Rigorous Machine Bounds for the Eigensystem of a General Complex Matrix*

By J. M. Varah

Introduction. We are concerned here with giving rigorous error bounds for the eigensystem of a general complex $n \times n$ matrix A, given an approximate eigensystem such as is furnished by [2]. In Section 1, we outline the technique in general terms and show that the bounds can be found in terms of computed quantities if $||E||_{\infty} = ||I - XY||_{\infty} < 1$, where X is the matrix of approximate eigenvectors and Y is an approximate inverse for X. Then in Sections 2, 3, and 4 we give the specific roundoff error bounds for these general error terms, which include all the rounding errors made during the computation. An Algol program using the method is given in the microfiche section, and the results for the matrix example given in [2] are presented in Section 5, using the results of [2] as the initial approximation.

1. Theoretical Bounds for the Eigensystem. We assume we have a complex matrix A° of order *n* represented for our calculation by the matrix *A*, with

$$A = A^{\circ} + \Delta, \qquad |\Delta_{ij}| \leq \delta \cdot \max_{1 \leq i, j \leq n} |A_{ij}|$$

and δ specified. We further assume a complete approximate eigensystem has been given for A, that is, a diagonal complex matrix Λ of eigenvalue approximations and a complex matrix X whose columns are approximations to the corresponding column eigenvectors of A, normalized in some way so that all components are less than or equal to 1.0 in modulus. We wish to give rigorous bounds for the true eigensystem of A° . In this section we outline the technique used, which follows Wilkinson [4, Chapter 9]. In Section 4, we will give the actual bounds used, which include bounds on the roundoff errors committed in the calculations.

We first perform a similarity transformation on A with X, assuming X can be inverted, giving

$$X^{-1}A^{\circ}X = \Lambda + P + Q,$$

where P is known exactly, and a bound for the modulus of each element of Q is known. To perform this similarity, let Y be a calculated approximate inverse for X, and define the following matrices:

$$F = AX - X\Lambda + Q_1$$
, $P = YF + Q_2$, $E = I - XY$,

where the elements of Q_1 and Q_2 are small. In the actual computation, we form

$$F = fl_2(AX - X\Lambda)$$
 and $P = fl_2(YF)$

using double-precision accumulation of inner products. Then we use Q_1 and Q_2 to

Received August 25, 1967. Revised January 24, 1968.

^{*} This work was supported in part by the Office of Naval Research under Contract No. Nonr-225(37) (NR-044-211).

denote the errors made in these machine computations. We will show that we can obtain bounds for $|Q_{ij}|$ assuming only that

$$||E||_{\infty} = \max_{1 \leq i \leq n} \left(\sum_{j=1}^{n} |E_{ij}| \right) < 1$$

Using these matrices above, we have

$$P = YF + Q_2$$

= $X^{-1}F + (Y - X^{-1})F + Q_2$
= $X^{-1}AX - \Lambda + X^{-1}Q_1 + (Y - X^{-1})F + Q_2$,

so that

 $X^{-1}A^{\circ}X = \Lambda + P + Q$, with $Q = -[X^{-1}\Delta X + X^{-1}Q_1 + (Y - X^{-1})F + Q_2]$. To bound Q, we note that if XY is invertible,

$$Y - X^{-1} = Y(I - (XY)^{-1})$$

= Y(I - (I - E)^{-1})
= -YE(I - E)^{-1}.

Also

$$X^{-1} = (Y^{-1} - EY^{-1})^{-1}$$

= Y(I - E)^{-1},

so that

$$Q = -[Y(I - E)^{-1}(\Delta X + Q_1) - YE(I - E)^{-1}F + Q_2].$$

Now, for $i = 1, 2, \dots, n$, let

$$\begin{split} \alpha_{i} &= \sum_{j=1}^{n} |Y_{ij}|, \quad \beta_{i} = \sum_{j=1}^{n} |X_{ji}|, \\ \gamma_{i} &= \max_{1 \le j \le n} |F_{ji}|, \quad \sigma_{i} = \max_{1 \le j \le n} |(Q_{1})_{ji}|, \\ \tau_{i} &= \sum_{j=1}^{n} |(YE)_{ij}|; \end{split}$$

and let $A_{\max} = \max_{1 \le i, j \le n} |A_{ij}|$. Then we have

$$\left| (\Delta X + Q_1)_{ij} \right| \leq \delta \cdot A_{\max} \cdot \beta_j + \sigma_j ,$$

so that

$$\left| [(I-E)^{-1}(\Delta X+Q_1)]_{ij} \right| \leq \left\| (I-E)^{-1} \right\|_{\infty} \cdot (\delta \cdot A_{\max} \cdot \beta_j + \sigma_j) .$$

