуль

$b_i$ :  $|b_i| > \sqrt{\frac{3\epsilon_2}{4\epsilon_1}}$ . Матрица  $\delta P$  отличается от нулевой только элементами с индексами  $(i, i), (i+1, i), (i, i+1)$  и  $(i+1, i+1)$ , кото-

рые образуют подматрицу  $\delta Q$ :

$$\delta Q = \begin{bmatrix} -\delta p_1 & -\delta q_1 \\ -\delta q_1 & \delta p_1 \end{bmatrix} + \begin{bmatrix} -u_1 & -v_1 \\ -v_1 & u_1 \end{bmatrix}.$$

Поэтому

$$\begin{aligned} \|\delta P\| &\leq \|\delta P\|_{\mathbb{E}} = \|\delta Q\|_{\mathbb{E}} = \sqrt{2(\delta p_1 + u_1)^2 + 2(\delta q_1 + v_1)^2} \leq \\ &\leq \sqrt{2(\delta p_1^2 + \delta q_1^2)} + \sqrt{2(u_1^2 + v_1^2)} \leq 4\sqrt{2}\epsilon_1 + \epsilon_2. \end{aligned}$$

Здесь норма  $\|\delta P\|$  оценивается сверху евклидовой нормой  $\|\delta P\|_{\mathbb{E}}$ , квадрат которой равен сумме квадратов всех элементов матрицы  $\delta P$ .

В случае "б" оценивается норма вектора  $\delta x = (\delta x_1, \delta x_2, \dots, \delta x_N)$ , причем  $\delta x_k = 0$ ,  $k \neq i$ ,  $k \neq i+1$ :

$$|\delta x_1| \leq 2\epsilon_1(|p_1||x_i| + |q_1||x_{i+1}|) + \frac{3}{2}\epsilon_2 \leq 2\epsilon_1\sqrt{x_i^2 + x_{i+1}^2} + \frac{3}{2}\epsilon_2,$$

$$|\delta x_2| \leq 2\epsilon_1(|q_1||x_i| + |p_1||x_{i+1}|) + \frac{3}{2}\epsilon_2 \leq 2\epsilon_1\sqrt{x_i^2 + x_{i+1}^2} + \frac{3}{2}\epsilon_2.$$

Отсюда следует, что  $\|\delta x\| = \sqrt{\delta x_1^2 + \delta x_2^2} \leq \sqrt{8\epsilon_1^2(x_i^2 + x_{i+1}^2)} + \sqrt{\frac{9}{2}\epsilon_2^2} \leq 2\sqrt{2}\epsilon_1\|x\| + \frac{3\sqrt{2}}{2}\epsilon_2 \leq 4\sqrt{2}\epsilon_1\|x\|$ , если выполнено неравенство:

$$\|x\| \geq \frac{3\epsilon_2}{4\epsilon_1}. \quad (8)$$

Пренебрегая в (?) отличием нормы  $\|x\|$  от нормы  $\|\tilde{x}\|$ , получим:

$$\|h\| \leq (8\sqrt{2}\epsilon_1 + \epsilon_2)\|\tilde{x}\| \leq 12\epsilon_1\|\tilde{x}\|, \quad (9)$$

поскольку заведомо  $\epsilon_2 < 0.6\epsilon_1$ .

4. Метод ортогональных встречных прогонки. Рассмотрим систему из  $N$  линейных алгебраических уравнений с трехдиагональной матрицей:

$$Gf = F, \quad (10)$$

или

$$\begin{bmatrix} C_{1,2} & C_{1,3} & & 0 \\ C_{2,1} & C_{2,2} & C_{2,3} & \\ & C_{3,1} & C_{3,2} & C_{3,3} \\ \hline & C_{N-1,1} & C_{N-1,2} & C_{N-1,3} \\ & & C_{N,1} & C_{N,2} \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ \vdots \\ f_{N-1} \\ f_N \end{bmatrix} = \begin{bmatrix} F_1 \\ F_2 \\ F_3 \\ \vdots \\ F_{N-1} \\ F_N \end{bmatrix}.$$

Покажем существование ортогональной матрицы  $U$  такой, что система  $UCf = UF$  (11)

будет содержать подсистему из двух уравнений относительно неизвестных  $f_i$  и  $f_{i+1}$  вида:

$$\begin{bmatrix} A_{1,1}^{(i)} & A_{1,2}^{(i)} \\ A_{2,1}^{(i)} & A_{2,2}^{(i)} \end{bmatrix} \begin{bmatrix} f_i \\ f_{i+1} \end{bmatrix} = \begin{bmatrix} G_1^{(i)} \\ G_2^{(i)} \end{bmatrix} \quad (12)$$

Для этого построим последовательность отражений типа  $P(k, s)$ ,  $|k-s|=1$ , причем в качестве вектора  $a$  на  $k$ -м шагу берется  $k$ -й вектор-столбец матрицы системы, полученной на  $(k-1)$ -м шагу. Если  $1 \leq i \leq N-1$ , то последовательность отражений состоит из правого и левого ходов прогонки. При  $i=1$  (т.е. требуется найти  $f_1$  и  $f_2$ ) достаточно выполнить только левый ход, а при  $i=N-1$  (неизвестные  $f_{N-1}$  и  $f_N$ ) - только правый ход.

На первом шагу правого хода прогонки матрица  $C$  и вектор  $F$  умножаются на  $P[1,2]$ . В силу определения  $P[1,2]$  преобразование сводится к вычислению произведения  $Q[1,2]$  и подматрицы, состоящей из элементов первой и второй строк системы (IO). Остальные элементы системы (IO) остаются без изменения. Результат представим в виде:

$$Q[1,2] \begin{bmatrix} C_{1,2} & C_{1,3} & 0 & \dots & 0 & F_1 \\ C_{2,1} & C_{2,2} & C_{2,3} & \dots & 0 & F_2 \end{bmatrix} = \begin{bmatrix} C_{1,1}^{(1)} & C_{1,2}^{(1)} & C_{1,3}^{(1)} & \dots & 0 & F_1^{(1)} \\ 0 & C_{2,2}^{(1)} & C_{2,3}^{(1)} & \dots & 0 & F_2^{(1)} \end{bmatrix},$$

где верхний индекс приписан новым значениям необязательно равных нулю элементов первой и второй строк системы:

$$C^{(1)}f = F^{(1)}, \quad C^{(1)} = P[1,2]C, \quad F^{(1)} = P[1,2]F. \quad (13)$$





Очевидно, (I7) и есть искомая система (II), включающая подсистему (I2) относительно  $f_i$  и  $f_{i+1}$ , где

$$\begin{aligned} A_{1,1}^{(i)} &= C_{i,2}^{(i-1)}, & A_{1,2}^{(i)} &= C_{i,3}^{(i-1)}, & G_1^{(i)} &= F_1^{(i-1)}, \\ A_{2,1}^{(i)} &= \tilde{C}_{i+1,1}^{(k)}, & A_{2,2}^{(i)} &= \tilde{C}_{i+1,2}^{(k)}, & G_2^{(i)} &= \tilde{F}_{i+1}^{(k)}, \quad k=N-i-1. \end{aligned}$$

Ортогональная матрица  $U$  является произведением  $(N-2)$ -х матриц отражений:  $U = P[i+1, i]P[i+2, i+1] \times \dots \times P[N, N-1] P[i-1, i] P[i-2, i-1] \times \dots \times P[1, 2]$ .

Обратим внимание на то, что правый ход и левый ход прогонки, приводящей к системе (I7), в действительности выполняются независимо. Поэтому систему (I7) можно составить из  $i$  первых строк системы (I5) и  $(N-i)$  последних строк системы  $\tilde{U}Cf = \tilde{U}F$ , где

$$\tilde{U} = P[i+1, i]P[i+2, i+1] \times \dots \times P[N, N-1].$$

При  $i = N-1$  коэффициенты подсистемы (I2) относительно  $f_{N-1}$  и  $f_N$  имеют вид:

$$\begin{aligned} A_{1,1}^{(i)} &= C_{i,2}^{(i-1)}, & A_{1,2}^{(i)} &= C_{i,3}^{(i-1)}, & G_1^{(i)} &= F_1^{(i-1)}, \\ A_{2,1}^{(i)} &= C_{N,1}, & A_{2,2}^{(i)} &= C_{N,2}, & G_2^{(i)} &= F_N. \end{aligned}$$

В этом случае  $U$  является произведением матриц отражений, определяемых правым ходом прогонки:  $U = P[N-2, N-1] P[N-3, N-2] \times \dots \times P[1, 2]$ .

