

Almost Diagonal Matrices with Multiple or Close Eigenvalues*

J. H. WILKINSON
*National Physical Laboratory
 Teddington, England*

1. INTRODUCTION

In a number of algorithms for finding eigenvalues of a matrix A_1 , the latter is reduced by an iterative sequence of similarity transformations to almost diagonal form. When A_1 has a multiple eigenvalue this is true of all the transforms (assuming exact computation). We are interested then in the nature of almost diagonal matrices with multiple eigenvalues. It turns out that such matrices have special characteristics which are of considerable interest as regards the convergence of iterative procedures for reducing a matrix to diagonal form.

2. THE HERMITIAN CASE

We first consider hermitian matrices with multiple eigenvalues. Let A be hermitian with eigenvalues $\lambda_1, \lambda_1, \dots, \lambda_1, \lambda_{r+1}, \lambda_{r+2}, \dots, \lambda_n$ the root λ_1 being precisely of multiplicity r . (A may have other multiple eigenvalues but this will not affect the argument.) Let δ be defined by the relation

$$3\delta = \min_{i=r+1}^n |\lambda_i - \lambda_1|, \quad (1)$$

and let

$$A = D + E, \quad (2)$$

* This work was done while the author was Visiting Professor in the Computer Science Department of Stanford University. It was supported by the National Science Foundation and the Office of Naval Research.

where D is the diagonal of A . Suppose we have

$$\|E\|_F = \varepsilon < \delta \quad (3)$$

where F denotes the Frobenius norm $(\sum \sum |e_{ij}|^2)^{1/2}$. When ε is small A may be regarded as almost diagonal. By the Wielandt-Hoffman theorem [5] the λ_i and a_{ii} may be ordered so that

$$\sum (\lambda_{p_i} - a_{ii})^2 \leq \varepsilon^2. \quad (4)$$

Let us permute the rows and columns of A similarly so that the a_{ii} associated with the λ_1 eigenvalues are the first r . Without loss of generality we can assume this was true originally and with appropriate numbering of the remaining $n - r$ eigenvalues inequality (4) becomes

$$\sum_1^n (\lambda_i - a_{ii})^2 \leq \varepsilon^2. \quad (5)$$

We write

$$A = \begin{bmatrix} F & G \\ G^T & H \end{bmatrix}, \quad (6)$$

where F is an $r \times r$ matrix.

If the eigenvalues of H are $\lambda'_{r+1}, \dots, \lambda'_n$, then since the off-diagonal elements of H are a subset of those of E , we have by the Wielandt-Hoffman theorem [5] with appropriate numbering of the λ'_i

$$\sum_{r+1}^n (\lambda'_i - a_{ii})^2 \leq \varepsilon^2. \quad (7)$$

Hence

$$\begin{aligned} |\lambda'_i - \lambda_i| &= |\lambda'_i - a_{ii} + a_{ii} - \lambda_i| \\ &\leq \varepsilon + \varepsilon < 2\delta \end{aligned} \quad (8)$$

and

$$\begin{aligned} |\lambda'_i - \lambda_1| &= |\lambda_i - \lambda_1 + \lambda'_i - \lambda_i| \\ &\geq |\lambda_i - \lambda_1| - |\lambda'_i - \lambda_i| \\ &\geq 3\delta - 2\delta \\ &= \delta. \end{aligned} \quad (9)$$

The matrix $H - \lambda_1 I$ is therefore nonsingular, i.e., it is of rank $n - r$. Now since A has λ_1 as a r -fold root it, too, is of rank $n - r$. We shall show that this means that F is especially related to G and H . We partition $A - \lambda_1 I$ in the form

$$A - \lambda_1 I = \begin{bmatrix} F - \lambda_1 I & G \\ G^T & H - \lambda_1 I \end{bmatrix}. \quad (10)$$

If we premultiply $A - \lambda_1 I$ by

$$\begin{bmatrix} I & -G(H - \lambda_1 I)^{-1} \\ \theta & I \end{bmatrix}, \quad (11)$$

its rank is unaltered and hence the derived matrix

$$\begin{bmatrix} F - \lambda_1 I - G(H - \lambda_1 I)^{-1}G^T & \theta \\ G^T & H - \lambda_1 I \end{bmatrix} \quad (12)$$

is also of rank $n - r$. Since $H - \lambda_1 I$ is already of rank $n - r$ this can be true only if

$$F - \lambda_1 I - G(H - \lambda_1 I)^{-1}G^T = \theta, \quad (13)$$

i.e.,

$$F = \lambda_1 I + G(H - \lambda_1 I)^{-1}G^T = \lambda_1 I + M \quad (\text{say}). \quad (14)$$