Also,

$$\left| [YE(I-E)^{-1}F]_{ij} \right| \leq \left\| (I-E)^{-1} \right\|_{\infty} \cdot \tau_i \cdot \gamma_j \right|$$

Thus, assuming $||E||_{\infty} < 1$,

$$|Q_{ij}| \leq \frac{\alpha_i (\delta \cdot A_{\max} \cdot \beta_j + \sigma_j) + \tau_i \cdot \gamma_j}{1 - \|E\|_{\infty}} + |(Q_2)_{ij}|.$$

We bound the eigenvalues of A° using Gerschgorin's theorems, which we now

794

state for reference. For a given matrix B, define the Gerschgorin disks

$$G_i = \left\{ \lambda : |\lambda - b_{ii}| \leq \sum_{j \neq i} |b_{ij}| \right\}, \quad i = 1, \cdots, n.$$

Then the first theorem of Gerschgorin states that all the eigenvalues of B are contained in the union of the n disks $\{G_i\}$. The second theorem states that if k of the disks are isolated from the others, then there are precisely k eigenvalues of B in the union of those k disks. For proofs of these theorems, see Marcus and Minc [1, p. 146].

We apply these theorems to the matrix $B = X^{-1}A^{\circ}X = \Lambda + P + Q$. If the elements of P and Q are small, and no other eigenvalue approximation Λ_{jj} is too close to Λ_{ii} , the *i*th disk will be isolated from the others, so that there is only one eigenvalue of B in the disk. For the *i*th disk to be isolated, we must have for all $k \neq i$,

$$|b_{kk} - b_{ii}| > \sum_{j \neq i} |b_{ij}| + \sum_{j \neq k} |b_{kj}|$$
.

For the *i*th eigenvalue, we can usually obtain a better bound by applying Gerschgorin's theorem to the matrix B modified by multiplying the *i*th row by β^{-m} and the *i*th column by β^{m} , where β is the number base of the machine used and m is a nonnegative integer, chosen as large as possible under the restriction that the *i*th Gerschgorin disk of this modified B matrix remain isolated. For such an m, the disk is defined by

$$|\lambda - (\Lambda_{ii} + P_{ii} + Q_{ii})| \leq \beta^{-m} \cdot \sum_{j \neq i} |P_{ij} + Q_{ij}|,$$

so that one eigenvalue λ_i of A° satisfies the inequality

$$|\lambda_{i} - (\Lambda_{ii} + P_{ii})| \leq r_{i} = |Q_{ii}| + \beta^{-m} \sum_{j \neq i} (|P_{ij}| + |Q_{ij}|).$$

The *i*th disk will be isolated if for all $k \neq i$,

$$\begin{aligned} |(\Lambda_{ii} + P_{ii} + Q_{ii}) - (\Lambda_{kk} + P_{kk} + Q_{kk})| &> \beta^{-m} \sum_{j \neq i} |P_{ij} + Q_{ij}| \\ &+ \beta^{m} |P_{ki} + Q_{ki}| + \sum_{j \neq k, i} |P_{kj} + Q_{kj}| , \end{aligned}$$

which holds if

$$|(\Lambda_{ii} + P_{ii}) - (\Lambda_{kk} + P_{kk})| > |Q_{ii}| + |Q_{kk}| + \beta^{-m} \sum_{j \neq i} (|P_{ij}| + |Q_{ij}|) + \beta^{m} (|P_{ki}| + |Q_{ki}|) + \sum_{j \neq k, i} (|P_{kj}| + |Q_{kj}|) .$$

If such a bound can be obtained for the *i*th eigenvalue, we can also bound the corresponding eigenvector of A° . We first bound the corresponding eigenvector u of $B = X^{-1}A^{\circ}X$. Since B is nearly diagonal, u is close to the unit vector e_i . Hence we can set $u_i = 1.0$ and bound the other components of u by using the relation $Bu = \lambda_i u$ and the bound for λ_i obtained above. In fact, the kth equation of $Bu = \lambda_i u$ gives

