

где $Q = Q^{(1)}Q^{(2)}\dots Q^{(k)}$. Для системы уравнений $Mx = 0$ полный ортогональный базис подпространства векторов, удовлетворяющих $Mx = 0$, составляют векторы e_{p+1}, \dots, e_n , где e_i — вектор, у которого все компоненты, кроме i -й, равны 0, а i -я компонента равна 1. Тогда система векторов $q_{p+1}, q_{p+2}, \dots, q_n$, где q_i — i -й столбец матрицы Q , составляет полный ортогональный базис подпространства векторов, удовлетворяющих системе уравнений $Bx = 0$. Это следует из того, что Q — ортогональная матрица и в силу (3.12).

Пусть \tilde{M} — квадратная $p \times p$ матрица такая, что $\tilde{m}_{ij} = m_{ij}$ при $i, j = 1, 2, \dots, k$, т. е. $\tilde{M} = (\tilde{M}^0)$. Решаем систему линейных уравнений $\tilde{M}y = \varphi$, где y — вектор размерности k . Так как $m_{ii} \neq 0$ для любого i , система имеет единственное решение y_0 . Пусть x_0 — вектор размерности n , первые k компонент которого равны компонентам вектора y_0 , остальные компоненты равны 0. Тогда ясно, что $Mx_0 = \varphi$. Рассмотрим вектор $x_1 = Qx_0$. Поскольку x_1 есть линейная комбинация векторов q_1, q_2, \dots, q_n , так как последние p компонент x_0 равны 0, а Q — ортогональная матрица, x_0 ортогонален векторам $q_{k+1}, q_{k+2}, \dots, q_n$. В итоге имеем $z_1(x_0) = q_{k+1}, \dots, z_p(x_0) = q_n, z_0(x_0) = x_1$. Рассуждая так же, как и в п. 1 § 3, можно получить оценки: $\|Q - Q^0\| \leq k\sqrt{k}\delta^{20}$, где $\delta = 20e_1$, $\|\tilde{M} - M\| \leq k\sqrt{k}\delta^{20}\|B\|$. Затем, если воспользоваться оценкой точности решения системы с возмущенной треугольной матрицей, которую можно найти в [6], и использовать оценку точности перемножения двух матриц из [3], то можно оценить точность нахождения вектора x_1 .

4. Некоторые замечания относительно погрешностей v_i и $\delta^{(s)}$: v_i представляет собой погрешность округления при сложении $p + 1$ векторов, которую легко оценить с помощью неравенств для сложения векторов; $\delta^{(s)}$ — погрешность решения системы линейных уравнений с треугольной матрицей на вычислительной машине. Ее тоже можно оценить, тем более что число обусловленности матрицы \tilde{R}_s системы для определения β_s ограничено для любого s : $\mu(\tilde{R}_s) \leq 4e^{20}$.

ЛИТЕРАТУРА

- Годунов С. К. О численном решении краевых задач для систем линейных обыкновенных дифференциальных уравнений. — Успехи мат. наук, 1961, т. 16, № 3, с. 171—175.
- Годунов С. К. Метод ортогональной прогонки для решения систем разностных уравнений. — Журн. вычисл. математики и мат. физики, 1962, т. 2, № 6, с. 972—983.
- Годунов С. К. Решение систем линейных уравнений. — Новосибирск: Наука. Сиб. отд-ние, 1980. — 177 с.
- Наймарк М. А. Линейные дифференциальные операторы. — М.: Наука, 1969. — 373 с.
- Кутателадзе С. С. Основы функционального анализа. — Новосибирск: Наука. Сиб. отд-ние, 1983. — 221 с.
- Уилкинсон Дж. Х. Алгебраическая проблема собственных значений. — М.: Наука, 1970. — 564 с.

АЛГОРИТМЫ ИСЧЕРПЫВАНИЯ ТРЕХДИАГОНАЛЬНЫХ СИММЕТРИЧЕСКИХ И ДВУХДИАГОНАЛЬНЫХ МАТРИЦ С ГАРАНТИРОВАННОЙ ОЦЕНКОЙ ТОЧНОСТИ

A. D. МИТЧЕНКО

ВВЕДЕНИЕ

Работа посвящена выводу нового варианта формул и анализу погрешностей при реализации так называемых алгоритмов ортогонального исчерпывания трехдиагональных симметрических и двухдиагональных

матриц, с помощью которых эти матрицы приводятся к диагональному виду. Если известны собственное значение λ симметрической трехдиагональной матрицы порядка m и соответствующий ему собственный вектор, то алгоритм исчерпывания позволяет с помощью подобных ортогональных преобразований вращения привести эту матрицу к клеточно-диагональной форме, в которой одна клетка онять является симметрической трехдиагональной матрицей порядка $(m-1)$, а вторая имеет порядок 1 и просто совпадает с λ . Циклическое повторение процесса исчерпывания дает возможность привести исходную матрицу к диагональной форме. Этот алгоритм хорошо известен и широко освещен в литературе (см., например, [1]). Вариант для сингулярного разложения двухдиагональной матрицы описан в [2].

При практическом применении известных реализаций алгоритмов исчерпывания иногда обнаруживается численная их неустойчивость, состоящая в том, что окончательная форма преобразованной матрицы отличается от желаемой. В работе выяснены причины этой неустойчивости и приведены способы их устранения, благодаря чему разработаны численно устойчивые варианты алгоритмов исчерпывания трехдиагональных симметрических и двухдиагональных матриц с гарантированной оценкой точности.

Продемонстрируем проблемы, возникающие при реализации алгоритмов исчерпывания, на примере симметрической трехдиагональной матрицы

$$A = \begin{bmatrix} d_1 & b_2 & & & 0 & & \\ b_2 & d_2 & b_3 & & & & \\ & & \ddots & & & & \\ & & & b_{m-1} & d_{m-1} & b_m & \\ 0 & & & & b_m & d_m & \end{bmatrix},$$

не имеющей нулевых элементов на побочных диагоналях, $b_i \neq 0$. Предположим, что известны собственное значение λ этой матрицы и отвечающий ему собственный вектор $u = (u_1, u_2, \dots, u_m)^T$, так что имеют место равенства

$$\begin{aligned} (d_1 - \lambda)u_1 + b_2u_2 &= 0; \\ b_iu_{i-1} + (d_i - \lambda)u_i + b_{i+1}u_{i+1} &= 0, \quad i = 2, \dots, m-1; \\ b_mu_{m-1} + (d_m - \lambda)u_m &= 0. \end{aligned} \quad (1)$$

Традиционно алгоритм исчерпывания состоит в следующем. Положим $u^{(1)} = u$ и подберем ортогональные матрицы вращения

$$C_i = \begin{bmatrix} 1 & & & & 0 & & \\ & 1 & & & & & \\ & & \ddots & & & & \\ & & & 1 & & & \\ & & & & c_i - s_i & & \\ & & & & s_i & c_i & \\ & & & & & & \\ 0 & & & & & & 1 \\ & & & & & & & \ddots \\ & & & & & & & & 1 \end{bmatrix}, \quad i = 2, \dots, m,$$

так, чтобы вектор $u^{(i)} = C_i u^{(i-1)}$ имел вид $u^{(i)} = (0, 0, \dots, 0, u_i^{(i)}, u_{i+1}, \dots, u_m)^T$. Другими словами, параметры $s_i, c_i (s_i^2 + c_i^2 = 1)$ находятся как решение уравнения $c_i u_{i-1}^{(i-1)} - s_i u_i = 0$. При этом компонента $u_i^{(i)} = s_i u_{i-1}^{(i-1)} + c_i u_i$. Одно из решений уравнения $c_i u_{i-1}^{(i-1)} - s_i u_i = 0$ есть

$$s_i = \frac{u_{i-1}^{(i-1)}}{\sqrt{[u_{i-1}^{(i-1)}]^2 + u_i^2}}, \quad c_i = \frac{u_i}{\sqrt{[u_{i-1}^{(i-1)}]^2 + u_i^2}}.$$

Тогда $u_i^{(i)} = \sqrt{[u_{i-1}^{(i-1)}]^2 + u_i^2}$ и выражения для s_i , c_i можно переписать так:

$$s_i = \frac{u_{i-1}^{(i-1)}}{u_i^{(i)}}, \quad c_i = \frac{u_i}{u_i^{(i)}}. \quad (2)$$

Отметим, что знаменатель в этих формулах не равен нулю, так как, очевидно, $u_i \neq 0$ и, следовательно, $u_i^{(i)} = \sqrt{u_1^2 + u_2^2 + \dots + u_i^2} \neq 0$. Поделим обе части i -го ($i = 1, 2, \dots, m-1$) из равенств (1) на $u_{i+1}^{(i+1)}$ и последнего — на $u_m^{(m)}$. В результате простых преобразований, использующих формулы (2), получим

$$\begin{aligned} s_2(d_1 - \lambda) + c_2 b_2 &= 0; \\ s_3 s_2 b_2 + s_3 c_2 (d_2 - \lambda) + c_3 b_3 &= 0; \\ s_{i+1} s_i c_{i-1} b_i + s_{i+1} c_i (d_i - \lambda) + c_{i+1} b_{i+1} &= 0, \quad i = 3, \dots, m-1; \\ s_m c_{m-1} b_m + c_m (d_m - \lambda) &= 0. \end{aligned} \tag{3}$$

С помощью этих равенств покажем, что матрица $\bar{A} = C_m C_{m-1} \cdots C_3 C_2 A C_2^* \times \cdots \times C_3^* \cdots C_{m-1}^* C_m^*$ имеет вид

$$\bar{A} = \begin{bmatrix} \bar{d}_1 & \bar{b}_2 & & & 0 \\ \bar{b}_2 & \bar{d}_2 & \bar{b}_3 & & \\ \vdots & \vdots & \ddots & & \\ 0 & \ddots & \ddots & \bar{d}_{m-2} & \bar{b}_{m-1} \\ & & & \bar{b}_{m-1} & \bar{d}_{m-1} \\ & & & 0 & 0 \end{bmatrix}, \quad (4)$$

т. е. матрица A в результате применения к ней $m - 1$ подобных преобразований вращения приводится к клеточно-диагональной форме, в которой одномерная клетка совпадает с собственным значением λ .

Обозначив $A_i = C_4 C_{4-1} \cdots C_3 C_2 A C_2^* C_3^* \cdots C_{i-1}^* C_i^*$, покажем сначала, что матрица A_i выглядит так:

$$A_i = \begin{bmatrix} \bar{d}_1 & \bar{b}_2 \\ \bar{b}_2 & \bar{d}_2 & \bar{b}_3 \\ & \ddots & \ddots & \ddots \\ & & \bar{b}_{i-1} & \bar{d}_{i-1} & u_i & -s_i b_{i+1} \\ & & & u_i & w_i & c_i b_{i+1} \\ & & -s_i b_{i+1} & c_i b_{i+1} & d_{i+1} & b_{i+2} \\ & & & & \ddots & \ddots \\ & & & & & b_{m-1} & d_{m-1} & b_m \\ & & & & & b_m & d_m & \end{bmatrix} \quad (5)$$

где для элементов u_i , w_i имеют место формулы

$$u_i = -s_i[s_i c_{i-1} b_i + c_i(d_i - \lambda)]; \quad w_i = c_i[s_i c_{i-1} b_i + c_i(d_i - \lambda)] + \lambda.$$

Для этого предположим, что матрица A_{i-1} имеет такой же вид, если значение индекса i заменить на $i-1$, и вычислим матрицу $A_i = C_i A_{i-1} C_i^*$. Нетрудно проверить, что вычисленная по этой формуле матрица A_i записывается следующим образом:

$$A_i = \begin{bmatrix} \bar{d}_1 & \bar{b}_1 \\ \bar{b}_2 & \bar{d}_2 & \bar{b}_3 \\ & \vdots & \vdots \\ & \bar{b}_{i-2} & \bar{d}_{i-3} & \bar{b}_{i-1} & z_{i-1} \\ & & \bar{b}_{i-1} & \bar{d}_{i-1} & u_i & -s_i b_{i+1} \\ & & z_{i-1} & u_i & w_i & c_i b_{i+1} \\ & & & -s_i b_{i+1} & c_i b_{i+1} & d_{i+1} & b_{i+2} \\ & & & & & \vdots & \vdots \\ & & & & & b_{m-1} & d_{m-1} & b_m \\ & & & & & b_m & d_m & \end{bmatrix} \quad 0$$

где использованы обозначения

$$\begin{aligned} \bar{b}_{i-1} &= c_i u_{i-1} + s_i s_{i-1} b_i; & z_{i-1} &= s_i u_{i-1} - c_i s_{i-1} b_i; \\ \bar{d}_{i-1} &= c_i (c_i w_{i-1} - s_i c_{i-1} b_i) - s_i (c_i c_{i-1} b_i - s_i d_i); \\ u_i &= c_i (s_i w_{i-1} + c_i c_{i-1} b_i) - s_i (s_i c_{i-1} b_i + c_i d_i); \\ w_i &= s_i (s_i w_{i-1} + c_i c_{i-1} b_i) + c_i (s_i c_{i-1} b_i + c_i d_i). \end{aligned} \quad (6)$$

Подставляя выражения для u_{i-1} , w_{i-1} в формулы для z_{i-1} , u_i , w_i , находим

$$\begin{aligned} z_{i-1} &= -s_{i-1} [s_i s_{i-1} c_{i-2} b_{i-1} + s_i c_{i-1} (d_{i-1} - \lambda) + c_i b_i]; \\ u_i &= c_i c_{i-1} [s_i s_{i-1} c_{i-2} b_{i-1} + s_i c_{i-1} (d_{i-1} - \lambda) + c_i b_i] - s_i [s_i c_{i-1} b_i + c_i (d_i - \lambda)]; \\ w_i &= s_i c_{i-1} [s_i s_{i-1} c_{i-2} b_{i-1} + s_i c_{i-1} (d_{i-1} - \lambda) + c_i b_i] + c_i [s_i c_{i-1} b_i + c_i (d_i - \lambda)] + \lambda. \end{aligned}$$

Используя теперь равенства (3), видим, что $z_{i-1} = 0$, а u_i , w_i имеют требуемые выражения

$$u_i = -s_i [s_i c_{i-1} b_i + c_i (d_i - \lambda)]; \quad w_i = c_i [s_i c_{i-1} b_i + c_i (d_i - \lambda)] + \lambda.$$

Таким образом, предположение, что матрица A_i имеет вид (5), обосновано.

Рассмотрим теперь матрицу A_i при $i = m$, т. е. матрицу $A_m = \bar{A}$. Чтобы убедиться в том, что она имеет вид (4), выпишем формулы для элементов u_m , w_m :

$$\begin{aligned} u_m &= -s_m [s_m c_{m-1} b_m + c_m (d_m - \lambda)]; \\ w_m &= c_m [s_m c_{m-1} b_m + c_m (d_m - \lambda)] + \lambda. \end{aligned}$$

Выражение, заключенное в квадратные скобки, представляет собой в точности левую часть последнего из равенств (3). Следовательно, $u_m = 0$, $w_m = \lambda$ и матрица \bar{A} действительно имеет вид (4). Переход от матрицы A к подобной ей матрице \bar{A} (4) и составляет сущность алгоритма исчерпывания.

Особо подчеркнем, что при получении матрицы \bar{A} использованы все равенства (3). В то же время последнее из них является в некотором смысле лишним. Конкретнее, для определения $(m-1)$ пар неизвестных s_i , c_i достаточно первых $(m-1)$ из этих равенств, а последнее должно выполняться автоматически, если λ — точное собственное значение. Если же λ мало отличается от точной величины собственного значения, то, казалось бы, это последнее уравнение должно быть почти выполненным. В действительности же бывают случаи, когда это не так, т. е. последнее уравнение имеет большую (по отношению к норме решения) невязку, хотя все остальные уравнения решены точно. По-видимому, первым обратил на это внимание Уилкинсон [см. [3], с. 286]. Один из наиболее показательных в этом отношении примеров рассмотрен в [4].

На практике параметры s_i , c_i подбираются именно из первых $(m-1)$ равенств (3). При этом из-за того, что последнее уравнение может иметь большую невязку, значения элементов u_m , w_m могут сильно отличаться от их точных значений, в частности значение элемента u_m может оказаться

ся недостаточно малым, чтобы им можно было пренебречь. Таким образом, первая причина неустойчивости алгоритма исчерпывания связана с недооценкой роли последнего уравнения, которое обычно считается выполненным. На первый взгляд, эту причину можно устранить, вычислив собственный вектор так, чтобы хорошо удовлетворялись все уравнения (1), — а значит, и уравнения (3), так как они являются следствием (1), — включая последнее, и определяя параметры s_i, c_i по формулам (2), т. е. непосредственно через компоненты собственного вектора. Однако простое преобразование этих формул к виду

$$s_i = \frac{1}{\sqrt{1 + c_{i-1}^2 \left(\frac{u_i}{u_{i-1}} \right)^2}}, \quad c_i = \frac{\frac{u_i}{u_{i-1}}}{\sqrt{1 + c_{i-1}^2 \left(\frac{u_i}{u_{i-1}} \right)^2}} \quad (7)$$

показывает, что параметры s_i, c_i определяются не столько самими компонентами собственного вектора, сколько их отношениями u_i/u_{i-1} . В таком случае ясно, что даже чрезвычайно точное определение собственного вектора в метрике евклидова пространства не гарантирует правильности выбора параметров s_i, c_i , поскольку оно все равно может приводить к неправильным отношениям компонент данного вектора. Это происходит, например, в том случае, когда среди компонент наряду с большими есть очень малые. Таким образом, вторая причина неустойчивости связана с недостаточным определением отношений компонент собственного вектора.

Задача вычисления отношений компонент собственного вектора с высокой точностью решена в [4], результаты которой использованы в данной работе. В частности, в главе 2 показано, что использование этих отношений позволяет добиться достаточно точного удовлетворения соотношений (3). При этом найдены некоторые «идеальные» параметры s_i, c_i ($s_i^2 + c_i^2 = 1$), которые удовлетворяют соотношениям (3) приближенно, т. е. с некоторыми возмущениями коэффициентов d_i, b_i . Приведена оценка этих возмущений. Параметры s_i, c_i названы идеальными по той причине, что они в точности удовлетворяют соотношению $s_i^2 + c_i^2 = 1$ и используются только при выводе формул алгоритма, т. е. только в теоретических рассмотрениях. Реально же в памяти машины хранятся лишь некоторые их «машинные» приближения s'_i, c'_i , мало (в смысле относительной погрешности) отличающиеся от s_i, c_i . Чтобы различие между параметрами s_i, c_i и s'_i, c'_i стало понятным, скажем, что s_i, c_i определяются в результате точных вычислений по формулам (7), в которых вместо c_{i-1} используется c'_{i-1} , а s'_i, c'_i представляют собой результаты машинной реализации этих формул. Обеспечение высокой относительной точности при вычислении s_i, c_i важно по нескольким причинам. Во-первых, от величины погрешности параметров c'_i зависит величина возмущений коэффициентов d_i, b_i , с которыми выполняются равенства (3) для параметров s_i, c_i . Во-вторых, параметры s_i, c_i используются при построении явных формул, задающих элементы \tilde{d}_i, \tilde{b}_i преобразованной матрицы \tilde{A} [см. (4)]. Понятно, что при реальных вычислениях в этих формулах можно использовать лишь приближения s'_i, c'_i . Поэтому от их точности зависит точность определения преобразованной матрицы \tilde{A} . Наконец, важно обеспечить, чтобы приближенная матрица C_i (см. с. 115) была близка к ортогональной, ведь результатами работы алгоритма, кроме преобразованной матрицы, являются и матрицы, задающие требуемые преобразования вращения. Алгоритм вычисления величин s_i, c_i , обеспечивающий их высокую относительную точность, описан в § 4 главы 2; § 3 той же главы посвящен непосредственно решению уравнений (3) относительно параметров s_i, c_i .

В главе 1 на основе соотношений (3) выводятся явные формулы для элементов преобразованной матрицы \tilde{A} . В традиционно используемом ал-

$$C'_i = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & c'_i - s'_i & \\ & & & s'_i & c'_i \\ 0 & & & & 1 \\ & & & & \\ & & & & 1 \end{bmatrix}$$

горитме исчертывания формулы для этих элементов имеют, как нетрудно видеть из (6), рекуррентный характер, что, конечно, затрудняет анализ погрешностей, возникающих при их реализации на вычислительной машине. В данной работе приведены окончательные формулы для этих элементов, так что сами преобразования вращения выполняются не в процессе вычислений, а аналитически. Поскольку уравнения (3) могут быть решены только приближенно, преобразованная матрица уже не будет иметь клеточно-диагональной формы (4), а окажется заполненной. В работе показано, что эта заполненная матрица может быть представлена в виде суммы двух матриц, одна из которых имеет требуемый вид (4), а другая, называемая матрицей погрешностей, — малую норму. Выводу формул для элементов преобразованной матрицы и оценке нормы матрицы погрешностей посвящен § 1.

При вычислении элементов преобразованной матрицы по полученным формулам неизбежно возникновение погрешностей. Эти погрешности имеют двоякую природу. Во-первых, как отмечено, вместо участвующих в этих формулах параметров s_i, c_i можно использовать только их машинные приближения s'_i, c'_i . Во-вторых, погрешности возникают и при непосредственном выполнении арифметических операций, предписываемых этими формулами. Оба источника погрешностей учитываются в § 5 главы 3.

Проблемы, аналогичные описанным, имеют место в алгоритме исчертывания двухдиагональной матрицы

$$A = \begin{bmatrix} a_1 & b_2 & & & 0 & & \\ & a_2 & b_3 & & & & \\ & & \ddots & & & & \\ 0 & & & a_{N-1} & b_N & & \\ & & & & a_N & & \end{bmatrix}$$

Кратко опишем его, предполагая элементы a_i, b_i отличными от нуля. Пусть известно сингулярное число σ этой матрицы. Напомним, что сингулярными числами матрицы A называются N наибольших собственных значений составной матрицы

$$S = \begin{bmatrix} 0 & A \\ A^* & 0 \end{bmatrix}.$$

Алгоритм состоит в определении ортогональных матриц вращения C_i, \bar{C}_i ($i = 2, \dots, N$)

$$C_i = \begin{bmatrix} 1 & & & & & & \\ & 1 & & & & & \\ & & 1 & & & & \\ & & & c_i - s_i & & & \\ & & & s_i & c_i & & \\ 0 & & & & & 1 & \\ & & & & & & \\ & & & & & & 1 \end{bmatrix};$$

$$\bar{C}_i = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & \\ & & & & c_i - \bar{s}_i \\ & & & & \bar{s}_i & c \\ 0 & & & & & 1 \\ & & & & & & \ddots \\ & & & & & & 1 \end{bmatrix}$$

так, чтобы матрица $\bar{A} = \bar{C}_N \bar{C}_{N-1} \cdots \bar{C}_3 \bar{C}_2 A C_2^* C_3^* \cdots C_{N-1}^* C_N^*$ имела вид

$$\bar{A} = \begin{bmatrix} \bar{a}_1 & \bar{b}_2 & & & & 0 & & \\ & \bar{a}_2 & \bar{b}_3 & & & & & \\ & & & \ddots & & & & \\ & & & & \bar{a}_{N-2} & \bar{b}_{N-1} & & \\ 0 & & & & & \bar{a}_{N-1} & 0 & \\ & & & & & & & \pm \sigma \end{bmatrix}. \quad (8)$$

Аналогично случаю трехдиагональной матрицы нетрудно показать, что параметры $s_i, c_i, \bar{s}_i, \bar{c}_i (s_i^2 + c_i^2 = 1, \bar{s}_i^2 + \bar{c}_i^2 = 1)$ должны удовлетворять равенствам

$$\begin{aligned} s_2 a_1 - s_2 \frac{\sigma^2}{a_1} + c_2 b_2 &= 0; \\ \bar{s}_2 b_2 - c_2 \frac{s_2}{\bar{s}_2} a_1 + \bar{c}_2 a_2 &= 0; \\ s_{i+1} c_i a_i - c_i \frac{s_{i+1} \dots s_2}{s_i \dots s_2} \frac{\sigma^2}{a_1} + c_{i+1} b_{i+1} &= 0, \\ \bar{s}_{i+1} \bar{c}_i b_{i+1} - c_{i+1} \frac{\bar{s}_{i+1} \dots \bar{s}_2}{\bar{s}_{i+1} \dots \bar{s}_2} a_1 + \bar{c}_{i+1} a_{i+1} &= 0 \\ c_N a_N - c_N \frac{s_{N-1} s_2 \sigma^2}{s_{N-1} \dots s_2} a_1 &= 0. \end{aligned} \quad i = 2, \dots, N-1, \quad (9)$$

Первое из этих равенств можно рассматривать как уравнение для определения параметров s_2, c_2 , второе — как уравнение для \bar{s}_2, \bar{c}_2 и т. д. При этом последнее равенство остается неиспользованным. Если σ известно с некоторой, пусть даже достаточно малой, погрешностью, то неиспользованное уравнение может иметь большую невязку. Это приводит к тому, что элемент матрицы \bar{A} , стоящий в позиции $(N-1, N)$, теоретически равный нулю, в практических расчетах оказывается недостаточно малым, чтобы им можно было пренебречь. Поэтому возникает задача научиться достаточно точно решать уравнения (9). При ее решении опять будем опираться на возможность предварительного определения отношений компонент собственного вектора с высокой точностью (см. [4]). В данном случае речь идет о собственном векторе трехдиагональной симметрической матрицы

$$T = \begin{bmatrix} 0 & a_1 & & & & & & & \\ a_1 & 0 & b_2 & & & & & & 0 \\ & b_2 & 0 & a_2 & & & & & \\ & & a_2 & 0 & b_3 & & & & \\ & & & & \ddots & & & & \\ & & & & & \ddots & & & \\ & & & & & & a_{N-1} & 0 & b_N \\ 0 & & & & & & b_N & 0 & a_N \\ & & & & & & & a_N & 0 \end{bmatrix}.$$

отвечающем ее почти собственному значению σ . Отметим, что матрица T получается из матрицы S с помощью одноименных перестановок строк и столбцов. Задачам определения параметров $s_i, c_i, \bar{s}_i, \bar{c}_i$, удовлетворяющих (приближенным) равенствам (9), и вычисления их достаточно точных машинных приближений $s'_i, c'_i, \bar{s}'_i, \bar{c}'_i$ посвящены § 3, 4 главы 2, где они решаются параллельно с аналогичными вопросами для случая трехдиагональной матрицы. Отметим, что с достаточной точностью удалось решить уравнения (9) только в случае, когда σ является наибольшим сингулярным числом матрицы A .

Поскольку уравнения (9) можно решить только приближенно, т. е. они выполняются с некоторыми возмущенными коэффициентами a_i, b_i , матрица \bar{A} уже не может иметь вида (8), а оказывается заполненной. Однако, как показано в данной работе, эту заполненную матрицу можно представить в виде суммы двух матриц, одна из которых имеет требуемую форму (8), а вторая (матрица погрешностей) — малую норму. Для элементов \bar{a}_i, \bar{b}_i первой матрицы получены простые явные формулы. Их вывод и оценка нормы матрицы погрешностей составляют содержание § 2 главы 1; § 6 главы 3 посвящен анализу погрешностей, возникающих при непосредственном вычислении элементов \bar{a}_i, \bar{b}_i по полученным формулам. В результате анализа находится оценка, показывающая, насколько вычисленная (т. е. рассматриваемая уже как таблица машинных чисел) преобразованная матрица отличается от некоторой матрицы, ортогонально эквивалентной исходной матрице A .

В § 7 и 8 главы 4 приведены общие схемы алгоритмов исчерпывания соответственно трехдиагональных симметрических и двухдиагональных матриц.

Автор считает приятным долгом отметить неизменную поддержку и внимание к его работе со стороны чл.-кор. АН СССР С. К. Годунова и выражает ему искреннюю признательность.

Глава 1

ФОРМУЛЫ АЛГОРИТМОВ ИСЧЕРПЫВАНИЯ.

