

## THE APPLICATION OF EIGENPAIR STABILITY TO BLOCK DIAGONALIZATION\*

NILOTPAL GHOSH<sup>†</sup>, WILLIAM W. HAGER<sup>‡</sup>, AND PURANDAR SARMAH<sup>‡</sup>

**Abstract.** An algorithm presented in Hager [*Comput. Math. Appl.*, 14 (1987), pp. 561–572] for diagonalizing a matrix is generalized to a block matrix setting. It is shown that the resulting algorithm is locally quadratically convergent. A global convergence proof is given for matrices with separated eigenvalues and with relatively small off-diagonal elements. Numerical examples along with comparisons to the QR method are presented.

**Key words.** block diagonalization, eigenpair stability, eigenvalues, eigenvectors

**AMS subject classifications.** 65F05, 65F15, 65F35

**PII.** S0036142995261252

**1. Introduction.** Let  $A$  denote an  $n \times n$  block matrix; that is, the  $(i, j)$  element  $A_{ij}$  of  $A$  is itself a matrix of dimension  $n_i \times n_j$ , where  $\sum_{i=1}^r n_i = n$  for some  $r \leq n$ . In this paper, we develop and analyze an algorithm for computing a block diagonalization  $X\Lambda X^{-1}$  of  $A$ , assuming one exists. Here  $\Lambda$  is a block diagonal matrix with diagonal blocks  $\Lambda_i, i = 1$  to  $r$ , and  $X$  is an invertible matrix whose  $i$ th block of columns is denoted by  $X_i$ . Hence, the equation  $A = X\Lambda X^{-1}$  is equivalent to the relation

$$AX_i = X_i\Lambda_i, \quad i = 1 \text{ to } r.$$

The algorithm developed in this paper is based on a stability result, Proposition 1, for a perturbation  $A(\varepsilon)X(\varepsilon) = X(\varepsilon)\Lambda(\varepsilon)$  of the original eigenequation. We show that if the spectrum of  $\Lambda_i$  and  $\Lambda_j$  are disjoint for each  $i \neq j$ , then there exist continuously differentiable solutions  $X(\varepsilon)$  and  $\Lambda(\varepsilon)$  to the perturbed equation. After differentiating the perturbed equation and applying Taylor's theorem, we obtain the following algorithm (throughout the paper, the subscript  $k$  denotes the iteration number while the subscripts  $i$  and  $j$  denote elements or submatrices of larger matrices).

**BLOCK DIAGONALIZATION ALGORITHM.** If  $X_k\Lambda_k X_k^{-1}$  is the current approximate diagonalization of  $A$ , then

$$(1) \quad \Lambda_{k+1} = \text{diag } X_k^{-1}AX_k \quad \text{and} \quad X_{k+1} = X_k(I + D),$$

where

$$(2) \quad \text{diag } D = 0 \quad \text{and} \quad D\Lambda_{k+1} - \Lambda_{k+1}D = \text{off } X_k^{-1}AX_k.$$

The notation “diag” and “off” above are defined in the following way: given a block matrix  $B$ , “diag  $B$ ” denotes the block diagonal matrix whose diagonal blocks coincide with the diagonal blocks of  $B$ , and “off  $B$ ” denotes the matrix that coincides with  $B$  except for the diagonal blocks which are replaced by blocks of zeros.

---

\*Received by the editors January 4, 1995; accepted for publication (in revised form) September 22, 1995. This research was supported by U.S. Army Research Office contract DAAL03-89-G-0082 and by the National Science Foundation.

<http://www.siam.org/journals/sinum/34-3/26125.html>

<sup>†</sup>Department of Mathematics, Catonsville Community College, Baltimore, MD 21228. Visiting Professor, Department of Mathematics, University of Florida, Gainesville, FL, 1990–1991.

<sup>‡</sup>Department of Mathematics, University of Florida, Gainesville, FL 32611 (hager@math.ufl.edu, puran@swamp.bellcore.com).

In section 3 we show that if the spectrum of  $\Lambda_i$  and  $\Lambda_j$  are disjoint for  $i \neq j$ , then this algorithm is locally quadratically convergent. Moreover, in the case of  $1 \times 1$  diagonal blocks, we establish a global convergence result in section 4 when the diagonal elements of  $A$  are separated and the off-diagonal elements are relatively small.

Some related Newton-type methods for computing or refining invariant spaces are presented in [2], [4], [5], [7], and [15]. The work of Anselone and Rall [2] justifies the application of Newton's method to the system

$$(3) \quad Ax - \lambda x = 0, \quad y^T x = 1,$$

where the vector  $x$  and the scalar  $\lambda$  are the unknowns and  $y$  is a normalization vector. This scheme, which typically converges locally quadratically for a simple eigenvalue, involves the inversion of a matrix of the following form in each iteration:

$$\begin{bmatrix} A - \lambda_k I & x_k \\ y^T & 0 \end{bmatrix}.$$

In [4] Chatelin examines the problem of approximating an invariant space spanned by the columns of some matrix  $Z$ . She applies Newton's method to the equation