(1)
$$[(\Lambda_{ii} + P_{ii}) - (\Lambda_{kk} + P_{kk})]u_k = P_{ki} + Q_{ki} + (Q_{kk} + \theta_1 r_i)u_k + \sum_{j \neq k, i} (P_{kj} + Q_{kj})u_j$$

where $|\theta_1| \leq 1$. We set $\rho_k = (\Lambda_{ii} + P_{ii}) - (\Lambda_{kk} + P_{kk})$ for convenience. We first obtain a crude bound for u_k by assuming $|u_j| \leq 1$ for $j \neq i$. This gives

$$|u_k| \leq {s_k}^{(1)} = \left[|Q_{kk}| + r_i + \sum_{j \neq k} (|P_{kj}| + |Q_{kj}|) \right] / |\rho_k|, \quad k = 1, \dots, n$$

Now we use this bound for all $|u_k|$ in (1), obtaining the more precise bounds

$$\left| u_{k} - \frac{P_{ki}}{\rho_{k}} \right| \leq \left[|Q_{ki}| + (|Q_{kk}| + r_{i})s_{k}^{(1)} + \sum_{j \neq k, i} (|P_{kj}| + |Q_{kj}|)s_{j}^{(1)} \right] / |\rho_{k}|,$$

or, using $U_k = P_{ki}/\rho_k$ as the approximation to u_k ,

$$|u_k - U_k| \leq s_k^{(2)}, \quad k \neq i \quad (\text{and } u_i = 1.0).$$

This bounds the eigenvector u of $B = X^{-1}A^{\circ}X$. To bound the corresponding eigenvector v = Xu of A° , we must transform the estimate and bound for u by multiplying by X. Thus $|v_k - (XU)_k| \leq \sum_{j=1}^{n} |X_{kj}| s_j^{(2)}$, $k = 1, \dots, n$. Finally, we can normalize the estimate so that its largest component in modulus is 1.0, obtaining

$$\left| v_k - \frac{(XU)_k}{(XU)_{\max}} \right| \leq \frac{\sum_{j=1}^n |X_{kj}| s_j^{(2)}}{|(XU)_{\max}|}, \quad k = 1, \dots, n.$$

2. Basic Roundoff Errors. To bound the errors in machine calculation, we let $\eta_1 = 2 \cdot (1.06) \cdot \beta^{1-t}$ as in [3, p. 19] where β is the floating-point number base of the machine and t, the number of base β digits carried in each single-precision floating-point number. Thus η_1 is an upper bound for the relative rounding error committed in each basic real single-precision floating-point operation. That is,

$$\begin{aligned} |fl(x+y) - (x+y)| &\leq \eta_1(|x|+|y|) ,\\ |fl(x\cdot y) - (x\cdot y)| &\leq \eta_1(|x\cdot y|) ,\\ |fl(x/y) - (x/y)| &\leq \eta_1(|x/y|) , \quad y \neq 0 . \end{aligned}$$

We also assume the square root routine on the machine gives answers of comparable accuracy, i.e.,

$$|fl(\sqrt{x}) - \sqrt{x}| \leq \eta_1 \sqrt{x} \, .$$

The factor 1.06 makes for easier accumulation of errors. η_1 could probably be taken smaller on most machines by a factor between 2 and 4, and the user may wish to adjust its value in the program. For further information, see Wilkinson [3].

To bound the errors in complex operations, let $z_1 = x_1 + iy_1$, $z_2 = x_2 + iy_2$. Then the following results are easily obtained.