При  $i = 1$  коэффициенты подсистемы относительно  $f_1$  и  $f_2$  имеют вид:

$$\begin{aligned} A_{1,1}^{(i)} &= C_{1,2}, & A_{1,2}^{(i)} &= C_{1,3}, & G_1^{(i)} &= F_1, \\ A_{2,1}^{(i)} &= \tilde{C}_{2,1}^{(N-2)}, & A_{2,2}^{(i)} &= \tilde{C}_{2,2}^{(N-2)}, & G_2^{(i)} &= \tilde{F}_2^{(N-2)}, \end{aligned}$$

а  $U$  - произведение матриц отражений, определяемых левым ходом прогонки:  $U = P[3, 2]P[4, 3] \times \dots \times P[N, N-1]$ .

5. Обусловленность матрицы подсистем. Пусть  $\max \alpha(C)$  и  $\min \alpha(C)$  - наибольшее и наименьшее сингулярные числа матрицы  $C$ ,  $0 < \min \alpha(C) \leq \max \alpha(C)$ , и пусть  $\mu(C)$  - обусловленность  $C$ :  $\mu(C) = \frac{\max \alpha(C)}{\min \alpha(C)}$ . Обозначим через  $A^{(i)}$  матрицу системы (I2). Для оценки  $\mu(A^{(i)})$  нам потребуется теорема отделимости Штурма [4]. Приведем ее формулировку.

Рассмотрим последовательность симметрических матриц  $A_p$  размерности  $r \times r$ , где  $r = 1, 2, \dots, N$ . Пусть  $\lambda_j(A_p)$ ,  $j = 1, 2, \dots, r$ , обозначает  $j$ -е собственное число матрицы  $A_p$ , причем  $\lambda_1(A_p) \leq \lambda_2(A_p) \leq \dots \leq \lambda_r(A_p)$ . Тогда  $\lambda_j(A_{i+1}) \leq \lambda_j(A_i) \leq \lambda_{j+1}(A_{i+1})$ .

Опираясь на теорему отделения, легко показать, что  $\mu(A^{(i)}) \leq \mu(C)$ . Действительно, матрица  $\tilde{C}^{(k)}$  системы (17) имеет следующую структуру:

$$\tilde{C}^{(k)} = \begin{bmatrix} B_1 & | & 0 \\ \hline -\frac{1}{B_2} & | & -\frac{1}{B_3} \end{bmatrix},$$

где  $B_1$  - подматрица, включающая  $A^{(i)}$  как блок в правом нижнем углу,  $0$  - нулевая подматрица. По этой причине матрица  $B_1 B_1^T$  имеет в правом нижнем углу подматрицу  $A^{(i)} [A^{(i)}]^T$ . Применяя теорему отделения к матрицам  $A^{(i)} [A^{(i)}]^T$  и  $B_1 B_1^T$ , получим:  $\min \sigma^2(B_1) \leq \min \sigma^2(A^{(i)}) \leq \max \sigma^2(A^{(i)}) \leq \max \sigma^2(B_1)$ , т.е.  $\mu(A^{(i)}) \leq \mu(B_1)$ . В свою очередь матрица  $\tilde{C}^{(k)} [\tilde{C}^{(k)}]^T$  имеет в левом верхнем углу подматрицу  $B_1 B_1^T$ . Следовательно,  $\mu(B_1) \leq \mu(\tilde{C}^{(k)})$ . Остается вспомнить, что  $\mu(\tilde{C}^{(k)}) = \mu(C)$ , так как  $\tilde{C}^{(k)} = UC$ , где  $U$  - ортогональная матрица.

6. Оценка погрешности решения подсистемы. Во многих отношениях целесообразно дополнить прогонку еще одним шагом. Умножим слева матрицу  $U$  на  $P[i, i+1]$ . Тогда система

$$Sf = V, \quad S = P[i, i+1]UC, \quad V = P[i, i+1]UF, \quad (18)$$

$$\|S\| = \|C\|, \quad \|V\| = \|F\|,$$

будет содержать подсистему относительно  $f_i$  и  $f_{i+1}$  с верхней треугольной матрицей:

$$L^{(i)} z_i = R^{(i)}, \quad (19)$$

$$L^{(i)} = Q[i, i+1]A^{(i)} = \begin{bmatrix} 1 & m_i \\ 0 & n_i \end{bmatrix}, \quad R^{(i)} = Q[i, i+1]G^{(i)} = \begin{bmatrix} r_i \\ r_{i+1} \end{bmatrix}, \quad z_i = \begin{bmatrix} f_i \\ f_{i+1} \end{bmatrix}.$$

Приведение системы (10) к (18) сопровождается накоплением ошибок округлений. В результате вместо (18) мы получим возмущенную сис-

тому относительно  $\tilde{r}$ ,  $\tilde{r} = (\tilde{r}_1, \tilde{r}_2, \dots, \tilde{r}_N)$ :

$$\tilde{S}\tilde{r} = \tilde{v}, \quad \tilde{S} = S + \Theta, \quad \tilde{v} = v + \Lambda, \quad (20)$$

где матрица  $\Theta$  и вектор  $\Lambda$  моделируют влияние ошибок округлений вследствие вычислений, связанных с  $(N-1)$  ортогональными преобразованиями отражений. Аналогично вместо (19) мы будем иметь систему с матрицей  $\tilde{L}^{(i)}$  и правой частью  $\tilde{R}^{(i)}$  относительно  $\tilde{z}_i$  и  $\tilde{r}_{i+1}$ :

$$\tilde{L}^{(i)} \tilde{z}_i = \tilde{R}^{(i)}, \quad (21)$$

$$\tilde{L}^{(i)} = \begin{bmatrix} \tilde{L}_i & \tilde{M}_i \\ 0 & \tilde{R}_{i+1} \end{bmatrix} = L^{(i)} + B^{(i)}, \quad \tilde{R}^{(i)} = \begin{bmatrix} \tilde{r}_i \\ \tilde{r}_{i+1} \end{bmatrix} = R^{(i)} + H^{(i)}, \quad \tilde{z}_i = \begin{bmatrix} \tilde{z}_i \\ \tilde{r}_{i+1} \end{bmatrix}.$$

Нашей ближайшей задачей является получение оценок для  $\|B^{(i)}\|$  и  $\|H^{(i)}\|$  - норм возмущений  $L^{(i)}$  и  $R^{(i)}$ .

Пусть  $c_k, s_k, \tilde{s}_k$  и  $e_k$ ,  $k = 1, 2, \dots, N$ , - вектор-столбцы матриц  $C, S, \tilde{S}$  и  $\Theta$  соответственно:  $\tilde{s}_k = s_k + e_k$ ,  $\|s_k\| = \|c_k\|$ . Согласно [3], оценка нормы вектора  $\Lambda$  с учетом погрешности машинного представления  $\Lambda$  имеет вид:  $\|\Lambda\| \leq N\alpha e^{N\alpha} \|\tilde{v}\|$ , где  $\alpha = 12\epsilon_1$  в соответствии с (9). Так как  $\|\tilde{v}\| \leq \|v\| + \|\Lambda\|$ , то  $\|\Lambda\|$  может быть выражена через  $\|F\|$ :  $\|\Lambda\| \leq N\alpha e^{N\alpha} \|F\| / (1 - N\alpha e^{N\alpha})$ ,  $N\alpha e^{N\alpha} < 1$ . Подобная оценка справедлива и для  $\|e_k\|$ . Если, однако, можно уточнить, если принять во внимание, что не более трех элементов вектора  $c_k$  различно от нуля. Тогда  $\|e_k\| \leq 5\alpha e^{5\alpha} \|\tilde{s}_k\|$  или  $\|e_k\| \leq 5\alpha e^{5\alpha} \|c_k\| / (1 - 5\alpha e^{5\alpha})$ ,  $5\alpha e^{5\alpha} < 1$ . В приведенных оценках, как следствие (8), пред-