$$(4) \quad AZ - Z(Y^T AZ) = 0,$$

where  $Y$  is a normalization matrix. Observe that at a solution to (4),  $Y^T Z = I$  assuming  $Y^T AZ$  is invertible. The condition  $Y^T Z = I$  is the generalization of the condition  $y^T x = 1$  in (3). Chatelin obtains a local quadratic convergence result for Newton's method applied to (4) assuming that zero is not an eigenvalue of  $A$  and that the number of eigenvalues of  $A$ , counting multiplicities, associated with the columns of  $Z$  is equal to the number of columns of  $Z$ . Each iteration of this scheme involves inverting a linear operator  $L$ , where  $L$  acting on a matrix  $M$  is defined by

$$L(M) = (I - ZY^T)AM - M(Y^T AZ).$$

In [7] Dongarra, Moler, and Wilkinson also consider Newton's method applied to (3); however, they focus on the special case where all the components of  $y$  are zero except for one component which is set to one. In [7] the authors consider both the modified Newton's method and the standard Newton scheme. In the modified Newton's method (see [12]), the Jacobian is evaluated at some given point, and in each Newton iteration, this fixed Jacobian is used instead of the Jacobian at the new iterate. Hence, each iteration of the modified scheme involves inverting a matrix of the form

$$\begin{bmatrix} A - \lambda I & x \\ y^T & 0 \end{bmatrix},$$

where  $\lambda$  and  $x$  are fixed approximations to an eigenpair. Generalizations of this idea for invariant spaces are also presented in [7].

In [5] Demmel examines the schemes [4] and [7], as well as a scheme that he attributes to Stewart [15], and he shows that in some sense, each of these schemes tries to solve an underlying Riccati equation. In comparing the schemes of Anselone and Rall, Chatelin, and Dongarra, Moler, and Wilkinson with the scheme proposed in this paper, we make the following observation: each iteration of any of these schemes

essentially involves the inversion of a matrix. With the block diagonalization scheme, one inversion yields a first-order correction to all the eigenvectors and eigenvalues while one inversion with any of the other methods yields a first-order correction to an invariant space. On the other hand, with these other schemes, one only needs an approximation to an invariant space to get an improved approximation while the scheme in this paper requires an approximation to a block diagonalization to get an improved approximation.

For another comparison between the block algorithm and Newton's method applied to (3), suppose that an approximate block diagonalization has been determined and we wish to obtain a first-order correction. For simplicity, suppose that all the eigenvalues are simple and separated so that Newton's method applied to (3) converges nicely. With the block algorithm, the flop count for 1 step is essentially  $\frac{10}{3}n^3$ —there are  $n^3$  flops involved in multiplying  $A$  and  $X_k$ ,  $\frac{4}{3}n^3$  flops to compute  $X_k^{-1}(AX_k)$ , and  $n^3$  flops to compute  $X_k(I + D)$ . Now consider Newton's method applied to (3); to be specific, let us consider the choice of Dongarra, Molder, and Wilkinson for  $y$ : each component is zero except for one component. To implement the Newton approach, let us initially reduce  $A$  to upper Hessenberg form by an orthogonal similarity transformation; without this reduction, the Newton approach would require on the order of  $n^4$  flops. The flop count (see [10]) for this reduction is about  $\frac{10}{3}n^3$  flops. As described in [7], the Newton system can be solved using a series of Givens rotations, followed by an LU factorization. The flop count for this process in the worst case is  $3n^2$ , and since there are  $n$  eigenvectors to update, the total flop count is  $3n^3$ . Finally, we need to expend  $n^3$  flops to multiply each eigenvector by the orthogonal factor to obtain an eigenvector of the original matrix. The total flop count for the Newton approach is  $\frac{22}{3}n^3$  flops. Hence, the block algorithm appears to be preferable relative to a flop count; in addition, the block algorithm is much easier to implement, especially when the blocks are larger than  $1 \times 1$ .

**2. The algorithm.** The algorithm developed in this paper is based on a stability result for the block diagonalization of a perturbation of the original matrix  $A$ .

PROPOSITION 1. *Let  $A(\varepsilon)$  denote an  $n \times n$  matrix which is continuously differentiable at  $\varepsilon = 0$  and which has the property  $A = A(0)$ . If for each  $i \neq j$  the spectrum of  $\Lambda_i$  and  $\Lambda_j$  is disjoint, then there exist continuously differentiable functions  $X(\varepsilon)$  and  $\Lambda(\varepsilon)$ , where  $\Lambda(\varepsilon)$  is block diagonal,  $X(0) = X$ ,  $\Lambda(0) = \Lambda$ , and for each  $\varepsilon$  near 0 we have*

$$(5) \quad A(\varepsilon)X(\varepsilon) = X(\varepsilon)\Lambda(\varepsilon).$$

*Proof.* Our proof of this result is based on the implicit function theorem. Instead of evaluating  $X(\varepsilon)$  directly, we introduce an auxiliary matrix  $C(\varepsilon)$  which is chosen such that  $X(\varepsilon) = XC(\varepsilon)$ . Since  $X$  is invertible, there is a one-to-one correspondence between  $C(\varepsilon)$  and  $X(\varepsilon)$ . It turns out that there are an infinite number of solutions to (5) with the desired properties. In order to achieve uniqueness, we impose the following constraints:  $C_{jj}(\varepsilon) = I$  for each  $j$ . With these side constraints, (5) takes the form

$$(6) \quad A(\varepsilon)XC(\varepsilon) = XC(\varepsilon)\Lambda(\varepsilon), \quad C_{jj}(\varepsilon) = I \quad \text{for each } j.$$