1.
$$|fl(z_1 + z_2) - (z_1 + z_2)| \le \eta_1(|z_1| + |z_2|).$$

2. $|fl(z_1 \cdot z_2) - (z_1 \cdot z_2)| \le (2 \sqrt{2})\eta_1|z_1 \cdot z_2|.$
3.(a) If

$$fl(z_1/z_2) \equiv fl\left(\frac{x_1x_2 + y_1y_2}{x_2^2 + y_2^2}\right) + i \cdot fl\left(\frac{x_2y_1 - x_1y_2}{x_2^2 + y_2^2}\right),$$

then $|fl(z_1/z_2) - (z_1/z_2)| \leq (5\sqrt{2})\eta_1 |z_1/z_2|.$

(b) Suppose $|x_2| \leq |y_2|$ and set $r = x_2/y_2$, $d = y_2 + rx_2$. If

$$fl(z_1/z_2) \equiv fl\left(\frac{x_1r+y_1}{d}\right) + i \cdot fl\left(\frac{y_1r-x_1}{d}\right),$$

then $|fl(z_1/z_2) - (z_1/z_2)| \leq (6 \sqrt{2})\eta_1 |z_1/z_2|.$

4.(a) If $fl(|z|) \equiv fl((x^2 + y^2)^{1/2})$, then $|fl(|z|) - |z|| \leq 2\eta_1 |z|$.

(b) If $|x| \leq |y|$ and $fl(|z|) \equiv fl(|y| \cdot (1 + |x/y|^2)^{1/2})$, then $|fl(|z|) - |z|| \leq 4\eta_1 |z|$. 5. In a real single-precision floating-point inner product,

$$\left| fl\left(\sum_{1}^{n} x_{i}y_{i}\right) - \left(\sum_{1}^{n} x_{i}y_{i}\right) \right| \leq n\eta_{1}\left(\sum_{1}^{n} |x_{i}y_{i}|\right) \equiv \epsilon_{1}\left(\sum_{1}^{n} |x_{i}y_{i}|\right).$$

6. In a complex floating-point inner product of single-precision factors accumulated in double precision and rounded to single precision,

$$\left| fl_2\left(\sum_{i=1}^n z_i w_i\right) - \left(\sum_{i=1}^n z_i w_i\right) \right| \leq \eta_1 \left| \sum_{i=1}^n z_i w_i \right| + \epsilon_2 \left(\sum_{i=1}^n |z_i w_i|\right)$$

where $\epsilon_2 = (n + 1/2)\eta_2$ and $\eta_2 = 2 \cdot (1.06)\beta^{1-2t}$ bounds the basic double-precision rounding error.

3. Machine Bounds for the Error Matrix Q. To bound the elements of Q, we need bounds for the quantities used to bound Q in Section 1. We use barred symbols to denote the machine bounds.

1. $\alpha_i \leq \bar{\alpha}_i = fl(\sum_{j=1}^n |Y_{ij}|) \cdot (1 + \epsilon_1 + 3\eta_1), \ i = 1, \dots, n.$ 2. $\beta_i \leq \bar{\beta}_i = fl(\sum_{j=1}^n |X_{ji}|) \cdot (1 + \epsilon_1 + 3\eta_1), \ i = 1, \dots, n.$ 3. $\gamma_i \leq \bar{\gamma}_i = \max_{1 \leq j \leq n} (fl(|F_{ji}|)) \cdot (1 + 5\eta_1), \ i = 1, \dots, n.$ 4. $A_{\max} \leq \bar{a} = \max_{1 \leq i, j \leq n} (fl(|A_{ij}|)) \cdot (1 + 5\eta_1).$ 5. To bound $||E||_{\infty}$, we have for $j \neq i$,

$$fl(E_{ij}) = -fl_2\left(\sum_{k=1}^n X_{ik}Y_{kj}\right)$$
$$= E_{ij}(1+\theta\eta_1) + \theta\epsilon_2\left(\sum_{k=1}^n |X_{ik}||Y_{kj}|\right)$$

and

$$fl(E_{ii}) = fl_2 \left(1 - \sum_{k=1}^n X_{ik} Y_{ki} \right)$$

= $E_{ii} (1 + \theta \eta_1) + \theta(\epsilon_2 + \eta_2) \left(\sum_{k=1}^n |X_{ik}| |Y_{ki}| \right) + \theta \eta_2.$

Here θ denotes a generic multiplier less than or equal to 1 in modulus. Thus if we set

$$\overline{\alpha} = fl\left(\sum_{i=1}^{n} \overline{\alpha}_{i}\right)(1+\epsilon_{1}),$$

we have, since $|X_{ik}| \leq 1$ for all *i* and *k*,

$$\|E\|_{\infty} \leq \overline{e} = \left(\max_{1 \leq i \leq n} \left[fl\left(\sum_{j=1}^{n} |E_{ij}|\right) \right] (1+\epsilon_1+4\eta_1) + (\epsilon_2+\eta_2)\overline{\alpha}+\eta_2 \right) (1+5\eta_1).$$