полагается выполненным условие:  $\min_k (\min \|c_k\|, \|F\|) > \frac{3\epsilon_2}{4\epsilon_1}$ . Мы выра-

ве считаем, что матрица  $\Theta$  имеет ту же структуру, что и  $S$ , и поэтому содержит подматрицу  $B^{(i)}$ , так же как элементы вектора  $H^{(i)}$  являются элементами  $\Lambda$ . В результате имеем следующие неравенства:

$$\|H^{(i)}\| \leq \|\Lambda\| \leq N\alpha e^{N\alpha} \|\tilde{v}\|,$$

$$\|B^{(i)}\| \leq \|B^{(i)}\|_F \leq 5\alpha e^{5\alpha} \sqrt{\|s_2\|^2 + \|s_{i+1}\|^2}$$

или

$$\|H^{(i)}\| \leq \frac{N\alpha e^{N\alpha}}{1 - N\alpha e^{N\alpha}} \|F\|, \quad \|B^{(i)}\| \leq \frac{5\alpha e^{5\alpha}}{1 - 5\alpha e^{5\alpha}} \sqrt{\|c_i\|^2 + \|c_{i+1}\|^2}. \quad (22)$$

При оценке близости решений систем (19) и (21) используется известное неравенство [3]. Отвлечемся на время от ранее принятых обозначений и сопоставим две системы линейных алгебраических уравнений:  $Ax = f$  и  $(A+B)y = f + g$ . При этом

$$\begin{aligned} \|y - x\| &\leq \|(A+B)^{-1}\| (\|B\| \|x\| + \|g\|) \leq \\ &\leq \frac{\|A^{-1}\|}{1 - \|B\| \|A^{-1}\|} (\|B\| \|x\| + \|g\|), \quad \|B\| \|A^{-1}\| < 1. \end{aligned} \quad (23)$$

Будем считать, что (19) является возмущенной системой (21) с оценками норм возмущений (22). Выпишем основные характеристики матрицы  $\tilde{L}^{(1)}$ , выражения которых очевидны:

$$\|L\| = \sqrt{\frac{\omega}{2} \left[ 1 + \sqrt{1 - \left( \frac{2\Delta}{\omega} \right)^2} \right]}, \quad \|L^{-1}\| = \frac{\|L\|}{\Delta}, \quad \mu = \frac{\|L\|^2}{\Delta}, \quad (24)$$

$$L \equiv \tilde{L}^{(1)}, \quad \omega = \tilde{l}_1^2 + \tilde{m}_1^2 + \tilde{n}_1^2, \quad \Delta = |\tilde{l}_1 \tilde{n}_1|, \quad \mu = \mu(L).$$

Рассмотрим погрешности двух типов. Обозначим через  $\tilde{r} = (\tilde{r}_1, \tilde{r}_2, \dots, \tilde{r}_N)$  вектор, элементы которого определяются как машинное решение системы (21). В этом случае элементы вектора  $\tilde{z}_1 + \kappa^{(1)}$ ,  $\tilde{z}_1 = (\tilde{r}_1, \tilde{r}_2, \dots, \tilde{r}_N)$ ,  $\kappa^{(1)} = (\kappa_1^{(1)}, \kappa_2^{(1)})$ , являются точным решением системы [3]:  $(L + B^{(1)})(\tilde{z}_1 + \kappa^{(1)}) = R + \tilde{H}^{(1)}$ ,  $R \equiv \tilde{r}^{(1)}$ , причем нормы возмущений оцениваются следующим образом:  $\|B^{(1)}\| \leq 2\sqrt{2}\epsilon_1 \|L\|$ ,  $\|\tilde{H}^{(1)}\| \leq \epsilon_1 \|R\| + \sqrt{2}\epsilon_2$ ,  $\|\kappa^{(1)}\| \leq \epsilon_2/\sqrt{2}$ . Полагая в (23)  $x = \tilde{z}_1$ ,  $y = \tilde{z}_1 + \kappa^{(1)}$ ,  $\|B\| = 2\sqrt{2}\epsilon_1 \|L\|$ ,  $\|g\| = \epsilon_1 \|R\| + \sqrt{2}\epsilon_2$ , получим:

$$\|\tilde{z}_1 - z_1\| \leq K_1 \|\tilde{z}_1\| + \phi_1, \quad (25)$$

где

$$K_1 = \frac{2\sqrt{2}\epsilon_1\mu}{1 - 2\sqrt{2}\epsilon_1\mu}, \quad \phi_1 = \frac{\|L^{-1}\|}{1 - 2\sqrt{2}\epsilon_1\mu} (\epsilon_1 \|R\| + \sqrt{2}\epsilon_2) + \frac{\epsilon_2}{\sqrt{2}}.$$

Из (25) следует, что  $\|\tilde{z}_1\| \leq K_1 \|\tilde{z}_1\| + \phi_1 + \|\tilde{z}_1\|$ . Поэтому

$$\|\tilde{z}_1\| \leq \frac{\|\tilde{z}_1\| + \phi_1}{1 - K_1}, \quad K_1 < 1.$$

С учетом последнего неравенства представим (25) в виде:

$$\|\tilde{z}_1 - z_1\| \leq \rho_1^{(1)} = \frac{K_1 \|\tilde{z}_1\| + \phi_1}{1 - K_1}. \quad (26)$$

Далее, пользуясь (23), сопоставим системы (19) и (21). В результате получим оценку нормы погрешности второго рода:

$$\|\tilde{z}_1 - z_1\| \leq \rho_2^{(1)} = \frac{\|L^{-1}\|}{1 - \xi \|L^{-1}\|} [\xi(\rho_1^{(1)} + \|\tilde{z}_1\|) + \eta], \quad \xi \|L^{-1}\| < 1, \quad (27)$$

где, согласно (22),

$$\xi = \frac{5\alpha e^{5\alpha}}{1 - 5\alpha e^{5\alpha}} \sqrt{\|C_1\|^2 + \|C_{i+1}\|^2}, \quad \eta = \frac{N\alpha e^{N\alpha}}{1 - N\alpha e^{N\alpha}} \|\tilde{x}\|.$$

Таким образом, значения  $\rho_1^{(1)}$  и  $\rho_2^{(1)}$  находятся после машинного вычисления  $\tilde{f}_1$  и  $\tilde{f}_{i+1}$ :  $\tilde{f}_{i+1} = \tilde{f}_{i+1}/\tilde{h}_i$ ,  $\tilde{f}_1 = (\tilde{f}_1 - \tilde{h}_1 \tilde{f}_{i+1})/\tilde{1}_1$ ,  $\|\tilde{z}_1\| = \sqrt{\tilde{f}_{i+1}^2 + \tilde{f}_1^2}$ . Вместе с ними вычисляются оценки абсолютной и относительной погрешностей;

$$\|\tilde{z}_1 - z_1\| \leq \rho_1 = \sqrt{[\rho_1^{(1)}]^2 + [\rho_2^{(1)}]^2}, \quad \|\tilde{z}_1 - z_1\| / \|\tilde{z}_1\| \leq \rho_1 / \|\tilde{z}_1\|, \quad (28)$$

7. Оценки погрешности решения системы. В этом пункте мы получим, используя (28), выражения вычисляемой (и поэтому наиболее эффективной) оценки погрешности решения системы (10), а затем установим связь оценки относительной погрешности решения системы с обусловленностью матрицы С. Заметим, что определение всех неизвестных  $f_j$  по методу встречных прогонок требует организации не менее  $N_0$  "встреч", где

$$N_0 = \begin{cases} N/2, & \text{если } N \text{ четное,} \\ (N+1)/2, & \text{если } N \text{ нечетное.} \end{cases}$$

В случае нечетного  $N$  одно из неизвестных вычисляется дважды. Пусть

это будет  $f_{N-1}$ . Следовательно,  $\|\tilde{f} - f\| = \left( \sum_{k=1}^{N_0} \|\tilde{z}_{2k-1} - z_{2k-1}\|^2 \right)^{1/2}$ ,