Abstractly, (6) has the form  $F(\varepsilon, C(\varepsilon), \Lambda(\varepsilon)) = 0$ , where the solutions  $C(\varepsilon)$  and  $\Lambda(\varepsilon)$  depend on  $\varepsilon$  and where the number of equations is equal to the number of unknown elements of  $C(\varepsilon)$  and  $\Lambda(\varepsilon)$ . At  $\varepsilon = 0$ , there is the trivial solution  $C(0) = I$ , and

$\Lambda(0) = \Lambda$ . By the implicit function theorem (see [1, Thm. 13.7] or [8, Thm. 10.8]), there exists a continuously differentiable solution for  $\varepsilon$  near zero if the derivative operator of the system (6) with respect to  $C$  and  $\Lambda$  is invertible at  $\varepsilon = 0$ . Establishing invertibility of the derivative is equivalent to showing that the only solution  $(\delta C, \delta \Lambda)$ , with  $\delta \Lambda$  block diagonal, to the following linear system is  $\delta C = 0$  and  $\delta \Lambda = 0$ :

$$\frac{\partial F}{\partial C}(0, C(0), \Lambda(0))\delta C + \frac{\partial F}{\partial \Lambda}(0, C(0), \Lambda(0))\delta \Lambda = \frac{\partial F}{\partial C}(0, I, \Lambda)\delta C + \frac{\partial F}{\partial \Lambda}(0, I, \Lambda)\delta \Lambda = 0.$$

In the context of (6), this reduces to

$$(7) \quad AX\delta C = X\delta C\Lambda + X\delta \Lambda, \quad \delta C_{jj} = 0, \quad j = 1 \text{ to } r.$$

Premultiplying (7) by  $X^{-1}$  yields

$$(8) \quad \Lambda\delta C = \delta C\Lambda + \delta \Lambda, \quad \delta C_{jj} = 0, \quad j = 1 \text{ to } r.$$

Since  $\Lambda$  and  $\delta \Lambda$  are block diagonal and  $\delta C_{jj} = 0$ , we deduce that  $\delta \Lambda = 0$ . For  $i \neq j$ , (8) implies that

$$(9) \quad \Lambda_i\delta C_{ij} - \delta C_{ij}\Lambda_j = 0.$$

Referring to [13, sect. 12.5, Thm. 1], we see that the unique solution to (9) is  $\delta C_{ij} = 0$  when the spectrum of  $\Lambda_i$  and  $\Lambda_j$  is disjoint. Since  $\delta C = \delta \Lambda = 0$ , we conclude that the derivative of (6) with respect to  $C$  and  $\Lambda$  at  $\varepsilon = 0$  is invertible. The implicit function theorem completes the proof.  $\square$

Proposition 1 is surprising since  $\Lambda(\varepsilon)$  is differentiable even though its eigenvalues may be nondifferentiable. To evaluate  $X'(0)$  and  $\Lambda'(0)$ , we differentiate (6) to obtain the relation

$$A'(0)X + AXC'(0) = XC'(0)\Lambda + X\Lambda'(0), \quad C'_{jj}(0) = 0, \quad j = 1 \text{ to } r.$$

Premultiplying by  $X^{-1}$  yields

$$(10) \quad X^{-1}A'(0)X + \Lambda C'(0) = C'(0)\Lambda + \Lambda'(0), \quad C'_{jj}(0) = 0, \quad j = 1 \text{ to } r.$$

If  $Z^T$  denotes  $X^{-1}$ , (10) is equivalent to the following relations:

$$(11) \quad \Lambda'_i(0) = Z_i^T A'(0)X_i, \quad i = 1 \text{ to } r,$$

and

$$(12) \quad C'_{ij}(0)\Lambda_j - \Lambda_i C'_{ij}(0) = Z_i^T A'(0)X_j \quad \text{for } i \neq j, \quad C'_{jj}(0) = 0, \quad j = 1 \text{ to } r.$$

With the notation “diag” and “off” of section 1, (11) and (12) can be expressed

$$(13) \quad \Lambda'(0) = \text{diag } X^{-1}A'(0)X, \quad \text{diag } C'(0) = 0, \quad C'(0)\Lambda - \Lambda C'(0) = \text{off } X^{-1}A'(0)X.$$

After obtaining  $C'(0)$  satisfying (13), we have  $X'(0) = XC'(0)$ .

When the spectrum of  $\Lambda_i$  and  $\Lambda_j$  is disjoint, (12) can be solved in the following way (see [3] or [10, p. 387] for details): replace  $\Lambda_i$  and  $\Lambda_j$  by their Schur decompositions and obtain an equivalent equation with upper triangular matrices in place of  $\Lambda_i$  and  $\Lambda_j$ . This upper triangular system can be solved directly by a substitution procedure.