АНАЛИЗ ПОГРЕШНОСТЕЙ, СВЯЗАННЫХ С НЕТОЧНОСТЬЮ РЕШЕНИЯ УРАВНЕНИЙ ОТНОСИТЕЛЬНО ПАРАМЕТРОВ ВРАЩЕНИЯ

§ 1. Алгоритм исчерпывания симметрической трехдиагональной матрицы

Рассмотрим симметрическую трехдиагональную матрицу

$$A = \begin{bmatrix} d_1 & b_2 & & & & \\ b_2 & d_2 & b_3 & & & 0 \\ & \ddots & \ddots & \ddots & & \\ & & & & \ddots & \\ 0 & & & b_{m-1} & d_{m-1} & b_m \\ & & & & b_m & d_m \end{bmatrix} \quad (1.1)$$

такую, что $b_i \neq 0$. Предположим, что η — некоторое приближение к собственному значению λ этой матрицы, удовлетворяющее неравенству

$$|\eta - \lambda| \leq \gamma M(A), \quad \gamma \ll 1, \quad (1.2)$$

и числа s_i, c_i удовлетворяют соотношениям

$$c_i^2 + s_i^2 = 1, \quad i = 2, \dots, m, \quad (1.3)$$

$$s_i > 0, \quad i = 2, \dots, m; \quad (1.4)$$

$$\begin{aligned}
& s_2(d_1 + \alpha_1 - \eta) + c_2 b_2 = 0; \\
& s_3 s_2 b_2 (1 + \beta_2) + s_3 c_2 (d_2 + \alpha_2 - \eta) + c_3 b_3 = 0; \\
& s_{i+1} s_i c_{i-1} b_i (1 + \beta_i) + s_{i+1} c_i (d_i + \alpha_i - \eta) + c_{i+1} b_{i+1} = 0, \quad i = 3, \dots, m-1; \\
& s_m c_{m-1} b_m (1 + \beta_m) + c_m (d_m + \alpha_m - \eta) = 0,
\end{aligned} \tag{1.5}$$

где возмущения α_i и β_i удовлетворяют оценкам

$$\begin{aligned}
|\alpha_i| &\leq \bar{\alpha} \mathcal{M}(A), \quad \bar{\alpha} \ll 1; \\
|\beta_i| &\leq \bar{\beta}, \quad \bar{\beta} \ll 1,
\end{aligned} \tag{1.6}$$

а $\mathcal{M}(A)$ обозначает одну из норм матрицы A :

$$\mathcal{M}(A) = \max \left\{ \begin{array}{l} |d_1| + |b_2|, \\ \max_{2 \leq i \leq m-1} (|d_i| + |b_i| + |b_{i+1}|), \\ |d_m| + |b_m|. \end{array} \right.$$

Отметим, что получение соотношений (1.5) подробно описано в § 3 главы 2. Там же приведены формулы для нахождения возмущений α_i и β_i .

Определим ортогональные матрицы

$$C_i = \begin{bmatrix} 1 & & & & & & & \\ & \ddots & & & & & & \\ & & 1 & & & & & \\ & & & \frac{c_i - s_i}{s_i} & & & & 0 \\ & & & s_i & c_i & & & \\ & & & & & \ddots & & \\ 0 & & & & & & 1 & \\ & & & & & & & \ddots \\ & & & & & & & & 1 \end{bmatrix}, \quad i = 2, \dots, m,$$

и их произведение $C = C_m C_{m-1} \cdots C_2 C_2$.

Теорема 1. Матрицу CAC^* , ортогонально подобную матрице A , можно представить в виде суммы двух матриц $CAC^* = \tilde{A} + R$, где матрица \tilde{A} имеет вид

$$\tilde{A}_i = \begin{bmatrix} \bar{d}_1 & \bar{b}_2 & & & & & & \\ \bar{b}_2 & \bar{d}_2 & \bar{b}_3 & & & & & 0 \\ & \vdots & \vdots & \ddots & & & & \\ & & & \bar{b}_{m-2} & \bar{d}_{m-2} & \bar{b}_{m-1} & & \\ 0 & & & & \bar{b}_{m-1} & \bar{d}_{m-1} & 0 & \\ & & & & & & 0 & \eta \end{bmatrix}, \tag{1.7}$$

а для спектральной нормы матрицы R имеет место оценка

$$\|R\| \leq 2(\tilde{\alpha} + \sqrt{m}\tilde{\beta})\mathcal{M}(A). \tag{1.8}$$

При этом для элементов \bar{d}_i , \bar{b}_i матрицы \tilde{A} справедливы формулы

$$\begin{aligned}
\bar{d}_1 &= d_2 - \frac{c_2 b_2}{s_2} + \frac{c_3 c_2 b_3}{s_3}; \\
\bar{d}_i &= d_{i+1} - \frac{c_{i+1} c_i b_{i+1}}{s_{i+1}} + \frac{c_{i+2} c_{i+1} b_{i+2}}{s_{i+2}}, \quad i = 2, \dots, m-2; \\
\bar{d}_{m-1} &= d_m - \frac{c_m c_{m-1} b_m}{s_m}; \\
\bar{b}_i &= \frac{s_i b_{i+1}}{s_{i+1}}, \quad i = 2, \dots, m-1.
\end{aligned} \tag{1.9}$$

Доказательство. Обозначив

$$d'_i = d_i + \alpha_i, \quad i = 1, \dots, m, \tag{1.10}$$

введем в рассмотрение вспомогательную матрицу

$$A' = \begin{bmatrix} d'_1 & b_2 & & & & \\ b_2 & d'_2 & b_3 & & & \\ & \ddots & \ddots & \ddots & & \\ & & & & 0 & \\ 0 & & & b_{m-1} & d'_{m-1} & b_m \\ & & & b_m & d'_m & \end{bmatrix}$$

и вычислим матрицу $CA'C^*$. Сформулируем рекуррентное правило для вычисления этой матрицы:

$$A_1 = A', \quad A_i = C_i A_{i-1} C_i^*, \quad i = 2, \dots, m, \quad A_m = CA'C^*.$$

Легко проверить, что матрица $A_2 = C_2 A_1 C_2^*$ имеет вид

$$A_2 = \begin{bmatrix} u_1 & v_2 & -s_2 b_3 & & & \\ v_2 & w_2 & c_2 b_3 & & & \\ -s_2 b_3 & c_2 b_3 & d'_3 & b_4 & & \\ & & \ddots & \ddots & \ddots & \\ 0 & & & b_{m-1} & d'_{m-1} & b_m \\ & & & b_m & d'_m & \end{bmatrix}$$

где

$$u_1 = c_2(c_2 d'_1 - s_2 b_2) - s_2(c_2 b_2 - s_2 d'_2);$$

$$v_2 = c_2(s_2 d'_1 + c_2 b_2) - s_2(s_2 b_2 + c_2 d'_2);$$

$$w_2 = s_2(s_2 d'_1 + c_2 b_2) + c_2(s_2 b_2 + c_2 d'_2).$$

Преобразуем выражения для u_1 , v_2 и w_2 :

$$u_1 = d'_2 + \frac{c_2^2}{s_2} [s_2(d'_1 - \eta) + c_2 b_2] - \frac{c_2 b_2}{s_2} - \frac{c_2}{s_2} [s_3 s_2 b_2 (1 + \beta_2) + s_3 c_2 (d'_2 - \eta)] + s_2 c_2 b_2 \beta_2;$$

$$v_2 = c_2[s_3(d'_1 - \eta) + c_2 b_2] - \frac{s_2}{s_3} [s_3 s_2 b_2 (1 + \beta_2) + s_3 c_2 (d'_2 - \eta)] + s_2^2 b_2 \beta_2;$$

$$w_2 = \eta + s_2[s_2(d'_1 - \eta) + c_2 b_2] + \frac{c_2}{s_3} [s_3 s_2 b_2 (1 + \beta_2) + s_3 c_2 (d'_2 - \eta)] - s_2 c_2 b_2 \beta_2.$$

Принимая во внимание первые два из соотношений (1.5) и обозначения (1.9), (1.10), получаем

$$u_1 = \bar{d}_1 + \alpha_3 + s_2 c_2 b_2 \beta_2; \quad v_2 = c_3 \bar{b}_2 + s_2^2 b_2 \beta_2;$$

$$w_2 = \eta - \frac{c_3 c_2 b_2}{s_3} - s_2 c_2 b_2 \beta_2.$$

Введем обозначения

$$b_{11} = s_2 c_2 b_2 \beta_2; \quad b'_{21} = b'_{12} = s_2^2 b_2 \beta_2; \quad b''_{22} = -s_2 c_2 b_2 \beta_2. \quad (1.11)$$

Тогда матрицу A_2 можно записать так:

$$A'_2 = \begin{bmatrix} \bar{d}_1 + \alpha_3 + b_{11} & c_3 \bar{b}_2 + b'_{12} & -s_2 b_3 & & & \\ c_3 \bar{b}_2 + b'_{21} & \eta - \frac{c_3 c_2 b_2}{s_3} + b''_{22} & c_2 b_3 & & & \\ -s_2 b_3 & c_2 b_3 & d'_3 & b_4 & & \\ & & \ddots & \ddots & \ddots & \\ 0 & & & b_{m-1} & d'_{m-1} & b_m \\ & & & b_m & d'_m & \end{bmatrix}$$

Предположим теперь, что матрица A_i выглядит следующим образом:

$$A_i = \begin{bmatrix} \bar{A}_{(i-1)} + X_{(i)} + B_{(i-1)} & b'_{1i} & 0 \\ \vdots & \vdots & 0 \\ c_{i+1}\bar{b}_i + b'_{i-1,i} - s_i b_{i+1} & b'_{i-2,i} & \\ b'_{i1} \dots b'_{i,i-2} c_{i+1}\bar{b}_i + b'_{i,i-1} & t_i + b'_{ii} & \frac{c_i b_{i+1}}{c_i b_{i+1}} \\ -s_i b_{i+1} & c_i b_{i+1} & A'_{(i+1)} \\ 0 & & \end{bmatrix}. \quad (I)$$

Здесь через $\bar{A}_{(i-1)}$, $X_{(i)}$, $B_{(i-1)}$ обозначены матрицы порядка $i-1$:

$$\bar{A}_{(i-1)} = \begin{bmatrix} \bar{d}_1 & \bar{b}_2 & 0 \\ \bar{b}_2 & \bar{d}_2 & \bar{b}_3 \\ \vdots & \vdots & \vdots \\ 0 & \bar{b}_{i-2} & \bar{d}_{i-2} & \bar{b}_{i-1} \\ & \bar{b}_{i-1} & d_{i-1} & \end{bmatrix},$$

$X_{(i)} = \text{diag}[a_2, a_3, \dots, a_i]$, $B_{(i-1)} = (b_{kl})$ ($k, l = 1, \dots, i-1$);
через $A'_{(i+1)}$ — трехдиагональная матрица порядка $m-i$:

$$A'_{(i+1)} = \begin{bmatrix} d'_{i+1} & b_{i+2} & 0 \\ b_{i+2} & d'_{i+2} & b_{i+3} \\ \vdots & \vdots & \vdots \\ 0 & \vdots & \vdots \\ b_{m-1} & d'_{m-1} & b_m \\ b_m & d'_m & \end{bmatrix},$$

а элемент t_i имеет выражение $t_i = \eta - c_{i+1}c_i b_{i+1}/s_{i+1}$. Ясно, что матрица A_i должна быть симметрической, поэтому $b_{kj} = b_{jk}$ ($k = 2, \dots, i-1$; $j = -1, \dots, k-1$); $b'_{ij} = b'_{ji}$ ($j = 1, \dots, i-1$). Легко проверить, что матрица $A_{i+1} = C_{i+1}A_iC_{i+1}^*$ имеет следующий вид:

$$A_{i+1} = \begin{bmatrix} \bar{A}_{(i-1)} + X_{(i)} + B_{(i-1)} & c_{i+1}b'_{1i} & s_{i+1}b'_{1i} & 0 \\ \vdots & \vdots & \vdots & 0 \\ c_{i+1}b'_{i-2,i} & s_{i+1}b'_{i-2,i} & & \\ z_i & y_i & & \\ c_{i+1}b'_{i1} \dots c_{i+1}b'_{i,i-2} & z_i & u_i & v_{i+1} - s_{i+1}b_{i+2} \\ s_{i+1}b'_{i1} \dots s_{i+1}b'_{i,i-2} & y_i & v_{i+1} & w_{i+1} & c_{i+1}b_{i+2} \\ 0 & & -s_{i+1}b_{i+2} & c_{i+1}b_{i+2} & A'_{(i+2)} \end{bmatrix},$$

где для элементов z_i , y_i , u_i , v_{i+1} и w_{i+1} принятые обозначения

$$z_i = c_{i+1}^2\bar{b}_i + s_{i+1}s_i b_{i+1} + c_{i+1}b'_{i-1,i};$$

$$y_i = s_{i+1}c_{i+1}\bar{b}_i - c_{i+1}s_i b_{i+1} + s_{i+1}b'_{i-1,i};$$

$$u_i = c_{i+1}[c_{i+1}(t_i + b'_{ii}) - s_{i+1}c_i b_{i+1}] - s_{i+1}(c_{i+1}c_i b_{i+1} - s_{i+1}d'_{i+1}); \quad (1.12)$$

$$v_{i+1} = c_{i+1}[s_{i+1}(t_i + b'_{ii}) + c_{i+1}c_i b_{i+1}] - s_{i+1}(s_{i+1}c_i b_{i+1} + c_{i+1}d'_{i+1});$$

$$w_{i+1} = s_{i+1}[s_{i+1}(t_i + b'_{ii}) + c_{i+1}c_i b_{i+1}] + c_{i+1}(s_{i+1}c_i b_{i+1} + c_{i+1}d'_{i+1}).$$

Из определения \bar{b}_i следует, что

$$z_i = \bar{b}_i + c_{i+1} b'_{i-1,i}; \quad y_i = s_{i+1} b'_{i-1,i}.$$

С учетом обозначения для t_i преобразуем выражения для u_i , v_{i+1} и w_{i+1} :

$$\begin{aligned} u_i &= d'_{i+1} - \frac{c_{i+1} c_i b_{i+1}}{s_{i+1}} - \frac{c_{i+1}}{s_{i+2}} [s_{i+2} s_{i+1} c_i b_{i+1} (1 + \beta_{i+1}) + s_{i+2} c_{i+1} (d'_{i+1} - \eta)] + \\ &\quad + c_{i+1}^2 b'_{ii} + c_{i+1} s_{i+1} c_i b_{i+1} \beta_{i+1}; \\ v_{i+1} &= -\frac{s_{i+1}}{s_{i+2}} [s_{i+2} s_{i+1} c_i b_{i+1} (1 + \beta_{i+1}) + s_{i+2} c_{i+1} (d'_{i+1} - \eta)] + \\ &\quad + c_{i+1} s_{i+1} b'_{ii} + s_{i+1}^2 c_i b_{i+1} \beta_{i+1}; \\ w_{i+1} &= \eta + \frac{c_{i+1}}{s_{i+2}} [s_{i+2} s_{i+1} c_i b'_{i+1} (1 + \beta_{i+1}) + s_{i+2} c_{i+1} (d'_{i+1} - \eta)] + \\ &\quad + s_{i+1}^2 b'_{ii} - c_{i+1} s_{i+1} c_i b_{i+1} \beta_{i+1}. \end{aligned}$$

Учитывая соотношения (1.5) и принимая во внимание обозначения (1.9) и (1.10), получаем

$$\begin{aligned} u_i &= \bar{d}_i + \alpha_{i+1} + c_{i+1}^2 b'_{ii} + c_{i+1} s_{i+1} c_i b_{i+1} \beta_{i+1}; \\ v_{i+1} &= c_{i+1} \bar{b}_{i+1} + c_{i+1} s_{i+1} b'_{ii} + s_{i+1}^2 c_i b_{i+1} \beta_{i+1}; \\ w_{i+1} &= t_{i+1} + s_{i+1}^2 b'_{ii} - c_{i+1} s_{i+1} c_i b_{i+1} \beta_{i+1}. \end{aligned}$$

После введения обозначений

$$\begin{aligned} b_{ii} &= c_{i+1}^2 b'_{ii} + c_{i+1} s_{i+1} c_i b_{i+1} \beta_{i+1}; \\ b'_{i,i+1} &= b'_{i+1,i} = c_{i+1} s_{i+1} b'_{ii} + s_{i+1}^2 c_i b_{i+1} \beta_{i+1}; \\ b'_{i+1,i+1} &= s_{i+1}^2 b'_{ii} - c_{i+1} s_{i+1} c_i b_{i+1} \beta_{i+1}; \\ b_{ij} &= b_{ji} = c_{i+1} b'_{ij}, \\ b'_{i+1,j} &= b'_{j,i+1} = s_{i+1} b_{ij} \end{aligned} \quad j = 1, \dots, i-1, \quad (1.13)$$

становится понятным, что матрица A_{i+1} записывается так же, как и матрица A_i , если только значение индекса i заменить на $i+1$. Таким образом, предположение, что матрица A_i ($i = 2, \dots, m-1$) имеет вид (I), обосновано.

Непосредственное вычисление матрицы A_m по формуле $A_m = C_m A_{m-1} C_m^*$ дает

$$A_m = \left[\begin{array}{cc|cc} & & c_m b'_{1,m-1} & s_m b'_{1,m-1} \\ & & \vdots & \vdots \\ & & c_m b'_{m-3,m-1} & s_m b'_{m-3,m-1} \\ \hline \bar{A}_{(m-2)} + X_{(m-1)} + B_{(m-2)} & & z_{m-1} & y_{m-1} \\ \hline c_m b'_{m-1,1} \dots c_m b'_{m-1,m-3} & z_{m-1} & u_{m-1} & v_m \\ s_m b'_{m-1,1} \dots s_m b'_{m-1,m-3} & y_{m-1} & v_m & w_m \end{array} \right],$$

где для элементов z_{m-1} , y_{m-1} , u_{m-1} , v_m и w_m справедливы формулы (1.12) при $i = m-1$. Легко видеть, что $z_{m-1} = \bar{b}_{m-1} + c_m b_{m-2,m-1}$, $y_{m-1} = s_m b_{m-2,m-1}$. Выражения для u_{m-1} , v_m и w_m после простых преобразований могут быть записаны в виде

$$\begin{aligned} u_{m-1} &= d'_m - \frac{c_m c_{m-1} b_m}{s_m} - c_m [s_m c_{m-1} b_m (1 + \beta_m) + c_m (d'_m - \eta)] + \\ &\quad + c_m^2 b'_{m-1,m-1} + c_m s_m c_{m-1} b_m \beta_m; \end{aligned}$$

$$v_m = -s_m [s_m c_{m-1} b_m (1 + \beta_m) + c_m (d'_m - \eta)] + c_m s_m b'_{m-1, m-1} + s_m^2 c_{m-1} b_m \beta_m;$$

$$w_m = \eta + c_m [s_m c_{m-1} b_m (1 + \beta_m) + c_m (d'_m - \eta)] + s_m^2 b'_{m-1, m-1} - c_m s_m c_{m-1} b_m \beta_m.$$

Учет последнего из соотношений (1.5), а также обозначений (1.9) и (1.10) позволяет заключить, что

$$u_{m-1} = \bar{d}_{m-1} + \alpha_m + c_m^2 b'_{m-1, m-1} + c_m s_m c_{m-1} b_m \beta_m;$$

$$v_m = c_m s_m b'_{m-1, m-1} + s_m^2 c_{m-1} b_m \beta_m;$$

$$w_m = \eta + s_m^2 b'_{m-1, m-1} - c_m s_m c_{m-1} b_m \beta_m.$$

Введем обозначения

$$\begin{aligned} b_{m-1, m-1} &= c_m^2 b'_{m-1, m-1} + c_m s_m c_{m-1} b_m \beta_m; \\ b'_{m-1, m} &\equiv b'_{m, m-1} = c_m s_m b'_{m-1, m-1} + s_m^2 c_{m-1} b_m \beta_m; \\ b'_{mm} &= s_m^2 b'_{m-1, m-1} - c_m s_m c_{m-1} b_m \beta_m; \\ b_{m-1, j} &\equiv b_{j, m-1} = c_m b'_{m-1, j}, \\ b'_{mj} &\equiv b'_{jm} = s_m b'_{m-1, j}, \end{aligned} \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} \quad j = 1, \dots, m-2,$$

которые, очевидно, совпадают с обозначениями (1.13) при $i = m-1$, а также положим

$$\begin{aligned} b_{m-1, m} &\equiv b_{m, m-1} = b'_{m-1, m}; \\ b_{mm} &= b'_{mm}; \\ b_{mj} &\equiv b'_{jm} = b'_{mj}, \quad j = 1, \dots, m-2. \end{aligned} \quad (1.14)$$

Ясно, что матрицу A_m можно теперь записать в виде

$$A_m = \bar{A} + X + B, \quad (1.15)$$

где $X = \text{diag} [\alpha_2, \alpha_3, \dots, \alpha_m, 0]$, а элементы b_{ij} ($i, j = 1, 2, \dots, m$) матрицы B вычисляются в соответствии с рекуррентными формулами (1.11), (1.13) (при $i = 2, \dots, m-1$) и (1.14). С другой стороны, в силу обозначений (1.10) понятно, что матрицу $A_m = CAC^*$ можно представить как

$$A_m = CAC^* + CYC^*, \quad (1.16)$$

если обозначить $Y = \text{diag} [\alpha_1, \alpha_2, \dots, \alpha_m]$. Сравнивая между собой (1.15) и (1.16), получим $CAC^* = \bar{A} + X - CYC^* + B$. Введем обозначение $R = X - CYC^* + B$. Так как

$$\|R\| \leq \|X\| + \|Y\| + \|B\|, \quad (1.17)$$

для оценки нормы матрицы R достаточно оценить нормы матриц X , Y и B . Очевидно, что

$$\begin{aligned} \|X\| &= \max_{2 \leq i \leq m} |\alpha_i| \leq \tilde{\alpha} \mathcal{M}(A); \\ \|Y\| &= \max_{1 \leq i \leq m} |\alpha_i| \leq \tilde{\alpha} \mathcal{M}(A). \end{aligned} \quad (1.18)$$

Для оценки нормы матрицы B введем в рассмотрение вспомогательную матрицу

$$D = \begin{bmatrix} 0 & -b_2 \beta_2 & & & & 0 \\ & 0 & -b_3 \beta_3 & & & \\ & & & \ddots & & \\ & & & & \ddots & \\ 0 & & & & & 0 - b_m \beta_m \end{bmatrix}$$

и построим матрицу CDC^* . Принимая во внимание определение матрицы C и выписывая для элементов матрицы CDC^* рекуррентные формулы аналогично тому, как это было сделано для элементов матрицы B , не-трудно убедиться, что элементы последней, стоящие на главной диагонали и выше ее, совпадают с соответствующими элементами матрицы CDC^* . Кроме того, в силу симметричности матрицы B ясно, что ее элементы, стоящие ниже главной диагонали, совпадают с соответствующими элементами матрицы CD^*C^* . Представим матрицу B в виде суммы B_1 и B_2 , где матрица B_1 — верхняя треугольная, содержащая на главной диагонали и выше ее соответствующие элементы матрицы B , а B_2 — нижняя треугольная с нулевой главной диагональю, содержащая ниже ее соответствующие элементы матрицы B . Очевидно, что для фробениусовых норм матриц B_1 и B_2 имеют место оценки $\|B_1\|_E \leq \|CDC^*\|_E = \|D\|_E$, $\|B_2\|_E \leq \|CD^*C^*\|_E = \|D\|_E$ и, следовательно, $\|B\| \leq \|B_1\|_E + \|B_2\|_E \leq 2\|D\|_E$. Так как

$$\begin{aligned} \|D\|_E &= \sqrt{\sum_{i=2}^m (b_i \beta_i)^2} \leq \max_{2 \leq i \leq m} |\beta_i| \sqrt{\sum_{i=2}^m b_i^2} \leq \\ &\leq \tilde{\beta} \sqrt{\sum_{i=2}^m b_i^2} \leq \sqrt{m} \tilde{\beta} \max_{2 \leq i \leq m} |\bar{b}_i| \leq \sqrt{m} \tilde{\beta} \mathcal{M}(A), \end{aligned}$$

оценка нормы матрицы B принимает вид $\|B\| \leq 2\sqrt{m}\tilde{\beta}\mathcal{M}(A)$. С помощью полученной оценки и оценки (1.18) из (1.17) следует окончательная оценка

$$\|R\| \leq 2(\tilde{\alpha} + \sqrt{m}\tilde{\beta})\mathcal{M}(A).$$

Теорема доказана.

§ 2. Алгоритм исчерпывания двухдиагональной матрицы

Рассмотрим двухдиагональную матрицу

$$A = \begin{bmatrix} a_1 & b_2 & & & & & 0 \\ a_2 & b_3 & & & & & \\ & \ddots & \ddots & & & & \\ & & & \ddots & & & \\ & & & & a_{N-1} & b_N & \\ 0 & & & & & & a_N \end{bmatrix}, \quad (2.1)$$

относительно которой предположим, что $a_i \neq 0$, $b_i \neq 0$. Пусть $\sigma > 0$ — некоторое достаточно хорошее приближение к наибольшему сингулярному числу $\sigma_N(A) = \|A\|$ матрицы A , т. е.

$$|\sigma - \|A\|| \leq \delta \|A\|, \quad \delta \ll 1, \quad (2.2)$$

и числа c_i , s_i , \bar{c}_i , \bar{s}_i удовлетворяют соотношениям

$$\left. \begin{array}{l} c_i^2 + s_i^2 = 1, \\ \bar{c}_i^2 + \bar{s}_i^2 = 1, \end{array} \right\} \quad i = 2, \dots, N, \quad (2.3)$$

$$s_i > 0, \quad \bar{s}_i > 0, \quad i = 2, \dots, N, \quad (2.4)$$

$$s_2 a_1 (1 + \hat{\alpha}_1) - s_2 \frac{\sigma^2}{a_1} + c_2 b_2 (1 + \check{\beta}_2) = 0;$$

$$\bar{s}_2 b_2 (1 + \hat{\beta}_2) - c_2 \frac{\bar{s}_2}{s_2} a_1 + \bar{c}_2 a_2 (1 + \check{\alpha}_2) = 0;$$

$$\left. \begin{aligned} s_{i+1}c_i a_i (1 + \hat{\alpha}_i) - \bar{c}_i \frac{s_{i+1} \cdots s_2 \sigma^2}{\bar{s}_i \cdots \bar{s}_2 \bar{a}_1} + c_{i+1} b_{i+1} (1 + \check{\beta}_{i+1}) = 0, \\ \bar{s}_{i+1} \bar{c}_i b_{i+1} (1 + \hat{\beta}_{i+1}) - c_{i+1} \frac{\bar{s}_{i+1} \cdots \bar{s}_2}{\bar{s}_{i+1} \cdots \bar{s}_2} a_1 + \bar{c}_{i+1} a_{i+1} (1 + \check{\alpha}_{i+1}) = 0, \end{aligned} \right\} i = 2, \dots, N-1; \quad (2.5)$$

$$c_N a_N (1 + \hat{\alpha}_N) - \bar{c}_N \frac{s_N \cdots s_2 \sigma^2}{\bar{s}_N \cdots \bar{s}_2 \bar{a}_1} = 0,$$

где возмущения $\hat{\alpha}_i$, $\check{\alpha}_i$, $\hat{\beta}_i$, $\check{\beta}_i$ удовлетворяют оценкам

$$|\hat{\alpha}_i| \leq \varepsilon, \quad |\check{\alpha}_i| \leq \varepsilon, \quad |\hat{\beta}_i| \leq \varepsilon, \quad |\check{\beta}_i| \leq \varepsilon. \quad (2.6)$$

Отметим, что получение соотношений (2.5) подробно описано в § 3 главы 2. Там же приведено значение величины ε , с которым удовлетворяются оценки (2.6) [см. (3.20)].