Now the elements of G are a subset of those of E and hence

$$\|G\|_E = \|G^T\|_E \leq \varepsilon, \quad (15)$$

while

$$(H - \lambda_1 I)^{-1} = R \operatorname{diag}(\lambda_i' - \lambda_1)^{-1} R^H, \quad (16)$$

where R is unitary. Hence from the unitary invariance of the Frobenius norm and from (9) and (15)

$$\begin{aligned} \|M\|_F &\leq \|G\|_E \max_i |\lambda_i' - \lambda_1|^{-1} \|G^T\|_E \\ &\leq \frac{\varepsilon^2}{\delta}. \end{aligned} \quad (17)$$

We see then that diagonal elements of F differ from λ_1 by quantities bounded by ε^2/δ and its off-diagonal elements are bounded by ε^2/δ .

When $\varepsilon \ll \delta$ this means that the largest off-diagonal element of A is never found in F , the matrix with the diagonal elements "associated" with the multiple root λ_1 . This has important consequences in connection with the classical Jacobi method [6, 10, 15] for diagonalizing hermitian matrices. At each stage in the reduction the largest off-diagonal element in the current matrix is annihilated but the theorem shows that after a certain stage such an off-diagonal element is never "associated" with two elements tending to the same multiple root.

This simple observation removes a difficulty in demonstrating that the classical Jacobi method is always ultimately quadratically convergent [7, 10, 11, 15]. A similar remark applies to the serial Jacobi method if a threshold strategy is used [9]. If at any stage the element which is annihilated is chosen to be one which is not small compared with the current norm of off-diagonal elements then this ensures that from a certain stage the annihilated element will not be associated with two diagonal elements tending to the same multiple root.

3. PATHOLOGICALLY CLOSE ROOTS

In practice when a transformation is made on a matrix having multiple roots, the transformed matrix merely has very close roots because of rounding errors. In discussing the convergence of the Jacobi method the quantity $\min |\lambda_i - \lambda_j|$ is of great importance and the presence of very close roots would appear to be serious. We now show that this is not so.

Suppose the roots of A are

$$\lambda_1, \lambda_2, \dots, \lambda_r; \lambda_{r+1}, \dots, \lambda_n, \quad (18)$$

where

$$\lambda_i = \lambda + \varepsilon_i, \quad i = 1, \dots, r, \quad (19)$$

and the ε_i are very small. The first r roots are therefore pathologically close. Define D and E as in (2), but δ by the relation

$$3\delta = \min_{i=r+1}^n |\lambda_i - \lambda|, \quad (20)$$

and assume that

$$\|E\|_F + \left(\sum \varepsilon_i^2\right)^{1/2} = \varepsilon < \delta. \quad (21)$$

Now A may be expressed in the form

$$A = R D_1 R^H, \quad (22)$$

where $D_1 = \text{diag}(\lambda_i)$ and D_1 can be separated into D_2 and D_3 where

$$D_2 = \text{diag}(\lambda, \lambda, \dots, \lambda, \lambda_{r+1}, \dots, \lambda_n), \quad (23)$$

$$D_3 = \text{diag}(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_r, 0, \dots, 0). \quad (24)$$

Hence

$$A = R(D_2 + D_3)R^H = R D_2 R^H + R D_3 R^H = B + C \quad (\text{say}). \quad (25)$$

The matrix B has λ as an r -fold root and to apply the result of the previous section we require only a bound for the Frobenius norm of its off-diagonal elements. Since $B = A - C$ such a bound is given by

$$\|E\|_F + \|C\|_F = \|E\|_F + \left(\sum \varepsilon_i^2\right)^{1/2} = \varepsilon < \delta.$$

The Frobenius norm of the off-diagonal elements of B "associated" with the multiple root is therefore bounded by ε^2/δ and hence that of the corresponding elements of A is bounded by $\varepsilon^2/\delta + \left(\sum \varepsilon_i^2\right)^{1/2}$.