As in [11], we can utilize these formulas for  $X'(0)$  and  $\Lambda'(0)$  in an iterative algorithm for block diagonalizing a matrix. In particular, given a matrix  $A$  and an approximate block diagonalization  $X\Lambda X^{-1}$ , let us consider the matrix

$$A(\varepsilon) = X\Lambda X^{-1} + \varepsilon(A - X\Lambda X^{-1}).$$

Observe that  $A(1) = A$ . Suppose that the spectrum of  $\Lambda_i$  and  $\Lambda_j$  is disjoint for  $i \neq j$  and that for  $\varepsilon$  between zero and one, the continuously differentiable functions  $X(\varepsilon)$  and  $\Lambda(\varepsilon)$  of Proposition 1 exist. Evaluating the first-order Taylor expansions of  $X(\varepsilon)$  and  $\Lambda(\varepsilon)$  around  $\varepsilon = 0$  and putting  $\varepsilon = 1$ , we obtain the following first-order approximations:

$$\Lambda(1) \approx \Lambda(0) + \Lambda'(0) = \Lambda + \text{diag}(X^{-1}(A - X\Lambda X^{-1})X) = \text{diag } X^{-1}AX,$$

and

$$X(1) \approx X(0) + X'(0) = X + XC'(0) = X(I + C'(0)),$$

where  $C'(0)$  satisfies (13). Based on these expansions, we arrive at the algorithm (1) and (2), where  $D$  corresponds to the matrix  $C'(0)$  above.

As another variation of (2), the matrix  $\Lambda_{k+1}$  can be replaced by the previous iterate  $\Lambda_k$  in the evaluation of  $D$ ; that is, the matrix  $D$  satisfies

$$(14) \quad \text{diag } D = 0 \quad \text{and} \quad D\Lambda_k - \Lambda_k D = \text{off } X_k^{-1}AX_k.$$

It turns out that with the modified formula (14) for  $D$ , the block diagonalization algorithm is 2-step quadratically convergent while the original formula (2) yields a 1-step quadratically convergent algorithm.

**3. Local quadratic convergence.** This section analyzes the local convergence of the block diagonalization algorithm while the next section analyzes one situation where convergence is guaranteed for the starting guess  $X_0 = I$ . Throughout this analysis,  $\|\cdot\|$  denotes any matrix norm with the following properties:

- P1.  $\|AB\| \leq \|A\| \|B\|$  for each  $A$  and  $B$ .
- P2.  $\|A\| \leq \|B\|$  whenever  $|a_{ij}| \leq |b_{ij}|$  for each  $i$  and  $j$ .
- P3.  $\|I\| = 1$ .

Below, subscripts  $i$  and  $j$  are used to denote elements of matrices while subscripts  $k$ ,  $l$ , and  $m$  are used to denote the iteration number. Finally, let  $B_\rho(X)$  denote the ball with center  $X$  and radius  $\rho$ :

$$B_\rho(X) = \{Y : \|Y - X\| \leq \rho\}.$$

Recall that there is an eigenspace associated with each eigenvalue of a matrix. At best, the block diagonalization algorithm converges to elements of the eigenspace associated with a block diagonalization of  $A$ . Given a block diagonalization  $X\Lambda X^{-1}$  of  $A$ , let  $S(X)$  denote the set of all matrices of the form  $XD$  for some invertible block diagonal matrix  $D$ .

**THEOREM 1.** *Suppose that  $A = X\Lambda X^{-1}$ , where  $\Lambda$  is block diagonal and the spectrum of  $\Lambda_i$  and  $\Lambda_j$  is disjoint for all  $i \neq j$ . Then for all  $X_0$  in a neighborhood of  $X$ , the iterates  $X_k$  and  $\Lambda_k$  generated by the block diagonalization algorithm approach limits  $X_\infty$  and  $\Lambda_\infty$  such that  $A = X_\infty \Lambda_\infty X_\infty^{-1}$ , where  $X_\infty$  lies in  $S(X)$  and where the*

root convergence order is at least two. That is, there exists a constant  $C$ , independent of  $k$ , such that

$$\|X_k - X_\infty\| + \|\Lambda_k - \Lambda_\infty\| \leq C2^{-2^k}.$$

*Proof.* Let  $F_\Lambda(Z)$  denote the linear operator defined by

$$F_\Lambda(Z) = \text{off}(Z\Lambda - \Lambda Z),$$

where the domain and range are the set of  $n \times n$  block matrices with diagonal blocks equal to zero. As noted previously,  $F_\Lambda$  is an invertible linear operator since the spectrum of  $\Lambda_i$  and  $\Lambda_j$  is disjoint for all  $i \neq j$ . Since  $F_\Lambda$  is a continuous function of  $\Lambda$ , there exist constants  $\gamma$  and  $\sigma$  such that  $\|F_M^{-1}\| \leq \gamma$  for every block diagonal  $M \in B_\sigma(\Lambda)$ . Choose  $\rho$  small enough that  $Y$  is invertible whenever  $Y \in B_\rho(X)$ , and let  $\beta$  be defined by

$$\beta = \max\{\|Y^{-1}\| : Y \in B_\rho(X)\}.$$

Decrease  $\rho$  further if necessary to ensure that

$$(15) \quad \rho\beta \leq 1 \quad \text{and} \quad \|\Lambda - Y^{-1}AY\| + \beta\rho\|Y^{-1}AY\| \leq \sigma$$

for every  $Y \in B_\rho(X)$ . We will show that there exist a constant  $c$  and a  $Y_{k+1} \in S(X)$  such that

$$(16) \quad \|Y_{k+1} - Y_k\| \leq c\|X_k - Y_k\| \quad \text{and} \quad \|Y_{k+1} - X_{k+1}\| \leq c\|X_k - Y_k\|^2,$$

where  $c$  is independent of  $X_k$  and  $Y_k$  in  $B_\rho(X)$ , and where  $X_{k+1}$  is generated by the block diagonalization algorithm.