Определим ортогональные матрицы C_i , \bar{C}_i ($i = 2, \dots, N$) и C_{N+1} :

$$C_i = \begin{bmatrix} 1 & & & & & & & \\ & \ddots & & & & & & \\ & & 0 & & & & & \\ & & & 1 & & & & \\ & & & & \bar{c}_i - s_i & & & \\ & & & & s_i & \bar{c}_i & & \\ & & & & & 1 & & \\ 0 & & & & & & \ddots & \\ & & & & & & & 1 \end{bmatrix}, \quad \bar{C}_i = \begin{bmatrix} 1 & & & & & & & \\ & \ddots & & & & & & \\ & & 0 & & & & & \\ & & & 1 & & & & \\ & & & & \bar{c}_i - \bar{s}_i & & & \\ & & & & \bar{s}_i & \bar{c}_i & & \\ & & & & & 1 & & \\ & & & & & & \ddots & \\ & & & & & & & 1 \end{bmatrix},$$

$$C_{N+1} = \text{diag}[1, 1, \dots, 1, \text{sign } a_1],$$

а также их произведения $C = C_{N+1} \cdot C_N \cdots C_3 \cdot C_2$, $\bar{C} = \bar{C}_N \cdots \bar{C}_3 \cdot \bar{C}_2$.

Теорема 2. Матрица $\bar{C}AC^*$, ортогонально эквивалентная матрице A , может быть представлена в виде суммы двух матриц \bar{A} и R :

$$\bar{C}AC^* = \bar{A} + R, \quad (2.7)$$

где матрица \bar{A} имеет вид

$$\bar{A} = \begin{bmatrix} \bar{a}_1 & \bar{b}_2 & & & & & & 0 \\ & \bar{a}_2 & \bar{b}_3 & & & & & \\ & & \ddots & & & & & \\ 0 & & & \bar{a}_{N-2} & \bar{b}_{N-1} & & & \\ & & & & \bar{a}_{N-1} & 0 & & \end{bmatrix}, \quad (2.8)$$

а для спектральной нормы матрицы R имеет место оценка

$$\|R\| \leq 2\sqrt{2}(\sqrt{N} + 2)\varepsilon\|A\|. \quad (2.9)$$

При этом элементы матрицы \bar{A} можно вычислить по формулам

$$\bar{a}_i = \frac{s_{i+1} a_{i+1}}{\bar{s}_{i+1}}, \quad i = 1, \dots, N-1; \quad (2.10)$$

$$\bar{b}_i = \frac{\bar{s}_i b_{i+1}}{s_{i+1}}, \quad i = 2, \dots, N-1.$$

Доказательство. Для доказательства потребуются равенства

$$\left. \begin{aligned} s_2 a_1 + c_2 b_2 &= s_2 \frac{\sigma^2}{a_1} + \bar{\Phi}_1; \\ \bar{s}_2 s_2 \frac{\sigma^2}{a_1} + \bar{c}_2 c_2 a_2 &= \bar{s}_2 a_1 + \Phi_2; \\ s_{i+1} \frac{s_i \dots s_2}{s_i \dots s_2} a_1 + c_{i+1} \bar{c}_i b_{i+1} &= \frac{s_{i+1} s_i \dots s_2 \sigma^2}{s_i \dots s_2} \frac{a_1}{a_1} + \bar{\Phi}_i; \\ \bar{s}_{i+1} \frac{s_{i+1} \dots s_i \dots s_2}{s_i \dots s_2} \frac{\sigma^2}{a_1} + \bar{c}_{i+1} c_{i+1} a_{i+1} &= \frac{\bar{s}_{i+1} \dots \bar{s}_2}{s_{i+1} \dots s_2} a_1 + \Phi_{i+1}, \end{aligned} \right\} \quad \begin{aligned} (2.11) \\ i = 2, \dots, N-1, \end{aligned}$$

которые достаточно просто следуют из соотношений (2.3) и (2.5), если для величин Φ_i , $\bar{\Phi}_i$ справедливы следующие рекуррентные формулы:

$$\left. \begin{aligned} \Phi_1 &= -s_2 a_1 \hat{a}_1 - c_2 b_2 \check{\beta}_2; \\ \Phi_2 &= -\bar{s}_2 \bar{\Phi}_1 - c_2 (\bar{s}_2 b_2 \hat{\beta}_2 + \bar{c}_2 a_2 \check{\alpha}_2); \\ \bar{\Phi}_i &= -s_{i+1} \Phi_i - \bar{c}_i (s_{i+1} c_i a_i \hat{\alpha}_i + c_{i+1} b_{i+1} \check{\beta}_{i+1}), \\ \Phi_{i+1} &= -\bar{s}_{i+1} \bar{\Phi}_i - c_{i+1} (\bar{s}_{i+1} \bar{c}_i b_{i+1} \hat{\beta}_{i+1} + \bar{c}_{i+1} a_{i+1} \check{\alpha}_{i+1}), \end{aligned} \right\} \quad i = 2, \dots, N-1.$$

Действительно, первое из этих равенств просто совпадает с первым из соотношений (2.5), если положить $\Phi_i = -s_2 a_1 \hat{a}_1 - c_2 b_2 \check{\beta}_2$. Используя теперь первое равенство (2.11) и второе равенство (2.5), а также (2.3), находим

$$\begin{aligned} \bar{s}_2 s_2 \frac{\sigma^2}{a_1} + \bar{c}_2 c_2 a_2 &= \bar{s}_2 (s_2 a_1 + c_2 b_2 - \bar{\Phi}_1) + \bar{c}_2 c_2 a_2 = \\ &= \bar{s}_2 s_2 a_1 + c_2 (\bar{s}_2 b_2 + \bar{c}_2 a_2) - \bar{s}_2 \bar{\Phi}_1 = \bar{s}_2 \frac{\bar{s}_2}{s_2} a_1 + c_2 \frac{\bar{s}_2}{s_2} a_1 - \\ &- \bar{s}_2 \bar{\Phi}_1 - c_2 (\bar{s}_2 b_2 \hat{\beta}_2 + \bar{c}_2 a_2 \check{\alpha}_2) = \frac{\bar{s}_2}{s_2} a_1 - \bar{s}_2 \bar{\Phi}_1 - c_2 (\bar{s}_2 b_2 \hat{\beta}_2 + \bar{c}_2 a_2 \check{\alpha}_2). \end{aligned}$$

Вводя обозначение $\Phi_2 = -\bar{s}_2 \bar{\Phi}_1 - c_2 (\bar{s}_2 b_2 \hat{\beta}_2 + \bar{c}_2 a_2 \check{\alpha}_2)$, будем иметь второе из равенств (2.11).

$$\bar{s}_2 s_2 \frac{\sigma^2}{a_1} + \bar{c}_2 c_2 a_2 = \frac{\bar{s}_2}{s_2} a_1 + \Phi_2.$$

Аналогичным образом получаются все остальные равенства (2.11), в которых использованы обозначения (2.12).

Переходя непосредственно к доказательству теоремы, введем обозначения

$$A_1 = A, \quad A'_i = A_{i-1} C_i^*, \quad A_i = \bar{C}_i A'_{i+1}, \quad i = 2, \dots, N,$$

из которых следует, что $\bar{C} A C^* = A_N C_{N+1}^*$. Матрица $A'_2 = A_1 C_2^*$, очевидно, имеет вид

$$A'_2 = \begin{bmatrix} c_2 a_1 - s_2 b_2 & s_2 a_1 + c_2 b_2 \\ -s_2 a_2 & c_2 a_2 & b_3 \\ & a_3 & b_4 \\ & & \ddots & \ddots \\ & & & \ddots & \ddots \\ & 0 & & & & 0 \\ & & & & & a_{N-1} & b_N \\ & & & & & & a_N \end{bmatrix}.$$

Используя соотношения (2.5) и (2.11), а также обозначения (2.10), преобразуем выражения для элементов матрицы A'_2 :

$$c_2 a_1 - s_2 b_2 = \frac{s_2}{s_2} \left[c_2 \frac{\bar{s}_2}{s_2} a_1 - \bar{s}_2 \bar{b}_2 (1 + \bar{\beta}_2) \right] + s_2 b_2 \bar{\beta}_2 = \bar{c}_2 \bar{a}_1 (1 + \bar{\alpha}_1) + s_2 b_2 \bar{\beta}_2,$$

$$s_2 a_1 + c_2 b_2 = s_2 \frac{\sigma^2}{a_1} + \bar{\Phi}_1.$$

Введем обозначения

$$b'_{11} = s_2 b_2 \bar{\beta}_2, \quad b''_{12} = \bar{\Phi}_1. \quad (2.13)$$

Матрица A'_2 принимает теперь вид

$$A'_2 = \begin{bmatrix} \bar{c}_2 \bar{a}_1 (1 + \bar{\alpha}_2) + b'_{11} & s_2 \frac{\sigma^2}{a_1} + b''_{12} & & \\ -s_2 a_2 & c_2 a_2 & b_3 & 0 \\ & & a_3 & b_4 \\ 0 & & & a_{N-1} & b_N \\ & & & & a_N \end{bmatrix}.$$

По индукции легко показать, что матрицы A'_i ($i = 3, \dots, N$) и A_i ($i = 2, \dots, N-1$) выглядят следующим образом:

$$A'_i = \begin{bmatrix} A'_{(i-1)} + X'_{(i)} + B'_{(i-1)} & b'_{1i} & 0 \\ b'_{i-1,1} \dots b'_{i-1,i-2} t_i + b'_{i-1,i-1} r_i + b''_{i-1,i} & b'_{2i} & \\ 0 & -s_i a_i & c_i a_i & \overline{A'_{(i+1)}} \end{bmatrix}; \quad (II)$$

$$A_i = \begin{bmatrix} \bar{A}_{(i-1)} + X_{(i)} + B_{(i-1)} & b'_{1i} & 0 \\ b'_{ii} \dots b'_{i,i-2} b'_{i,i-1} & b'_{i-2,i} & \\ 0 & x_i + b'_{i-1,i} & -s_i b_{i+1} & \overline{\tilde{A}_{(i+1)}} \\ & y_i + b''_{ii} & c_i b_{i+1} & \end{bmatrix}; \quad (III)$$

где приняты обозначения

$$\bar{A}'_{(i-1)} = \begin{bmatrix} \bar{a}_1 \bar{b}_2 & & 0 \\ \bar{a}_2 \bar{b}_3 & & \\ 0 & \ddots & \\ & \bar{a}_{i-2} \bar{b}_{i-1} & \end{bmatrix}, \quad \tilde{A}_{(i+1)} = \begin{bmatrix} a_{i+1} b_{i+2} & & 0 \\ \vdots & \ddots & \\ 0 & a_{N-1} b_N & \\ & & a_N \end{bmatrix};$$

$$X'_{(i)} = \begin{bmatrix} \bar{a}_1 \bar{\alpha}_2 & \bar{b}_2 \bar{\beta}_3 & & 0 \\ \bar{a}_2 \bar{\alpha}_3 & \bar{b}_3 \bar{\beta}_4 & & \\ 0 & & \ddots & \\ & & & \bar{a}_{i-2} \bar{\alpha}_{i-1} & \bar{b}_{i-1} \bar{\beta}_{i-1} \end{bmatrix};$$

$$B'_{(i-1)} = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1,i-2} & b_{1,i-1} \\ b_{21} & b_{22} & \dots & b_{2,i-2} & b_{2,i-1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ b_{i-2,1} & b_{i-2,2} & \dots & b_{i-2,i-2} & b_{i-2,i-1} \end{bmatrix};$$

$$\bar{A}_{(i-1)} = \begin{bmatrix} \bar{A}'_{(i-1)} \\ 0 \dots 0 \bar{a}_{i-1} \end{bmatrix}, \quad \bar{A}'_{(i+1)} = \begin{bmatrix} b_{i+1} 0 \dots 0 \\ \bar{A}_{(i+1)} \end{bmatrix};$$

$$X_{(i)} = \begin{bmatrix} X'_{(i)} \\ 0 \dots 0 \bar{a}_{i-1} \bar{\alpha}_i \end{bmatrix}, \quad B_{(i-1)} = \begin{bmatrix} B'_{(i-1)} \\ b_{i-1,1} b_{i-1,2} \dots b_{i-1,i-1} \end{bmatrix},$$

$$t_i = \bar{c}_i \bar{a}_{i-1} (1 + \check{\alpha}_i); \quad r_i = \frac{s_i s_{i-1} \dots s_2 \sigma^2}{s_{i-1} \dots s_2 a_1};$$

$$x_i = c_{i+1} \bar{b}_i (1 + \check{\beta}_{i+1}); \quad y_i = \frac{\bar{s}_i \dots \bar{s}_2}{s_i \dots s_2} a_1.$$

Продемонстрируем только переход от матрицы A'_i к матрице A_i . Вычисления по формуле $A_i = \bar{C}_i A'_i$ дают

$$A_i = \left[\begin{array}{c|c} \bar{A}'_{(i-1)} + X'_{(i)} + B'_{(i-1)} & \begin{matrix} b'_{1i} \\ b'_{2i} \\ \vdots \\ b'_{i-2,i} \end{matrix} \\ \hline \bar{c}_i b'_{i-1,1} \dots \bar{c}_i b'_{i-1,i-2} u_{i-1} & 0 \\ \bar{s}_i b'_{i-1,1} \dots \bar{s}_i b'_{i-1,i-2} z_i & \bar{c}_i b_{i+1} \\ \hline 0 & \bar{A}_{(i+1)} \end{array} \right]$$

где элементы u_{i-1} , z_i , v_i , w_i имеют следующие выражения:

$$\begin{aligned} u_{i-1} &= \bar{c}_i (t_i + b'_{i-1,i-1}) + \bar{s}_i s_i a_i; \\ z_i &= \bar{s}_i (t_i + b'_{i-1,i-1}) - \bar{c}_i s_i a_i; \\ v_i &= \bar{c}_i (r_i + b''_{i-1,i}) - \bar{s}_i c_i a_i; \\ w_i &= \bar{s}_i (r_i + b''_{i-1,i}) + \bar{c}_i c_i a_i. \end{aligned} \tag{2.14}$$

Преобразуем эти выражения. В силу определения \bar{a}_{i-1} и обозначения для t_i легко видеть, что $u_{i-1} = \bar{a}_{i-1} (1 + \check{\alpha}_i) + \bar{c}_i b'_{i-1,i-1} - \bar{s}_i s_i a_i \check{\alpha}_i$, $z_i = \bar{s}_i b'_{i-1,i-1} + \bar{c}_i s_i a_i \check{\alpha}_i$. Выражение для v_i с учетом обозначения для r_i преобразуем к виду

$$v_i = \frac{\bar{s}_i}{s_{i+1}} \left[\bar{c}_i \frac{s_{i+1} s_i \dots s_2 \sigma^2}{s_i \dots s_2 a_1} - s_{i+1} c_i a_i (1 + \hat{\alpha}_i) \right] + \bar{c}_i b''_{i-1,i} + \bar{s}_i c_i a_i \hat{\alpha}_i,$$

откуда, принимая во внимание (2.5), определение \bar{b}_{i-1} и обозначение для x_i , получим $v_i = x_i + \bar{c}_i b''_{i-1,i} + s_i c_i a_i \hat{\alpha}_i$. Наконец, из (2.11) и обозначения для y_i следует, что $w_i = y_i + s_i b'_{i-1,i} + \varphi_i$. После введения обозначений

$$\begin{aligned} b'_{i-1,i-1} &= \bar{c}_i b'_{i-1,i-1} - \bar{s}_i s_i a_i \check{\alpha}_i, \quad b'_{i,i-1} = \bar{s}_i b'_{i-1,i-1} + \bar{c}_i s_i \check{\alpha}_i; \\ b'_{i-1,i} &= \bar{c}_i b''_{i-1,i} + \bar{s}_i c_i a_i \hat{\alpha}_i, \quad b''_{ii} = \bar{s}_i b''_{i-1,i} + \varphi_i; \\ b'_{i-1,j} &= \bar{c}_i b'_{i-1,j}, \\ b'_{ij} &= \bar{s}_i b'_{i-1,j}, \end{aligned} \quad \left. \begin{array}{l} j = 1, \dots, i-2, \\ \end{array} \right\} \tag{2.15}$$

становится понятным, что матрица A_i действительно принимает вид (III).

Аналогичным образом выполняется переход от матрицы A_i к матрице A'_{i+1} по формуле $A'_{i+1} = A_i C^*_{i+1}$. Нетрудно проверить, что A'_{i+1} записывается в виде (II), если только значение индекса i заменено на $i+1$. При этом должны быть сделаны следующие обозначения:

$$\begin{aligned} b_{i-1,i} &= c_{i+1} b'_{i-1,i} - s_{i+1} \bar{s}_i b_{i+1} \bar{b}_{i+1}; \\ b'_{i-1,i+1} &= s_{i+1} b'_{i-1,i} + c_{i+1} \bar{s}_i b_{i+1} \bar{b}_{i+1}; \\ b''_{ii} &= c_{i+1} b''_{ii} + s_{i+1} \bar{c}_i b_{i+1} \bar{\beta}_{i+1}, \quad b''_{i,i+1} = s_{i+1} b''_{ii} + \bar{\Phi}_i; \\ b_{ji} &= c_{i+1} b'_{ji}, \\ b_{j,i+1} &= s_{i+1} b'_{ji}, \end{aligned} \quad j = 1, \dots, i-2. \quad (2.16)$$

Рассмотрим еще переход от матрицы A'_N к матрице A_N по формуле $A_N = \bar{C}_N A'_N$. Непосредственно проверяется, что A_N имеет вид

$$A_N = \left[\begin{array}{c|c} \bar{A}'_{(N-1)} + X'_{(N)} + B'_{(N-1)} & b'_{1N} \\ \hline c_N b'_{N-1,1} \dots c_N b'_{N-1,N-2} u_{N-1} & v_N \\ s_N b'_{N-1,1} \dots s_N b'_{N-1,N-2} z_N & w_N \end{array} \right],$$

где для элементов u_{N-1} , z_N , v_N и w_N справедливы формулы (2.14) при $i=N$. Из определения \bar{a}_{N-1} легко заключить, что $u_{N-1} = a_{N-1}(1 + \alpha_N) + c_N b'_{N-1,N-1} - s_N s_N a_N \alpha_N$, $z_N = s_N b'_{N-1,N-1} + c_N s_N a_N \alpha_N$. Выражение для v_N с учетом последнего из соотношений (2.5) преобразуется в $v_N = c_N b''_{N-1,N} + s_N c_N a_N \hat{\alpha}_N$. Использование того же соотношения позволяет получить выражение для w_N :

$$w_N = \frac{s_N \dots s_2 \sigma^2}{s_N \dots s_2 a_1} + s_N b''_{N-1,N} - c_N c_N a_N \hat{\alpha}_N.$$

С другой стороны, из последнего из соотношений (2.11) непосредственно вытекает, что

$$w_N = \frac{s_N \dots s_2}{s_N \dots s_2} a_1 + s_N b''_{N-1,N} + \Phi_N.$$

Вводя обозначение

$$p = \frac{s_N \dots s_2}{s_N \dots s_2} a_1 \quad (2.17)$$

и сравнивая два различных выражения для w_N , получим уравнение для определения величины p :

$$(\sigma^2 - p^2)/p = \Phi_N + c_N c_N a_N \hat{\alpha}_N, \quad (2.18)$$

которое исследуем далее.

Воспользуемся обозначениями (2.15) при значении индекса $i=N$ и введем новые обозначения

$$\begin{aligned} b_{Nj} &= b'_{Nj}, \\ b_{jN} &= b'_{jN}, \\ b_{NN} &= b''_{NN}. \end{aligned} \quad j = 1, \dots, N-1, \quad (2.19)$$

Выражения для элементов u_{N-1} , z_N , v_N и w_N теперь можно переписать следующим образом: $u_{N-1} = a_{N-1}(1 + \alpha_N) + b_{N-1,N-1}$, $z_N = b_{N,N-1}$, $v_N =$

$= b_{N-1, N}$, $w_N = p + b_{NN}$. Легко понять, что матрица A_N принимает вид

$$A_N = \left[\begin{array}{c|c} \bar{A}_{(N-1)} + X_{(N)} + B_{(N-1)} & \begin{matrix} b_{1N} \\ \vdots \\ b_{N-1, N} \end{matrix} \\ \hline b_{N1} \dots b_{N, N-1} & p + b_{NN} \end{array} \right].$$

Отсюда следует, что окончательную матрицу $\bar{C}AC^* = A_N C_{N+1}^*$ можно представить в виде суммы $\bar{C}AC^* = \bar{A} + X + BC_{N+1}^* + Y$, где \bar{A} имеет требуемый вид (2.8)

$$X = \left[\begin{array}{c|c} X_{(N)} & \begin{matrix} 0 \\ \vdots \\ 0 \end{matrix} \\ \hline 0 \dots 0 & 0 \end{array} \right]; \quad B = \left[\begin{array}{c|c} B_{(N-1)} & \begin{matrix} b_{1N} \\ \vdots \\ b_{N-1, N} \end{matrix} \\ \hline b_{N1} \dots b_{N, N-1} & b_{NN} \end{array} \right],$$

а единственный ненулевой элемент матрицы Y расположен в ее нижнем правом углу и равен $p \operatorname{sign} a_1 - \sigma$. Обозначив $R = X + BC_{N+1}^* + Y$, оценим норму матрицы R . В силу неравенства

$$\|R\| \leq \|X\| + \|B\| + \|Y\| \quad (2.20)$$

достаточно найти оценки норм матриц X , B и Y .

Сначала оценим норму матрицы B , предварительно напомнив, что ее элементы b_{ij} ($i, j = 1, \dots, N$) вычисляются в соответствии с рекуррентными формулами (2.13), (2.15) при $i = 2, \dots, N$, (2.16) при $i = 2, \dots, N-1$ и (2.19). Введем в рассмотрение две вспомогательные двухдиагональные матрицы

$$D_1 = \begin{bmatrix} -a_1 \hat{\alpha}_1 - d_2 \check{\beta}_2 & & & & 0 \\ -a_2 \hat{\alpha}_2 - b_3 \check{\beta}_3 & \ddots & & & \\ & \ddots & \ddots & & \\ 0 & & -a_{N-1} \hat{\alpha}_{N-1} - b_N \check{\beta}_N & & \\ & & & -a_N \hat{\alpha}_N & \end{bmatrix},$$

$$D_2 = \begin{bmatrix} 0 & -b_2 \hat{\beta}_2 & & & 0 \\ -a_2 \check{\alpha}_2 - b_3 \hat{\beta}_3 & \ddots & & & \\ & \ddots & \ddots & & \\ 0 & & -a_{N-1} \check{\alpha}_{N-1} - b_N \hat{\beta}_N & & \\ & & & -a_N \hat{\alpha}_N & \end{bmatrix}.$$

и положим

$$E = \bar{C}_N \cdots \bar{C}_3 \bar{C}_2 D_1 C_2 C_3 \cdots C_N, \quad F = \bar{C}_N \cdots \bar{C}_3 \bar{C}_2 D_2 C_2 C_3 \cdots C_N, \quad (2.21)$$

Если для элементов матриц E и F выписать рекуррентные формулы аналогично тому, как это сделано для элементов матрицы B , и сравнить получающиеся формулы с (2.13), (2.15), (2.16) и (2.19), то нетрудно убедиться, что имеют место равенства

$$b_{ij} = \begin{cases} e_{ij}, & \text{если } j > i, \\ f_{ij}, & \text{если } j \leq i. \end{cases} \quad (2.22)$$

Здесь e_{ij} и f_{ij} — элементы матриц E и F соответственно. Разобьем теперь матрицу B на сумму двух треугольных матриц B_1 и B_2 , отнеся в верхнюю треугольную B_1 , имеющую нулевую главную диагональ, элементы

матрицы B , расположенные выше ее главной диагонали, а в нижнюю треугольную B_2 — остальные элементы матрицы B . Ясно, что $\|B\| \leq \|B_1\| + \|B_2\|$. В силу определения матриц B_1 и B_2 и равенств (2.22) для фробениусовых норм этих матриц имеют место оценки $\|B_1\|_F \leq \|E\|_F$, $\|B_2\|_F \leq \|F\|_F$. Отсюда, учитывая определение (2.21) матриц E и F и ортогональность матриц C_i , \bar{C}_i , а также связь между фробениусовой и спектральной нормами, находим $\|B_1\| \leq \|B_1\|_F \leq \|E\|_F \leq \sqrt{N}\|E\| = \sqrt{N}\|D_1\|$ и аналогично $\|B_2\| \leq \sqrt{N}\|D_2\|$. Получим еще оценки норм матриц D_1 и D_2 через норму исходной матрицы A . Для нормы двухдиагональной матрицы D_1 имеет место оценка (см., например, (2.30) из [4]):

$$\|D_1\| \leq \max \begin{cases} \max_{1 \leq i \leq N-1} (|a_i \hat{\alpha}_i| + |b_{i+1} \check{\beta}_{i+1}|), \\ \max_{2 \leq i \leq N} (|a_i \hat{\alpha}_i| + |b_i \check{\beta}_i|), \end{cases}$$

из которой в силу неравенств (2.6) следует

$$\|D_1\| \leq \varepsilon \max \begin{cases} \max_{1 \leq i \leq N-1} (|a_i| + |b_{i+1}|), \\ \max_{2 \leq i \leq N} (|a_i| + |b_i|). \end{cases} \quad (2.23)$$

Далее, известно (см., например, [5], § 3), что для нормы двухдиагональной матрицы A справедливы оценки снизу $\sqrt{a_i^2 + b_i^2} \leq \|A\|$, $\sqrt{a_i^2 + b_{i+1}^2} \leq \|A\|$, из которых очевидным образом следуют неравенства

$$|a_i| + |b_i| \leq \sqrt{2}\|A\|, \quad |a_i| + |b_{i+1}| \leq \sqrt{2}\|A\|. \quad (2.24)$$

Огрубляя с помощью этих неравенств оценку (2.23), получим $\|D_1\| \leq \sqrt{2}\varepsilon\|A\|$. Аналогичным образом запишем оценку нормы матрицы D_2 : $\|D_2\| \leq \sqrt{2}\varepsilon\|A\|$. В результате окончательная оценка нормы матрицы B принимает вид $\|B\| \leq 2\sqrt{2}\sqrt{N}\varepsilon\|A\|$.

Для оценки нормы матрицы X воспользуемся упомянутой оценкой (2.30) из [4]:

$$\|X\| \leq \max \left\{ \begin{array}{l} \max_{2 \leq i \leq N-1} (|\bar{a}_{i-1} \check{\alpha}_i| + |\bar{b}_i \check{\beta}_{i+1}|), \\ \max_{2 \leq i \leq N-1} (|\bar{a}_i \check{\alpha}_{i+1}| + |\bar{b}_i \check{\beta}_{i+1}|) \end{array} \right\} \leq \sqrt{2}\varepsilon\|A\|.$$

Для оценки нормы матрицы Y потребуется некоторое представление элемента p , который, как уже отмечалось, удовлетворяет квадратному уравнению (2.18). Прежде чем приступить к исследованию решений этого уравнения, получим удобное выражение и оценку для его правой части $\varphi_N + c_N c_N a_N \hat{\alpha}_N$. Оказывается, и это нетрудно проверить непосредственными вычислениями, что величина $\varphi_N + c_N c_N a_N \hat{\alpha}_N$ совпадает с элементом матрицы $G = \bar{C}_N \cdots \bar{C}_3 \bar{C}_2 (-D_1 + D_2) C_2 C_3 \cdots C_N$, стоящим в ее нижнем правом углу. Отсюда сразу следует оценка $|\varphi_N + c_N c_N a_N \hat{\alpha}_N| \leq \|G\| = \| -D_1 + D_2 \| \leq 2\sqrt{2}\varepsilon\|A\|$, очевидно, эквивалентная равенству $\varphi_N + c_N c_N a_N \hat{\alpha}_N = \gamma_1 \|A\|$, которое выполнено при некотором γ_1 таком, что $|\gamma_1| \leq 2\sqrt{2}\varepsilon$. Заметим, что оценку (2.2) также можно записать в виде равенства $\sigma = (1 + \gamma_2)\|A\|$, справедливого при некотором γ_2 , удовлетворяющем оценке $|\gamma_2| \leq \delta$. Уравнение (2.18) теперь перепишем в виде $(\sigma^2 - p^2)/p = \sigma\gamma_1/(1 + \gamma_2)$ или $(\sigma^2 - p^2)/p = \gamma_3\sigma$, где $|\gamma_3| \leq 2\sqrt{2}\varepsilon/(1 - \delta)$.