6. We also need a bound for $\tau_i = \sum_{j=1}^n |(YE)_{ij}|$. Now

$$fl(YE)_{ij} = fl_2 \left(\sum_{k=1}^n Y_{ik} \cdot fl(E_{kj}) \right)$$
$$= \left(\sum_{k=1}^n Y_{ik} E_{kj} \right) (1 + \theta \eta_1) + \sum_{k=1}^n Y_{ik} (E_{kj} - fl(E_{kj}))$$
$$+ \theta \epsilon_2 \cdot \sum_{k=1}^n \left(|Y_{ik}| \cdot |fl(E_{kj})| \right).$$

Thus

$$\tau_i \leq \bar{\tau}_i = \left[\left(\sum_{j=1}^n |fl_2(YE)_{ij}| \right) (1 + \epsilon_1 + 4\eta_1) + \bar{\alpha}_i (\epsilon_2 + \eta_2) \bar{\alpha} + (\eta_1 + \epsilon_2) \bar{e} + \eta_2 \right] (1 + 6\eta_1) .$$

7. To bound σ_j , recall

$$fl(F_{ij}) = fl_2 \left(\left(\sum_{k=1}^n A_{ik} X_{kj} \right) - X_{ij} \Lambda_{jj} \right)$$

= $(AX - X\Lambda)_{ij} + \theta \eta_1 |F_{ij}| + \theta(\epsilon_2 + \eta_2) \left(\sum_{k=1}^n |A_{ik|} X_{kj}| \right)$
+ $\theta((2\sqrt{2})\eta_2) |X_{ij}| |\Lambda_{jj}|$
= $(AX - X\Lambda)_{ij} + (Q_1)_{ij}$,

and thus

$$\sigma_j \leq \overline{\sigma}_j = (\epsilon_2 + \eta_2) \overline{a} \overline{\beta}_j + 3\eta_2 |\Lambda_{jj}| + \eta_1 \overline{\gamma}_j.$$

8. Finally to bound $|(Q_2)_{ij}|$, we have

,

$$fl(P_{ij}) = fl_2 \left(\sum_{k=1}^{n} Y_{ik} F_{kj} \right) \\ = (YF)_{ij} + \theta \eta_1 |P_{ij}| + \theta \epsilon_2 \left(\sum_{k=1}^{n} |Y_{ik}| |F_{kj}| \right),$$

so that

$$|(Q_2)_{ij}| \leq \eta_1 |P_{ij}| + \epsilon_2 \overline{\alpha}_i \overline{\gamma}_j.$$

To give a rigorous machine bound for Q, we have to account for the errors made in computing the above bounds as well. Thus, for example,

$$|P_{ij}| \leq fl(|P_{ij}|) \cdot (1 + 4\eta_1)$$
.

Finally, we obtain

$$|Q_{ij}| \leq \overline{Q}_{ij} = fl \left\{ \eta_1 |P_{ij}| + \epsilon_2 \overline{\alpha}_i \overline{\gamma}_j + \frac{\overline{\alpha}_i (\delta \overline{a} \overline{\beta}_j + \overline{\sigma}_j) + \overline{\tau}_i \overline{\gamma}_j}{1 - \overline{e} - \eta_1 (1 + \overline{e})} \right\} (1 + 12\eta_1)$$

For this to exist, we must have the denominator positive. Thus corresponding to the theoretical condition $||E||_{\infty} < 1$ of Section 1, we have the machine condition $\bar{e} < (1 - \eta_1)/(1 + \eta_1)$.