Suppose for example a matrix A has the roots

$$1 - 2(10^{-10}), \quad 1 - 10^{-10}, \quad 1, \quad 4, \quad 5$$

and

$$\|E\|_F + 2^{1/2}(10^{-10}) = 10^{-5}.$$

The off-diagonal elements of A associated with the close roots will then have a Frobenius norm bounded by

$$10^{-10} + 2^{1/2}(10^{-10})$$

and therefore they will all be far smaller than the largest off-diagonal element of A . Hence at such a stage in the classical Jacobi method or the threshold serial Jacobi method with a matrix having the root distribution above, the current rotation will not be in a plane associated with the close roots. In fact, with the above example one sweep of the threshold serial Jacobi method will reduce the norm of off-diagonal elements from 10^{-5} to 10^{-10} . Provided we do not wish to reduce the norm below this level the presence of the close roots has no adverse influence. (In fact, it is beneficial since it ensures that the main weight in the off-diagonal positions is concentrated on fewer elements.)

4. NONHERMITIAN MATRICES

The above proofs may give the impression that the result above is associated specifically with hermitian matrices. In fact a closely related result is true for any diagonalizable matrix having an r -fold root.

Again let A have the roots $\lambda_1, \dots, \lambda_1, \lambda_{r+1}, \dots, \lambda_n$. Let

$$A = D + E \quad (D \text{ diagonal}), \tag{26}$$

$$\min_{\lambda_i \neq \lambda_j} |\lambda_i - \lambda_j| = 3\delta, \tag{27}$$

$$\|E\|_\infty = \varepsilon < \delta. \tag{28}$$

Suppose the λ_i ($i = r + 1, \dots, n$) include a root λ_k of multiplicity s so that $A - \lambda_k I$ is of rank $n - s$; then by a theorem due to Fan and Hoffman [2]* there are at least s indices i for which

$$|a_{ii} - \lambda_k| \leq \sum_{j \neq i} |a_{ij}| \leq \varepsilon < \delta, \tag{29}$$

the last two inequalities following from (28). This means that there are at least s of the a_{ii} lying in the disk with center λ_k and radius ε .

At the moment it appears possible that there may be more than s ; if so let us associate the first s qualifying a_{ii} with λ_k . The a_{ii} associated with λ_k and λ_l ($k, l > r, \lambda_k \neq \lambda_l$) must be different since the disks with centers λ_k and λ_l are disjoint, from (27) and (28). Hence we have associated precisely $n - r$ diagonal elements with $\lambda_{r+1}, \dots, \lambda_n$. Now permute rows and columns of A similarly so that these $n - r$ diagonal elements are in the southeast corner. We can assume that A was in this form originally.

As in the symmetric case, $A - \lambda_1 I$ is of rank $n - r$ and partitioning $A - \lambda_1 I$ in the form

$$A - \lambda_1 I = \begin{bmatrix} F - \lambda_1 I & G \\ K & H - \lambda_1 I \end{bmatrix} \tag{30}$$

we have

$$F - \lambda_1 I - G(H - \lambda_1 I)^{-1}K = \theta, \tag{31}$$

* I would like to thank Dr. B. Levinger of Case Institute for drawing my attention to this result while we were both enjoying the hospitality of Applied Mathematics Division, Argonne National Laboratory.

provided $H - \lambda_1 I$ is nonsingular. Now

$$H - \lambda_1 I = \text{diag}(a_{ii} - \lambda_1) + L = D_1 + L \quad (\text{say}), \quad (32)$$

where L is the matrix of off-diagonal elements of H . (These elements are a subset of those of E .) Since

$$\begin{aligned} |a_{ii} - \lambda_1| &\geq |\lambda_i - \lambda_1| - |a_{ii} - \lambda_i| \\ &> 2\delta \quad (i = r + 1, \dots, n), \end{aligned} \quad (33)$$