Решая это уравнение, получим $p_{1,2} = \pm\sigma[\sqrt{1 + \gamma_3^2/4} \mp \gamma_3/2]$, откуда следует, что $p_1 = \sigma(1 + \gamma_4)$, $p_2 = -\sigma(1 + \gamma_5)$, где $|\gamma_4| \leq |\gamma_3| \leq 2\sqrt{2}\varepsilon/(1 - \delta)$, $|\gamma_5| \leq 2\sqrt{2}\varepsilon/(1 - \delta)$. Из определения (2.17) p и положительности чисел s_i , \bar{s}_i вытекает, что знак корня следует выбирать равным знаку элемента

a_1 , вследствие чего

$$p = \begin{cases} \sigma(1 + \gamma_4), & \text{если } a_1 > 0, \\ -\sigma(1 + \gamma_5), & \text{если } a_1 < 0. \end{cases}$$

Последнее равенство, очевидно, можно записать в виде $p = \sigma \operatorname{sign} a_1 (1 + \gamma_6)$, где $|\gamma_6| \leq 2\sqrt{2}\varepsilon/(1-\delta)$. Из этого равенства следует, что $p \operatorname{sign} a_1 - \sigma = \sigma \gamma_6$. Помня, что $p \operatorname{sign} a_1 - \sigma$ есть единственный ненулевой элемент матрицы Y , легко получить оценку ее нормы

$$\|Y\| = |p \operatorname{sign} a_1 - \sigma| = \sigma |\gamma_6| \leq \frac{2\sqrt{2}\varepsilon(1+\delta)}{1-\delta} \|A\|.$$

Подставляя полученные оценки норм матриц B , X и Y в оценку (2.20), находим

$$\|R\| \leq 2\sqrt{2}\sqrt{N}\varepsilon\|A\| + \sqrt{2}\varepsilon\|\bar{A}\| + \frac{2\sqrt{2}\varepsilon(1+\delta)}{1-\delta} \|A\|, \quad (2.25)$$

С помощью этого неравенства и равенства $R = \bar{C}AC^* - \bar{A}$, следующего из (2.7), можно найти оценку нормы матрицы \bar{A} в терминах нормы исходной матрицы A . Действительно, $\|R\| \geq \|\bar{A}\| - \|\bar{C}AC^*\| = \|\bar{A}\| - \|A\|$, следовательно,

$$\|\bar{A}\| \leq \frac{1 + 2\sqrt{2}\sqrt{N}\varepsilon + [2\sqrt{2}\varepsilon(1+\delta)]/(1-\delta)}{1 - \sqrt{2}\varepsilon} \|A\|.$$

Подстановка полученной оценки в (2.25) дает

$$\|R\| \leq \frac{\sqrt{2}[(2\sqrt{N}+1)(1-\delta) + 2(1+\delta)]}{(1-\delta)(1-\sqrt{2}\varepsilon)} \varepsilon \|A\|.$$

Предполагая величины ε и δ настолько малыми, что выполнено неравенство $5\delta + 2\sqrt{2}(\sqrt{N}+2)(1-\delta)\varepsilon \leq 1$, полученную оценку огрубим следующим образом:

$$\|R\| \leq 2\sqrt{2}(\sqrt{N}+2)\varepsilon\|A\|.$$

Теорема доказана.

Глава 2

РЕШЕНИЕ УРАВНЕНИЙ ДЛЯ ПАРАМЕТРОВ, ЗАДАЮЩИХ ОРТОГОНАЛЬНЫЕ ПРЕОБРАЗОВАНИЯ ВРАЩЕНИЯ, И ИХ МАШИННОЕ ВЫЧИСЛЕНИЕ

§ 3. Решение уравнений для параметров, определяющих преобразования вращения

В этом параграфе показано, как найти параметры c_i, s_i , удовлетворяющие соотношениям (1.3) — (1.5), а также параметры $c_i, s_i, \bar{c}_i, \bar{s}_i$, удовлетворяющие соотношениям (2.3) — (2.5), т. е. рассмотрены случаи трехдиагональной симметрической и двухдиагональной матриц. Получены оценки возмущений α_i, β_i , с которыми удовлетворяются уравнения (1.5), и возмущений $\hat{\alpha}_i, \hat{\alpha}_i, \hat{\beta}_i, \hat{\beta}_i$, с которыми выполняются уравнения (2.5).

В случае трехдиагональной матрицы A (1.1) будем предполагать, что известны достаточно хорошее приближение η к собственному значению λ матрицы A и отношения $\bar{\mathcal{P}}_2, \bar{\mathcal{P}}_3, \dots, \bar{\mathcal{P}}_m$ последовательных компонент почти собственного вектора $\bar{u} = (\bar{u}_1, \bar{u}_2, \dots, \bar{u}_m)^T$ этой матрицы, отвечающего ее почти собственному значению η . Конкретнее, предполо-

жим, что имеют место неравенство

$$|\eta - \lambda| \leq \gamma \mathcal{M}(A), \quad \gamma \ll 1, \quad (3.1)$$

и соотношения

$$\begin{aligned} d_i + \bar{\alpha}_i - \eta + b_2 \bar{\beta}_2 &= 0; \\ \frac{b_i(1 + \bar{\beta}_i)}{\bar{\beta}_i} + d_i + \bar{\alpha}_i - \eta + b_{i+1} \bar{\beta}_{i+1} &= 0, \quad i = 2, \dots, m-1; \\ \frac{b_m(1 + \bar{\beta}_m)}{\bar{\beta}_m} + d_m + \bar{\alpha}_m - \eta &= 0, \end{aligned} \quad (3.2)$$

где возмущения $\bar{\alpha}_i$ и $\bar{\beta}_i$ удовлетворяют оценкам

$$\begin{aligned} |\bar{\alpha}_i| &\leq \bar{\alpha} \mathcal{M}(A), \quad \bar{\alpha} \ll 1; \\ |\bar{\beta}_i| &\leq \bar{\beta}, \quad \bar{\beta} \ll 1; \end{aligned} \quad (3.3)$$

а $\mathcal{M}(A)$ обозначает одну из норм матрицы A :

$$\mathcal{M}(A) = \max \begin{cases} |d_1| + |b_2|, \\ \max_{2 \leq i \leq m-1} (|d_i| + |b_i| + |b_{i+1}|), \\ |d_m| + |b_m|. \end{cases}$$

Вычисление η — стандартная задача, оно подробно описано, например, в [5]. Возможность же получения чисел $\bar{\beta}_2, \bar{\beta}_3, \dots, \bar{\beta}_m$, удовлетворяющих соотношениям (3.2), является в точности утверждением теоремы 2, доказанной в [4]. Отметим для дальнейшего, что из [4] следует, что числа $\bar{\beta}_i$ удовлетворяют неравенствам

$$\frac{\varepsilon_2}{\alpha^2} < \frac{\varepsilon_1^2}{2} \leq |\bar{\beta}_i| \leq \frac{2}{\varepsilon_1^2} < \frac{1}{\varepsilon_2}, \quad (3.4)$$

гарантирующим отличие их от нуля, а также возможность размещения в ячейке ЭВМ ЕС, не вызывая переполнения разрядной сетки. Здесь α — основание системы счисления, принятой в ЭВМ, а ε_1 и ε_2 — характеристики ее разрядной сетки (об их определении см., например, [4], § 4). Определение этих характеристик напомним в следующем параграфе, где они существенно используются при анализе выполнения арифметических операций на вычислительной машине. А пока приведем их значения для ЭВМ ЕС в случае использования чисел с «двойной точностью»: $\varepsilon_1 \approx 0,22 \cdot 10^{-15}$, $1/\varepsilon_2 \approx 0,72 \cdot 10^{16}$. Основание α системы счисления, принятой в этих машинах, равно 16. Напомним также одно обозначение, введенное в § 3 работы [4]. Если a — арифметическое выражение, то под $(a)_m$ понимается результат вычисления значения этого выражения в вычислительной машине.

Введем теперь в рассмотрение последовательности чисел s_i, c_i, c'_i , определив их следующим образом:

$$\begin{aligned} s_2 &= \frac{1}{\sqrt{1 + \bar{\beta}_2^2}}, \quad c_2 = \frac{\bar{\beta}_2}{\sqrt{1 + \bar{\beta}_2^2}}, \quad c'_2 = (c_2)_m; \\ s_i &= \frac{1}{\sqrt{1 + (c'_{i-1} \bar{\beta}_i)^2}}, \quad c_i = \frac{c'_{i-1} \bar{\beta}_i}{\sqrt{1 + (c'_{i-1} \bar{\beta}_i)^2}}, \quad c'_i = (c_i)_m, \quad i = 3, \dots, m. \end{aligned} \quad (3.5)$$

Предположим, что c'_i можно вычислить так, что

$$c'_i = c_i (1 + \kappa_i), \quad |\kappa_i| \leq \kappa \ll 1. \quad (3.6)$$

Это предположение обосновано в § 4, где приведен алгоритм вычисле-

ния c_i и анализируются погрешности, возникающие при его реализации на вычислительной машине.

Из (3.5) непосредственно следуют равенства

$$c_2 = s_2 \bar{\mathcal{P}}_2, \quad c_i = s_i c_{i-1} \bar{\mathcal{P}}_i, \quad i = 3, \dots, m, \quad (3.7)$$

которые с учетом (3.6) легко приводят к равенствам

$$\begin{aligned} \frac{c'_2}{\bar{\mathcal{P}}_2} &= s_2(1 + \kappa_2); \quad \frac{c'_i}{\bar{\mathcal{P}}_i} = s_i c_{i-1} (1 + \kappa_{i-1})(1 + \kappa_i), \quad i = 3, \dots, m-1; \\ \frac{c'_m}{\bar{\mathcal{P}}_m} &= s_m c_{m-1} (1 + \kappa_{m-1}). \end{aligned} \quad (3.8)$$

Умножая первое из соотношений (3.2) на s_2 , i -е — на $s_{i+1} c'_i$, последнее — на s_m и учитывая равенства (3.6) — (3.8), получим

$$\begin{aligned} s_2(d_1 + \bar{\alpha}_1 - \eta) + c_2 b_2 &= 0; \\ s_3 s_2 b_2 (1 + \bar{\beta}_2) (1 + \kappa_2) + s_3 c_2 (d_2 + \bar{\alpha}_2 - \eta) (1 + \kappa_2) + c_3 b_3 &= 0; \\ s_{i+1} s_i c_{i-1} b_i (1 + \bar{\beta}_i) (1 + \kappa_{i-1}) (1 + \kappa_i) + s_{i+1} c_i (d_i + \bar{\alpha}_i - \eta) (1 + \kappa_i) + \\ + c_{i+1} b_{i+1} &= 0, \quad i = 3, \dots, m-1; \\ s_m c_{m-1} b_m (1 + \bar{\beta}_m) (1 + \kappa_{m-1}) + c_m (d_m + \bar{\alpha}_m - \eta) &= 0. \end{aligned}$$

Вводя обозначения

$$\begin{aligned} \alpha_1 &= \bar{\alpha}_1; \quad \beta_2 = (1 + \bar{\beta}_2) (1 + \kappa_2) - 1; \\ \alpha_i &= (d_i - \eta) \kappa_i + \bar{\alpha}_i (1 + \kappa_i), \quad i = 2, \dots, m-1; \\ \beta_i &= (1 + \bar{\beta}_i) (1 + \kappa_{i-1}) (1 + \kappa_i) - 1, \quad i = 3, \dots, m-1; \\ \alpha_m &= \bar{\alpha}_m; \quad \beta_m = (1 + \bar{\beta}_m) (1 + \kappa_{m-1}) - 1, \end{aligned} \quad (3.9)$$

полученные соотношения можно записать в виде

$$\begin{aligned} s_2(d_1 + \alpha_1 - \eta) + c_2 b_2 &= 0; \\ s_3 s_2 b_2 (1 + \beta_2) + s_3 c_2 (d_2 + \alpha_2 - \eta) + c_3 b_3 &= 0; \\ s_{i+1} s_i c_{i-1} b_i (1 + \beta_i) + s_{i+1} c_i (d_i + \alpha_i - \eta) + c_{i+1} b_{i+1} &= 0, \quad i = 3, \dots, m-1; \\ s_m c_{m-1} b_m (1 + \beta_m) + c_m (d_m + \alpha_m - \eta) &= 0. \end{aligned}$$

Они и являются основным результатом данного параграфа и, очевидно, полностью совпадают с соотношениями (1.5). Что касается соотношений (1.3) и неравенств (1.4), то их выполнение непосредственно вытекает из определения (3.5).

Обратимся теперь к случаю двухдиагональной матрицы A (2.1). Предположим, что нам известны некоторое достаточно хорошее приближение σ к наибольшему сингулярному числу $\sigma_N(A) = \|A\|$ матрицы A , а также отношения $\mathcal{P}_1, \mathcal{R}_2, \mathcal{P}_2, \dots, \mathcal{R}_N, \mathcal{P}_N$ последовательных компонент почти собственного вектора $w = (u_1, v_1, u_2, v_2, \dots, u_N, v_N)^T$ трехдиагональной симметрической матрицы

$$S = \begin{bmatrix} 0 & a_1 & & & & & & & \\ a_1 & 0 & b_2 & & & & & & 0 \\ & b_2 & 0 & a_2 & & & & & \\ & & a_2 & 0 & b_3 & & & & \\ & & & \ddots & \ddots & \ddots & & & \\ & & & & \ddots & \ddots & \ddots & & \\ 0 & & & & & a_{N-1} & 0 & b_N & \\ & & & & & b_N & 0 & a_N & \\ & & & & & & a_N & 0 & \end{bmatrix},$$

отвечающего ее почти собственному значению σ . Более точно, предполагается, что имеют место неравенство

$$|\sigma - \|A\|| \leq \delta \|A\|, \quad \delta \ll 1, \quad (3.10)$$

и соотношения

$$\left. \begin{aligned} -\sigma + a_1(1 + \xi_1)\mathcal{P}_1 &= 0; \\ \frac{a_{i-1}(1 + \xi_{i-1})}{\mathcal{P}_{i-1}} - \sigma + b_i(1 + \xi_i)\mathcal{R}_i &= 0, \\ \frac{b_i(1 + \xi_i)}{\mathcal{R}_i} - \sigma + a_i(1 + \xi_i)\mathcal{P}_i &= 0, \\ \frac{a_N(1 + \xi_N)}{\mathcal{P}_N} - \sigma &= 0, \end{aligned} \right\} i = 2, \dots, N; \quad (3.11)$$

где возмущения $\xi_1, \xi_2, \xi_3, \xi_4, \xi_5$ удовлетворяют оценкам

$$|\xi_1| \leq \rho, \quad |\xi_2| \leq \rho, \quad |\xi_3| \leq \rho, \quad |\xi_4| \leq \rho, \quad \rho \ll 1. \quad (3.12)$$

Заметим, что основанием для соотношений (3.11) служит теорема 1 из [4], которую достаточно просто переформулировать на данный случай. Подробно получение этих соотношений описано в § 8 главы 4, где приводится общая схема предлагаемого алгоритма исчерпывания. Особо подчеркнем, что под σ понимается приближение именно к наибольшему сингулярному числу, в противном случае нельзя гарантировать малости возмущений $\xi_1, \xi_2, \xi_3, \xi_4, \xi_5$.

Умножая второе из соотношений (3.11) на \mathcal{P}_1 , третье — на \mathcal{R}_2 и вообще $(2i-1)$ -е — на \mathcal{R}_i , $2i$ -е — на \mathcal{P}_i , преобразуем их к виду

$$\left. \begin{aligned} -\sigma + a_1(1 + \xi_1)\mathcal{P}_1 &= 0; \\ a_{i-1}(1 + \xi_{i-1}) - \sigma \mathcal{P}_{i-1} + b_i(1 + \xi_i)\mathcal{P}_{i-1}\mathcal{R}_i &= 0, \\ b_i(1 + \xi_i) - \sigma \mathcal{R}_i + a_i(1 + \xi_i)\mathcal{R}_i \mathcal{P}_i &= 0, \\ a_N(1 + \xi_N) - \sigma \mathcal{P}_N &= 0. \end{aligned} \right\} i = 2, \dots, N; \quad (3.13)$$

Введем в рассмотрение последовательности чисел $s_i, c_i, c'_i, \bar{s}_i, \bar{c}_i, \bar{c}'_i$, определив их следующим образом:

$$\left. \begin{aligned} s_2 &= \frac{1}{\sqrt{1 + (\mathcal{P}_1 \mathcal{R}_2)^2}}, \quad c_2 = \frac{\mathcal{P}_1 \mathcal{R}_2}{\sqrt{1 + (\mathcal{P}_1 \mathcal{R}_2)^2}}, \quad c'_2 = (c_2)_M; \\ \bar{s}_2 &= \frac{1}{\sqrt{1 + (\mathcal{R}_2 \mathcal{P}_2)^2}}, \quad \bar{c}_2 = \frac{\mathcal{R}_2 \mathcal{P}_2}{\sqrt{1 + (\mathcal{R}_2 \mathcal{P}_2)^2}}, \quad \bar{c}'_2 = (\bar{c}_2)_M; \\ s_i &= \frac{1}{\sqrt{1 + (c'_{i-1} \mathcal{P}_{i-1} \mathcal{R}_i)^2}}, \quad c_i = \frac{c'_{i-1} \mathcal{P}_{i-1} \mathcal{R}_i}{\sqrt{1 + (c'_{i-1} \mathcal{P}_{i-1} \mathcal{R}_i)^2}}, \quad c'_i = (c_i)_M, \\ \bar{s}_i &= \frac{1}{\sqrt{1 + (\bar{c}'_{i-1} \mathcal{R}_i \mathcal{P}_i)^2}}, \quad \bar{c}_i = \frac{\bar{c}'_{i-1} \mathcal{R}_i \mathcal{P}_i}{\sqrt{1 + (\bar{c}'_{i-1} \mathcal{R}_i \mathcal{P}_i)^2}}, \quad \bar{c}'_i = (\bar{c}_i)_M, \end{aligned} \right\} i = 3, \dots, N. \quad (3.14)$$

Предположим, что величины c'_i и \bar{c}'_i можно вычислить так, что

$$c'_i = c_i(1 + \tau_i), \quad \bar{c}'_i = \bar{c}_i(1 + \bar{\tau}_i), \quad (3.15)$$

где $|\tau_i| \leq \tau, |\bar{\tau}_i| \leq \bar{\tau}, \tau \ll 1$. Из определения (3.14) ясно, что $s_i > 0, \bar{s}_i > 0$. А из того, что числа $\mathcal{P}_i, \mathcal{R}_i$ отличны от нуля, и из предположений (3.15) легко следует, что и все числа $c_i, c'_i, \bar{c}_i, \bar{c}'_i$ отличны от нуля.

Оставляя первое и последнее из соотношений (3.13) без изменений, а второе, третье и вообще $2i$ -е и $(2i+1)$ -е умножая соответственно на $s_2, \bar{s}_2, s_{i+1}c_i$ и $\bar{s}_{i+1}c_i$, преобразуем их к виду

$$\begin{aligned} -\sigma + a_1(1 + \xi_1)\mathcal{P}_1 &= 0; \\ s_2a_1(1 + \hat{\xi}_1) - s_2\sigma\mathcal{P}_1 + c_2b_2(1 + \xi_2) &= 0; \\ -\bar{s}_2b_2(1 + \hat{\xi}_2) - \bar{s}_2\sigma\mathcal{R}_2 + \bar{c}_2a_2(1 + \xi_2) &= 0; \\ s_{i+1}c_i a_i(1 + \xi_i) - s_{i+1}c_i \sigma \mathcal{P}_i + c_{i+1}b_{i+1} \frac{1 + \xi_{i+1}}{1 + \tau_i} &= 0, \\ s_{i+1}\bar{c}_i b_{i+1}(1 + \hat{\xi}_{i+1}) - \bar{s}_{i+1}\bar{c}_i \sigma \mathcal{R}_{i+1} + \bar{c}_{i+1}a_{i+1} \frac{1 + \xi_{i+1}}{1 + \tau_i} &= 0, \\ a_N(1 + \hat{\xi}_N) - \sigma\mathcal{P}_N &= 0. \end{aligned} \quad (3.16)$$

При получении этих соотношений учтены равенства

$$\begin{aligned} c_2 &= s_2\mathcal{P}_1\mathcal{R}_2; \quad \bar{c}_2 = \bar{s}_2\mathcal{R}_2\mathcal{P}_2; \\ c_{i+1} &= s_{i+1}c_i' \mathcal{P}_i\mathcal{R}_{i+1}; \quad \bar{c}_{i+1} = \bar{s}_{i+1}\bar{c}_i' \mathcal{R}_{i+1}\mathcal{P}_{i+1}, \quad i = 2, \dots, N-1, \end{aligned} \quad (3.17)$$

непосредственно следующие из (3.14), а также предположения (3.15).

Покажем теперь, как, пользуясь первым из полученных соотношений, исключить $\mathcal{P}_1, \mathcal{R}_2, \mathcal{P}_2, \dots, \mathcal{R}_N, \mathcal{P}_N$ из остальных соотношений (3.16). Из формул (3.17) легко следуют рекуррентные соотношения для \mathcal{P}_i и \mathcal{R}_i :

$$\begin{aligned} \mathcal{P}_2 &= \frac{\bar{c}_2 s_2}{c_2 \bar{s}_2} \mathcal{P}_1, \quad \mathcal{P}_i = \frac{\bar{c}_i c_{i-1}' s_i}{c_i' c_{i-1} s_i} \mathcal{P}_{i-1}, \quad i = 3, \dots, N; \\ \mathcal{R}_2 &= \frac{c_2}{s_2} \frac{1}{\mathcal{P}_1}, \quad \mathcal{R}_2 = \frac{c_3}{c_2} \frac{1}{c_2} \frac{\bar{s}_2}{s_3} \mathcal{R}_2, \quad \mathcal{R}_i = \frac{c_i}{c_{i-1}} \frac{\bar{c}_{i-2} \bar{s}_{i-1}}{c_{i-1} c_{i-2}} \mathcal{R}_{i-1}, \quad i = 4, \dots, N, \end{aligned}$$

последовательное применение которых дает

$$\begin{aligned} \mathcal{P}_i &= \frac{\bar{c}_i c_{i-1}' \bar{c}_{i-1} c_{i-2}' \bar{c}_{i-2} \bar{c}_{i-2}}{c_i' c_{i-1} c_{i-1} c_{i-2} c_{i-2}'} \dots \frac{\bar{c}_2' \bar{c}_2 s_i \dots s_2}{c_2' c_2 s_i \dots s_2} \mathcal{P}_1, \quad i = 3, \dots, N; \\ \mathcal{R}_i &= \frac{c_i}{c_{i-1}} \frac{c_{i-1}' \bar{c}_{i-2} c_{i-2}'}{c_{i-1}' c_{i-1} c_{i-2} c_{i-2}'} \dots \frac{c_3 c_2'}{c_3' c_2} \frac{c_2 \bar{s}_{i-1} \dots \bar{s}_2}{c_2' c_2 s_{i-1} \dots s_2} \frac{1}{\mathcal{P}_1}, \quad i = 3, \dots, N. \end{aligned}$$

Из первого соотношения (3.16) следует формула $\mathcal{P}_1 = \sigma/[a_1(1 + \xi_1)]$. Подставляя ее в выражения для \mathcal{P}_i и \mathcal{R}_i и учитывая (3.15), получим

$$\begin{aligned} \mathcal{R}_2 &= c_2 \frac{1}{s_2} \frac{a_1(1 + \xi_1)}{\sigma}; \\ \mathcal{P}_2 &= \frac{\bar{c}_2 s_2}{c_2 \bar{s}_2} \frac{\sigma}{a_1(1 + \xi_1)}; \\ \mathcal{R}_3 &= \frac{c_3}{c_2} \frac{1}{1 + \tau_2} \frac{\bar{s}_2}{s_3 s_2} \frac{a_1(1 + \xi_1)}{\sigma}; \\ \mathcal{P}_i &= \frac{c_i}{c_i} \frac{(1 + \tau_{i-1}) \dots (1 + \tau_2)}{(1 + \tau_{i-1}) \dots (1 + \tau_2)} \frac{s_i \dots s_2}{s_i \dots s_2} \frac{\sigma}{a_1(1 + \xi_1)}, \quad i = 3, \dots, N; \\ \mathcal{R}_i &= \frac{c_i}{c_{i-1}} \frac{(1 + \bar{\tau}_{i-2}) \dots (1 + \bar{\tau}_2)}{(1 + \bar{\tau}_{i-1})(1 + \bar{\tau}_{i-2}) \dots (1 + \bar{\tau}_2)} \frac{\bar{s}_{i-1} \dots \bar{s}_2}{s_i s_{i-1} \dots s_2} \frac{a_1(1 + \xi_1)}{\sigma}, \quad i = 4, \dots, N. \end{aligned} \quad (3.18)$$

Введем обозначения

$$\begin{aligned}
 \hat{\alpha}_1 &= (1 + \hat{\zeta}_1)(1 + \check{\zeta}_1) - 1; \\
 \check{\beta}_2 &= (1 + \check{\xi}_2)(1 + \check{\zeta}_1) - 1, \quad \hat{\beta}_2 = \frac{1 + \hat{\xi}_2}{1 + \check{\zeta}_1} - 1; \\
 \check{\alpha}_2 &= \frac{1 + \check{\xi}_2}{1 + \check{\zeta}_1} - 1, \quad \hat{\alpha}_2 = (1 + \hat{\zeta}_2)(1 + \check{\zeta}_1) - 1; \\
 \check{\beta}_3 &= \frac{(1 + \check{\xi}_3)(1 + \check{\zeta}_1)}{1 + \tau_2} - 1, \quad \hat{\beta}_3 = \frac{(1 + \hat{\xi}_3)(1 + \tau_2)}{1 + \check{\zeta}_1} - 1; \\
 \check{\alpha}_i &= \frac{(1 + \check{\xi}_i)(1 + \tau_{i-1}) \dots (1 + \tau_2)}{(1 + \check{\zeta}_1)(1 + \bar{\tau}_{i-1}) \dots (1 + \bar{\tau}_2)} - 1, \\
 \hat{\alpha}_i &= \frac{(1 + \hat{\xi}_i)(1 + \check{\zeta}_1)(1 + \bar{\tau}_{i-1}) \dots (1 + \bar{\tau}_2)}{(1 + \tau_{i-1}) \dots (1 + \tau_2)} - 1, \\
 \check{\beta}_i &= \frac{(1 + \check{\xi}_i)(1 + \check{\zeta}_1)(1 + \bar{\tau}_{i-2}) \dots (1 + \bar{\tau}_2)}{(1 + \tau_{i-1})(1 + \tau_{i-2}) \dots (1 + \tau_2)} - 1, \\
 \hat{\beta}_i &= \frac{(1 + \hat{\xi}_i)(1 + \tau_{i-1})(1 + \tau_{i-2}) \dots (1 + \tau_2)}{(1 + \check{\zeta}_1)(1 + \bar{\tau}_{i-2}) \dots (1 + \bar{\tau}_2)} - 1,
 \end{aligned} \tag{3.19}
 \left. \begin{array}{l} i = 3, \dots, N; \\ i = 4, \dots, N. \end{array} \right\}$$

С учетом этих обозначений и формул (3.18) второе, третье и следующие соотношения (3.16) перепишем в виде

$$\begin{aligned}
 s_2 a_1 (1 + \hat{\alpha}_1) - s_2 \frac{\sigma^2}{a_1} + c_2 b_2 (1 + \check{\beta}_2) &= 0; \\
 \bar{s}_2 b_2 (1 + \hat{\beta}_2) - c_2 \frac{\bar{s}_2}{s_2} a_1 + \bar{c}_2 a_2 (1 + \check{\alpha}_2) &= 0; \\
 s_{i+1} c_i a_i (1 + \hat{\alpha}_i) - \bar{c}_i \frac{s_{i+1} s_i \dots s_2}{\bar{s}_i \dots \bar{s}_2} \frac{\sigma^2}{a_1} + c_{i+1} b_{i+1} (1 + \check{\beta}_{i+1}) &= 0, \\
 \bar{s}_{i+1} \bar{c}_i b_{i+1} (1 + \hat{\beta}_{i+1}) - \bar{c}_{i+1} \frac{\bar{s}_{i+1} \dots \bar{s}_2}{\bar{s}_{i+1} \dots \bar{s}_2} a_1 + \bar{c}_{i+1} a_{i+1} (1 + \check{\alpha}_{i+1}) &= 0, \\
 c_N a_N (1 + \hat{\alpha}_N) - \bar{c}_N \frac{s_N \dots s_2}{\bar{s}_N \dots \bar{s}_2} \frac{\sigma^2}{a_1} &= 0.
 \end{aligned} \tag{3.14}
 \left. \begin{array}{l} i = 2, \dots, N-1; \\ \end{array} \right.$$

Они, очевидно, полностью совпадают с соотношениями (2.5). Выполнение же соотношений (2.3) непосредственно следует из определения (3.14) чисел c_i , s_i , \bar{c}_i , \bar{s}_i .