EIGENSYSTEM BOUNDS: FOURTH ITERATION

ABS(LAMBDA	[10]=(8.1	2276,59240	,40516,95387,4180-02))	≤ 7,08773666250-13
EIGENVECTO ABS(X[1,1 ABS(X[2,1 ABS(X[3,1 ABS(X[3,1 ABS(X[4,1 ABS(X[4,1 ABS(X[6,1 ABS(X[7,1 ABS(X[7,1 ABS(X[9,1 ABS(X[0,1 ABS(X[10,1 ABS(X[12,1	R BOUNDS A 0] - (-9.8 0] - (1.1 0] - (7.3 0] - (-1.4 0] - (1.5 0] - (-3.3 0] - (-9.6 0] - (1.3 0] - (-9.1 0] - (3.8 0] - (-9.1 0] - (1.0	RE: 0976,06019 0958,96145 3905,41289 4903,69720 5770,04376 4790,51622 3703,80170 3735,07467 9476,1761 41457,47745 8772,34075 0000,00000	,70599,72072,0420-10)) ,28815,97919,5130-08)) ,82543,43573,1710-09)) ,68494,60904,9790-06)) ,90727,26666,8130-05)) ,02285,25975,7050-05)) ,00350,37346,5620-04)) ,49867,89746,5830-02)) ,76852,50508,6800-02)) ,22297,68340,1980-01)) ,95947,12999,6050-01)) ,00000,00000,0000+00))	<pre>≤ 1.1068600432e=14 ≤ 1.0313422383e=14 ≤ 5.2100951406e=15 ≤ 6.8429513045e=15 ≤ 2.7379035124e=14 ≤ 1.9301862161e=13 ≤ 1.2611504385e=12 ≤ 6.4441614122e=12 ≤ 2.3154536600e=11 ≤ 5.1921389637e=11 ≤ 5.4842010918e=11</pre>
ABS(LAMBDA	[11]-(1.4	3646,51976	,92204,86973,645 0-01))	≤ 1,54164296650-13
EIGENVECTO ABS(X[1,1 ABS(X[2,1 ABS(X[3,1 ABS(X[3,1 ABS(X[3,1 ABS(X[4,1 ABS(X[5,1 ABS(X[6,1 ABS(X[6,1 ABS(X[9,1 ABS(X[10,1 ABS(X[12,1	R BOUNDS A 1] - (6.2 1] - (-3.7 1] - (-2.5 1] - (4.0 1] - (-8.4 1] - (-1.7 1] - (1.4 1] - (-1.3 1] - (2.9 1] - (2.9 1] - (2.6 5 1] - (1.0	RE : 0702,69189 0033,96340 4717,68428 9451,15813 5053,59816 7654,47289 7089,38187 4829,72206 1604,61433 4847,38166 6353,48023 0000,00000	<pre>,10568,26641,5050-09)) ,09112,79938,7790-08)) ,92117,74995,8880-07)) ,68153,62233,4870-06)) ,8153,6520,5520-06)) ,81438,28545,1170-04)) ,69934,90432,8240-03)) ,76638,38936,0380-03)) ,37340,37369,1660-02)) ,70738,33609,4350-01)) ,07795,40335,0460-01)) ,00000,00000,0000000000000000000000</pre>	<pre>≤ 1.7434299658@=14 ≤ 1.6308351611@=14 ≤ 6.7782059504@=15 ≤ 9.3479387732@=15 ≤ 8.2431766586@=15 ≤ 1.3763628955@=14 ≤ 7.4496003126@=14 ≤ 5.5882308860@=13 ≤ 3.0773949604@=12 ≤ 1.1497666843@=11 ≤ 2.6381681437@=11 ≤ 2.8151569027@=11</pre>
ABS(LAMBDA	[12]=(2.8	4749,72055	,84781,88282,617 0-01))	≤ 1,38734395220-14
EIGENVECTO ABS(X[1,1 ABS(X[2,1 ABS(X[2,1 ABS(X[3,1 ABS(X[3,1 ABS(X[5,1 ABS(X[5,1 ABS(X[6,1 ABS(X[6,1 ABS(X[9,1 ABS(X[1,1] ABS(X[1,1]	R BOUNDS A 2] - (-1.1 2] - (2.9 2] - (3.8 2] - (-2.0 2] - (-7.0 2] - (-7.0 2] - (-1.2 2] - (-1.5 2] - (4.0 2] - (1.1 2] - (-1.1 2] - (-1.0 2] - (-1.0	RE: 9061,05278 9064,21369 4817,17071 1687,57372 9670,09125 4897,06845 6492,74040 3780,95663 8486,82483 3416,62084 5250,2/944	<pre>,84764,66934,149e-07)) ,10219,92917,680e-07)) ,07064,87942,594e-06)) ,24147,30515,215e-05)) ,93027,69031,150e-05)) ,14587,61341,644e-04)) ,69843,68850,891e-04)) ,08615,36267,732e-02)) ,41952,45413,714e-02)) ,13488,33735,374e-01)) ,15221,20985,641e-01))</pre>	<pre> 4 6.03021981460-15 5.64086595480-15 2 2.37149954760-15 3 2.38238313000-15 4 2.86864505230-15 5 2.68463512930-14 1.78391958520-13 1.00323886670-12 3.86573279990-12 9.00723779610-12 9.068012599880-12 9.48012599880-12 </pre>