Assuming, for the moment, the existence of the constant  $c$  in (16), the proof is completed in the following way: define  $Y_0 = X$  and choose  $X_0$  close enough to  $X$  such that

$$(17) \quad \|X_0 - X\| \leq \text{minimum} \left\{ \frac{1}{(4c)}, \frac{\rho}{(4c)}, \frac{\rho}{4} \right\}.$$

Since  $Y_0 = X$ , both  $X_0$  and  $Y_0$  lie in  $B_\rho(X)$ . Proceeding by induction, suppose that  $X_1, Y_1, \dots, X_k, Y_k$  all lie in  $B_\rho(X)$ . We use (16) to show that  $X_{k+1}$  and  $Y_{k+1}$  also lie in  $B_\rho(X)$ . The second inequality in (16) leads to the relation

$$c\|Y_k - X_k\| \leq (c\|X - X_0\|)^{2^k} \quad \text{or} \quad \|Y_k - X_k\| \leq \|X - X_0\|(c\|X - X_0\|)^{2^k - 1}.$$

Combining this with the bound (17), we have

$$(18) \quad \|Y_k - X_k\| \leq \|X - X_0\| \frac{1}{4^{2^k - 1}} = \frac{4\|X - X_0\|}{4^{2^k}} \leq \frac{\rho}{4^{2^k}}.$$

Utilizing the first inequalities in (16) and (18), and the bound in (17), we get

$$\begin{aligned} \|Y_{k+1} - X\| &\leq \|Y_{k+1} - Y_k\| + \|Y_k - X\| \leq c\|X_k - Y_k\| + \|Y_k - X\| \\ &\leq \frac{4c}{4^{2^k}}\|X - X_0\| + \|Y_k - X\| \leq \frac{\rho}{4^{2^k}} + \|Y_k - X\|. \end{aligned}$$

Repeated application of this inequality yields

$$\|Y_{k+1} - X\| \leq \rho \sum_{l=0}^k \frac{1}{4^{2^l}} < \frac{\rho}{2}.$$

Thus,  $Y_{k+1} \in B_{\rho/2}(X)$ . By (18), with  $k$  replaced by  $k + 1$ , we have

$$\|Y_{k+1} - X_{k+1}\| \leq \frac{\rho}{4^{2^{k+1}}}.$$

Hence,  $X_{k+1}$  lies in  $B_\rho(X)$  and the induction step is complete.

By the first inequality in (16) and by (18) we have

$$\begin{aligned} (19) \quad \|Y_k - Y_l\| &\leq \sum_{m=l}^{k-1} \|Y_{m+1} - Y_m\| \leq c \sum_{m=l}^{k-1} \|Y_m - X_m\| \leq c\rho \sum_{m=l}^{k-1} \frac{1}{4^{2^m}} \\ &\leq c\rho \sum_{m=l}^{\infty} \frac{1}{4^{2^m}} \leq \frac{2c\rho}{4^{2^l}} \end{aligned}$$

for any  $k \geq l$ . Thus, for  $X_0$  sufficiently close to  $X$ , the  $Y_k$  form a Cauchy sequence approaching some limit  $X_\infty$ , and by (18) the  $X_k$  approach  $X_\infty$  as well. The triangle inequality

$$\|X_k - X_\infty\| \leq \|X_k - Y_k\| + \|Y_k - X_\infty\|$$

combined with (18) and (19) complete the proof of quadratic convergence for the  $X_k$ . Since  $\Lambda_{k+1} = \text{diag } X_k^{-1}AX_k$ , we conclude that the  $\Lambda_k$  approach  $\Lambda_\infty = X_\infty^{-1}AX_\infty$  at the same rate that the  $X_k$  approach  $X_\infty$ .

To establish (16), we proceed in the following way. Let  $E$  be a matrix chosen so that  $X_k = Y_k(I + E)$ ; that is,  $E = Y_k^{-1}X_k - I = Y_k^{-1}(X_k - Y_k)$ . By (15) and (18)

$$(20) \quad \|E\| = \|Y_k^{-1}(X_k - Y_k)\| \leq \|Y_k^{-1}\| \|X_k - Y_k\| \leq \beta \|X_k - Y_k\| \leq \frac{\beta\rho}{4^{2^k}} \leq \frac{1}{4^{2^k}} \leq \frac{1}{2}.$$