Для завершения рассмотрений, связанных с соотношениями (2.5), дадим оценки возмущений $\check{\alpha}_i$, $\hat{\alpha}_i$, $\check{\beta}_i$ и $\hat{\beta}_i$. Для примера оценим величину α_i . Очевидно, что

$$\check{\alpha}_i \leq \frac{(1 + \rho)(1 + \tau)^{i-2}}{(1 - \rho)(1 - \tau)^{i-2}} - 1 \leq \frac{(1 + \rho)(1 + \tau)^{N-2}}{(1 - \rho)(1 - \tau)^{N-2}} - 1.$$

В предположениях $N^2\tau \leq 1$, $2(2N\tau + \rho)\rho \leq \tau$, которые при всех разумных порядках N совершенно необременительны, легко получить оценку $\alpha_i \leq 2(N\tau + \rho)$. Подобным же образом оценивается α_i снизу: $\alpha_i \geq -2(N\tau + \rho)$. Следовательно, $|\alpha_i| \leq 2(N\tau + \rho)$. Аналогичные оценки имеют место и для остальных возмущений. Таким образом, обозначив

$$\varepsilon = 2(N\tau + \rho), \tag{3.20}$$

можно записать $|\hat{\alpha}_i| \leq \varepsilon$, $|\check{\alpha}_i| \leq \varepsilon$, $|\hat{\beta}_i| \leq \varepsilon$, $|\check{\beta}_i| \leq \varepsilon$.

Как уже отмечалось, значение ρ оценивает величину возмущений $\hat{\zeta}_i, \tilde{\zeta}_i, \hat{\xi}_i, \tilde{\xi}_i$, с которыми удовлетворены уравнения (3.11). Получение этих уравнений и оценка величины ρ — самостоятельная задача, решению которой посвящена работа [4]. Величина же τ , оценивающая погрешность машинного вычисления чисел c_i и \bar{c}_i , оценивается аналогично тому, как в следующем параграфе определена погрешность вычисления чисел c_i в случае трехдиагональной матрицы.

§ 4. Машинное вычисление параметров, задающих преобразования вращения

Основной целью данного параграфа является обоснование предположений (3.6), обеспечивающих высокую относительную точность вычисления величин c_i [см. формулы (3.5)], а также получение аналогичного результата относительно вычисления величин s_i . Подобные результаты сформулированы и для параметров $c_i, s_i, \bar{c}_i, \bar{s}_i$, используемых в алгоритме исчерпывания двухдиагональной матрицы [см. формулы (3.14)].

Предварительно напомним определение параметров ε_1 и ε_2 , которыми будем характеризовать разрядную сетку вычислительной машины [см. [4], § 4]. При этом будем рассматривать машины, использующие представление чисел с «плавающей точкой». Каждое число $x \neq 0$ в памяти таких машин хранится в виде $x = \alpha^t s$, где α — основание системы счисления, принятой в машине, t — целочисленный порядок, $s (1/\alpha \leq |s| < 1)$ — мантисса числа x . Число разрядов, отводимых для размещения мантиссы, ограничивает относительную точность представления чисел в машине. В качестве характеристики этой точности используют число ε_1 , определяемое как наименьшее положительное машинное число, для которого $(1 + \varepsilon_1)_m > 1$. Число разрядов, отводимых для размещения порядка, задает величину наибольшего и наименьшего по модулю чисел, которые могут быть представлены в машине. Обозначим наибольшее по модулю машинное число через $1/\varepsilon_2$. Тогда наименьшим по модулю отличным от нуля машинным числом будет

$$\frac{\varepsilon_2}{\alpha^2} \left(1 - \frac{\varepsilon_1}{\alpha}\right) < \frac{\varepsilon_2}{\alpha^2}.$$

Погрешность выполнения арифметических операций определяется следующим неравенством:

$$|(a * b)_m - (a * b)| \leq \varepsilon_1 |a * b| + \varepsilon_2 / \alpha^2, \quad (4.1)$$

где под знаком «*» подразумевается один из знаков арифметических операций сложения, вычитания, умножения или деления. Первое слагаемое $\varepsilon_1 |a * b|$ в правой части (4.1) определяет погрешность, возникающую из-за ограниченности числа разрядов, отводимых для размещения мантиссы результата операции. Второе слагаемое описывает погрешность, возникающую в том случае, когда модуль результата операции меньше, чем наименьшее по модулю отличное от нуля машинное число. В этом случае $(a * b)_m$ принудительно полагается равным нулю, а возникающая из-за этого погрешность, очевидно, не превосходит ε_2 / α^2 . В дальнейшем будем называть ее погрешностью исчезновения порядка.

Перейдем непосредственно к машинному вычислению параметров, определяющих преобразования вращения. Представим числа c'_{i-1} и \bar{P}_i ($c'_{i-1} \neq 0, \bar{P}_i \neq 0$) в виде

$$\begin{aligned} c'_{i-1} &= a^{k_{i-1}} q_{i-1}, \quad \frac{1}{\alpha} \leq |q_{i-1}| < 1; \\ \bar{P}_i &= \alpha^{m_i} p_i, \quad \frac{1}{\alpha} \leq |p_i| < 1 \end{aligned} \quad (4.2)$$

и сформулируем алгоритм вычисления c_i по формуле

$$c_i = \frac{c'_{i-1} \bar{\varphi}_i}{\sqrt{1 + (c'_{i-1} \bar{\varphi}_i)^2}}$$

следующим образом:

- 1) $x = (q_{i-1} p_i)_m$;
- 2) $j = k_{i-1} + m$;
- 3) $y = (x x)_m$;
- 4) представим y в виде

$$y = \alpha^t z, \quad 1/\alpha \leq |z| < 1; \quad (4.3)$$

- 5) $t = t + 2j$;
- 6) $g = \max(1, l)$;
- 7) если g нечетно, положим $g = g + 1$;
- 8) $u = (\alpha^{1-g}/\alpha)_m, \quad v = (\alpha^{l-g} z)_m$;
- 9) $w = (u + v)_m$;
- 10) $\mu = (\sqrt{w})_m$;
- 11) $v = (x/\mu)_m$;
- 12) представим v в виде $v = \alpha^h q_i, \quad 1/\alpha \leq |q_i| < 1$;
- 13) $k_i = h + j - g/2$.

Оценим погрешности, возникающие при реализации описанного алгоритма, предварительно заметив, что п. 2, 4—7, 12 и 13 выполняются на машине абсолютно точно. Кроме того, использование арифметики вынесенных порядков позволяет при выполнении п. 4, 3, 9—11 избежать погрешности исчезновения порядка, в результате чего оценки погрешностей, возникающих при выполнении этих пунктов, упрощаются. Учитывая высказанные замечания, получим

$$\begin{aligned} |x - q_{i-1} p_i| &\leq \varepsilon_1 |q_{i-1} p_i|, \\ |y - (q_{i-1} p_i)^2| &\leq |y - x^2| + |x^2 - (q_{i-1} p_i)^2| \leq \\ &\leq \varepsilon_1 x^2 + |x - q_{i-1} p_i| |x + q_{i-1} p_i| \leq [(1 + \varepsilon_1)^3 - 1] (q_{i-1} p_i)^2 < \\ &< \varepsilon_1 (3 + 4\varepsilon_1) (q_{i-1} p_i)^2. \end{aligned} \quad (4.4)$$

Остановимся на выполнении п. 8 и 9. При вычислении одной из величин u или v возможно появление погрешности исчезновения порядка. Эта погрешность возникает всякий раз, когда порядки 1 и l отличаются на достаточно большую величину (например, для машин серии ЕС ЭВМ эта величина равна 64, если число $\max(1, l)$ четно, и 63 — в противном случае). В силу определения g ясно, что хотя бы одна из величин u или v вычисляется точно. Предположим для определенности, что точно вычисляется u , а при вычислении v происходит исчезновение порядка, т. е. $\alpha^{l-g} z < \varepsilon_2/\alpha^2$. В этом случае $v = 0$, $u = \alpha^{1-g}/\alpha$, операция сложения u и v осуществляется, очевидно, точно и, следовательно,

$$\begin{aligned} \left| w - \left(\alpha^{1-g} \frac{1}{\alpha} + \alpha^{l-g} z \right) \right| &= \alpha^{l-g} z < \frac{\varepsilon_2}{\alpha^2} = \frac{\varepsilon_2}{\alpha^{2-g}} \alpha^{1-g} \frac{1}{\alpha} < \\ &< \frac{\varepsilon_2}{\alpha^{2-g}} \left(\alpha^{1-g} \frac{1}{\alpha} + \alpha^{l-g} z \right) < \varepsilon_1 \left(\alpha^{1-g} \frac{1}{\alpha} + \alpha^{l-g} z \right). \end{aligned}$$

При выводе последнего неравенства использовано то обстоятельство, что в данном случае $\max(1, l) = 1$, вследствие чего $2 - g = 0$, и неравенство $\varepsilon_2 \ll \varepsilon_1$. Случай, когда v вычисляется точно, а исчезновение порядка происходит при вычислении u , анализируется аналогично. Наконец, если вычисление u и v осуществляется без исчезновения порядка, т. е. точно, погрешность выполнения п. 9 оценивается стандартным образом:

$$\left| w - \left(\alpha^{1-g} \frac{1}{\alpha} + \alpha^{l-g} z \right) \right| = |w - (u + v)| \leq \varepsilon_1 (u + v) = \varepsilon_1 \left(\alpha^{1-g} \frac{1}{\alpha} + \alpha^{l-g} z \right).$$

Итак, во всех рассматриваемых случаях имеет место оценка

$$\left| w - \left(\alpha^{1-g} \frac{1}{\alpha} + \alpha^{l-g} z \right) \right| \leq \varepsilon_1 \left(\alpha^{1-g} \frac{1}{\alpha} + \alpha^{l-g} z \right),$$

которая очевидным образом преобразуется к виду

$$|\alpha^g w - (1 + \alpha^l z)| \leq \varepsilon_1 (1 + \alpha^l z).$$

Учитывая определение l и z и оценки погрешностей выполнения предыдущих пунктов, теперь нетрудно получить оценку

$$|\alpha^g w - [1 + \alpha^{2j} (q_{i-1} p_i)^2]| \leq \varepsilon_1 (4 + 8\varepsilon_1) [1 + \alpha^{2j} (q_{i-1} p_i)^2]. \quad (4.5)$$

Перейдем к оценкам погрешностей, допускаемых на следующих пунктах алгоритма:

$$|\mu - Vw| \leq \varepsilon_1 Vw;$$

$$\left| v - \frac{x}{Vw} \right| \leq \left| v - \frac{x}{\mu} \right| + \left| \frac{x}{\mu} - \frac{x}{Vw} \right| \leq \varepsilon_1 \left| \frac{x}{\mu} \right| + \frac{|x| |\mu - Vw|}{\mu Vw} \leq \frac{2\varepsilon_1}{1 - \varepsilon_1} \frac{|x|}{Vw}.$$

Последнее неравенство после умножения обеих его частей на $\alpha^{-g/2}$ принимает вид

$$\left| \alpha^{-g/2} v - \frac{x}{V\alpha^g w} \right| \leq \frac{2\varepsilon_1}{1 - \varepsilon_1} \frac{|x|}{V\alpha^g w}. \quad (4.6)$$

Наконец, на основе неравенств (4.4) — (4.6) с помощью несложных, но довольно громоздких выкладок можно получить оценку

$$\left| \alpha^{-g/2} v - \frac{q_{i-1} p_i}{\sqrt{1 + \alpha^{2j} (q_{i-1} p_i)^2}} \right| \leq \varepsilon_1 (5 + 35\varepsilon_1) \frac{q_{i-1} p_i}{\sqrt{1 + \alpha^{2j} (q_{i-1} p_i)^2}}.$$

Умножая обе части этой оценки на α^l и учитывая формулы (4.2), а также формулы п. 2, 12 и 13 алгоритма (4.3), ее легко преобразовать к виду

$$\left| \alpha^{k_i} q_i - \frac{c'_{i-1} \bar{\varphi}_i}{\sqrt{1 + (c'_{i-1} \bar{\varphi}_i)^2}} \right| \leq \varepsilon_1 (5 + 35\varepsilon_1) \frac{c'_{i-1} \bar{\varphi}_i}{\sqrt{1 + (c'_{i-1} \bar{\varphi}_i)^2}}.$$

Обозначив теперь величину $\alpha^{k_i} q_i$ через c'_i , получим

$$|c'_i - c_i| \leq \varepsilon_1 (5 + 35\varepsilon_1) |c_i|. \quad (4.7)$$

Особо подчеркнем, что при определении числа $c'_i \equiv (c_i)_m$ по формуле $c'_i = \alpha^{k_i} q_i$ операция умножения на машине не выполняется, т. е. подразумевается, что c_i хранится в памяти машины в виде двух чисел — порядка k_i и мантиссы q_i .

Пользуясь малостью ε_1 , оценку (4.7) можно огрубить, например, следующим образом: $|c'_i - c_i| \leq 5,001\varepsilon_1 |c_i|$. Очевидно, что эта оценка эквивалентна равенству $c'_i = c_i (1 + \varepsilon_i)$, которое имеет место при некотором ε_i таком, что $|\varepsilon_i| \leq 5,001\varepsilon_1$. Следовательно, в качестве значения ε_i , оценивающего сверху величину погрешности ε_i [см. формулы (3.6)], можно взять $\varepsilon = 5,001\varepsilon_1$.

Для вычисления s_i по формуле [см. (3.5)]

$$s_i = \frac{1}{\sqrt{1 + (c'_{i-1} \bar{\varphi}_i)^2}}$$

можно сформулировать алгоритм, аналогичный (4.3). Этот алгоритм совпадает с (4.3) по всем пунктам, за исключением п. 11—13. Эти пункты формулируются теперь следующим образом:

11') $v = (1/\mu)_m$;

12') представим v в виде $v = \alpha^l r_i$, $1/\alpha \leq |r_i| < 1$;

13') $n_i = h - g/2$.

В результате применения этого алгоритма находятся числа n_i и r_i , которые и определяют величину $s'_i = (s_i)_M : s'_i = \alpha^{n_i} r_i$. Анализ погрешностей, проводимый аналогично предыдущему, позволяет утверждать, что $s'_i = s_i(1 + \omega_i)$, где $|\omega_i| \leq \omega$, а в качестве значения ω можно взять $\omega = 4,001\epsilon_1$.

Обратимся теперь к вычислению чисел c_i , s_i , \bar{c}_i , \bar{s}_i , используемых в алгоритме исчерпывания двухдиагональной матрицы [см. формулы (3.14)]. Представив числа c'_{i-1} , \mathcal{P}_{i-1} и \mathcal{R}_i в виде

$$c'_{i-1} = \alpha^{h_{i-1}} q_{i-1}, \quad \frac{1}{\alpha} \leq |q_{i-1}| < 1;$$

$$\mathcal{P}_{i-1} = \alpha^{m_{i-1}} p_{i-1}, \quad \frac{1}{\alpha} \leq |p_{i-1}| < 1;$$

$$\mathcal{R}_i = \alpha^{n_i} r_i, \quad \frac{1}{\alpha} \leq |r_i| < 1;$$

алгоритм вычисления c_i по формуле

$$c_i = \frac{c'_{i-1} \mathcal{P}_{i-1} \mathcal{R}_i}{\sqrt{1 + (c'_{i-1} \mathcal{P}_{i-1} \mathcal{R}_i)^2}}$$

сформулируем следующим образом:

- 1) $x = (q_{i-1} p_{i-1} r_i)_M$;
- 2) $j = k_{i-1} + m_{i-1} + n_i$;
- 3) $y = (xx)_M$;
- 4) представим y в виде $y = \alpha^t z$, $1/\alpha \leq |z| < 1$;
- 5) $l = t + 2j$;
- 6) $g = \max(1, l)$;
- 7) если g нечетно, положим $g = g + 1$;
- 8) $u = \left(\alpha^{1-g} \frac{1}{\alpha}\right)_M$, $v = (\alpha^{l-g} z)_M$;
- 9) $w = (u + v)_M$;
- 10) $\mu = (\sqrt{w})_M$;
- 11) $v = (x/\mu)_M$;
- 12) представим v в виде $v = \alpha^h q_i$, $1/\alpha \leq |q_i| < 1$;
- 13) $k_i = h + j - g/2$.

Легко показать, что величина $c'_i = \alpha^{h_i} q_i$, получаемая в результате работы этого алгоритма, связана с c_i соотношением $c'_i = c_i(1 + \tau_i)$, где $|\tau_i| \leq \tau$, а в качестве τ можно взять $\tau = 7,001\epsilon_1$. Совершенно аналогичный алгоритм можно предложить для вычислений величины

$$\bar{c}_i = \frac{\bar{c}'_{i-1} \mathcal{R}_i \mathcal{P}_i}{\sqrt{1 + (\bar{c}'_{i-1} \mathcal{R}_i \mathcal{P}_i)^2}}.$$

При этом вычисленная величина $\bar{c}'_i = (\bar{c}_i)_M$ оказывается связанный с величиной \bar{c}_i соотношением $\bar{c}'_i = \bar{c}_i(1 + \bar{\tau}_i)$, где $|\bar{\tau}_i| \leq \tau$.

Отметим, наконец, что и для вычисления величин s_i и \bar{s}_i можно сформулировать соответствующие алгоритмы, результатами выполнения которых являются числа s'_i и \bar{s}'_i , удовлетворяющие равенствам $s'_i = s_i(1 + \pi_i)$, $\bar{s}'_i = \bar{s}_i(1 + \bar{\pi}_i)$, где $|\pi_i| \leq \pi$, $|\bar{\pi}_i| \leq \pi$, а в качестве π можно взять $\pi = 5,001\epsilon_1$.

Глава 3

МАШИННОЕ ВЫЧИСЛЕНИЕ ПРЕОБРАЗОВАННЫХ МАТРИЦ

§ 5. Анализ погрешностей, возникающих при вычислении элементов преобразованной трехдиагональной матрицы

В этом параграфе описано вычисление элементов \bar{d}_i , \bar{b}_i матрицы (1.7) \bar{A} , получаемой в результате применения алгоритма исчерпывания к трехдиагональной симметрической матрице A , и проанализированы возникающие при этом погрешности. Напомним формулы, по которым определяются эти элементы [см. (1.9)],

$$\begin{aligned}\bar{d}_1 &= d_2 - \frac{c_2 b_2}{s_2} + \frac{c_3 c_2 b_3}{s_3}; \\ \bar{d}_i &= d_{i+1} - \frac{c_{i+1} c_i b_{i+1}}{s_{i+1}} + \frac{c_{i+2} c_{i+1} b_{i+2}}{s_{i+2}}, \quad i = 2, \dots, m-2; \\ \bar{d}_{m-1} &= d_m - \frac{c_m c_{m-1} b_m}{s_m}; \\ \bar{b}_i &= \frac{s_i b_{i+1}}{s_{i+1}}, \quad i = 2, \dots, m-1.\end{aligned}$$

Как уже отмечалось, погрешности, возникающие при вычислении по этим формулам, имеют двоякую природу. Прежде всего, они определяются тем, что вместо точных значений параметров s_i и c_i , фигурирующих в формулах, в памяти машины имеются лишь некоторые их приближения s'_i и c'_i , так что реальные вычисления могут использовать только эти приближенные значения. В таком случае ясно, что вместо точных элементов \bar{d}_i и \bar{b}_i можно вычислить только элементы \bar{d}'_i и \bar{b}'_i , определяемые по формулам

$$\begin{aligned}\bar{d}'_i &= d_{i+1} - \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} + \frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}}; \\ \bar{b}'_i &= \frac{s'_i b_{i+1}}{s'_{i+1}}.\end{aligned}\tag{5.1}$$

Второй источник погрешностей связан с непосредственной реализацией на вычислительной машине операций, предписываемых формулами (5.1).

Приступая к анализу погрешностей, напомним, что в § 4 разработан алгоритм вычисления величин s'_i и c'_i , гарантирующий выполнение равенств

$$s'_i = s_i (1 + \omega_i), \quad c'_i = c_i (1 + \kappa_i).\tag{5.2}$$

Эти равенства имеют место при некоторых ω_i и κ_i таких, что $|\omega_i| \leq \omega = 4,001\epsilon_1$, $|\kappa_i| \leq \kappa = 5,001\epsilon_1$. Здесь ϵ_1 , а также используемая ниже величина ϵ_2 — характеристики разрядной сетки используемой ЭВМ (см. § 4). Важно отметить, что упомянутый алгоритм позволяет определить каждую из величин s'_i и c'_i с помощью двух чисел: порядка и мантиссы, так что

$$s'_i = m(s'_i) \alpha^{p(s'_i)}, \quad c'_i = m(c'_i) \alpha^{p(c'_i)},\tag{5.3}$$

где $p(s'_i)$ и $p(c'_i)$ — порядки чисел s'_i и c'_i соответственно, а $m(s'_i)$ и $m(c'_i)$ — их мантиссы, удовлетворяющие неравенствам

$$\frac{1}{\alpha} \leq |m(s'_i)| < 1, \quad \frac{1}{\alpha} \leq |m(c'_i)| < 1.\tag{5.4}$$

Под α здесь понимается основание системы счисления, принятой в вычислительной машине. Отметим также, что в силу (5.3) и (5.4), $s'_i \neq 0$, $c'_i \neq 0$.

Получим оценки величин $\bar{d}'_i - \bar{d}_i$ и $\bar{b}'_i - b_i$. В силу равенств (5.2) ясно, что для этих величин имеют место формулы

$$\begin{aligned}\bar{d}'_i - \bar{d}_i &= -\frac{c_{i+1}c_i b_{i+1}}{s_{i+1}} \left[\frac{(1+\kappa_{i+1})(1+\kappa_i)}{1+\omega_{i+1}} - 1 \right] + \\ &+ \frac{c_{i+2}c_{i+1}b_{i+2}}{s_{i+2}} \left[\frac{(1+\kappa_{i+2})(1+\kappa_{i+1})}{1+\omega_{i+2}} - 1 \right]; \\ \bar{b}'_i - \bar{b}_i &= \frac{s_i b_{i+1}}{s_{i+1}} \left(\frac{1+\omega_i}{1+\omega_{i+1}} - 1 \right),\end{aligned}$$

из которых легко следуют оценки

$$\begin{aligned}|\bar{d}'_i - \bar{d}_i| &\leq \left| \frac{(1+\kappa)^2}{1-\omega} - 1 \right| \left(\left| \frac{c_{i+1}c_i b_{i+1}}{s_{i+1}} \right| + \left| \frac{c_{i+2}c_{i+1}b_{i+2}}{s_{i+2}} \right| \right); \\ |\bar{b}'_i - \bar{b}_i| &\leq \frac{2\omega}{1-\omega} \left| \frac{s_i b_{i+1}}{s_{i+1}} \right|.\end{aligned}$$

Вводя обозначения

$$\hat{\kappa} = \frac{(1+\kappa)^2}{1-\omega} - 1, \quad \hat{\omega} = \frac{2\omega}{1-\omega} \quad (5.5)$$

и учитывая формулу для \bar{b}_i , перепишем полученные оценки в виде

$$\begin{aligned}|\bar{d}'_i - \bar{d}_i| &\leq \hat{\kappa} \left(\left| \frac{c_{i+1}c_i b_{i+1}}{s_{i+1}} \right| + \left| \frac{c_{i+2}c_{i+1}b_{i+2}}{s_{i+2}} \right| \right); \\ |\bar{b}'_i - \bar{b}_i| &\leq \hat{\omega} |\bar{b}_i|.\end{aligned} \quad (5.6)$$

Перейдем теперь к анализу погрешностей, связанных с вычислением величин \bar{d}'_i и \bar{b}'_i по формулам (5.1). Представив числа b_{i+1} и b_{i+2} как

$$b_{i+1} = m(b_{i+1})\alpha^{p(b_{i+1})}, \quad \frac{1}{\alpha} \leq |m(b_{i+1})| < 1;$$

$$b_{i+2} = m(b_{i+2})\alpha^{p(b_{i+2})}, \quad \frac{1}{\alpha} \leq |m(b_{i+2})| < 1$$

и обозначив через $(\bar{d}'_i)_M$ результат машинного вычисления величины \bar{d}'_i , сформулируем алгоритм вычисления:

- 1) $x = (m(c'_{i+1})m(c'_i)m(b_{i+1})/m(s'_{i+1}))_M;$
- 2) $j = p(c'_{i+1}) + p(c'_i) + p(b_{i+1}) - p(s'_{i+1});$
- 3) $u = (x\alpha^j)_M;$
- 4) $y = (m(c'_{i+2})m(c'_{i+1})m(b_{i+2})/m(s'_{i+2}))_M;$
- 5) $k = p(c'_{i+2}) + p(c'_{i+1}) + p(b_{i+2}) - p(s'_{i+2});$
- 6) $v = (y\alpha^k)_M;$
- 7) $w = (v - u)_M;$
- 8) $(\bar{d}'_i)_M = (d_{i+1} + w)_M.$

Оценим погрешности, возникающие при реализации этого алгоритма. Заметим, что вынесение порядков, используемое в алгоритме, позволяет при выполнении п. 1 и 4 избежать исчезновения порядка (см. § 4). В таком случае нетрудно проверить, что оценки погрешностей, допускаемых при выполнении этих пунктов, имеют вид

$$\begin{aligned}\left| x - \frac{m(c'_{i+1})m(c'_i)m(b_{i+1})}{m(s'_{i+1})} \right| &\leq [(1+\varepsilon_1)^3 - 1] \left| \frac{m(c'_{i+1})m(c'_i)m(b_{i+1})}{m(s'_{i+1})} \right|; \\ \left| y - \frac{m(c'_{i+2})m(c'_{i+1})m(b_{i+2})}{m(s'_{i+2})} \right| &\leq [(1+\varepsilon_1)^3 - 1] \left| \frac{m(c'_{i+2})m(c'_{i+1})m(b_{i+2})}{m(s'_{i+2})} \right|.\end{aligned}$$

Умножая обе части первого из полученных неравенств на α^j , а обе части второго — на α^k , перепишем их в виде

$$\begin{aligned} \left| \alpha^j x - \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} \right| &\leq [(1 + \varepsilon_1)^j - 1] \left| \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} \right|; \\ \left| \alpha^k y - \frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}} \right| &\leq [(1 + \varepsilon_1)^k - 1] \left| \frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}} \right|. \end{aligned} \quad (5.7)$$