если  $N$  четное, и  $\|\tilde{f} - f\| \leq \left( \sum_{k=1} \|\tilde{z}_{2k-1} - z_{2k-1}\|^2 + \|\tilde{z}_{N-1} - z_{N-1}\|^2 \right)^{1/2}$ ,

если  $N$  нечетное. Условимся обозначать знаком  $\Sigma$  сумму  $N_0$  слагаемых как в том, так и в другом случае, объединив их одним неравенством:  $\|\tilde{f} - f\| \leq \left( \Sigma \|\tilde{z}_i - z_i\|^2 \right)^{1/2}$ . Каждое из слагаемых суммы, согласно (28), меньше  $\rho_1^2$ . Поэтому

$$\|\tilde{f} - f\| \leq \rho = \left( \Sigma \rho_1^2 \right)^{1/2}, \quad \frac{\|\tilde{f} - f\|}{\|\tilde{f}\|} \leq \delta_0 = \frac{\rho}{\|\tilde{f}\|}. \quad (29)$$

Представляет интерес зависимость относительной погрешности  $\|\tilde{f} - f\| / \|\tilde{f}\|$  от обусловленности матрицы С или обусловленности матрицы S,  $\mu(C) = \mu(S)$ , свойственная рассматриваемому методу орто-

гональных встречных прогонки. Для этого нам потребуется равномерная оценка характеристик матрицы  $\tilde{S}$ ,  $\tilde{S} = S + \theta$ , системы (20), чтобы затем оценить характеристики матрицы  $L \equiv \tilde{L}^{(1)}$  системы (21). Имеем:

$$\begin{aligned}\|\tilde{S}\| &\leq \|S\| + \|\theta\| = \|C\| + \|\theta\|, \\ \|\tilde{S}^{-1}\| &\leq \frac{\|S^{-1}\|}{1 - \|\theta\| \|S^{-1}\|} = \frac{\|C^{-1}\|}{1 - \|\theta\| \|C^{-1}\|}, \\ \mu(\tilde{S}) &= \|\tilde{S}\| \|\tilde{S}^{-1}\| \leq \frac{\|C^{-1}\| (\|C\| + \|\theta\|)}{1 - \|\theta\| \|C^{-1}\|},\end{aligned}$$

где

$$\|\theta\| \leq \|\theta\|_E \leq K_C \|C\|, \quad K_C = \frac{5 N \alpha e^{5\alpha}}{1 - 5 \alpha e^{5\alpha}}. \quad (30)$$

Согласно п.5, наибольшие и наименьшие сингулярные числа матриц  $L$  и  $\tilde{S}$  удовлетворяют неравенствам:  $1/\|\tilde{S}^{-1}\| \leq 1/\|L^{-1}\| \leq \|L\| \leq \|\tilde{S}\|$ . Отсюда следуют оценки характеристик матрицы  $L$ :

$$\left. \begin{aligned}\|L\| &\leq \|\tilde{S}\| \leq (1 + K_C) \|C\|, \\ \|L^{-1}\| &\leq \|\tilde{S}^{-1}\| \leq \frac{\|C^{-1}\|}{1 - K_C \mu(C)}, \\ \mu(L) &\leq \mu(\tilde{S}) \leq \mu_0, \quad K_C \mu(C) < 1, \\ \mu_0 &= \frac{1 + K_C}{1 - K_C \mu(C)} \mu(C).\end{aligned}\right\} \quad (31)$$

Обратимся вначале к оценке  $\|\tilde{z}_i - z_i\|$ , преобразовав неравенство (25) к виду:  $\|\tilde{z}_i - z_i\| \leq K_\mu \|\tilde{z}_i\| + \phi$ ,  $K_i \leq K_\mu$ ,  $\phi_i \leq \phi$ , где

$$K_\mu = \frac{2\sqrt{2} \epsilon_1 \mu_0}{1 - 2\sqrt{2} \epsilon_1 \mu_0}, \quad 2\sqrt{2} \epsilon_1 \mu_0 < 1. \quad (32)$$

Определим  $\phi$ . Из выражения  $\phi_i$  и (31) следует, что

$$\phi_i \leq \frac{\|C^{-1}\| (\epsilon_1 \|F\| + \sqrt{2} \epsilon_2)}{[1 - K_C \mu(C)] [1 - 2\sqrt{2} \epsilon_1 \mu_0]} + \frac{\epsilon_2}{\sqrt{2}}.$$

Так как  $\|F\| \leq \|C\| \|f\|$ , то

$$\phi_i \leq \left\{ \frac{\mu(C) \left( \epsilon_1 + \frac{\sqrt{2} \epsilon_2}{\|F\|} \right)}{[1-K_C \mu(C)] [1-2\sqrt{2} \epsilon_1 \mu_0]} + \frac{\epsilon_2 \|C\|}{\sqrt{2} \|F\|} \right\} \|f\|.$$

Потребуем, чтобы норма  $F$  была ограничена снизу (для системы (19) аналогичное ограничение на  $\|R^{(1)}\|$  было бы обременительным):  $\|F\| \geq \frac{\sqrt{2} \epsilon_2}{\epsilon_1} \left(1 + \frac{1}{2} \|C\|\right)$ . При этом в качестве  $\phi$  можно взять выражение:

$$\phi = \epsilon_1 \left[ \frac{2\mu_0}{(1+K_C)(1-2\sqrt{2}\epsilon_1\mu_0)} + 1 \right] \|f\|.$$

Для оценки  $\|\tilde{x} - \tilde{f}\|$  просуммируем все  $\|\tilde{z}_i - z_i\|^2$ :  $\|\tilde{x} - \tilde{f}\| \leq \sqrt{\sum \|\tilde{z}_i - z_i\|^2} \leq K_\mu \sqrt{\sum \|\tilde{z}_i\|^2} + \sqrt{\frac{N+1}{2}} \phi \leq \sqrt{2} K_\mu \|\tilde{x}\| + \frac{1}{2} \sqrt{N+1} (K_\mu + \sqrt{2} \epsilon_1) \|f\|$ . Пусть известно, что  $\|\tilde{x} - f\| \leq \delta_1' \|f\|$ . Тогда  $\|\tilde{x}\| \leq (1 + \delta_1') \|f\|$ , и, следовательно -

$$\frac{\|\tilde{x} - \tilde{f}\|}{\|f\|} \leq \delta_1'' = \sqrt{2} (1 + \delta_1') K_\mu + \frac{\sqrt{N+1}}{2} (K_\mu + \sqrt{2} \epsilon_1).$$

Вычисление  $\delta_1'$  связано с использованием неравенства (27), которое запишем в виде:

$$\|\tilde{z}_i - z_i\| \leq \frac{\|L^{-1}\|}{1 - \xi \|L^{-1}\|} (\xi \|z_i\| + \eta), \quad \xi < K_C \|C\|.$$

С учетом (31) и неравенства

$$\frac{\eta}{\|C\|} = K_F \frac{\|F\|}{\|C\|} \leq K_F \|f\|, \quad K_F = \frac{N \alpha \epsilon^{N\alpha}}{1 - N \alpha \epsilon^{N\alpha}}, \quad (33)$$

получим  $\|\tilde{z}_i - z_i\| \leq K_F (K_C \|z_i\| + K_F \|f\|)$ , где

$$K_F = \frac{\mu_0}{1 - K_C (\mu_0 - 1)}, \quad K_C (\mu_0 - 1) < 1. \quad (34)$$

Для оценки  $\|\tilde{x} - f\|$  просуммируем все  $\|\tilde{z}_i - z_i\|^2$ :

$$\begin{aligned} \|\tilde{x} - f\| &\leq \sqrt{\sum \|\tilde{z}_i - z_i\|^2} \leq K_F \left[ K_C \sqrt{\sum \|\tilde{z}_i\|^2} + \sqrt{\frac{N+1}{2}} K_F \|f\| \right] \leq \\ &\leq K_F \left[ \sqrt{2} K_C \|\tilde{x}\| + \sqrt{\frac{N+1}{2}} K_F \|f\| \right] \leq K_F \left[ \sqrt{2} K_C (1 + \rho_1) + \sqrt{\frac{N+1}{2}} K_F \right] \|f\|. \end{aligned}$$