It follows that  $I + E$  is invertible and

$$(21) \quad \|(I + E)^{-1}\| \leq \frac{1}{1 - \|E\|} \leq 2.$$

Defining  $M_k = Y_k^{-1}AY_k$ , which is block diagonal since  $Y_k \in S(X)$ , observe that

$$(22) \quad \begin{cases} X_k^{-1}AX_k = (I + E)^{-1}Y_k^{-1}AY_k(I + E) \\ \quad = [I - E + E^2(I + E)^{-1}]M_k(I + E) \\ \quad = M_k - EM_k + M_kE - EM_kE + E^2(I + E)^{-1}M_k(I + E) \\ \quad = M_k - EM_k + M_kE + L, \end{cases}$$

where

$$L = E^2(I + E)^{-1}M_k(I + E) - EM_kE.$$

By (15) with  $Y = Y_k$  we have

$$(23) \quad \|\Lambda - M_k\| = \|\Lambda - Y_k^{-1}AY_k\| \leq \sigma,$$

and combining this with (21) we have

$$(24) \quad \|L\| \leq \|E\|^2 \|M_k\| \left(1 + \frac{1 + \|E\|}{1 - \|E\|}\right) = \frac{2\|M_k\| \|E\|^2}{1 - \|E\|} \leq 4(\sigma + \|\Lambda\|)\|E\|^2.$$

Defining  $\bar{\Lambda} = \Lambda_{k+1} = \text{diag } X_k^{-1}AX_k$  and referring to (22), the matrix  $D$  involved in the evaluation of  $X_{k+1}$  satisfies

$$\text{off } (D\bar{\Lambda} - \bar{\Lambda}D) = D\bar{\Lambda} - \bar{\Lambda}D = \text{off } X_k^{-1}AX_k = \text{off } (M_kE - EM_k + L),$$

which implies that

$$(25) \quad \text{off } ((D + E)\bar{\Lambda} - \bar{\Lambda}(D + E)) = \text{off } ((M_k - \bar{\Lambda})E - E(M_k - \bar{\Lambda}) + L).$$

By (21), (22), (23), and the next-to-last inequality in (24) we have

$$(26) \quad \begin{aligned} \|\bar{\Lambda} - M_k\| &= \|\text{diag } X_k^{-1}AX_k - M_k\| \\ &\leq \|X_k^{-1}AX_k - M_k\| = \|M_kE - EM_k + L\| \\ &\leq 2\|M_k\| \|E\| + \|L\| \leq \frac{2\|M_k\| \|E\|}{1 - \|E\|} \leq 4(\sigma + \|\Lambda\|)\|E\|, \end{aligned}$$

from which it follows that

$$(27) \quad \begin{aligned} \|\text{off } ((M_k - \bar{\Lambda})E - E(M_k - \bar{\Lambda}) + L)\| &\leq \|(M_k - \bar{\Lambda})E - E(M_k - \bar{\Lambda}) + L\| \\ &\leq 2\|E\| \|\bar{\Lambda} - M_k\| + \|L\| \leq 12(\sigma + \|\Lambda\|)\|E\|^2. \end{aligned}$$

By (26) we have

$$(28) \quad \|\Lambda - \bar{\Lambda}\| \leq \|\bar{\Lambda} - M_k\| + \|M_k - \Lambda\| \leq \frac{2\|M_k\| \|E\|}{1 - \|E\|} + \|M_k - \Lambda\|.$$

Referring to (20), we see that

$$(29) \quad \|E\| \leq \frac{1}{2} \quad \text{and} \quad \|E\| \leq \frac{\beta\rho}{4}.$$

Combining (28) and (29) gives

$$\|\Lambda - \bar{\Lambda}\| \leq \|M_k - \Lambda\| + \beta\rho\|M_k\|.$$

By (15)  $\|\Lambda - \bar{\Lambda}\| \leq \sigma$ , which ensures that  $\|F_{\bar{\Lambda}}^{-1}\| \leq \gamma$ . This bound in conjunction with (25) and (27) implies that

$$(30) \quad \|\text{off } (D + E)\| \leq 12\gamma(\sigma + \|\Lambda\|)\|E\|^2.$$

Let  $G$  be defined by  $G = \text{off } (D + E)$ . Substituting  $D = \text{off } D = G - \text{off } E$  in (1) and utilizing the identity  $X_k = Y_k(I + E)$  yields

$$X_{k+1} = X_k(I + D) = Y_k(I + E)(I + G - \text{off } E).$$

Hence, we have

$$(31) \quad X_{k+1} = Y_k(I + \text{diag } E) + Y_k(I + E)G - Y_kE\text{off } E.$$



Defining  $Y_{k+1} = Y_k(I + \text{diag } E)$ , observe that

$$(32) \quad \|Y_{k+1} - Y_k\| \leq \|Y_k\| \|\text{diag } E\| \leq (\rho + \|X\|)\|E\|$$

since  $Y_k$  lies in  $B_\rho(X)$ . Referring to (20),  $\|E\| \leq \beta\|X_k - Y_k\|$ , and combining this with (32) gives

$$(33) \quad \|Y_{k+1} - Y_k\| \leq \beta(\|X\| + \rho)\|X_k - Y_k\|,$$

which yields the first inequality in (16). Moreover, by (31), we have

$$\|X_{k+1} - Y_{k+1}\| \leq \|Y_k\|(\|I + E\| \|G\| + \|E\|^2) \leq (\rho + \|X\|)(\|I + E\| \|G\| + \|E\|^2).$$