D_1 is nonsingular and

$$H - \lambda_1 I = D_1 [I + D_1^{-1} L]. \quad (34)$$

Now

$$\begin{aligned} \|D_1^{-1} L\|_\infty &\leq \|D_1^{-1}\|_\infty \|L\|_\infty \\ &\leq \max_{i=r+1}^n |a_{ii} - \lambda_1|^{-1} \|L\|_\infty \\ &\leq \frac{\varepsilon}{2\delta} < \frac{1}{2}, \end{aligned} \quad (35)$$

and hence

$$\begin{aligned} \|(H - \lambda_1 I)^{-1}\|_\infty &\leq \frac{\|D_1^{-1}\|_\infty}{(1 - \|D_1^{-1} L\|_\infty)} \\ &\leq \frac{1/2\delta}{1/2} \\ &= \frac{1}{\delta}. \end{aligned} \quad (36)$$

Equation (31) therefore gives

$$F = \lambda_1 I + G(H - \lambda_1 I)^{-1} K = \lambda_1 I + M \quad (\text{say}),$$

where

$$\begin{aligned} \|M\|_\infty &\leq \|G\|_\infty \|(H - \lambda_1 I)^{-1}\|_\infty \|K\|_\infty \\ &\leq \frac{\varepsilon^2}{\delta}. \end{aligned} \quad (37)$$

This now shows that precisely $n - r$ of the a_{jj} were associated with the λ_i ($i = r + 1, \dots, n$) and the remaining r diagonal elements are all in a disk of radius ε^2/δ centered on λ_1 . Again off-diagonal elements "associated" with the multiple eigenvalue are bounded by ε^2/δ and are therefore well below the level of the largest off-diagonal elements when $\varepsilon \ll \delta$.

The result is at first sight surprising since the condition of the *eigenvalue* problem of A seems not to be involved. Indeed a result may be proved which is only marginally weaker even when A is defective (though not as far as λ_1 is concerned). In this respect it is the hypothesis $\|E\|_\infty \leq \varepsilon$ which is deceptive. If B has an ill-conditioned eigenvalue problem then, in order to derive a similarity transformation $X^{-1}BX = A$ such that A is almost diagonal with $\|E\|_\infty$ less than a prescribed quantity, we shall, in general, have to work to higher precision if B is ill-conditioned than if it is well-conditioned. In the hermitian case the hypothesis does not have this deceptive feature.

5. PATHOLOGICALLY CLOSE ROOTS IN NONHERMITIAN CASE

The deceptive nature of the result becomes apparent as soon as we consider the effect of very close roots. Assume now that A is nondefective and let X be a matrix having as its columns n independent eigenvectors of A . Then we have

$$A = X \operatorname{diag}(\lambda_i) X^{-1}. \quad (38)$$

Using a similar notation to that in Section 3 we have in the case of r very close roots

$$A = XD_2X^{-1} + XD_3X^{-1} = B + C, \quad (39)$$

where B now has an r -fold root. In the hermitian case X is unitary and $\|C\|_F = \|D_3\|_F$, but now all we can say is

$$\|C\| \leq \|X\| \|D_3\| \|X^{-1}\|, \quad (40)$$

and we see that the condition number κ of X with respect to inversion is inevitably involved. It is clear that it is the minimum value of $\|X\| \|X^{-1}\|$ for all permissible X that is relevant [1]. It should be emphasized, though, that the possession of a multiple root or of a set of very close roots does not imply that $\|X\| \|X^{-1}\|$ is necessarily large. Provided the close roots are well conditioned the fact that the eigenvector problem is ill conditioned is irrelevant.

6. ITERATIVE REFINEMENT OF AN EIGENSYSTEM

The above results have important consequences in connection with procedures for the refinement of a computed eigensystem of a matrix [12, 13, 15]. In such procedures one starts with a computed set of eigenvalues and eigenvectors μ_i and x_i . Let X be the matrix having columns x_i and define R and S by the relation

$$AX - X \operatorname{diag}(\mu_i) = R, \quad (41)$$

$$X^{-1}AX - \operatorname{diag}(\mu_i) = X^{-1}R = S. \quad (42)$$

If the system were exact both R and S would be null. In practice neither R nor S can be computed exactly with the given X because of rounding errors but with well-designed procedures an \bar{S} is determined with a low relative error. Hence we have

$$X^{-1}AX = \operatorname{diag}(\mu_i) + \bar{S} + (S - \bar{S}). \quad (43)$$

If the computed system is accurate \bar{S} is small, and with good procedures for calculating R and $X^{-1}R$ a bound is obtained for $\|S - \bar{S}\|$ which is small compared with $\|\bar{S}\|$. (Note \bar{S} is computed explicitly but a bound for the norm only is determined for $S - \bar{S}$.) The matrix sum on the right of (43) is therefore an almost diagonal matrix which is exactly similar to A .