Отсюда следует выражение  $\delta'_1$ :

$$\frac{\|\tilde{f}-f\|}{\|f\|} \leq \delta'_1 = K_F \frac{2K_C + \sqrt{N+1} K_F}{\sqrt{2}-2K_C K_F}, \quad \sqrt{2} K_C K_F < 1.$$

Теперь мы можем выписать искомую оценку  $\delta_1$  относительной погрешности:

$$\frac{\|\tilde{f}-f\|}{\|f\|} \leq \delta'_1 + \delta'_2, \quad \frac{\|\tilde{f}-f\|}{\|f\|} \leq \delta_1 = \frac{\delta'_1 + \delta'_2}{1 - \delta'_1 - \delta'_2}, \quad (35)$$

$$\delta'_1 = K_F \frac{2K_C + \sqrt{N+1} K_F}{\sqrt{2} - 2K_C K_F},$$

$$\delta'_2 = \sqrt{2} K_\mu (1 + \delta'_1) + \frac{\sqrt{N+1}}{2} (K_\mu + \sqrt{2} \epsilon_1),$$

где  $K_C, K_\mu, K_F$  и  $K_F$  определены равенствами (30)-(34).

Формула (35), дающая оценку относительной погрешности решения системы (10), свидетельствует о том, что методу ортогональных встречных прогонок свойственна повышенная чувствительность к возмущениям правых частей. Действительно, если обусловленность матрицы  $C$  и размерность системы  $N$  "не слишком велики", то числитель и знаменатель дроби  $(\delta'_1 + \delta'_2)/(1 - \delta'_1 - \delta'_2)$  содержат слагаемое порядка  $\epsilon_1 N \sqrt{N} \mu(C)$ . Отчасти этот вывод основывается на оценке  $\|H^{(1)}\|$  (см. (22)), которая пропорциональна  $\|F\|$  и поэтому, как правило, пессимистична. Но вместе с тем вывод отражает и суть явления.

**Описание алгоритма.** Хотя вычисление всех неизвестных, определяемых системой (10), связано с  $N_0$  встречными прогонами, совсем не обязательно всякий раз начинать прогоны заново. Чтобы сформировать  $N_0$  систем (12), достаточно один раз проделать правый ход до  $i=N-2$  и левый до  $i=3$ , как это предложено в п.4. При этом требуется дополнительный массив  $W$  с элементами  $w_{i,j}$  размерности  $2N_0 \times 3$ . Кроме того, текущие значения прогоночных коэффициентов мы будем присваивать коэффициентам  $\alpha_i$ ,  $\beta_i$  и  $\gamma_i$ . Напомним, что при нечетном  $N$  мы условились вычислять дважды  $f_{N-1}$ . Поэтому удобнее начинать прогонку с левого хода.

Левому ходу прогонки предшествует заполнение строки с номером  $2N_0$  массива  $W$ :  $w_{2N_0,1} = \beta_0 = C_{N,1}$ ,  $w_{2N_0,2} = \alpha_0 = C_{N,2}$ ,

$W_{2N_0,3} = \gamma_0 = F_N$ . На первом шагу вычисляются  $P_{N-1}$  и  $Q_{N-1}$ -элементы матрицы  $Q[N, N-1]$  (см. (2), (4)):

$$P_{N-1} = \frac{\alpha_0}{\sqrt{\alpha_0^2 + C_{N-1,3}^2}}, \quad Q_{N-1} = \frac{\sigma C_{N-1,3}}{\sqrt{\alpha_0^2 + C_{N-1,3}^2}}, \quad \sigma = \begin{cases} +1, & \alpha_0 \geq 0, \\ -1, & \alpha_0 < 0, \end{cases}$$

а затем параметры  $\alpha_1, \beta_1$  и  $\gamma_1$  по формулам:

$$\alpha_1 = \tilde{\alpha}_{N-1,2}^{(1)} = P_{N-1} C_{N-1,2} - Q_{N-1} \beta_0,$$

$$\beta_1 = \tilde{C}_{N-1,1}^{(1)} = P_{N-1} C_{N-1,1},$$

$$\gamma_1 = \tilde{F}_{N-1}^{(1)} = P_{N-1} F_{N-1} - Q_{N-1} \gamma_0.$$

На втором шагу вычисляются  $P_{N-2}$  и  $Q_{N-2}$  - элементы матрицы  $Q[N-1, N-2]$ :

$$P_{N-2} = \frac{\alpha_1}{\sqrt{\alpha_1^2 + C_{N-2,3}^2}}, \quad Q_{N-2} = \frac{\sigma C_{N-2,3}}{\sqrt{\alpha_1^2 + C_{N-2,3}^2}}, \quad \sigma = \begin{cases} +1, & \alpha_1 \geq 0, \\ -1, & \alpha_1 < 0, \end{cases}$$

а затем параметры  $\alpha_2, \beta_2$  и  $\gamma_2$  по формулам:

$$\alpha_2 = \tilde{\alpha}_{N-2,2}^{(2)} = P_{N-2} C_{N-2,2} - Q_{N-2} \beta_1,$$

$$\beta_2 = \tilde{C}_{N-2,1}^{(2)} = P_{N-2} C_{N-2,1},$$

$$\gamma_2 = \tilde{F}_{N-2}^{(2)} = P_{N-2} F_{N-2} - Q_{N-2} \gamma_1$$

и так далее:  $i = 1, 2, \dots, N-2$ ,

$$\alpha_i = \tilde{\alpha}_{N-i,2}^{(i)} = P_{N-i} C_{N-i,2} - Q_{N-i} \beta_{i-1},$$

$$\beta_i = \tilde{C}_{N-i,1}^{(i)} = P_{N-i} C_{N-i,1},$$

$$\gamma_i = \tilde{F}_{N-i}^{(i)} = P_{N-i} F_{N-i} - Q_{N-i} \gamma_{i-1},$$

где

$$P_{N-i} = \frac{\alpha_{i-1}}{\sqrt{\alpha_{i-1}^2 + C_{N-i,3}^2}}, \quad Q_{N-i} = \frac{\sigma C_{N-i,3}}{\sqrt{\alpha_{i-1}^2 + C_{N-i,3}^2}}, \quad \sigma = \begin{cases} +1, & \alpha_{i-1} \geq 0, \\ -1, & \alpha_{i-1} < 0. \end{cases}$$

При четном  $i$  заполняются четные строки массива  $W$ :  $W_{i,1} = \beta_i$ ,  $W_{i,2} = -\alpha_i$ ,  $W_{i,3} = \gamma_i$ ,  $i = 2(N_0 - k)$ ,  $k = 0, 1, \dots, N_0 - 1$ , на этом левый ход прогонки завершается.

Правому ходу прогонки предшествует заполнение первой строки массива  $W$ :  $W_{1,1} = \alpha_0 = C_{1,2}$ ,  $W_{1,2} = \beta_0 = C_{1,3}$ ,  $W_{1,3} = \gamma_0 = F_1$ . В результате имеем систему (I2) (в других обозначениях) относительно  $f_1$  и  $f_2$ :  $W_{1,1}f_1 + W_{1,2}f_2 = W_{1,3}$ ;  $W_{2,1}f_1 + W_{2,2}f_2 = W_{2,3}$ . Здесь элементы  $W_{2,1}$ ,  $W_{2,2}$  и  $W_{2,3}$  были определены на  $(N-2)$ -м шагу левого хода прогонки.