Combining this with the bound for  $G = \text{off } (D + E)$  in (30) and the relation  $\|I + E\| \leq 1 + \|E\| \leq \frac{3}{2}$  from (29) gives

$$(34) \quad \begin{aligned} \|X_{k+1} - Y_{k+1}\| &\leq (\rho + \|X\|)(18\gamma(\sigma + \|\Lambda\|) + 1)\|E\|^2 \\ &\leq \beta^2(\rho + \|X\|)(18\gamma(\sigma + \|\Lambda\|) + 1)\|X_k - Y_k\|^2. \end{aligned}$$

The constant  $c$  in (16) is the maximum of the constants in (33) and (34) □

**4. A priori convergence.** This section analyzes one situation where the block diagonalization algorithm is guaranteed to converge. Define  $A_0 = A$  and for  $k > 0$ , let  $A_k$  denote  $X_k^{-1}AX_k$ . With this notation, the block diagonalization algorithm, starting from  $X_0 = I$ , can be expressed in the following way:

$$(35) \quad \Lambda_{k+1} = \text{diag } A_k, \quad A_{k+1} = (I + D_k)^{-1}A_k(I + D_k), \quad X_{k+1} = X_k(I + D_k),$$

where

$$\text{diag } D_k = 0 \quad \text{and} \quad D_k\Lambda_{k+1} - \Lambda_{k+1}D_k = \text{off } A_k.$$

Throughout this section, we assume that the diagonal blocks of  $A$  are all  $1 \times 1$ . In addition to the three properties imposed for the norm in section 3, we assume the following:

P4.  $\|\text{diag } A\| = \text{maximum } \{|A_{ii}| : i = 1 \text{ to } n\}$  for each  $A$ .

Let  $\sigma(A)$  be the minimum separation of diagonal elements defined by

$$\sigma(A) = \text{minimum}_{i \neq j} |A_{jj} - A_{ii}|.$$

In the case of  $1 \times 1$  blocks, the equation for  $D_k$  can be solved explicitly. In particular, we have

$$(36) \quad D_{k,ij} = \frac{A_{k,ij}}{A_{k,jj} - A_{k,ii}}$$

for all  $i \neq j$ . Taking absolute values yields

$$|D_{k,ij}| \leq |A_{k,ij}|/\sigma(A_k).$$

It follows that

$$(37) \quad \|D_k\| \leq \frac{\|\text{off } A_k\|}{\sigma(A_k)}.$$

These observations are the basis for the following theorem.

THEOREM 2. *If  $\|\text{off } A\| < \sigma(A)(\sqrt{3} - 1)/2$ , then  $A$  is diagonalizable, and the iterates  $\Lambda_k$  and  $X_k$  generated by the block diagonalization algorithm, with the starting guess  $X_0 = I$ , approach limits  $\Lambda_\infty$  and  $X_\infty$ , respectively, where  $A = X_\infty \Lambda_\infty X_\infty^{-1}$ .*

*Proof.* Let  $\alpha$  denote the ratio  $\|\text{off } A\|/\sigma(A)$ , which is less than or equal to  $(\sqrt{3} - 1)/2$  by assumption. Hence, the inequality

$$(38) \quad \|\text{off } A_k\| \leq \alpha\sigma(A_k)$$

holds trivially at  $k = 0$ . Proceeding by induction, assume that (38) holds at iteration  $k$ . By (37) and (38),

$$(39) \quad \|D_k\| \leq \frac{\|\text{off } A_k\|}{\sigma(A_k)} \leq \alpha.$$

Letting  $B$  denote  $I + D_k$ , it follows that

$$(40) \quad \|B^{-1}\| = \|(I + D_k)^{-1}\| \leq \frac{1}{1 - \|D_k\|} \leq \frac{1}{1 - \alpha}.$$

By (35) we have

$$\begin{aligned} A_{k+1} - \Lambda_{k+1} &= B^{-1}A_k B - \Lambda_{k+1} \\ &= B^{-1}(\text{off } A_k + \Lambda_{k+1})B - \Lambda_{k+1} \\ &= B^{-1}\text{off } A_k B + B^{-1}\Lambda_{k+1}B - \Lambda_{k+1} \\ &= B^{-1}\text{off } A_k B + B^{-1}(\Lambda_{k+1}B - B\Lambda_{k+1}) \\ &= B^{-1}\text{off } A_k B - B^{-1}\text{off } A_k \\ &= B^{-1}\text{off } A_k(B - I) \\ &= B^{-1}\text{off } A_k D_k. \end{aligned}$$

Combining this with (38), (39), and (40) yields

$$(41) \quad \begin{aligned} \|\text{off } A_{k+1}\| &\leq \|A_{k+1} - \Lambda_{k+1}\| \leq \|B^{-1}\| \|\text{off } A_k\| \|D_k\| \\ &\leq \frac{\alpha}{1 - \alpha} \|\text{off } A_k\| \leq \frac{\alpha^2}{1 - \alpha} \sigma(A_k). \end{aligned}$$

Also note that

$$(42) \quad \begin{aligned} \|\text{diag } (A_{k+1} - A_k)\| &\leq \|A_{k+1} - \text{diag } A_k\| = \|A_{k+1} - \Lambda_{k+1}\| \\ &\leq \frac{\alpha}{1 - \alpha} \|\text{off } A_k\| \leq \frac{\alpha^2}{1 - \alpha} \sigma(A_k). \end{aligned}$$