Now when A has a multiple root corresponding to a linear divisor our result shows that provided \bar{S} is small (and hence $S - \bar{S}$ is *very* small), the off-diagonal elements of \bar{S} associated with the multiple roots will be far smaller than the largest off-diagonal elements of \bar{S} . When none of the roots of A is ill conditioned we shall find typically that if $\|\bar{S}\|_\infty = \epsilon$ then the bound for $\|S - \bar{S}\|_\infty$ will be approximately $2^{-t}\epsilon$ (with a t -digit mantissa binary computer). The diagonal elements of $\operatorname{diag}(\mu_i) + \bar{S}$ associated with the multiple roots will differ by quantities of the order of ϵ^2 and the associated off-diagonal elements will be of order ϵ^2 .

Hence after suitable permutations of rows and columns the right-hand side of (43) will have the form

$$\text{diagonal} + \begin{matrix} r \\ \left\{ \begin{matrix} \overbrace{\epsilon^2 L} \\ \epsilon N & \begin{matrix} \epsilon M \\ \epsilon P \end{matrix} \end{matrix} \right\} + (S - S) \end{matrix} \quad (44)$$

and the bound for $\|\bar{S} - S\|$ will usually be of order at least as small as ε^2 . Premultiplication of the first r rows by $k\varepsilon$ and the first r columns by $(1/k)\varepsilon$ then modifies the second matrix to the form

$$\begin{bmatrix} \varepsilon^2 L & k\varepsilon^2 M \\ \frac{1}{k} N & \varepsilon P \end{bmatrix} \quad (45)$$

and because of this Gerschgorin's theorem gives just as fine bounds for multiple roots as for well-separated roots.

Forgetting rounding errors for the moment it is interesting to consider what can be achieved with an approximate matrix X of eigenvectors which can be expressed in the form

$$\bar{X} = X(I + \varepsilon E), \quad (46)$$

where $\|E\|_\infty = 1$ and X is a matrix of exact normalized eigenvectors. We have

$$\begin{aligned} \bar{X}^{-1} A \bar{X} &= (I + \varepsilon E)^{-1} X^{-1} A X (I + \varepsilon E) \\ &= (I - \varepsilon E + \varepsilon^2 E^2 - \dots) \text{diag}(\lambda_i) (I + \varepsilon E) \\ &= \text{diag}(\lambda_i) + \varepsilon F + \text{terms in } \varepsilon^2, \text{ etc.}, \end{aligned} \quad (47)$$

where

$$f_{ij} = -\lambda_j e_{ij} + \lambda_i e_{ij}. \quad (48)$$

We see that the elements f_{ij} are zero whenever $\lambda_i = \lambda_j$. Hence the off-diagonal elements associated with multiple eigenvalues are of order ε^2 .

Notice that when A has eigenvalues which, while not being truly coincident, have separations which are appreciably smaller than ε , (48) shows that the associated off-diagonal elements are again appreciably smaller than ε and a simple application of Gerschgorin's theorem using diagonal similarity transformations gives bounds for the relevant eigenvalues which are of the order of ε^2 or of the separations, whichever is the larger. The weakest bounds arise when the separations are themselves of order ε . The bounds are then of order ε and cannot be improved *merely* by diagonal similarity transformations.

When the procedure for refining an eigensystem is used iteratively, then provided the system is not too ill conditioned the final eigensystem

is "correct to working accuracy." Generally we can assume that the final computed system of vectors satisfies a relation of the form

$$\bar{X} = X + E \quad \text{where} \quad \|E\|_\infty \leq n \cdot 2^{-t} \|X\|_\infty. \quad (49)$$

Hence we have

$$\begin{aligned} \bar{X}^{-1}A\bar{X} &= (X^{-1} - X^{-1}EX^{-1} - \dots)AX(I + X^{-1}E) \\ &= \text{diag}(\lambda_i) - X^{-1}E \text{diag}(\lambda_i) + \text{diag}(\lambda_i)X^{-1}E + \dots \end{aligned} \quad (50)$$