Остановимся на п. 3 и 6. При выполнении каждого из них возможно появление погрешности исчезновения порядка. Например, при выполнении п. 3 она возникает в том случае, если $|x\alpha^j| < \varepsilon_2/\alpha^2$. При этом величина u полагается равной нулю и, следовательно, имеет место оценка

$$|u - x\alpha^j| < \varepsilon_2/\alpha^2. \quad (5.8)$$

Если при выполнении п. 3 порядок не исчезает, то величина u вычисляется точно и оценка (5.8) остается справедливой. Аналогично оценивается погрешность выполнения п. 6:

$$|v - y\alpha^k| < \varepsilon_2/\alpha^2. \quad (5.9)$$

Погрешность выполнения п. 7 определяется стандартным образом. $|w - (v - u)| \leq \varepsilon_1 |v - u| + \varepsilon_2/\alpha^2$. Из этой оценки с помощью (5.8) и (5.9) легко получить $|w - (y\alpha^k - x\alpha^j)| \leq \varepsilon_1 |x\alpha^j| + \varepsilon_1 |y\alpha^k| + (3 + 2\varepsilon_1)\varepsilon_2/\alpha^2$. Отсюда, в свою очередь, в силу (5.7) следует

$$\begin{aligned} \left| w - \left(\frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}} - \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} \right) \right| &\leq [(1 + \varepsilon_1)^4 - 1] \left(\left| \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} \right| + \right. \\ &\quad \left. + \left| \frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}} \right| \right) + \frac{(3 + 2\varepsilon_1)\varepsilon_2}{\alpha^2}. \end{aligned} \quad (5.10)$$

Оценим, наконец, погрешность выполнения п. 8: $|(\bar{d}'_i)_m - (d_{i+1} + w)| \leq \varepsilon_1 |d_{i+1} + w| + \varepsilon_2/\alpha^2$. Из этой оценки с помощью (5.10) нетрудно получить

$$\begin{aligned} |(\bar{d}'_i)_m - \bar{d}'_i| &\leq \varepsilon_1 |d_{i+1}| + [(1 + \varepsilon_1)^5 - 1] \left(\left| \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} \right| + \left| \frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}} \right| \right) + \\ &\quad + \frac{[4 + \varepsilon_1(5 + 2\varepsilon_1)]\varepsilon_2}{\alpha^2}. \end{aligned}$$

Учитывая равенства (5.2) и обозначение (5.5) для $\hat{\kappa}$, полученную оценку преобразуем к виду

$$\begin{aligned} |(\bar{d}'_i)_m - \bar{d}'_i| &\leq \varepsilon_1 |d_{i+1}| + [(1 + \varepsilon_1)^5 - 1](1 + \hat{\kappa}) \left(\left| \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} \right| + \right. \\ &\quad \left. + \left| \frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}} \right| \right) + \frac{[4 + \varepsilon_1(5 + 2\varepsilon_1)]\varepsilon_2}{\alpha^2}. \end{aligned} \quad (5.11)$$

Теперь можно оценить отличие вычисленной величины $(\bar{d}'_i)_m$ от \bar{d}_i , а именно из (5.6) и (5.11) следует оценка

$$\begin{aligned} |(\bar{d}'_i)_m - \bar{d}_i| &\leq \varepsilon_1 |d_{i+1}| + \Delta \left(\left| \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} \right| + \left| \frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}} \right| \right) + \\ &\quad + \frac{[4 + \varepsilon_1(5 + 2\varepsilon_1)]\varepsilon_2}{\alpha^2}, \end{aligned} \quad (5.12)$$

где через Δ обозначена величина

$$\Delta = [(1 + \varepsilon_1)^5 - 1](1 + \hat{\kappa}) + \hat{\kappa}. \quad (5.13)$$

Покажем, что выражения $c_{i+1}c_i b_{i+1}/s_{i+1}$, $c_{i+2}c_{i+1}b_{i+2}/s_{i+2}$ ограничены. Точнее говоря, попытаемся оценить эти выражения через величину элементов исходной матрицы A . Для этого рассмотрим соотношение $s_{i+1}s_i c_{i-1}b_i(1 + \beta_i) + s_{i+1}c_i(d_i + \alpha_i - \eta) + c_{i+1}b_{i+1} = 0$, использованное в § 1. Умножая обе части соотношения на c_i/s_{i+1} и перенося первое и второе слагаемые в правую часть, получим

$$\frac{c_{i+1}c_i b_{i+1}}{s_{i+1}} = -c_i s_i c_{i-1} b_i (1 + \beta_i) - c_i^2 (d_i + \alpha_i - \eta).$$

Отсюда с учетом равенства $c_i^2 + s_i^2 = 1$ следует, что

$$\left| \frac{c_{i+1}c_i b_{i+1}}{s_{i+1}} \right| \leq \frac{1}{2} |b_i|(1 + \beta_i) + |d_i| + |\eta| + |\alpha_i|.$$

В силу неравенств (1.6) эту оценку перепишем в виде

$$\left| \frac{c_{i+1}c_i b_{i+1}}{s_{i+1}} \right| \leq \frac{1}{2} |b_i|(1 + \tilde{\beta}) + |d_i| + |\eta| + \tilde{\alpha} \mathcal{M}(A).$$

Аналогичным образом будем иметь

$$\left| \frac{c_{i+2}c_{i+1}b_{i+2}}{s_{i+2}} \right| \leq \frac{1}{2} |b_{i+1}|(1 + \tilde{\beta}) + |d_{i+1}| + |\eta| + \tilde{\alpha} \mathcal{M}(A).$$

С помощью полученных оценок из (5.12) нетрудно вывести

$$\begin{aligned} |(\bar{d}_i')_m - \bar{d}_i| &\leq \varepsilon_1 |d_{i+1}| + \frac{1}{2} \Delta (|b_i| + |d_i| + |b_{i+1}|)(1 + \tilde{\beta}) + \\ &+ \Delta \left(\frac{1}{2} |d_i| + |d_{i+1}| + 2|\eta| \right) + 2\Delta \tilde{\alpha} \mathcal{M}(A) + \frac{[4 + \varepsilon_1(5 + 2\varepsilon_1)]\varepsilon_2}{\alpha^2}. \end{aligned} \quad (5.14)$$

Из определения величины $\mathcal{M}(A)$ легко следует, что $|d_i| \leq \mathcal{M}(A)$, $|d_{i+1}| \leq \mathcal{M}(A)$, $|b_i| + |d_i| + |b_{i+1}| \leq \mathcal{M}(A)$. Кроме того, поскольку $\mathcal{M}(A)$ — одна из норм матрицы A , а модуль любого собственного значения матрицы не превосходит любой ее нормы, из неравенства (1.2) следует $|\eta| \leq (1 + \gamma)\mathcal{M}(A)$. Прежде чем огрубить оценку (5.14) с помощью полученных неравенств, оценим последнее слагаемое в ее правой части. Предполагая, что матрица A нормирована, т. е. $\mathcal{M}(A) = O(1)$, и для характеристик вычислительной машины ε_1 и ε_2 выполнено соотношение $\varepsilon_2 \ll \varepsilon_1$, оценим это слагаемое следующим образом: $[4 + \varepsilon_1(5 + 2\varepsilon_1)]\varepsilon_2/\alpha^2 \leq \varepsilon_1 \mathcal{M}(A)$. Из (5.14) теперь следует оценка $|(\bar{d}_i')_m - \bar{d}_i| \leq 2(2\Delta + \varepsilon_1)(1 + \tilde{\beta})(1 + \tilde{\alpha} + \gamma)\mathcal{M}(A)$. Вводя обозначение

$$\tilde{\Delta} = 2(2\Delta + \varepsilon_1)(1 + \tilde{\beta})(1 + \tilde{\alpha} + \gamma), \quad (5.15)$$

эту оценку окончательно запишем в виде

$$|(\bar{d}_i')_m - \bar{d}_i| \leq \tilde{\Delta} \mathcal{M}(A). \quad (5.16)$$

Обозначим теперь через $(\bar{b}_i')_m$ результат машинного вычисления величины \bar{b}_i' по формуле (5.1) и сформулируем алгоритм этого вычисления следующим образом:

- 1) $x = (m(s'_i)m(b_{i+1})/m(s'_{i+1}))_m$;
- 2) $j = p(s'_i) + p(b_{i+1}) - p(s'_{i+1})$;
- 3) $(\bar{b}_i')_m = (x\alpha^j)_m$.

Учтем погрешности, возникающие при реализации этого алгоритма на вычислительной машине. Поскольку вынесение порядков, используемое в алгоритме, позволяет при выполнении п. 1 избежать исчезновения порядка, погрешность выполнения этого пункта оценивается совсем просто:

$$\left| x - \frac{m(s'_i)m(b_{i+1})}{m(s'_{i+1})} \right| \leq \varepsilon_1(2 + \varepsilon_1) \left| \frac{m(s'_i)m(b_{i+1})}{m(s'_{i+1})} \right|.$$

Принимая во внимание формулу (5.1) для \bar{b}'_i , легко проверить, что полученная оценка после умножения обеих ее частей на α^j принимает вид

$$|x\alpha^j - \bar{b}'_i| \leq \varepsilon_1(2 + \varepsilon_1)|\bar{b}'_i|. \quad (5.17)$$

Погрешность выполнения п. 3 оценивается точно так же, как в алгоритме вычисления величины \bar{d}'_i . В результате получим оценку

$$|(\bar{b}'_i)_m - x\alpha^j| < \frac{\varepsilon_2}{\alpha^2}. \quad (5.18)$$

Из (5.17) и (5.18) легко следует $|(\bar{b}'_i)_m - \bar{b}'_i| \leq \varepsilon_1(2 + \varepsilon_1)|\bar{b}'_i| + \varepsilon_2/\alpha^2$. Из этой оценки, в свою очередь, с помощью (5.6) нетрудно получить

$$|(\bar{b}'_i)_m - \bar{b}'_i| \leq \varepsilon_1(2 + \varepsilon_1)(1 + \hat{\omega})|\bar{b}'_i| + \varepsilon_2/\alpha^2. \quad (5.19)$$

Теперь можно оценить отличие вычисленной величины $(\bar{b}'_i)_m$ от \bar{b}'_i .

Используя еще раз оценки (5.6) и (5.19), получаем $|(\bar{b}'_i)_m - \bar{b}'_i| \leq [\varepsilon_1(2 + \varepsilon_1)(1 + \hat{\omega}) + \hat{\omega}]|\bar{b}'_i| + \varepsilon_2/\alpha^2$. Вводя обозначение

$$\tilde{\omega} = \varepsilon_1(2 + \varepsilon_1)(1 + \hat{\omega}) + \hat{\omega}, \quad (5.20)$$

полученную оценку окончательно запишем в виде

$$|(\bar{b}'_i)_m - \bar{b}'_i| \leq \tilde{\omega}|\bar{b}'_i| + \varepsilon_2/\alpha^2. \quad (5.21)$$

Имея оценки погрешностей, допускаемых при вычислении каждого из элементов \bar{d}'_i , \bar{b}'_i матрицы \bar{A} , можно оценить погрешность всей вычисленной матрицы, т. е. матрицы с элементами $(\bar{d}'_i)_m$ и $(\bar{b}'_i)_m$. Обозначим эту матрицу через $(\bar{A})_m$ и оценим норму разности $(\bar{A})_m - \bar{A}$. Для этого матрицу $(\bar{A})_m - \bar{A}$ представим в виде суммы двух матриц R_1 и R_2 :

$$(\bar{A})_m - \bar{A} = R_1 + R_2, \quad (5.22)$$

где $R_1 = \text{diag}[(\bar{d}'_1)_m - \bar{d}_1, (\bar{d}'_2)_m - \bar{d}_2, \dots, (\bar{d}'_{m-1})_m - \bar{d}_{m-1}, 0]$, а матрица R_2 имеет вид

$$R_2 = \begin{bmatrix} 0 & (\bar{b}'_2)_m - \bar{b}_2 & & & & & 0 \\ (\bar{b}'_2)_m - \bar{b}_2 & 0 & (\bar{b}'_3)_m - \bar{b}_3 & & & & \\ & & & \ddots & & & \\ & & & & 0 & & (\bar{b}'_{m-1})_m - \bar{b}_{m-1} \\ 0 & & & & (\bar{b}'_{m-1})_m - \bar{b}_{m-1} & 0 & 0 \\ & & & & & 0 & 0 \end{bmatrix}.$$

В силу (5.22),

$$\|(\bar{A})_m - \bar{A}\| \leq \|R_1\| + \|R_2\|, \quad (5.23)$$

следовательно, достаточно оценить нормы матриц R_1 и R_2 . Учитывая (5.16), оценим норму матрицы R_1 :

$$\|R_1\| = \max_{1 \leq i \leq m-1} |(\bar{d}'_i)_m - \bar{d}_i| \leq \tilde{\Delta}M(A).$$

Для оценки нормы матрицы R_2 воспользуемся оценками (2.30) из [4] и (5.21) из настоящего параграфа:

$$\begin{aligned} \|R_2\| &\leq \max_{2 \leq i \leq m-2} (|(\bar{b}'_i)_m - \bar{b}'_i| + |(\bar{b}'_{i+1})_m - \bar{b}'_{i+1}|) \leq \\ &\leq \tilde{\omega} \max_{2 \leq i \leq m-2} (|\bar{b}'_i| + |\bar{b}'_{i+1}|) + \frac{2\varepsilon_2}{\alpha^2}. \end{aligned}$$

Введем в рассмотрение норму

$$\mathcal{M}(\bar{A}) = \max \begin{cases} |\bar{d}_1| + |\bar{b}_2|, \\ \max_{2 \leq i \leq m-2} (|\bar{b}_i| + |\bar{d}_i| + |\bar{b}_{i+1}|), \\ |\bar{d}_{m-1}| + |\bar{b}_{m-1}|, \\ |\eta| \end{cases}$$

матрицы \bar{A} , аналогичную норме $\mathcal{M}(A)$ матрицы A . Используя эту норму, оценку нормы матрицы R_2 можно огрубить следующим образом: $\|R_2\| \leq \omega \mathcal{M}(\bar{A}) + 2e_2/\alpha^2$. Предполагая выполненным неравенство $2e_2/\alpha^2 \leq \varepsilon_1 \mathcal{M}(\bar{A})$, еще раз огрубим полученную оценку

$$\|R_2\| \leq (\tilde{\omega} + \varepsilon_1) \mathcal{M}(\bar{A}). \quad (5.24)$$

Здесь использовано уже упоминавшееся предположение о нормированности матрицы A , $\mathcal{M}(A) = O(1)$, вследствие которого норма $\mathcal{M}(\bar{A})$ также не может быть слишком малой.

Легко показать, что спектральная норма матрицы \bar{A} и норма $\mathcal{M}(\bar{A})$ связаны соотношением

$$\mathcal{M}(\bar{A}) \leq 2\|\bar{A}\|. \quad (5.25)$$

Действительно, очевидно, что $\|\bar{A}\| \geq \sqrt{\bar{b}_i^2 + \bar{d}_i^2 + \bar{b}_{i+1}^2}$. Но $\sqrt{\bar{b}_i^2 + \bar{d}_i^2 + \bar{b}_{i+1}^2} \geq (1/2)(|\bar{b}_i| + |\bar{d}_i| + |\bar{b}_{i+1}|)$. Следовательно, $|\bar{b}_i| + |\bar{d}_i| + |\bar{b}_{i+1}| \leq 2\|\bar{A}\|$. Аналогично можно показать, что $|\bar{d}_1| + |\bar{b}_2| \leq 2\|\bar{A}\|$, $|\bar{d}_{m-1}| + |\bar{b}_{m-1}| \leq 2\|\bar{A}\|$. Кроме того, очевидно, что $|\eta| \leq 2\|\bar{A}\|$. Следовательно, неравенство (5.25) выполнено. Огрубляя с его помощью оценку (5.24), получим

$$\|R_2\| \leq 2(\tilde{\omega} + \varepsilon_1)\|\bar{A}\|. \quad (5.26)$$

Оценим теперь величину $\|\bar{A}\|$ через величину $\mathcal{M}(A)$. Для этого воспользуемся оценкой (1.8), которая в силу равенства $CAC^* = \bar{A} + K$ записывается в виде

$$\|\bar{A} - CAC^*\| \leq 2(\tilde{\alpha} + \sqrt{m}\tilde{\beta})\mathcal{M}(A). \quad (5.27)$$

Отсюда следует, что $\|\bar{A}\| \leq \|A\| + 2(\tilde{\alpha} + \sqrt{m}\tilde{\beta})\mathcal{M}(A)$. Но спектральная норма симметрической матрицы не превосходит любой другой ее нормы. Следовательно, $\|\bar{A}\| \leq [1 + 2(\tilde{\alpha} + \sqrt{m}\tilde{\beta})]\mathcal{M}(A)$. Оценку (5.26) теперь можно огрубить следующим образом: $\|R_2\| \leq 2(\tilde{\omega} + \varepsilon_1)[1 + 2(\tilde{\alpha} + \sqrt{m}\tilde{\beta})]\mathcal{M}(A)$. Вводя обозначение

$$\tilde{\Delta} = 2(\tilde{\omega} + \varepsilon_1)[1 + 2(\tilde{\alpha} + \sqrt{m}\tilde{\beta})], \quad (5.28)$$

полученную оценку окончательно запишем в виде $\|R_2\| \leq \tilde{\Delta}\mathcal{M}(A)$. Подставляя полученные оценки норм матриц R_1 и R_2 в (5.23), находим

$$\|(\bar{A})_n - \bar{A}\| \leq (\tilde{\Delta} + \tilde{\Delta})\mathcal{M}(A). \quad (5.29)$$

Это и есть требуемая оценка.

В заключение приведем элементарно следующую из (5.27) и (5.29) оценку $\|(\bar{A})_n - CAC^*\| \leq [2(\tilde{\alpha} + \sqrt{m}\tilde{\beta}) + \tilde{\Delta} + \tilde{\Delta}]\mathcal{M}(A)$, которая показывает, насколько отличается вычисленная матрица $(\bar{A})_n$ от матрицы CAC^* , ортогонально подобной исходной матрице A . Таким образом, доказана

Теорема 3. Пусть параметры c_i и s_i , определяющие ортогональные преобразования вращения, и задаваемая с их помощью матрица C опре-

делены так же, как в § 1. Пусть, далее, $(\bar{A})_m$ обозначает матрицу

$$(\bar{A})_m = \begin{bmatrix} (\bar{d}'_1)_m & (\bar{b}'_2)_m & & & & & \\ (\bar{b}'_2)_m & (\bar{d}'_2)_m & (\bar{b}'_3)_m & & & & \\ & & & \ddots & & & \\ & & & & (\bar{b}'_{m-2})_m & (\bar{d}'_{m-2})_m & (\bar{b}'_{m-1})_m \\ 0 & & & & (\bar{b}'_{m-1})_m & (\bar{d}'_{m-1})_m & 0 \\ & & & & & & 0 \end{bmatrix},$$

элементы которой получены в результате машинного вычисления по формулам

$$\bar{d}'_i = d_{i+1} - \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} + \frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}},$$

$$\bar{b}'_i = \frac{s'_i b_{i+1}}{s'_{i+1}}.$$

Здесь c'_i и s'_i — имеющиеся в памяти машины приближения к точным значениям параметров c_i и s_i , связанные с ними равенствами $c'_i = c_i(1 + \kappa_i)$, $s'_i = s_i(1 + \omega_i)$, где $|\kappa_i| \leq \kappa \ll 1$, $|\omega_i| \leq \omega \ll 1$. Тогда имеет место оценка $\|(\bar{A})_m - CA^*C^*\| \leq [2(\tilde{\alpha} + \sqrt{m}\tilde{\beta}) + \tilde{\Delta} + \tilde{\tilde{\Delta}}] \mathcal{M}(A)$.

В этой оценке константы α и β определяются точностью решения уравнений (1.5) для параметров c_i и s_i , а константы $\tilde{\Delta}$ и $\tilde{\tilde{\Delta}}$ — в соответствии с (5.15) и (5.28) и зависят в основном от величины значений κ и ω , а также от точности вычисления элементов \bar{d}'_i и \bar{b}'_i по выписанным формулам.

§ 6. Машинное вычисление преобразованной двухдиагональной матрицы

В этом параграфе описано вычисление элементов \bar{a}_i , \bar{b}_i матрицы (2.8) \bar{A} по формулам

$$\bar{a}_i = \frac{s_{i+1} a_{i+1}}{s_{i+1}}, \quad i = 1, \dots, N-1;$$

$$\bar{b}_i = \frac{s'_i b_{i+1}}{s_{i+1}}, \quad i = 2, \dots, N-1,$$
(6.1)

а также анализируются возникающие при этом погрешности. В первую очередь, они определяются тем, что вместо точных значений параметров s_i и \bar{s}_i , используемых в формулах (6.1), в памяти машины имеются лишь некоторые их приближения s'_i и \bar{s}'_i . Поэтому реальные вычисления можно вести только по формулам

$$\bar{a}'_i = \frac{s'_{i+1} a_{i+1}}{\bar{s}'_{i+1}}, \quad i = 1, \dots, N-1;$$

$$\bar{b}'_i = \frac{\bar{s}'_i b_{i+1}}{s_{i+1}}, \quad i = 2, \dots, N-1,$$
(6.2)

а не по формулам (6.1).

Чтобы оценить возникающие при таком изменении формул погрешности, вспомним, что в § 4 разработан простой алгоритм вычисления величин s'_i и \bar{s}'_i , гарантирующий выполнение равенств, связывающих эти

вычисленные значения с точными значениями s_i и \bar{s}_i ,

$$s'_i = s_i(1 + \pi_i), \quad \bar{s}'_i = \bar{s}_i(1 + \bar{\pi}_i). \quad (6.3)$$

Там же показано, что данные равенства имеют место при некоторых значениях величин π_i и $\bar{\pi}_i$ таких, что

$$|\pi_i| \leq \pi, \quad |\bar{\pi}_i| \leq \bar{\pi}, \quad (6.4)$$

где в качестве значений π можно взять $\pi = 5,001\epsilon_1$. Для дальнейшего отметим также, что в упомянутом алгоритме каждая из величин s'_i и \bar{s}'_i определяется с помощью двух чисел: целочисленного порядка и мантиссы, так что

$$s'_i = m(s'_i) \alpha^{p(s'_i)}, \quad \bar{s}'_i = m(\bar{s}'_i) \alpha^{p(\bar{s}'_i)}, \quad (6.5)$$

где $p(s'_i)$ и $p(\bar{s}'_i)$ — порядки чисел s'_i и \bar{s}'_i соответственно, а $m(s'_i)$ и $m(\bar{s}'_i)$ — их мантиссы, удовлетворяющие неравенствам

$$1/\alpha \leq |m(s'_i)| < 1, \quad 1/\alpha \leq |m(\bar{s}'_i)| < 1. \quad (6.6)$$

Наконец, отметим, что из (6.5) и (6.6) следует отличие от нуля чисел s'_i и \bar{s}'_i .

Теперь уже можно оценить погрешности, обусловленные заменой точных формул (6.1) формулами (6.2). В силу равенств (6.3) ясно, что для величин $a'_i - a_i$ и $b'_i - b_i$ имеют место выражения

$$a'_i - a_i = \left(\frac{1 + \pi_{i+1}}{1 + \bar{\pi}_{i+1}} - 1 \right) a_i, \quad b'_i - b_i = \left(\frac{1 + \bar{\pi}_i}{1 + \pi_{i+1}} - 1 \right) b_i.$$

С учетом неравенств (6.4) из этих формул следуют оценки $|a'_i - a_i| \leq \leq [2\pi/(1 - \pi)] |a_i|$, $|b'_i - b_i| \leq [2\bar{\pi}/(1 - \bar{\pi})] |\bar{b}_i|$. Вводя обозначение

$$\hat{\pi} = \frac{2\pi}{1 - \pi}, \quad (6.7)$$

перепишем полученные оценки в виде

$$|a'_i - a_i| \leq \hat{\pi} |a_i|, \quad |b'_i - b_i| \leq \hat{\pi} |\bar{b}_i|. \quad (6.8)$$

Вторым источником погрешностей, возникающих при вычислении элементов матрицы \bar{A} , является непосредственное выполнение арифметических операций, предписываемых формулами (6.2). Во избежание переполнения, а также исчезновения порядка, возможных при выполнении промежуточных операций, предлагается вести вычисления по специальным алгоритмам. Представив числа a_{i+1} и b_{i+1} в виде

$$a_{i+1} = m(a_{i+1}) \alpha^{p(a_{i+1})}, \quad \frac{1}{\alpha} \leq |m(a_{i+1})| < 1;$$

$$b_{i+1} = m(b_{i+1}) \alpha^{p(b_{i+1})}, \quad \frac{1}{\alpha} \leq |m(b_{i+1})| < 1,$$

алгоритмы вычисления \bar{a}'_i и \bar{b}'_i сформулируем следующим образом:

- 1) $x = (m(s'_{i+1}) m(a_{i+1}) / m(\bar{s}'_{i+1}))_m;$
- 2) $j = p(s'_{i+1}) + p(a_{i+1}) - p(\bar{s}'_{i+1});$
- 3) $(\bar{a}'_i)_m = (x \alpha^j)_m;$

- 1) $y = (m(\bar{s}'_i) m(b_{i+1}) / m(s'_{i+1}))_m;$
- 2) $k = p(\bar{s}'_i) + p(b_{i+1}) - p(s'_{i+1});$
- 3) $(\bar{b}'_i)_m = (y \alpha^k)_m,$

Здесь через $(\bar{a}'_i)_m$ и $(\bar{b}'_i)_m$ обозначены результаты машинного вычисления величин \bar{a}_i и \bar{b}_i по сформулированным алгоритмам.