На первом шагу вычисляются  $p_1$  и  $q_1$  - элементы матрицы  $Q[1,2]$  (см. (I), (3)):

$$p_1 = \frac{\sigma \alpha_0}{\sqrt{\alpha_0^2 + C_{2,1}^2}}, \quad q_1 = \frac{\sigma C_{2,1}}{\sqrt{\alpha_0^2 + C_{2,1}^2}}, \quad \sigma = \begin{cases} +1, & \alpha_0 \geq 0, \\ -1, & \alpha_0 < 0, \end{cases}$$

а затем параметры  $\alpha_1$ ,  $\beta_1$  и  $\gamma_1$  по формулам:

$$\begin{aligned} \alpha_1 &= C_{2,2}^{(1)} = p_1 C_{2,2} - q_1 \beta_0, \\ \beta_1 &= C_{2,3}^{(1)} = p_1 C_{2,3}, \\ \gamma_1 &= F_2^{(1)} = p_1 F_2 - q_1 \gamma_0. \end{aligned}$$

На втором шагу вычисляются  $p_2$  и  $q_2$  - элементы матрицы  $Q[2,3]$ :

$$p_2 = \frac{\sigma \alpha_1}{\sqrt{\alpha_1^2 + C_{3,1}^2}}, \quad q_2 = \frac{\sigma C_{3,1}}{\sqrt{\alpha_1^2 + C_{3,1}^2}}, \quad \sigma = \begin{cases} +1, & \alpha_1 \geq 0, \\ -1, & \alpha_1 < 0, \end{cases}$$

а затем параметры  $\alpha_2$ ,  $\beta_2$  и  $\gamma_2$  по формулам:

$$\begin{aligned} \alpha_2 &= C_{3,2}^{(2)} = p_2 C_{3,2} - q_2 \beta_1, \\ \beta_2 &= C_{3,3}^{(2)} = p_2 C_{3,3}, \\ \gamma_2 &= F_3^{(2)} = p_2 F_3 - q_2 \gamma_1. \end{aligned}$$

Кроме того, заполняется третья строка массива  $W$ :  $W_{3,1} = \alpha_2$ ,  $W_{3,2} = \beta_2$ ,  $W_{3,3} = \gamma_2$ . Вместе с этим оказывается сформированной система (I2) относительно  $f_3$  и  $f_4$ :  $W_{3,1}f_3 + W_{3,2}f_4 = W_{3,3}$ ,  $W_{4,1}f_3 + W_{4,2}f_4 = W_{4,3}$ . Здесь элементы  $W_{4,1}$ ,  $W_{4,2}$  и  $W_{4,3}$  были определены на  $(N-4)$ -м шагу левого хода прогонки. Формулы  $i$ -го шага имеют вид:  $i = 1, 2, \dots, N-2$ ,

$$\alpha_i = C_{i+1,2}^{(i)} = p_i C_{i+1,2} - q_i \beta_{i-1},$$

$$\beta_i = C_{i+1,3}^{(1)} = p_i C_{i+1,3},$$

$$\gamma_i = F_{i+1}^{(1)} = p_i F_{i+1} - q_i \gamma_{i-1},$$

где

$$p_i = \frac{\sigma \alpha_{i-1}}{\sqrt{\alpha_{i-1}^2 + C_{i+1,1}^2}}, \quad q_i = \frac{\sigma C_{i+1,1}}{\sqrt{\alpha_{i-1}^2 + C_{i+1,1}^2}}, \quad \sigma = \begin{cases} +1, & \alpha_{i-1} \geq 0, \\ -1, & \alpha_{i-1} < 0. \end{cases}$$

При четном  $i$  заполняются нечетные строки массива  $W$ :  $W_{i+1,1} = \alpha_i$ ,  $W_{i+1,2} = \beta_i$ ,  $W_{i+1,3} = \gamma_i$ ,  $i = 2K$ ,  $K = 0, 1, \dots$ ,  $i \leq N-2$ , и, таким образом, формируются системы (I2) относительно  $f_i$  и  $f_{i+1}$ :

$$W_{i,1} f_i + W_{i,2} f_{i+1} = W_{i,3}, \quad (36)$$

$$W_{i+1,1} f_i + W_{i+1,2} f_{i+1} = W_{i+1,3}.$$

Наконец, если  $N$  нечетное, то значения  $\alpha_{N-2}$ ,  $\beta_{N-2}$  и  $\gamma_{N-2}$  присваиваются элементам  $(2N_0 - 1)$ -й строки массива  $W$ :  $W_{2N_0-1,1} = \alpha_{N-2}$ ,  $W_{2N_0-1,2} = \beta_{N-2}$ ,  $W_{2N_0-1,3} = \gamma_{N-2}$ , чтобы образовать систему (I2) относительно  $f_{N-1}$  и  $f_N$ :  $W_{2N_0-1,1} f_{N-1} + W_{2N_0-1,2} f_N = W_{2N_0-1,3}$ ,  $W_{2N_0,1} f_{N-1} + W_{2N_0,2} f_N = W_{2N_0,3}$ . Правый ход прогонки, а вместе с ним и формирование  $N_0$  систем (I2), из которых могут быть найдены все неизвестные  $f_j$ , завершен.

Далее, каждая из систем (36) приводится к виду (I9) при помощи ортогонального преобразования отражения. Коэффициенты  $l_i$ ,  $m_i$ ,  $r_i$  и  $r_{i+1}$  вычисляются по формулам:

$$l_i = -\sigma \sqrt{W_{i,1}^2 + W_{i+1,1}^2},$$

$$m_i = -pW_{i,2} - qW_{i+1,2}, \quad r_i = -pW_{i,3} - qW_{i+1,3},$$

$$n_i = -qW_{i,2} + pW_{i+1,2}, \quad r_{i+1} = pW_{i+1,3} - qW_{i,3},$$

где

$$p = \frac{\sigma W_{i,1}}{\sqrt{W_{i,1}^2 + W_{i+1,1}^2}}, \quad q = \frac{\sigma W_{i+1,1}}{\sqrt{W_{i,1}^2 + W_{i+1,1}^2}}, \quad \sigma = \begin{cases} +1, & W_{i,1} \geq 0, \\ -1, & W_{i,1} < 0. \end{cases}$$

Вместе с  $f_i$  и  $f_{i+1}$ , которые рассчитываются по формулам  $f_{i+1} = r_{i+1}/n_i$ ,  $f_i = (r_i - m_i/f_{i+1})/l_i$ , метод позволяет найти оценки аб-

солютной и относительной погрешности решения (см.(28)) каждой из систем и тем самым системы (10) (см.(29)). Обратим также внимание на возможность вычисления оценки относительной погрешности решения в норме  $\|f\|_{\infty} = \max_k |f_k|$ .

**ЗАМЕЧАНИЕ.** Чтобы перейти к следующему шагу прогонки, требуется построить матрицу отражения  $Q[i, i+1]$ , или  $Q[i+1, i]$ , и три отраженных вектора. Однако при этом используется только один из двух элементов отраженного вектора, и, следовательно, достаточно учесть только погрешность его вычисления. При оценке нормы возмущения отраженного вектора в ходе прогонки это обстоятельство позволяет в  $\sqrt{2}$  раз уменьшить прежнее значение  $\alpha$ , определенное неравенством (9). Нетрудно показать, что в оценках погрешности решения системы (10) можно считать  $\alpha$  равным  $8\epsilon_1$ .

Обратимся к характеристике трудоемкости метода. На каждом шагу прогонки, всего  $2(N-2)$  шагов, выполняются следующие арифметические операции: 3 сложения + 7 умножений + 2 деления + I извлечение квадратного корня. Кроме того,  $N_0$  раз повторяется 5 сложений + II умножений + 3 деления + I извлечение квадратного корня. На практике мы должны также учесть затраты на вычисление коэффициентов системы (10), которые во многих случаях являются определяющими. Поэтому основное качество трудоемкости метода выражается в линейной зависимости от  $N$  числа затрачиваемых арифметических операций.