Equation (50) shows the real limitation on the attainable accuracy with computation of a prescribed precision. The off-diagonal elements of $\bar{X}^{-1}A\bar{X}$ are certainly bounded by $2n \cdot 2^{-t} \|X\|_\infty \|X^{-1}\|_\infty \max|\lambda_i|$, if we ignore the quadratic and higher-order terms in E . Writing

$$2n \cdot 2^{-t} \|X\|_\infty \|X^{-1}\|_\infty \max|\lambda_i| = \beta, \quad (51)$$

the bounds attainable for the eigenvalues using Gerschgorin's theorem and diagonal transformations can be expressed in the following form.

Let the eigenvalues be divided into three groups. The first group consists of multiple eigenvalues; the second group consists of eigenvalues with a minimum separation which is less than β , and the third group consists of the remainder. For an eigenvalue in the first group having a minimum separation of δ_1 from all other eigenvalues the bound is of the order of β^2/δ_1 . For a member of the second group having separations of up to s from its close neighbors and a minimum separation of order δ_2 from all others the bound is of the order of $s + (\beta^2/\delta_2)$. For a member of the third group having a minimum separation from *all* other eigenvalues of δ_3 the bound is of the order of β^2/δ_3 . In general unless $\|X\|_\infty \|X^{-1}\|_\infty$ is quite large the bounds are all appreciably better than $2^{-t} \max|\lambda_i|$ except when s is of the order of magnitude of β .

This result has been amply confirmed in practice—multiple eigenvalues being found, in general, to the same high precision as well-separated eigenvalues.

REFERENCES

- 1 F. L. Bauer, Optimally scaled matrices, *Numer. Math.* **5**(1963), 73–87.
- 2 K. Fan and A. J. Hoffman, Lower bounds for the rank and location of the eigenvalues of a matrix, *Nat. Bureau of Standards Appl. Math. Series* **39**(1954), 117–130.
- 3 S. Gerschgorin, Über die Abgrenzung der Eigenwerte einer Matrix, *Izv. Akad. Nauk. SSSR, Ser. fiz-mat.* **6**(1931), 749–754.

- 4 P. Henrici, On the speed of convergence of cyclic and quasi-cyclic Jacobi methods for computing eigenvalues of Hermitian matrices, *J. Soc. Industr. Appl. Math.* **6**(1958), 144–162.
- 5 A. J. Hoffman and H. W. Wielandt, The variation of the spectrum of a normal matrix, *Duke Math. J.* **20**(1953), 37–39.
- 6 C. G. J. Jacobi, Über ein leichtes Verfahren, die in der Theorie der Säculärstörungen vorkommenden Gleichungen numerisch aufzulösen, *Crelle's J.* **30**(1846), 51–94.
- 7 H. P. M. van Kempken, On the convergence of the classical Jacobi method for real symmetric matrices with non-distinct eigenvalues, *Numer. Math.* **9**(1966), 11–18.
- 8 H. P. M. van Kempken, On the quadratic convergence of the special cyclic Jacobi method, *Numer. Math.* **9**(1966), 19–22.
- 9 D. A. Pope and C. Tompkins, Maximizing functions of rotations, *J. Ass. Comp. Mach.* **4**(1957), 459–466.
- 10 A. Schönhage, Zur Konvergenz des Jacobi-Verfahrens, *Numer. Math.* **3**(1961), 374–380.
- 11 A. Schönhage, Zur quadratischen Konvergenz des Jacobi-Verfahrens, *Numer. Math.* **6**(1964), 410–412.
- 12 J. Varah, Ph. D. Thesis, Stanford (1967).
- 13 J. H. Wilkinson, Rigorous error bounds for computed eigensystems, *Computer J.* **4**(1961), 230–241.
- 14 J. H. Wilkinson, Note on the quadratic convergence of the cyclic Jacobi process, *Numer. Math.* **4**(1962), 296–300.
- 15 J. H. Wilkinson, *The Algebraic Eigenvalue Problem*, Oxford Univ. Press, 1965.
- 16 J. H. Wilkinson, The QR algorithm for real symmetric matrices with multiple eigenvalues, *Computer J.* **8**(1965), 85–87.

Received September 14, 1967