Учтем погрешности, возникающие при реализации алгоритма (6.9) на вычислительной машине. Нетрудно проверить, что погрешность выполнения п. 1 оценивается так:

$$\left| x - \frac{m(s'_{i+1}) m(a_{i+1})}{m(s'_{i+1})} \right| \leq \varepsilon_1 (2 + \varepsilon_1) \left| \frac{m(s'_{i+1}) m(a_{i+1})}{m(s'_{i+1})} \right|.$$

Принимая во внимание формулы (6.2) и умножая обе части полученной оценки на α^j , перепишем ее в виде $|x\alpha^j - \bar{a}'_i| \leq \varepsilon_1 (2 + \varepsilon_1) |\bar{a}'_i|$. Пункт 2, очевидно, выполняется точно. При выполнении же п. 3 может возникнуть только погрешность исчезновения порядка. Эта погрешность появляется в том случае, если $|x\alpha^j| < \varepsilon_2/\alpha^2$. При этом величина $(\bar{a}'_i)_m$ получается равной нулю, а возникающая из-за этого погрешность, очевидно, не превосходит ε_2/α^2 : $|(\bar{a}'_i)_m - x\alpha^j| \leq \varepsilon_2/\alpha^2$. Окончательно, оценка погрешности, возникающей при выполнении алгоритма (6.9), примет вид

$$|(\bar{a}'_i)_m - \bar{a}'_i| \leq \varepsilon_1 (2 + \varepsilon_1) |\bar{a}'_i| + \varepsilon_2/\alpha^2. \quad (6.11)$$

Для погрешности выполнения алгоритма (6.10), очевидно, имеет место аналогичная оценка

$$|(\bar{b}'_i)_m - \bar{b}'_i| \leq \varepsilon_1 (2 + \varepsilon_1) |\bar{b}'_i| + \varepsilon_2/\alpha^2. \quad (6.12)$$

В силу (6.8) из (6.11) и (6.12) легко следуют оценки

$$\begin{aligned} |(\bar{a}'_i)_m - \bar{a}'_i| &\leq \varepsilon_1 (2 + \varepsilon_1) (1 + \hat{\pi}) |\bar{a}'_i| + \varepsilon_2/\alpha^2; \\ |(\bar{b}'_i)_m - \bar{b}'_i| &\leq \varepsilon_1 (2 + \varepsilon_1) (1 + \hat{\pi}) |\bar{b}'_i| + \varepsilon_2/\alpha^2. \end{aligned} \quad (6.13)$$

Теперь уже можно достаточно просто оценить отличие вычисленных элементов $(\bar{a}'_i)_m$ и $(\bar{b}'_i)_m$ от \bar{a}_i и \bar{b}_i соответственно. Действительно, из (6.8) и (6.13) легкоходим

$$\begin{aligned} |(\bar{a}'_i)_m - \bar{a}_i| &\leq [\varepsilon_1 (2 + \varepsilon_1) (1 + \hat{\pi}) + \hat{\pi}] |\bar{a}_i| + \varepsilon_2/\alpha^2; \\ |(\bar{b}'_i)_m - \bar{b}_i| &\leq [\varepsilon_1 (2 + \varepsilon_1) (1 + \hat{\pi}) + \hat{\pi}] |\bar{b}_i| + \varepsilon_2/\alpha^2. \end{aligned}$$

После введения обозначения

$$\tilde{\pi} = \varepsilon_1 (2 + \varepsilon_1) (1 + \hat{\pi}) + \hat{\pi} \quad (6.14)$$

перепишем полученные оценки

$$|(\bar{a}'_i)_m - \bar{a}_i| \leq \tilde{\pi} |\bar{a}_i| + \frac{\varepsilon_2}{\alpha^2}, \quad |(\bar{b}'_i)_m - \bar{b}_i| \leq \tilde{\pi} |\bar{b}_i| + \frac{\varepsilon_2}{\alpha^2}. \quad (6.15)$$

Обозначим теперь через $(\bar{A})_m$ вычисленную матрицу, т. е. матрицу с элементами $(\bar{a}'_i)_m$, $(\bar{b}'_i)_m$, и оценим ее отличие от матрицы \bar{A} . Для этого воспользуемся оценкой (2.30) из [4] и только что полученными оценками (6.15):

$$\begin{aligned} \|(\bar{A})_m - \bar{A}\| &\leq \max \left\{ \max_{1 \leq i \leq N-2} (|(\bar{a}'_i)_m - \bar{a}_i| + |(\bar{b}'_{i+1})_m - \bar{b}_{i+1}|), \right. \\ &\quad \left. \max_{2 \leq i \leq N-1} (|(\bar{a}'_i)_m - \bar{a}_i| + |(\bar{b}'_i)_m - \bar{b}_i|) \right\} \leq \\ &\leq \tilde{\pi} \max \left\{ \max_{1 \leq i \leq N-2} (|\bar{a}_i| + |\bar{b}_{i+1}|), \right. \\ &\quad \left. \max_{2 \leq i \leq N-1} (|\bar{a}_i| + |\bar{b}_i|) \right\} + \frac{2\varepsilon_2}{\alpha^2}. \end{aligned} \quad (6.16)$$

На основе таких же рассуждений, с помощью которых получены неравенства (2.24), нетрудно показать, что $|\bar{a}_i| + |\bar{b}_{i+1}| \leq \sqrt{2} \|\bar{A}\|$, $|\bar{a}_i| + |\bar{b}_i| \leq$

$\leq \sqrt{2}\|A\|$. Тогда из (6.16) легко следует

$$\|(\bar{A})_m - \bar{A}\| \leq \sqrt{2}\pi\|\bar{A}\| + \frac{2\varepsilon_2}{\alpha^2}. \quad (6.17)$$

Из оценки (2.9), которую можно переписать в виде

$$\|\bar{A} - \bar{C}AC^*\| \leq 2\sqrt{2}(\sqrt{N} + 2)\varepsilon\|A\|, \quad (6.18)$$

следует, что $\|\bar{A}\| \leq [1 + 2\sqrt{2}(\sqrt{N} + 2)\varepsilon]\|A\|$. Подставляя последнюю оценку в (6.17), получим $\|(\bar{A})_m - \bar{A}\| \leq \sqrt{2}\pi[1 + 2\sqrt{2}(\sqrt{N} + 2)\varepsilon]\|A\| + 2\varepsilon_2/\alpha^2$. Наконец, предположим, что выполнено неравенство $2\varepsilon_2/\alpha^2 \leq \varepsilon_1\|A\|$. Это предположение опирается на условие нормированности матрицы A и неравенство $\varepsilon_2 \ll \varepsilon_1$. Тогда оценку нормы матрицы $(\bar{A})_m - \bar{A}$ можно записать так:

$$\|(\bar{A})_m - \bar{A}\| \leq (\sqrt{2}\pi + \varepsilon_1)[1 + 2\sqrt{2}(\sqrt{N} + 2)\varepsilon]\|A\|. \quad (6.19)$$

В заключение приведем оценку

$$\|(\bar{A})_m - \bar{C}AC^*\| \leq \{2\sqrt{2}(\sqrt{N} + 2)\varepsilon + (\sqrt{2}\pi + \varepsilon_1)[1 + 2\sqrt{2}(\sqrt{N} + 2)\varepsilon]\}\|A\|,$$

которая элементарно следует из (6.18) и (6.19) и показывает, насколько вычисленная матрица $(\bar{A})_m$ отличается от матрицы $\bar{C}AC^*$, ортогонально эквивалентной исходной матрице A . Таким образом, нами доказана

Теорема 4. Пусть параметры c_i, s_i, \bar{c}_i и \bar{s}_i , определяющие преобразования вращения, и задаваемые с их помощью ортогональные матрицы C и \bar{C} определены так же, как в § 2. Пусть, далее, $(\bar{A})_m$ обозначает матрицу

$$(\bar{A})_m = \begin{bmatrix} (\bar{a}'_1)_m & (\bar{b}'_2)_m \\ (\bar{a}'_2)_m & (\bar{b}'_3)_m \\ \vdots & \vdots \\ (\bar{a}'_{N-2})_m & (\bar{b}'_{N-1})_m \\ 0 & (\bar{a}'_{N-1})_m \\ & 0 \\ & \sigma \end{bmatrix},$$

элементы $(\bar{a}'_i)_m, (\bar{b}'_i)_m$ которой получены в результате машинного вычисления по формулам

$$\bar{a}'_i = \frac{\bar{s}'_{i+1}a_{i+1}}{\bar{s}'_{i+1}}, \quad \bar{b}'_i = \frac{\bar{s}'_i b_{i+1}}{\bar{s}'_{i+1}}.$$

Здесь \bar{s}'_i и \bar{s}'_i — имеющиеся в памяти машины приближения к точным значениям параметров s_i и \bar{s}_i , связанные с ними равенствами $\bar{s}'_i = s_i(1 + \pi_i)$, $\bar{s}'_i = \bar{s}_i(1 + \bar{\pi}_i)$, которые имеют место при некоторых значениях π_i и $\bar{\pi}_i$ таких, что

$$|\pi_i| \leq \pi \ll 1, \quad |\bar{\pi}_i| \leq \pi. \quad (6.20)$$

Тогда справедлива оценка

$$\|(\bar{A})_m - \bar{C}AC^*\| \leq \{2\sqrt{2}(\sqrt{N} + 2)\varepsilon + (\sqrt{2}\pi + \varepsilon_1)[1 + 2\sqrt{2}(\sqrt{N} + 2)\varepsilon]\}\|A\|.$$

В этой оценке константа ε определяется точностью решения уравнений (2.5) относительно параметров c_i, s_i, \bar{c}_i и \bar{s}_i , а константа π определяется по формуле (6.14) и зависит от значения π , с которым удовлетворяются неравенства (6.20), а также от точности вычисления элементов \bar{a}'_i, \bar{b}'_i по приведенным формулам.

Г л а в а 4

ОБЩИЕ СХЕМЫ АЛГОРИТМОВ ИСЧЕРПЫВАНИЯ ТРЕХДИАГОНАЛЬНЫХ СИММЕТРИЧЕСКИХ И ДВУХДИАГОНАЛЬНЫХ МАТРИЦ

§ 7. Общая схема исчерпывания симметрической трехдиагональной матрицы

Рассмотрим симметрическую трехдиагональную матрицу

$$\tilde{A} = \begin{bmatrix} \tilde{d}_1 & \tilde{b}_2 & & & & \\ \tilde{b}_2 & \tilde{d}_2 & \tilde{b}_3 & & & \\ & \ddots & \ddots & \ddots & & \\ & & & \tilde{b}_{m-1} & \tilde{d}_{m-1} & \tilde{b}_m \\ 0 & & & & \tilde{b}_m & \tilde{d}_m \end{bmatrix}$$

такую, что

$$\mathcal{M}(\tilde{A}) = \max \left\{ \max_{2 \leq i \leq m-1} (|\tilde{b}_i| + |\tilde{d}_i| + |\tilde{b}_{i+1}|), \frac{|\tilde{d}_1| + |\tilde{b}_2|}{|\tilde{d}_m| + |\tilde{b}_m|} \right\} = O(1),$$

т. е. матрица \tilde{A} предполагается нормированной. Отметим, что относительно величины элементов \tilde{b}_i ее побочных диагоналей никаких предложений не делается.

Пусть $\lambda(\tilde{A})$ — одно из собственных значений матрицы \tilde{A} . В этом параграфе опишем общую схему исчерпывания собственного значения $\lambda(\tilde{A})$ и приведем оценку погрешности, возникающей при этом исчерпывании. Предположим, что η — некоторое приближение к $\lambda(\tilde{A})$ такое, что

$$|\eta - \lambda(\tilde{A})| \leq \gamma \mathcal{M}(\tilde{A}). \quad (7.1)$$

Отметим, что приближенное собственное значение η можно получить, например, с помощью алгоритма бисекций, основанного на применении теоремы Штурма. В [5], где этот алгоритм подробно описан и исследован, показано (см. [5], § 10), что при $\mathcal{M}(\tilde{A}) = O(1)$ значение η , получаемое в результате работы алгоритма, удовлетворяет оценке (7.1) с $\gamma = 6\varepsilon_1$.

Наряду с матрицей \tilde{A} введем в рассмотрение

$$A = \begin{bmatrix} d_1 & b_2 & & & & & 0 \\ b_2 & d_2 & b_3 & & & & \\ & \ddots & \ddots & \ddots & & & \\ & & & b_{m-1} & d_{m-1} & b_m & \\ 0 & & & & b_m & d_m & \end{bmatrix},$$

положив $d_i = \tilde{d}_i$, а элементы b_i ее побочных диагоналей определив следующим образом:

$$b_i = \begin{cases} \tilde{b}_i, & \text{если } |\tilde{b}_i| \geq \varepsilon_1 \mathcal{M}(\tilde{A}); \\ \text{иначе } \begin{cases} \varepsilon_1 \mathcal{M}(\tilde{A}), & \text{если } \tilde{b}_i > 0; \\ -\varepsilon_1 \mathcal{M}(\tilde{A}), & \text{если } \tilde{b}_i \leq 0. \end{cases} & \end{cases}$$

Легко видеть, что

$$(1 - 2\epsilon_1)\mathcal{M}(\tilde{A}) \leq \mathcal{M}(A) \leq (1 + 2\epsilon_1)\mathcal{M}(\tilde{A}). \quad (7.2)$$

Кроме того, нетрудно показать, что соответствующие собственные значения матриц A и \tilde{A} различаются не более чем на $2\epsilon_1\mathcal{M}(\tilde{A})$. Отсюда из оценки (7.1) следует, что η является приближением к некоторому собственному значению $\lambda(A)$ матрицы A , удовлетворяющим оценке

$$|\eta - \lambda(A)| \leq \gamma \mathcal{M}(\tilde{A}), \quad (7.3)$$

где

$$\gamma = \tilde{\gamma} + 2\epsilon_1. \quad (7.4)$$

С помощью алгоритма, детально описанного в § 3 работы [4], определим виртуальный почти собственный вектор $u = (u_1, u_2, \dots, u_m)^T$ матрицы A , отвечающий ее почти собственному значению η . Напомним, что вектор называется виртуальным, если он задается с помощью не своих компонент, а их последовательных отношений. В используемом алгоритме виртуальный почти собственный вектор u определяется двумя последовательностями $\{\mathcal{R}_i\}$, $\{Q_i\}$ и целым числом i_0 так, что отношения $u_2/u_1, u_3/u_2, \dots, u_{i_0-1}/u_{i_0-2}$ равны соответственно числам $\mathcal{R}_2, \mathcal{R}_3, \dots, \mathcal{R}_{i_0-1}$, а отношения $u_{i_0}/u_{i_0-1}, u_{i_0+1}/u_{i_0}, \dots, u_m/u_{m-1}$ — числам $1/Q_{i_0}, 1/Q_{i_0+1}, \dots, 1/Q_m$.

Определим $\bar{\mathcal{P}}_2, \bar{\mathcal{P}}_3, \dots, \bar{\mathcal{P}}_m$ в соответствии со следующим правилом: $\bar{\mathcal{P}}_i = \mathcal{R}_i, i = 2, \dots, i_0 - 1; \bar{\mathcal{P}}_i = (1 + \psi_i)/Q_i, i = i_0, \dots, m$. Здесь ψ_i моделируют погрешности делений $1/Q_i$, поэтому, очевидно, $|\psi_i| \leq \epsilon_1$. Согласно теореме 2 из § 2 работы [4], числа $\bar{\mathcal{P}}_i$ удовлетворяют соотношениям

$$\begin{aligned} d_1 - \eta + \bar{\alpha}_1 + b_2 \bar{\mathcal{P}}_2 &= 0; \\ \frac{b_i(1 + \bar{\beta}_i)}{\bar{\mathcal{P}}_i} + d_i - \eta + \bar{\alpha}_i + b_{i+1} \bar{\mathcal{P}}_{i+1} &= 0, \quad i = 2, \dots, m-1; \\ \frac{b_m(1 + \bar{\beta}_m)}{\bar{\mathcal{P}}_m} + d_m - \eta + \bar{\alpha}_m &= 0, \end{aligned}$$

где для возмущений $\bar{\alpha}_i$ и $\bar{\beta}_i$ имеют место оценки

$$\begin{aligned} |\bar{\alpha}_i| &\leq \alpha'_i(1 + \epsilon_1) + \epsilon_1(|d_i| + |\eta|), \quad |\bar{\alpha}_m| \leq \alpha''_m, \quad |\bar{\beta}_m| \leq \epsilon_1; \\ |\bar{\alpha}_i| &\leq \max \left\{ \frac{\alpha'_i(1 + \epsilon_1) + \epsilon_1(|d_i| + |\eta|)}{1 - |\beta''_{i+1}|}, \frac{\alpha''_i(1 + \epsilon_1) + (|\beta''_{i+1}| + \epsilon_1)(|d_i| + |\eta|)}{1 - |\beta''_{i+1}|} \right\}, \quad i = 2, \dots, m-1. \quad (7.5) \\ |\bar{\beta}_i| &\leq \max \left\{ |\beta'_i|(1 + \epsilon_1) + \epsilon_1, \frac{|\beta''_{i+1}| + \epsilon_1(2 + \epsilon_1)}{1 - |\beta''_{i+1}|} \right\}. \end{aligned}$$

Используемые здесь величины $\alpha'_i, \alpha''_i, \beta'_i, \beta''_i$ подчиняются, в свою очередь, следующим оценкам (см. [4], § 4):

$$\gamma \mathcal{M}(\tilde{A}) < \alpha'_i \leq \gamma \mathcal{M}(\tilde{A}) + 16\epsilon_1 \mathcal{M}(A), \quad |\beta'_i| \leq 10\epsilon_1;$$

$$\gamma \mathcal{M}(\tilde{A}) < \alpha''_i \leq \gamma \mathcal{M}(\tilde{A}) + 16\epsilon_1 \mathcal{M}(A), \quad |\beta''_i| \leq 10\epsilon_1.$$

Упростим оценки (7.5). Согласно сделанному ранее замечанию относительно возможного значения величины $\tilde{\gamma}$ в оценке (7.1) примем $\tilde{\gamma} = 6\epsilon_1$. Тогда из (7.1) следует, что в качестве значения γ можно взять $\gamma = 8\epsilon_1$. Принимая теперь во внимание неравенства (7.2) и пренебрегая членами второго порядка малости, оценки величин α_i и β_i запишем в виде $|\alpha_i| \leq \gamma \mathcal{M}(A)$, $|\beta_i| \leq \beta$, где в качестве значений α и β можно взять $\alpha = 46\epsilon_1$, $\beta = 12\epsilon_1$.

С помощью чисел $\bar{\mathcal{P}}_i$ определим числа c_i , s_i , c'_i и s'_i следующим образом:

$$\begin{aligned} s_2 &= \frac{1}{\sqrt{1 + \bar{\mathcal{P}}_2^2}}, \quad c_2 = \frac{\bar{\mathcal{P}}_2}{\sqrt{1 + \bar{\mathcal{P}}_2^2}}, \quad c'_2 = (c_2)_M, \quad s'_2 = (s_2)_M, \\ s_i &= \frac{1}{\sqrt{1 + (c'_{i-1}\bar{\mathcal{P}}_i)^2}}, \quad c_i = \frac{c'_{i-1}\bar{\mathcal{P}}_i}{\sqrt{1 + (c'_{i-1}\bar{\mathcal{P}}_i)^2}}, \quad c'_i = (c_i)_M, \\ s'_i &= (s_i)_M, \quad i = 3, \dots, m. \end{aligned} \tag{7.6}$$

В § 4 главы 2 разработан специальный алгоритм вычисления величин c_i и s_i по приведенным формулам. Получаемые в результате применения этого алгоритма машинные числа c'_i и s'_i связаны с c_i и s_i равенствами $c'_i = c_i(1 + \kappa_i)$, $s'_i = s_i(1 + \omega_i)$, которые имеют место при некоторых κ_i и ω_i таких, что $|\kappa_i| \leq \kappa = 5,001\epsilon_1$, $|\omega_i| \leq \omega = 4,001\epsilon_1$. Отметим, что в используемом алгоритме каждая из величин c'_i и s'_i определяется с помощью двух чисел — мантиссы и целочисленного порядка, а именно c'_i определяется числами $m(c'_i)$ и $p(c'_i)$, а s'_i определяется числами $m(s'_i)$ и $p(s'_i)$ так, что

$$c'_i = m(c'_i) \alpha^{p(c'_i)}, \quad 1/\alpha \leq |m(c'_i)| < 1;$$

$$s'_i = m(s'_i) \alpha^{p(s'_i)}, \quad 1/\alpha \leq |m(s'_i)| < 1.$$

Здесь α — основание системы счисления, принятой в вычислительной машине. Подчеркнем также, что алгоритм гарантирует отличие от нуля чисел c'_i и s'_i .

В § 3 главы 2 показано, что числа c_i и s_i удовлетворяют в таком случае соотношениям

$$s_2(d_1 + \alpha_1 - \eta) + c_2 b_2 = 0;$$

$$s_3 s_2 b_2 (1 + \beta_2) + s_3 c_2 (d_2 + \alpha_2 - \eta) + c_3 b_3 = 0;$$

$$s_{i+1} s_i c_{i-1} b_i (1 + \beta_i) + s_{i+1} c_i (d_i + \alpha_i - \eta) + c_{i+1} b_{i+1} = 0, \quad i = 3, \dots, m-1;$$

$$s_m c_{m-1} b_m (1 + \beta_m) + c_m (d_m + \alpha_m - \eta) = 0,$$

где возмущения α_i и β_i определяются формулами

$$\alpha_1 = \bar{\alpha}_1, \quad \beta_2 = (1 + \bar{\beta}_2)(1 + \kappa_2) - 1;$$

$$\alpha_i = (d_i - \eta) \kappa_i + \bar{\alpha}_i (1 + \kappa_i), \quad i = 2, \dots, m-1;$$

$$\beta_i = (1 + \bar{\beta}_i)(1 + \kappa_{i-1})(1 + \kappa_i) - 1, \quad i = 3, \dots, m-1;$$

$$\alpha_m = \bar{\alpha}_m, \quad \beta_m = (1 + \bar{\beta}_m)(1 + \kappa_{m-1}) - 1.$$

Из этих равенств следуют оценки $|\alpha_i| \leq \tilde{\alpha} M(A)$, $|\beta_i| \leq \tilde{\beta}$, в которых в качестве значений $\tilde{\alpha}$ и $\tilde{\beta}$ с точностью до членов первого порядка малости по ϵ_1 можно взять

$$\tilde{\alpha} = \bar{\alpha} + 2\kappa = 56,002\epsilon_1, \quad \tilde{\beta} = \bar{\beta} + 2\kappa = 22,002\epsilon_1. \tag{7.7}$$

Определим теперь матрицу

$$\begin{aligned} (\bar{A})_M &= \begin{bmatrix} (\bar{d}'_1)_M & (\bar{b}'_2)_M & & & & 0 \\ (\bar{b}'_2)_M & (\bar{d}'_2)_M & (\bar{b}'_3)_M & & & \\ & \cdot & \cdot & \cdot & \cdot & \\ & & & & & \\ & & & & & \\ 0 & & & -(\bar{b}'_{m-2})_M & (\bar{d}'_{m-2})_M & (\bar{b}'_{m-1})_M \\ & & & & (\bar{b}'_{m-1})_M & (\bar{d}'_{m-1})_M & 0 \\ & & & & & & 0 & \eta \end{bmatrix}, \end{aligned}$$

элементы $(\bar{d}'_i)_m$, $(\bar{b}'_i)_m$, которой получаются в результате машинного вычисления по формулам

$$\bar{d}'_1 = d_2 - \frac{c'_2 b_2}{s'_2} + \frac{c'_3 c'_2 b_3}{s'_3};$$

$$\bar{d}'_i = d_{i+1} - \frac{c'_{i+1} c'_i b_{i+1}}{s'_{i+1}} + \frac{c'_{i+2} c'_{i+1} b_{i+2}}{s'_{i+2}}, \quad i = 2, \dots, m-2;$$

$$\bar{d}'_{m-1} = d_m - \frac{c'_m c'_{m-1} b_m}{s'_m},$$

$$\bar{b}'_i = \frac{s'_i b_{i+1}}{s'_{i+1}}, \quad i = 2, \dots, m-1.$$

В § 5 главы 3 детально описаны алгоритмы, с помощью которых проводятся вычисления по этим формулам.

Из рассмотрений § 5 следует, что определенная таким образом матрица $(\bar{A})_m$ по норме мало отличается от некоторой матрицы, подобной матрице A . Точнее говоря, из теоремы 3 следует оценка

$$\|(\bar{A})_m - CAC^*\| \leq [2(\tilde{\alpha} + \sqrt{m}\tilde{\beta}) + \tilde{\Delta} + \tilde{\tilde{\Delta}}]M(A). \quad (7.8)$$

Напомним, что матрица C представляет собой произведение $C = C_m \cdots C_3 \cdot C_2$ ортогональных матриц вращения

$$C_i = \begin{bmatrix} 1 & & & & \\ & \ddots & & & \\ & & 0 & & \\ & & & 1 & \\ & & & & c_i - s_i \\ & & & & s_i & c_i \\ & & & & & 1 \\ & & & & & & \ddots \\ & & & & & & & 0 \\ & & & & & & & & 1 \end{bmatrix}$$

которые конструируются с помощью параметров c_i и s_i ($c_i^2 + s_i^2 = 1$), определенных формулами (7.6). Возможные значения констант α и β определяются формулами (7.7), а величины $\tilde{\Delta}$ и $\tilde{\tilde{\Delta}}$ — соответственно формулами (5.15) и (5.28). Напомним эти формулы:

$$\tilde{\Delta} = 2(2\Delta + \varepsilon_1)(1 + \tilde{\beta})(1 + \tilde{\alpha} + \gamma), \quad \tilde{\tilde{\Delta}} = 2(\tilde{\omega} + \varepsilon_1)[1 + \frac{1}{2}(\tilde{\alpha} + \sqrt{m}\tilde{\beta})].$$

Значения величин Δ и $\tilde{\omega}$, в свою очередь, определяются формулами (5.13) и (5.20):

$$\Delta = [(1 + \varepsilon_1)^2 - 1](1 + \hat{\alpha}) + \hat{\alpha}, \quad \tilde{\omega} = \varepsilon_1(2 + \varepsilon_1)(1 + \hat{\omega}) + \hat{\omega}.$$

Наконец, выпишем еще формулы (5.5), согласно которым определяются величины $\hat{\alpha}$ и $\hat{\omega}$,

$$\hat{\alpha} = (1 + \alpha)^2 / (1 - \omega) - 1, \quad \hat{\omega} = 2\omega / (1 - \omega).$$

Пренебрегая в последних формулах членами второго порядка малости по ε_1 , в качестве значений $\hat{\alpha}$ и $\hat{\omega}$ примем

$$\hat{\alpha} = 2\alpha + \omega = 14,003\varepsilon_1, \quad \hat{\omega} = 2\omega = 8,002\varepsilon_1.$$

Аналогичным образом упростим выражения для величин Δ и $\tilde{\omega}$:

$$\Delta = 5\epsilon_1 + \hat{\kappa} = 19,003\epsilon_1, \quad \tilde{\omega} = 2\epsilon_1 + \hat{\omega} = 10,002\epsilon_1.$$

Наконец, пренебрегая членами второго порядка в формулах для $\tilde{\Delta}$ и $\tilde{\tilde{\Delta}}$, получим

$$\tilde{\Delta} = 2(2\Delta + \epsilon_1) = 78,012\epsilon_1, \quad \tilde{\tilde{\Delta}} = 2(\tilde{\omega} + \epsilon_1) = 22,004\epsilon_1. \quad (7.9)$$

Итак, оценка (7.8) имеет место при значениях констант $\tilde{\alpha}$ и $\tilde{\beta}$, определяемых формулами (7.7), и значениях констант $\tilde{\Delta}$ и $\tilde{\tilde{\Delta}}$, определяемых формулами (7.9).

В заключение оценим отличие вычисленной матрицы $(\bar{A})_m$ от CAC^* , подобной исходной матрице \bar{A} . Соответствующая оценка легко следует из оценки (7.8) и очевидной оценки $\|A - \bar{A}\| \leq 2\epsilon_1 \mathcal{M}(\bar{A})$. Действительно,

$$\begin{aligned} \|(\bar{A})_m - CAC^*\| &\leq \|(\bar{A})_m - CAC^*\| + \|A - \bar{A}\| \leq \\ &\leq [2(\tilde{\alpha} + \sqrt{m}\tilde{\beta}) + \tilde{\Delta} + \tilde{\tilde{\Delta}}] \mathcal{M}(A) + 2\epsilon_1 \mathcal{M}(\bar{A}). \end{aligned}$$

Воспользовавшись еще раз неравенствами (7.2) и пренебрегая членами второго порядка, запишем полученную оценку в терминах нормы $\mathcal{M}(\bar{A})$ исходной матрицы \bar{A}

$$\|(\bar{A})_m - CAC^*\| \leq [2(\tilde{\alpha} + \sqrt{m}\tilde{\beta}) + \tilde{\Delta} + \tilde{\tilde{\Delta}} + 2\epsilon_1] \mathcal{M}(\bar{A}).$$

Подчеркнем, что к результатам работы описанного алгоритма кроме матрицы $(\bar{A})_m$ относится и значение выражения, стоящего в правой части полученной оценки. Кроме того, результатами работы алгоритма следует считать и числа c_i и s_i , а также величины κ и ω , оценивающие отличие этих чисел от точных значений параметров c_i и s_i , с помощью которых конструируется ортогональная матрица C .

§ 8. Общая схема исчерпывания двухдиагональной матрицы

Рассмотрим двухдиагональную матрицу

$$\bar{A} = \begin{bmatrix} \tilde{a}_1 & \tilde{b}_2 & & & 0 & & \\ & \tilde{a}_2 & \tilde{b}_3 & & & & \\ & & \ddots & & & & \\ & & & \ddots & & & \\ 0 & & & & \tilde{a}_{N-1} & \tilde{b}_N & \\ & & & & & \tilde{a}_N & \end{bmatrix},$$

относительно которой предположим, что

$$\mathcal{K}(\bar{A}) = \max \left\{ \begin{array}{l} \max_{1 \leq i \leq N-1} (|\tilde{a}_i| + |\tilde{b}_{i+1}|), \\ \max_{2 \leq i \leq N} (|\tilde{a}_i| + |\tilde{b}_i|) \end{array} \right\} = O(1). \quad (8.1)$$

Пусть $\sigma_N(\bar{A})$ — наибольшее сингулярное число матрицы \bar{A} , т. е. $\sigma_N(\bar{A}) = \|\bar{A}\|$. В этом параграфе описана общая схема исчерпывания сингулярного числа $\sigma_N(\bar{A})$ и приведена оценка погрешности, возникающей при этом исчерпывании. Заметим, что предположение (8.1) означает, что первым этапом работы алгоритма должно быть масштабирование исходной матрицы. Предположим, что σ — приближение к $\sigma_N(\bar{A}) = \|\bar{A}\|$ такое, что

$$|\sigma - \|\bar{A}\|| \leq \delta \mathcal{K}(\bar{A}). \quad (8.2)$$

Отметим, что если для получения σ воспользоваться алгоритмом бисек-

ций, подробно описанным в § 10 работы [5], то в качестве значения величины $\tilde{\delta}$ в оценке (8.2) можно взять $\tilde{\delta} = 4\epsilon_1$.