4. Machine Bounds for the Eigenvalues and Eigenvectors. We first form

$$\overline{\lambda}_i = fll(\Lambda_{ii} + P_{ii}), \qquad i = 1, \dots, n$$

as improved estimates for the eigenvalues, where fll denotes double-precision add and store. We include roundoff here by adding $\eta_2(|\Lambda_{ii}| + |P_{ii}|)$ to \overline{Q}_{ii} . We also form

$$t_i = fl\left(\sum_{k\neq i} \left(|P_{ik}| + \overline{Q}_{ik}\right)\right) \cdot \left(1 + \epsilon_1 + 5\eta_1\right), \qquad i = 1, \dots, n,$$

bounding the off-diagonal row sums of (P + Q).

To bound the error in $\overline{\lambda}_i$, we use the Gerschgorin bounds obtained in Section 1. We form for $j \neq i$,

$$\begin{split} \overline{\rho}_j &= f l_2 (\overline{\lambda}_i - \overline{\lambda}_j) , \\ \mu_j &= f l(|\overline{\rho}_j|) , \end{split}$$

and

$$\mathbf{v}_j = fl((\eta_1 + \eta_2)(|\overline{\lambda}_i| + |\overline{\lambda}_j|)) \cdot (1 + 7\eta_1)$$

so that $|\bar{\rho}_j - (\bar{\lambda}_i - \bar{\lambda}_j)| \leq \nu_j$. Then we can say rigorously that

$$|\lambda_i - \bar{\lambda}_i| \leq \bar{r}_i = fl(Q_{ii} + \beta^{-m}t_i) \cdot (1 + \eta_1)$$

where m is the largest nonnegative integer such that the *i*th Gerschgorin disk is isolated, i.e. so that for all $j \neq i$,

$$\mu_j > fl(4\eta_1\mu_j + \nu_j + \bar{r}_i + Q_{jj} + t_j + \beta^m(|P_{ji}| + \bar{Q}_{ji})) \cdot (1 + 7\eta_1)$$

To choose m initially, note that the largest term in the above expression is usually the last, so we pick the largest m such that

$$\beta^m(|P_{ji}| + \overline{Q}_{ji}) < \mu_j, \qquad j \neq i.$$

Then we test each of the above more stringent requirements, decreasing m until they all hold or until m < 0. If the latter is true, we conclude that the *i*th eigenvalue cannot be isolated. Otherwise, we proceed to bound the corresponding eigenvector.

We first bound the eigenvector u of $B = X^{-1}A^{\circ}X$. Corresponding to identity (1) of Section 1, we have in terms of machine quantities, for $k \neq i$,

(2)
$$\overline{\rho}_k u_k = P_{ki} + \theta \overline{Q}_{ki} + \theta (\overline{Q}_{kk} + \overline{r}_i + \nu_k) u_k + \sum_{j \neq k, i} (P_{kj} + \theta Q_{kj}) u_j.$$

Again we first obtain a crude bound for $|u_k|$ by taking moduli and replacing $|u_j|$, $j = 1, \dots, n$ on the right-hand side by the upper bound 1.0, obtaining for $k \neq i$,

$$|u_k| \leq \bar{s}_k^{(1)} = fl\left(\frac{\bar{Q}_{kk} + \bar{r}_i + \nu_k + t_k}{\mu_k}\right) \cdot (1 + 8\eta_1) .$$