**9. Обобщение метода.** Рассмотрим систему линейных алгебраических уравнений с ленточной матрицей:

$$Cf = F, \quad (37)$$

или

$$\sum_{j=1}^{M_0} C_{1,j} f_j = F_1, \quad i = 1, 2, \dots, M_2,$$

$$\sum_{j=1}^{M_1} C_{1,j} f_j = F_1, \quad v = i+j-M_2-1, \quad i = M_2+1, \dots, N-M_3,$$

$$\sum_{j=1}^{M_0} C_{1,j} f_j = F_1, \quad v = N-M_0+j, \quad i = N-M_3+1, \dots, N.$$

Здесь  $M_0, M_1, M_2$  и  $M_3$  - параметры ленточного массива размерности  $N \times M_1$ ;  $M_0$  - общее число "граничных" уравнений, или порядок разностной краевой задачи, представленной системой (37),  $M_1 = M_0 + 1$  -



$b = P[i, i+m] a$ , если

$$b_i = -\sigma \sqrt{a_1^2 + \dots + a_{i+m}^2}, \quad |b_i| \neq 0, \quad \sigma = \begin{cases} +1, & a_i \geq 0, \\ -1, & a_i < 0, \end{cases}$$

$$b_{i+s} = 0, \quad s = 1, 2, \dots, m, \quad b_k = a_k, \quad k \neq i, \quad k \neq i+s;$$

и  $b = P[i+m, i] a$ , если

$$b_{i+m} = -\sigma \sqrt{a_i^2 + \dots + a_{i+m}^2}, \quad |b_{i+m}| \neq 0, \quad \sigma = \begin{cases} +1, & a_{i+m} \geq 0, \\ -1, & a_{i+m} < 0, \end{cases}$$

$$b_{i+s} = 0, \quad s = 0, 1, \dots, m-1, \quad b_k = a_k, \quad k \neq i+m, \quad k \neq i+s.$$

В этих обозначениях мы опишем правый и левый ходы прогонки как способ построения матрицы  $U$ .

Правый ход состоит из  $i-1$  шагов:

$$c^{(j)}_f = F^{(j)}, \quad j = 1, 2, \dots, i-1,$$

$$c^{(j)} = P[j, j+M_2] c^{(j-1)}, \quad c^{(0)} \equiv c,$$

$$F^{(j)} = P[j, j+M_2] F^{(j-1)}, \quad F^{(0)} \equiv F.$$

Здесь на каждом шагу в качестве вектора  $a$  берется  $j$ -й вектор-столбец матрицы  $c^{(j-1)}$ . Из определения  $P[j, j+M]$  следует, что строки с номерами  $i, i+1, \dots, i+M_2-1$  системы

$$c^{(i-1)}_f = F^{(i-1)} \quad (40)$$

содержат только элементы вектора  $z_1$ . На рис. 2 схематично изображен результат правого хода прогонки применительно к системе с матрицей  $C$ , представленной на рис. 1. Целью прогонки является формирование подсистемы относительно  $f_3, f_4, f_5$  и  $f_6$ . Звездочкой отмечены возмущенные элементы матрицы  $C$ . В прямоугольнике заключены строки системы с матрицей  $C^{(2)}$ , содержащие  $f_3, f_4, f_5$  и  $f_6$ . Кроме того, отмечены элементы  $C$ , ставшие нулями.

$$C^{(2)} = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & \cdot & \cdot \\ & & x & x & x & x & x \\ & & & x & x & x & x \\ & & & & x & x & x & x \end{bmatrix}$$

Рис. 2

Левый ход прогонки состоит из  $k$  шагов,  $k = N+1 - M_0 - i$ :

$$\tilde{c}^{(j)}_f = \tilde{F}^{(j)}, \quad j = 1, 2, \dots, k,$$

$$\tilde{c}^{(j)} = P[N+1-j, N+1-M_0-j] \tilde{c}^{(j-1)}, \quad \tilde{c}^{(0)} \equiv c^{(i-1)},$$

$$\tilde{F}^{(j)} = P[N+1-j, N+1-M_2-j]F^{(j-1)}, \quad \tilde{F}^{(0)} \equiv F^{(1-1)}$$

Здесь на каждом шагу в качестве вектора  $a$  берется  $(N+1-j)$ -й вектор-столбец матрицы  $\tilde{C}^{(j-1)}$ . Очевидно, строки с номерами  $i+M_2, i+M_2+1, \dots, i+M_0-1$  системы

$$\tilde{C}^{(k)}x = \tilde{F}^{(k)} \quad (41)$$

содержат только элементы вектора  $z_1$ . Приведение к (41) таково, что строки с номерами  $1, 2, \dots, i+M_2-1$  систем (40) и (41) соответственно совпадают. Поэтому система (41) содержит подсистему (39). Прогонка завершена, и мы можем выписать представление матрицы  $U$  в виде произведения  $N-M_0$  матриц отражений:

$$U = P[M_0+1, M_2+1]x \dots x P[N, N-M_2]x P[i-1, i-1+M_2]x \dots x P[1, 1+M_2]. \quad (42)$$

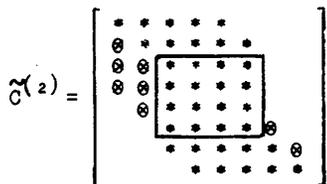


Рис. 3

На рис. 3 схематично представлен результат прогонки в примере, рассматриваемом на рис. 1-2. В прямоугольник заключены элементы матрицы подсистемы относительно  $f_3, f_4, f_5$  и  $f_6$ .

Заметим, что система (39) может быть сформирована из  $M_2$ -х строк системы (40) с номерами  $i, i+1, \dots, i+M_2-1$  и  $M_3$ -х строк системы  $\tilde{U}x = \tilde{F}$ , где  $\tilde{U} = P[M_0+1, M_2+1]x \dots x P[N, N-M_2]$ , т.е. правый ход и левый ход прогонки независимы.

10. Тестовые примеры. Рассмотрим численное решение 3-х примеров иллюстративного характера с целью продемонстрировать универсальность предлагаемого метода. Каждый из примеров - система из  $N$  линейных алгебраических уравнений с трехдиагональной матрицей  $C$ , размерность и обусловленность которой не слишком велики. Вычисления были выполнены на системе HP-2000 с параметрами:  $\epsilon_1 = 1,1921 \cdot 10^{-7}$ ,  $\epsilon_2 = 5,8775 \cdot 10^{-39}$ . Результаты представлены в виде таблиц, выданных на печать ЭВМ, с обозначениями (см. п.п. 6 и 7):

$F(I)$  - значение  $\tilde{F}_1$ ;

$ERR[F(I)]$  - оценка  $|\tilde{F}_1 - f_1|$ , равная  $\rho_1$ ,  $\rho_1 = \|\tilde{F}_1 - z_1\|$  (см. (28));

$COND(L)$  - значение  $\mu(\tilde{L}^{(1)})$ , при котором вычислялось  $\tilde{F}_1$ ;

$F(I) - T(I)$  - реальная разность численного и точного решений,  
т.е.  $\tilde{f}_i - f_i$ .

Кроме того, вычисляются следующие характеристики решения:

$\text{MAX COND}(L)$  - максимальное значение  $\mu(\tilde{L}^1)$ ;

$\text{ERR}(T\emptyset)$  - реальное значение  $\|\tilde{f} - f\|_\infty / \|\tilde{f}\|_\infty$ ;

$\text{ERR}(T2)$  - реальное значение  $\|\tilde{f}_2 - f\| / \|\tilde{f}\|$ ;

$\text{ERR}(F\emptyset)$  - оценка  $\|\tilde{f} - f\|_\infty / \|\tilde{f}\|_\infty$ , равная  $\max \rho_i / \|\tilde{f}\|_\infty$ ;

$\text{ERR}(F2)$  - оценка  $\|\tilde{f} - f\| / \|\tilde{f}\|$ , равная  $\delta_0$  (см. (29)).

При формировании выдачи результатов мы не заботились о том, чтобы печатались только верные знаки.