Наряду с матрицей \tilde{A} будем рассматривать матрицу

$$A = \begin{bmatrix} a_1 & b_2 & & & & \\ & a_2 & b_3 & & & \\ & & \ddots & \ddots & & \\ & & & 0 & b_N & \\ & & & & a_{N-1} & b_N \\ & & & & & a_N \end{bmatrix},$$

элементы которой определим следующим образом:

$$a_i = \begin{cases} \tilde{a}_i, & \text{если } |\tilde{a}_i| \geq \epsilon_1 \mathcal{K}(\tilde{A}); \\ \text{иначе} & \begin{cases} \epsilon_1 \mathcal{K}(\tilde{A}), & \text{если } \tilde{a}_i > 0; \\ -\epsilon_1 \mathcal{K}(\tilde{A}), & \text{если } \tilde{a}_i \leq 0; \end{cases} \end{cases} \quad (8.3)$$

$$b_i = \begin{cases} \tilde{b}_i, & \text{если } |\tilde{b}_i| \geq \epsilon_1 \mathcal{K}(\tilde{A}); \\ \text{иначе} & \begin{cases} \epsilon_1 \mathcal{K}(\tilde{A}), & \text{если } \tilde{b}_i > 0; \\ -\epsilon_1 \mathcal{K}(\tilde{A}), & \text{если } \tilde{b}_i \leq 0. \end{cases} \end{cases}$$

Введем в рассмотрение величину

$$\mathcal{K}(A) = \max \begin{cases} \max_{1 \leq i \leq N-1} (|a_i| + |b_{i+1}|), \\ \max_{2 \leq i \leq N} (|a_i| + |b_i|). \end{cases}$$

Легко видеть, что

$$(1 - 2\epsilon_1) \mathcal{K}(\tilde{A}) \leq \mathcal{K}(A) \leq (1 + 2\epsilon_1) \mathcal{K}(\tilde{A}). \quad (8.4)$$

Так как сингулярные числа матрицы A являются в то же время собственными значениями симметрической трехдиагональной матрицы

$$S = \begin{bmatrix} 0 & a_1 & & & & & \\ a_1 & 0 & b_2 & & & & \\ & b_2 & 0 & a_2 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & 0 & 0 & a & \\ & & & & & a_N & 0 \end{bmatrix},$$

имеющей порядок $2N$, то нетрудно показать, что сингулярные числа матриц A и \tilde{A} различаются не более чем на $2\epsilon_1 \mathcal{K}(\tilde{A})$. Отсюда следует, что σ является приближением к наибольшему сингулярному числу $\sigma_n(A) = \|A\|$ матрицы A , удовлетворяющим оценке

$$|\sigma - \|A\|| \leq \delta K(\tilde{A}), \quad (8.5)$$

где

$$\delta = \tilde{\delta} + 2\epsilon_1. \quad (8.6)$$

С помощью алгоритма, описанного в § 3 работы [4], определим виртуальный почти собственный вектор $w = (w_1, w_2, \dots, w_{2N-1}, w_{2N})^T$ матрицы S , отвечающий ее почти собственному значению σ . Для применения этого алгоритма необходимо заметить только, что используемая в его формулировке величина $\mathcal{M}(S)$, представляющая собой одну из норм матрицы S , совпадает с величиной $\mathcal{K}(A)$. В результате применения этого алгоритма находятся две последовательности чисел $Q_2, Q_3, \dots, Q_{2N-1}, Q_{2N}; T_2, T_3, \dots, T_{2N-1}, T_{2N}$ и целое число j_0 такие, что отно-

шения $w_1/w_1, w_3/w_2, \dots, w_{j_0-1}/w_{j_0-2}$ компонент вектора w равны соответственно числам $\mathcal{Q}_2, \mathcal{Q}_3, \dots, \mathcal{Q}_{j_0-1}$, а отношения $w_{j_0}/w_{j_0-1}, w_{j_0+1}/w_{j_0}, \dots, w_{2N}/w_{2N-1}$ — числам $1/T_{j_0}, 1/T_{j_0+1}, \dots, 1/T_{2N}$.

Определим числа $\mathcal{P}_1, \mathcal{R}_2, \mathcal{P}_2, \dots, \mathcal{R}_N, \mathcal{P}_N$ по следующему правилу. Если j_0 чётно, то представим j_0 в виде $j_0 = 2i_0$ и положим

$$\begin{aligned}\mathcal{P}_i &= \mathcal{Q}_{2i}, \quad \mathcal{R}_{i+1} = \mathcal{Q}_{2i+1}, \quad i = 1, \dots, i_0 - 1; \\ \mathcal{P}_i &= (1 + \psi_i)/T_{2i}, \quad \mathcal{R}_{i+1} = (1 + \varphi_{i+1})/T_{2i+1}, \quad i = i_0, \dots, N - 1; \\ \mathcal{P}_N &= (1 + \psi_N)/T_{2N}.\end{aligned}$$

Если же j_0 нечетно, то представим его в виде $j_0 = 2i_0 + 1$ и положим

$$\begin{aligned}\mathcal{P}_1 &= \mathcal{Q}_2, \quad \mathcal{R}_i = \mathcal{Q}_{2i-1}, \quad \mathcal{P}_i = \mathcal{Q}_{2i}, \quad i = 2, \dots, i_0; \\ \mathcal{R}_i &= (1 + \varphi_i)/T_{2i-1}, \quad \mathcal{P}_i = (1 + \psi_i)/T_{2i}, \quad i = i_0 + 1, \dots, N.\end{aligned}$$

Здесь возмущения φ_i и ψ_i моделируют погрешности, возникающие при выполнении операций деления $1/T_{2i-1}$ и $1/T_{2i}$ соответственно так, что $|\varphi_i| \leq \varepsilon_1, |\psi_i| \leq \varepsilon_1$. Из доказательства теоремы 1 работы [4] легко получить следующее утверждение. Числа $\mathcal{P}_1, \mathcal{R}_2, \mathcal{P}_2, \dots, \mathcal{R}_N, \mathcal{P}_N$ удовлетворяют соотношениям

$$\left. \begin{aligned} -\sigma + \delta''_1 + a_1(1 + \alpha''_1)(1 + \varphi_1)\mathcal{P}_1 &= 0; \\ \frac{a_{i-1}(1 + \alpha'_{i-1})}{(1 + \varphi_{i-1})\mathcal{P}_{i-1}} - \sigma + \delta'_{i-1} + b_i(1 + \beta''_i)(1 + \psi_i)\mathcal{R}_i &= 0, \\ \frac{b_i(1 + \beta'_i)}{(1 + \psi_i)\mathcal{R}_i} - \sigma + \delta''_i + a_i(1 + \alpha''_i)(1 + \varphi_i)\mathcal{P}_i &= 0, \\ \frac{a_N(1 + \alpha'_N)}{(1 + \psi_N)\mathcal{P}_N} - \sigma + \delta'_N &= 0. \end{aligned} \right\} i = 2, \dots, N; \quad (8.7)$$

Согласно теореме 3 из работы [4] возмущения $\delta'_i, \delta''_i, \alpha'_i, \alpha''_i, \beta'_i, \beta''_i$ подчиняются оценкам

$$\begin{aligned} |\delta'_i| &\leq \delta \mathcal{K}(\tilde{A}) + 16\varepsilon_1 \mathcal{K}(A), \quad |\delta''_i| \leq \delta \mathcal{K}(\tilde{A}) + 16\varepsilon_1 \mathcal{K}(A); \\ |\alpha'_i| &\leq 10\varepsilon_1, \quad |\alpha''_i| \leq 10\varepsilon_1, \quad |\beta'_i| \leq 10\varepsilon_1, \quad |\beta''_i| \leq 10\varepsilon_1.\end{aligned}$$

Заметим, что если учесть специфику матрицы S , а именно наличие у нее нулевой главной диагонали, то можно получить несколько лучшие оценки. Здесь же в целях простоты изложения использована непосредственно упомянутая теорема 3.

Для дальнейшего получим оценки величин δ'_i и δ''_i в терминах величины σ . Из неравенств (2.24) легко следует оценка

$$\mathcal{K}(\tilde{A}) \leq \sqrt{2} \|\tilde{A}\|. \quad (8.8)$$

Принимая теперь во внимание неравенства (8.4), имеем

$$|\delta'_i| \leq \sqrt{2} [\delta + 16\varepsilon_1(1 + 2\varepsilon_1)] \|\tilde{A}\|. \quad (8.9)$$

Из (8.2) и (8.8) нетрудно получить оценку $\|\tilde{A}\| \leq \sigma/(1 - \sqrt{2}\delta)$. Огрубляя с ее помощью оценку (8.9), запишем

$$|\delta'_i| \leq \frac{\sqrt{2} [\delta + 16\varepsilon_1(1 + 2\varepsilon_1)]}{1 - \sqrt{2}\delta} \sigma.$$

Наконец, учитывая замечание о возможном значении константы δ в оценке (8.2) ($\delta = 4\varepsilon_1$) и определение (8.6) величины δ , полученную оценку можно записать в виде

$$|\delta'_i| \leq 10\varepsilon_1 \sigma, \quad (8.10)$$

где в качестве значения μ принято $\mu = 31,5\varepsilon_1$. Оценка (8.10), очевидно, эквивалентна равенству $\delta'_i = \mu'_i\sigma$, которое имеет место при некотором μ'_i таком, что $|\mu'_i| \leq \mu$. Аналогичным образом можно получить равенство $\delta''_i = \mu''_i\sigma$, справедливое при некотором значении μ''_i , удовлетворяющем оценке $|\mu''_i| \leq \mu$.

Преобразуем теперь соотношения (8.7). Необходимые преобразования продемонстрируем на примере соотношения

$$\frac{b_i(1+\beta'_i)}{(1+\varphi_i)\mathcal{R}_i} - \sigma + \delta''_i + a_i(1+\alpha''_i)(1+\varphi_i)\mathcal{P}_i = 0.$$

Поскольку выражение $-\sigma + \delta''_i$ можно записать в виде $-\sigma + \delta''_i = -\sigma(1 - \mu''_i)$, это соотношение после деления обеих его частей на $1 - \mu''_i$ принимает вид

$$\frac{b_i(1+\beta'_i)}{(1+\varphi_i)(1-\mu''_i)\mathcal{R}_i} - \sigma + \frac{a_i(1+\alpha''_i)(1+\varphi_i)}{1-\mu''_i}\mathcal{P}_i = 0.$$

Вводя обозначения

$$\xi_i = \frac{1+\beta'_i}{(1+\varphi_i)(1-\mu''_i)} - 1, \quad \tilde{\xi}_i = \frac{(1+\alpha''_i)(1+\varphi_i)}{1-\mu''_i} - 1, \quad (8.11)$$

перепишем его в виде

$$\frac{b_i(1+\xi_i)}{\mathcal{R}_i} - \sigma + a_i(1+\tilde{\xi}_i)\mathcal{P}_i = 0. \quad (8.12)$$

Из формул (8.11) и оценок величин $\beta'_i, \alpha''_i, \varphi_i, \varphi'_i, \mu''_i$ легко получить оценки величин ξ_i и $\tilde{\xi}_i$:

$$|\xi_i| \leq \rho, \quad |\tilde{\xi}_i| \leq \rho, \quad (8.13)$$

справедливые при значении $\rho = 43\varepsilon_1$. Итак, соотношение (8.12) имеет место при некоторых значениях величин ξ_i и $\tilde{\xi}_i$, удовлетворяющих оценкам (8.13). Аналогично преобразуются остальные соотношения (8.7). Таким образом, мы показали, что числа $\mathcal{P}_1, \mathcal{R}_2, \mathcal{P}_2, \dots, \mathcal{R}_N, \mathcal{P}_N$ удовлетворяют соотношениям

$$\left. \begin{aligned} -\sigma + a_1(1+\tilde{\xi}_1)\mathcal{P}_1 &= 0; \\ \frac{a_{i-1}(1+\tilde{\xi}_{i-1})}{\mathcal{P}_{i-1}} - \sigma + b_i(1+\tilde{\xi}_i)\mathcal{R}_i &= 0, \\ \frac{b_i(1+\tilde{\xi}_i)}{\mathcal{R}_i} - \sigma + a_i(1+\tilde{\xi}_i)\mathcal{P}_i &= 0, \\ \frac{a_N(1+\tilde{\xi}_N)}{\mathcal{P}_N} - \sigma &= 0 \end{aligned} \right\} i = 2, \dots, N; \quad (8.14)$$

при некоторых $\tilde{\xi}_1, \tilde{\xi}_2, \tilde{\xi}_3, \tilde{\xi}_4$ таких, что $|\tilde{\xi}_1| \leq \rho, |\tilde{\xi}_2| \leq \rho, |\tilde{\xi}_3| \leq \rho, |\tilde{\xi}_4| \leq \rho$, где $\rho = 43\varepsilon_1$.

Определим теперь числа $c_i, s_i, \bar{c}_i, \bar{s}_i$ и $c'_{i_1} s'_{i_1}, \bar{c}'_{i_2} \bar{s}'_{i_2}$ следующим образом:

$$\begin{aligned} s_2 &= \frac{1}{\sqrt{1+(\mathcal{P}_1\mathcal{R}_2)^2}}, & c_2 &= \frac{\mathcal{P}_1\mathcal{R}_2}{\sqrt{1+(\mathcal{P}_1\mathcal{R}_2)^2}}, & c'_2 &= (c_2)_M, & s'_2 &= (s_2)_M; \\ \bar{s}_2 &= \frac{1}{\sqrt{1+(\mathcal{R}_2\mathcal{P}_2)^2}}, & \bar{c}_2 &= \frac{\mathcal{R}_2\mathcal{P}_2}{\sqrt{1+(\mathcal{R}_2\mathcal{P}_2)^2}}, & \bar{c}'_2 &= (\bar{c}_2)_M, & \bar{s}'_2 &= (\bar{s}_2)_M; \end{aligned}$$

$$\left. \begin{array}{l} s_i = \frac{1}{\sqrt{1 + (c'_{i-1}\varphi_{i-1}\varrho_i)^2}}, \quad c_i = \frac{c'_{i-1}\varphi_{i-1}\varrho_i}{\sqrt{1 + (c'_{i-1}\varphi_{i-1}\varrho_i)^2}}, \\ c'_i = (c_i)_M, \quad s'_i = (s_i)_M, \\ \bar{s}_i = \frac{1}{\sqrt{1 + (\bar{c}'_{i-1}\varrho_i\varphi_i)^2}}, \quad \bar{c}_i = \frac{\bar{c}'_{i-1}\varrho_i\varphi_i}{\sqrt{1 + (\bar{c}'_{i-1}\varrho_i\varphi_i)^2}}, \\ \bar{c}'_i = (\bar{c}_i)_M, \quad \bar{s}'_i = (\bar{s}_i)_M, \end{array} \right\} i = 3, \dots, N. \quad (8.15)$$

Формулы $c'_i = (c_i)_M$, $s'_i = (s_i)_M$, $\bar{c}'_i = (\bar{c}_i)_M$, $\bar{s}'_i = (\bar{s}_i)_M$ означают, что числа c'_i , s'_i , \bar{c}'_i , \bar{s}'_i являются результатами машинного вычисления величин c_i , s_i , \bar{c}_i , \bar{s}_i . При этом предполагается, что вычисления ведутся по специальным алгоритмам, подробно описанным и исследованным в § 4 главы 2. Анализ этих алгоритмов позволяет гарантировать выполнение равенств.

$$c'_i = c_i(1 + \tau_i), \quad s'_i = s_i(1 + \pi_i), \quad \bar{c}'_i = \bar{c}_i(1 + \bar{\tau}_i), \quad \bar{s}'_i = \bar{s}_i(1 + \bar{\pi}_i), \quad (8.16)$$

где величины τ_i , π_i , $\bar{\tau}_i$, $\bar{\pi}_i$ удовлетворяют оценкам

$$|\tau_i| \leq \tau, \quad |\pi_i| \leq \pi, \quad |\bar{\tau}_i| \leq \bar{\tau}, \quad |\bar{\pi}_i| \leq \bar{\pi}.$$

при $\tau = 7,001\epsilon_1$, $\pi = 5,001\epsilon_4$. Напомним также, что в указанных алгоритмах каждая из величин c_i , s_i , \bar{c}_i , \bar{s}_i определяется с помощью двух чисел: порядка и мантиссы, хранящихся в разных ячейках памяти машины. Значения мантисс этих величин, как обычно, заключены в промежутке от $1/\alpha$ до 1, где α — основание системы счисления, принятой в машине. В силу этого все числа c'_i , s'_i , \bar{c}'_i , \bar{s}'_i отличны от нуля.

В § 3 главы 2 показано, что из соотношений (8.14) и равенств (8.16) следуют соотношения

$$s_2 a_1 (1 + \hat{\alpha}_1) - s_2 \frac{\sigma^2}{a_1} + c_2 b_2 (1 + \check{\beta}_2) = 0;$$

$$\bar{s}_2 b_2 (1 + \hat{\beta}_2) - c_2 \frac{\bar{s}_2}{s_2} a_1 + \bar{c}_2 a_2 (1 + \check{\alpha}_2) = 0;$$

$$\left. \begin{array}{l} s_{i+1} c_i a_i (1 + \hat{\alpha}_i) - \bar{c}_i \frac{s_{i+1} s_i \dots s_2 \sigma^2}{\bar{s}_i \dots \bar{s}_2} \frac{a_1}{a_1} + c_{i+1} b_{i+1} (1 + \check{\beta}_{i+1}) = 0, \\ s_{i+1} \bar{c}_i b_{i+1} (1 + \hat{\beta}_{i+1}) - c_{i+1} \frac{s_{i+1} \dots \bar{s}_2}{s_{i+1} \dots s_2} a_1 + \bar{c}_{i+1} a_{i+1} (1 + \check{\alpha}_{i+1}) = 0, \end{array} \right\} i = 2, \dots, N-1;$$

$$c_N a_N (1 + \hat{\alpha}_N) - \bar{c}_N \frac{s_N \dots s_2 \sigma^2}{\bar{s}_N \dots \bar{s}_2} \frac{a_1}{a_1} = 0,$$

которые имеют место при некоторых $\hat{\alpha}_i$, $\check{\alpha}_i$, $\hat{\beta}_i$, $\check{\beta}_i$ таких, что

$$|\hat{\alpha}_i| \leq \varepsilon, \quad |\check{\alpha}_i| \leq \varepsilon, \quad |\hat{\beta}_i| \leq \varepsilon, \quad |\check{\beta}_i| \leq \varepsilon,$$

где в качестве значения константы ε можно взять $\varepsilon = 2(N\tau + \rho)$.

Дальнейшие рассмотрения непосредственно опираются на результаты § 6 главы 3. Определим матрицу

$$(\bar{A})_M = \begin{bmatrix} (\bar{a}_1)_M & (\bar{b}_2)_M & & & \\ & (\bar{a}'_2)_M & (\bar{b}'_3)_M & 0 & \\ & & \ddots & \ddots & \\ & & & (\bar{a}'_{N-2})_M & (\bar{b}'_{N-1})_M \\ 0 & & & & (\bar{a}'_{N-1})_M & 0 \\ & & & & & \sigma \end{bmatrix}.$$

элементы $(\bar{a}'_i)_m$, $(\bar{b}'_i)_m$, которой получаются в результате машинного вычисления по формулам

$$\bar{a}'_i = \frac{s'_{i+1} a_{i+1}}{s'_{i+1}}, \quad i = 1, \dots, N-1;$$

$$\bar{b}'_i = \frac{\bar{s}'_i b_{i+1}}{s'_{i+1}}, \quad i = 2, \dots, N-1.$$

В § 6 детально описаны алгоритмы, с помощью которых проводятся вычисления по этим формулам.

Из рассмотрений § 6 следует, что вычисленная таким образом матрица $(\bar{A})_m$ по норме мало отличается от некоторой матрицы, ортогонально эквивалентной матрице A , а именно из теоремы 4 следует оценка

$$\|(\bar{A})_m - \bar{C}AC^*\| \leq \{2\sqrt{2}(\sqrt{N} + 2)\varepsilon + (\sqrt{2}\tilde{\pi} + \varepsilon_1)[1 + 2\sqrt{2}(\sqrt{N} + 2)\varepsilon]\}\|A\|. \quad (8.17)$$

Напомним, что матрицы C и \bar{C} представляют собой произведения ортогональных матриц

$$C = C_{N+1}C_N \cdots C_3C_2, \quad \bar{C} = \bar{C}_N \cdots \bar{C}_3\bar{C}_2,$$

где C_i , \bar{C}_i ($i = 2, \dots, N$) — элементарные матрицы вращения, конструируемые с помощью параметров c_i , s_i ($c_i^2 + s_i^2 = 1$) и \bar{c}_i , \bar{s}_i ($\bar{c}_i^2 + \bar{s}_i^2 = 1$) соответственно, определенных формулами (8.15), а матрица C_{N+1} имеет вид

$$C_{N+1} = \text{diag}[1, 1, \dots, 1, \text{sign } a_1].$$

Характер константы ε в этой оценке можно представить себе с помощью формулы $\varepsilon = 2(N\tau + \rho)$, в которой значения величин τ и ρ равны соответственно $7,001\varepsilon_1$ и $43\varepsilon_1$. Значение константы $\tilde{\pi}$ определяется формулой (6.14) $\tilde{\pi} = \varepsilon_1(2 + \varepsilon_1)(1 + \hat{\pi}) + \hat{\pi}$. В свою очередь, величина $\hat{\pi}$ вычисляется по формуле (6.7): $\hat{\pi} = 2\pi/(1 - \pi)$, где $\pi = 5,001\varepsilon_1$. Пренебрегая в этих формулах членами второго порядка малости по ε_1 , находим, что в качестве значения $\tilde{\pi}$ можно взять $\tilde{\pi} = 12,002\varepsilon_1$.

Итак, оценка (8.17) имеет место при следующих значениях, входящих в нее констант ε и $\tilde{\pi}$: $\varepsilon = 2(7,001N + 43)\varepsilon_1$, $\tilde{\pi} = 12,002\varepsilon_1$.

В заключение оценим отличие вычисленной матрицы $(\bar{A})_m$ от $\bar{C}AC^*$, ортогонально эквивалентной исходной матрице A . Из определения (8.3) элементов a_i , b_i матрицы A и неравенства (8.8) с помощью неоднократно упоминавшейся оценки (2.30) из [4] легко получить оценку $\|A - \bar{A}\| \leq 2\sqrt{2}\varepsilon_1\|\bar{A}\|$. Из этой оценки следует, что $\|A\| \leq (1 + 2\sqrt{2}\varepsilon_1)\|\bar{A}\|$. Используя (8.17) и полученные оценки, находим

$$\begin{aligned} \|(\bar{A})_m - \bar{C}AC^*\| &\leq \|(\bar{A})_m - \bar{C}AC^*\| + \|A - \bar{A}\| \leq \\ &\leq \{2\sqrt{2}(\sqrt{N} + 2)\varepsilon + (\sqrt{2}\tilde{\pi} + \varepsilon_1)[1 + 2\sqrt{2}(\sqrt{N} + 2)\varepsilon]\}(1 + 2\sqrt{2}\varepsilon_1)\|\bar{A}\| + \\ &\quad + 2\sqrt{2}\varepsilon_1\|\bar{A}\|. \end{aligned}$$

Пренебрегая в этой оценке членами второго порядка, окончательно получим

$$\|(\bar{A})_m - \bar{C}AC^*\| \leq \{2\sqrt{2}[(\sqrt{N} + 2)\varepsilon + \varepsilon_1] + \sqrt{2}\tilde{\pi} + \varepsilon_1\}\|\bar{A}\|.$$

Подчеркнем, что к результатам работы описанного алгоритма кроме матрицы $(\bar{A})_m$ относятся значение выражения, стоящего в правой части полученной оценки, а также числа c_i , s_i , \bar{c}_i , \bar{s}_i и значения констант τ и $\tilde{\pi}$, оценивающих отличие этих чисел от точных значений параметров c_i , s_i , \bar{c}_i , \bar{s}_i , с помощью которых конструируются ортогональные матрицы C и \bar{C} .

ЛИТЕРАТУРА

1. Rutishauser H. On Jacobi rotation patterns.— Proceedings A. M. S. Symposia in Applied Mathematics, 1963, v. 15, p. 219—239.
2. Golub G., Kahan W. Calculating the singular values and pseudoinverse of a matrix.— J. SIAM Numer. Anal. Ser. B, 1965, v. 2, N 2, p. 205—224.
3. Уилкинсон Дж. Х. Алгебраическая проблема собственных значений.— М.: Наука, 1970.— 564 с.
4. Годунов С. К., Костин В. И., Митченко А. Д. Вычисление собственного вектора симметрической трехдиагональной матрицы.— Сиб. мат. журн., 1985, т. XXVI, № 5.
5. Годунов С. К. Решение систем линейных уравнений.— Новосибирск: Наука. Сиб. отд-ние, 1980.— 177 с.

ОБОБЩЕННЫЕ РЕШЕНИЯ И РЕГУЛЯРИЗАЦИЯ СИСТЕМ НЕРАВЕНСТВ

B. A. БУЛАВСКИЙ

Рассматривается подход к решению несовместных систем неравенств, включающий построение их обобщенных решений и регуляризацию. Этот подход представляет естественное обобщение развитого А. Н. Тихоновым [1, 2] способа регуляризации систем линейных алгебраических уравнений и обладает следующими свойствами. Во-первых, он не выводит за рамки линейных задач; во-вторых, позволяет построить класс релаксационных алгоритмов, сходящихся независимо от совместности решаемой задачи, который включает в себя, например, методы последовательного проектирования [3—6] и градиентный метод. Наконец, предлагаемый подход применим для коррекции несовместных (включая двойственную несовместность) задач линейного программирования, которым в последнее время уделяется внимание [7]. Материал статьи является развитием ранее опубликованных результатов автора [5, 8, 9] и основан на использовании аппарата задач с условиями дополнительности. Укажем также, что вводимые понятия и значительная часть результатов могут быть перенесены на нелинейный случай.

В заключение введение оговорим некоторые обозначения и терминологию. Все изложение ведется для конечномерных пространств, и векторы трактуются как векторы-столбцы, так что при умножении на матрицу они стоят справа. Транспонирование обозначается штрихом в качестве верхнего индекса, например, если x — вектор-столбец, то x' — вектор-строка. Для положительной и отрицательной частей числа α приняты обозначения α_+ и α_- . Эти же обозначения используются для векторов: символ x_+ обозначает вектор, составленный из положительных частей компонент вектора x . Покомпонентно применяется и знак неравенства между векторами. Наконец, термины «положительно определенная матрица Q » или «положительно полуопределенная матрица Q » не предполагают симметрию матрицы Q , а лишь означают, что выполняются неравенства $z'Qz > 0$ при $z \neq 0$ или $z'Qz \geq 0$ при всех z .

§ 1. ЗАДАЧИ С УСЛОВИЯМИ ДОПОЛНИТЕЛЬНОСТИ

Линейной задачей с условиями дополнительности называют задачу

$$Mz + q \geq 0, z \geq 0, z'(Mz + q) = 0, \quad (1.1)$$

где $z \in R^n$ — искомый вектор, а матрица $M \in R^{n \times n}$ и свободный член $q \in R^n$ заданы. В последние годы задачам такого типа посвящено много исследований. Обзор результатов и библиографию можно найти, напри-