Now we use this bound for $|u_j|$ $(j \neq i)$ in the right-hand side of (2), obtaining the estimate $\overline{u}_k = fl(P_k_i/\overline{\rho}_k)$ for u_k and the bound, for $k \neq i$,

$$\begin{aligned} |u_{k} - \overline{\mu}_{k}| &\leq \bar{s}_{k}^{(2)} \\ &= fl \bigg\{ \frac{\overline{Q}_{ki} + \bar{s}_{k}^{(1)} (\overline{Q}_{kk} + \overline{r}_{i} + \nu_{k}) + \sum_{j \neq k, i} (|P_{kj}| + Q_{kj}) \overline{s}_{j}^{(1)}}{\mu_{k}} + 9\eta_{1} |\overline{u}_{k}| \bigg\} \\ &\cdot (1 + \epsilon_{1} + 11\eta_{1}) \,. \end{aligned}$$

To bound the corresponding eigenvector v = Xu of A° , we multiply through by X. Thus our estimate for v is $\overline{v} = X\overline{u}$, which we form and store in double precision. So

$$\bar{v}_j = fll\left(\sum_{k=1}^n X_{jk}\overline{u}_k\right)$$

and $|\bar{v}_j - (X\bar{u})_j| \leq \epsilon_2 \xi$, where

800

$$\xi = fl\left(\sum_{k=1}^{n} |\overline{u}_k|\right) \cdot (1 + \epsilon_1 + 3\eta_1) .$$

Also, from the above bound for $|u_k - \overline{u}_k|$, we have

$$|v_{j} - (X\overline{u})_{j}| \leq \sum_{k=1}^{n} |X_{jk}| \bar{s}_{k}^{(2)}$$

where we know $u_i = \overline{u}_i = 1.0$ and $\overline{s}_i^{(2)} = 0$. This gives for $j = 1, \dots, n$,

$$|v_{j} - \bar{v}_{j}| \leq \bar{s}_{j}^{(3)} = fl\left(\sum_{k=1}^{n} |X_{jk}| \bar{s}_{k}^{(2)} + \epsilon_{2} \xi\right) \cdot (1 + \epsilon_{1} + 5\eta_{1}) .$$

Finally we normalize \bar{v} so its largest component in modulus is 1.0 by dividing by the largest component \bar{v}_{jmax} in double precision, obtaining

$$\left|v_{j} - fll\left(\frac{\bar{v}_{j}}{\bar{v}_{j\max}}\right)\right| \leq fl\left(\frac{\bar{s}_{j}^{(3)} + 9\eta_{2}|\bar{v}_{j}|}{|\bar{v}_{j\max}|}\right) \cdot (1 + 8\eta_{1}) .$$

5. Use of the Program. The program, M-10, given in the microfiche section herein, is coded in standard Algol 60 except for the addition of *complex* and *long* (double-precision) declarations. Arithmetic operations between two long variables is assumed to be done in double precision, and we assume the *abs* function is defined for a complex argument and gives the modulus. It is important to note that the technique can be applied repeatedly, using the output improved eigensystem as the input for the next iteration.

The program, translated into Burroughs Algol for real matrices, has been tested on scores of matrices using the Burroughs B5500 at Stanford University. The one example given here is the 12×12 Frank matrix for which the approximate eigensystem is given in [2]. Because the approximations to the smallest eigenvalues were so poor, on the first iteration we could not isolate the smallest three eigenvalues. But the improved eigensystem was more accurate, and on the fourth iteration, we obtained estimates for the eigenvalues which agreed with those published by Wilkinson [3, p. 152] to all 15 decimal places he gives, and which were guaranteed to at least 12 decimal places, as were the eigenvectors. We list the fourth iteration results for the smallest three eigenvalues.

Acknowledgment. The author would like to thank Professor G. E. Forsythe for his supervision and guidance, and Dr. J. H. Wilkinson for his many helpful suggestions.

Computer Science Department Stanford University Stanford, California 94305

801

^{1.} M. MARCUS & H. MINC, A Survey of Matrix Theory and Matrix Inequalities, Allyn & Bacon, Boston, Mass., 1964. MR 29 #112. 2. J. M. VARAH, "The calculation of the eigenvectors of a general complex matrix by inverse

iteration," Math. Comp., v. 22, 1968, pp. 785-791.
 3. J. H. WILKINSON, Rounding Errors in Algebraic Processes, Notes on Applied Science No.
 32, HMSO, London; Prentice-Hall, Englewood Cliffs, N. J., 1963. MR 28 #4661.

^{4.} J. H. WILKINSON, The Algebraic Eigenvalue Problem, Clarendon Press, Oxford, 1965. MR **32** #1894.