Приведем еще одну оценку относительной погрешности решения системы (10), характеризующую другой вариант метода ортогональной прогонки. В результате прямого хода прогонки, состоящего из  $N-1$  шагов правого хода в методе ортогональных встречных прогонок, система (10) принимает вид:  $C^{(N-1)}f = F^{(N-1)}$ . Очевидно, согласно (15),  $C^{(N-1)}$  - верхняя треугольная, ленточная матрица с тремя диагоналями [5, 6]. Под обратным ходом прогонки понимается последовательное определение элементов вектора  $f$ , начиная с  $f_N$ . Обусловленность соответствующей  $C^{(N-1)}$  возмущенной матрицы, учитывающей погрешности ортогональных преобразований отражений, ограничена сверху величиной  $\mu_0$ :

$$\mu_0 = \frac{1 + 5\sqrt{N} \alpha e^{5\alpha}}{1 - 5\sqrt{N} \alpha e^{5\alpha}} \mu(C), \quad \alpha = 12\epsilon_1.$$

Следуя [3], выпишем оценку относительной погрешности:

$$\frac{\|\tilde{f} - f\|}{\|\tilde{f}\|} \leq \delta_2 = \frac{\tau_1' + \tau_1''}{1 - \tau_1''}, \quad (43)$$

$$\tau_1' = \frac{5\sqrt{N} \alpha e^{5\alpha} + N \alpha e^{N\alpha}}{1 - 5\sqrt{N} \alpha e^{5\alpha} \mu_0}, \quad \tau_1'' = \epsilon_1 \left( \frac{11 \mu_0}{1 - 9 \mu_0} + 1 \right).$$

В случае, когда обусловленность матрицы  $C$  известна, мы можем сравнить метод ортогональных встречных прогонок и метод ортогональной прогонки при помощи оценок  $\delta_0$ ,  $\delta_1$  (см. (29) и (35)) и  $\delta_2$ . Кроме этих характеристик, мы приведем значения:  $\Delta_1$  - реальная относительная погрешность метода ортогональных встречных прогонок и  $\Delta_2$  - реальная относительная погрешность метода ортогональной прогонки.





В табл. 3 численное решение сопоставляется с точным.

Т а б л и ц а 3

MAX COND(L) = 17.5813

F(I)	ERR(F(I))	COND(L)	F(I) - T(I)
.333333	3.05107E-04	17.5813	-1.79814E-07
.2	3.05107E-04	17.5813	-2.98023E-08
.142857	2.79826E-04	13.0373	-3.94070E-08
.111111	2.79826E-04	13.0373	2.99003E-08
9.09091E-02	2.69702E-04	11.2588	-1.49012E-08
7.69231E-02	2.69702E-04	11.2588	-1.49012E-08
6.66667E-02	2.64724E-04	10.4054	1.49012E-08
5.88235E-02	2.64724E-04	10.4054	7.45058E-09
5.26315E-02	2.61786E-04	10.0433	-5.21541E-08
.047619	2.61786E-04	10.0433	-3.72529E-08
4.34783E-02	2.59895E-04	10.0417	0
.04	2.59895E-04	10.0417	4.47035E-08
3.70371E-02	2.58634E-04	10.3998	7.45058E-08
3.44827E-02	2.58634E-04	10.3998	-2.99003E-08
.032258	2.57797E-04	11.247	-6.70552E-08
.030303	2.57797E-04	11.247	1.86265E-08
2.85716E-02	2.57268E-04	13.0118	1.60187E-07
.027027	2.57268E-04	13.0118	-2.23517E-08
2.56409E-02	2.56978E-04	17.5822	-9.31323E-08
2.43903E-02	2.56978E-04	17.5813	1.49012E-08
2.32559E-02	2.55995E-04	17.5813	1.26660E-07

ERR(T0) = 5.36442E-07      ERR(T2) = 1.49415E-26  
ERR(F0) = 9.15320E-04      ERR(F2) = 1.37468E-03

Параметры  $\delta_0$ ,  $\delta_1$ ,  $\delta_2$ ,  $\Delta_1$  и  $\Delta_2$  имеют следующие значения:  $\delta_0 = 1,87 \cdot 10^{-3}$ ,  $\delta_1 = 1,10 \cdot 10^{-2}$ ,  $\delta_2 = 7,05 \cdot 10^{-3}$ ,  $\Delta_1 = 1,49 \cdot 10^{-6}$ ,  $\Delta_2 = 1,99 \cdot 10^{-7}$ .

ПРИМЕР 3.  $N = 6$ ,  $\mu(C) = 509,1$ ,

$$(C) = \begin{bmatrix} 7/5 & II/3 & & & & \\ 0 & 7/5 & II/3 & & & \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ & & 0 & 7/5 & II/3 & \\ & & & 0 & 7/5 & \end{bmatrix}.$$

В табл.4 численное решение сопоставляется с точным. Параметры  $\delta_0$ ,  $\delta_1$ ,  $\delta_2$ ,  $\Delta_1$  и  $\Delta_2$  имеют следующие значения:  $\delta_0 = 4,03 \cdot 10^{-3}$ ,  $\delta_1 = 1,48 \cdot 10^{-2}$ ,  $\delta_2 = 1,42 \cdot 10^{-2}$ ,  $\Delta_1 = 3,69 \cdot 10^{-6}$ ,  $\Delta_2 = 2,91 \cdot 10^{-6}$ .

Использование ЭМ с меньшими значениями  $\epsilon_1$  уменьшило бы "запас прочности" оценок, но не радикальным образом.

Т а б л и ц а 4

MAX COND(L) = 400.114			
F(I)	ERR(F(I))	COND(L)	F(I)-T(I)
.333333	1.79119E-03	400.114	-6.55651E-07
.2	1.79119E-03	400.114	2.93023E-07
.142857	2.15359E-04	59.3575	-1.19209E-07
.111111	2.15359E-04	59.3575	5.96046E-09
9.09091E-02	2.71510E-05	9.74506	0
7.69231E-02	2.71510E-05	9.74506	0
ERR(T0) = 1.96696E-06		ERR(T2) = 3.69830E-06	
ERR(F0) = 5.34359E-03		ERR(F2) = 4.03207E-03	

II. **З а к л ю ч е н и е.** В сравнении с другими универсальными методами решения системы линейных алгебраических уравнений с трехдиагональной матрицей предлагаемый вариант имеет, на наш взгляд, достоинство, которое выражается в возможностях контроля погрешности в различных нормах. При этом затраты, связанные с контролем, вполне соизмеримы с затратами на само решение. В частности, контролем не отвергается случай, когда, как в примере 3, часть элементов вектора  $f$  оказывается вычисленной с приемлемой точностью, но погрешность решения системы может не удовлетворять заданным ограничениям.

Для определения оценки относительной погрешности нам не потребовалось привлекать трудоемкую процедуру вычисления обусловленности матрицы системы. Тем не менее приведенные примеры показывают, что найденная относительная погрешность  $\delta_0$  (формула (29)), определяемая из решения системы, эффективнее оценки  $\delta_2$  (формула (43)), и, видимо, это следует ожидать вообще ( $\rho_0$  меньше  $\rho_1$  (формула (35)) по построению).

В данной работе не ставилась задача исчерпывающего анализа предлагаемого метода. Мы не можем также сказать в настоящее время, что метод опробован на большом числе практических примеров. Но, как нам кажется, указанных обстоятельств достаточно, чтобы оправдать наш интерес к ортогональным встречным прогонкам.

Автор выражает благодарность С.К.Годунову за внимание и интерес к данной работе, содержание которой во многом отражает попытку автора ответить на его вопросы и замечания, а также участникам семинара, руководимого Д.С.Завьяловым, и семинара, руководимого С.К.Годуновым, за обсуждение работы.

## Л и т е р а т у р а

1. Вычислительные методы линейной алгебры. Библиографический указатель 1928-1974 г. /Под редакцией В.В.Воеводина.- Новосибирск, 1976. - 418 с.
2. САМАРСКИЙ А.А., НИКОЛАЕВ Е.С. Методы решения сеточных уравнений. М.: Наука, 1978. - 592 с.
3. ГОДУНОВ С.К. Решение систем линейных уравнений. Новосибирск: Наука, 1980. - 177 с.
4. БЕЛЛМАН Р. Введение в теорию матриц.-М.: Наука, 1976.-352с.
5. УИЛКИНСОН Дж. РАЙНС С. Справочник алгоритмов на языке АЛГОЛ. Линейная алгебра. -М.: Машиностроение, 1976. - 389 с.
6. ВОЕВОДИН В.В. Численные методы алгебры (теория и алгоритмы).-М.: Наука, 1966.

Поступила в ред.-изд.отд.  
15 июня 1983 года