Applying the “diag” operator to the identity  $A_{k+1} = A_k + (A_{k+1} - A_k)$  and referring to (42) we have

$$(43) \quad \begin{aligned} \sigma(A_{k+1}) &\geq \sigma(A_k) - 2\|\text{diag } (A_{k+1} - A_k)\| \geq \sigma(A_k) - \frac{2\alpha^2\sigma(A_k)}{1 - \alpha} \\ &= \frac{(1 - \alpha - 2\alpha^2)\sigma(A_k)}{1 - \alpha}. \end{aligned}$$

Combining this with (41) gives

$$\frac{\|\text{off } A_{k+1}\|}{\sigma(A_{k+1})} \leq \frac{\alpha^2}{1 - \alpha - 2\alpha^2} \leq \alpha$$

for  $\alpha \leq (\sqrt{3} - 1)/2$ . This completes the induction step, and (38) holds for all  $k$ .

Observe that in (41) we establish the relation

$$\|\text{off } A_{k+1}\| \leq \frac{\alpha}{1-\alpha} \|\text{off } A_k\|.$$

Repeated application of this inequality gives

$$(44) \quad \|\text{off } A_k\| \leq \left(\frac{\alpha}{1-\alpha}\right)^k \|\text{off } A\|.$$

Since  $\alpha < \frac{1}{2}$ , the off-diagonal elements of  $A_k$  all tend to zero. Also, by (42) and (44),

$$(45) \quad \|\text{diag } A_{k+1} - \text{diag } A_k\| \leq \left(\frac{\alpha}{1-\alpha}\right)^{k+1} \|\text{off } A\|.$$

Thus, the diagonal of  $A_k$  forms a Cauchy sequence approaching some limit  $\Lambda_\infty$ , while the off-diagonal elements tend to zero.

We now show that no two diagonal elements of  $\Lambda_\infty$  are equal. By the first inequality in (43) and by (45) we have

$$\sigma(A_{k+1}) \geq \sigma(A_k) - 2 \left(\frac{\alpha}{1-\alpha}\right)^{k+1} \|\text{off } A\| = \sigma(A_k) - 2\alpha\sigma(A) \left(\frac{\alpha}{1-\alpha}\right)^{k+1}.$$

Repeated application of this inequality yields

$$\begin{aligned} \sigma(A_k) &\geq \sigma(A) \left(1 - 2\alpha \sum_{i=1}^k \left(\frac{\alpha}{1-\alpha}\right)^i\right) \geq \sigma(A) \left(1 - 2\alpha \sum_{i=1}^\infty \left(\frac{\alpha}{1-\alpha}\right)^i\right) \\ &= \sigma(A) \left(1 - 2\alpha \left(\frac{\alpha/(1-\alpha)}{1-\alpha/(1-\alpha)}\right)\right) = \sigma(A) \left(\frac{1-2\alpha-2\alpha^2}{1-2\alpha}\right). \end{aligned}$$

Since  $1 - 2\alpha - 2\alpha^2 > 0$  for  $\alpha < (\sqrt{3} - 1)/2$ ,  $\sigma(A_k)$  is uniformly bounded away from zero.

Combining the lower bound for  $\sigma(A_k)$  with (37) and (44) yields

$$(46) \quad \begin{aligned} \|D_k\| &\leq cr^k, \quad \text{where } r = \frac{\alpha}{1-\alpha} \leq \frac{1}{\sqrt{3}} \quad \text{and} \\ c &= \frac{\|\text{off } A\|(1-2\alpha)}{\sigma(A)(1-2\alpha-2\alpha^2)} = \frac{\alpha(1-2\alpha)}{1-2\alpha-2\alpha^2}. \end{aligned}$$

By the equation for  $X_{k+1}$

$$\|X_{k+1}\| = \|X_k(I + D_k)\| \leq (1 + \|D_k\|)\|X_k\| \leq (1 + cr^k)\|X_k\| \leq \|X_0\| \prod_{l=0}^k (1 + cr^l).$$

Since

$$\log \left(\prod_{l=0}^\infty (1 + cr^l)\right) \leq \frac{c}{(1-r)},$$

$\|X_k\|$  is uniformly bounded by some constant  $\beta$  independent of  $k$  (see [1, p. 209, Thm. 8.55]). The estimate

$$\|X_{k+1} - X_k\| = \|X_k D_k\| \leq \|X_k\| \|D_k\| \leq \beta cr^k$$

TABLE 1

$n$	$k$	$\ \text{off } A_k\ _\infty$
10	4	$2.0 \times 10^{-9}$
40	4	$2.7 \times 10^{-9}$
160	4	$2.7 \times 10^{-9}$
640	4	$2.7 \times 10^{-9}$

TABLE 2

$k$	$\ \text{off } A_k\ _\infty$
1	$4 \times 10^{-1}$
2	$3 \times 10^{-2}$
3	$1 \times 10^{-4}$
4	$2 \times 10^{-9}$

implies that the  $X_k$  form a Cauchy sequence that approaches some limit  $X_\infty$ . Since  $\|D_k\| < 1$  for each  $k$  by (39), the  $X_k$  are invertible and

$$(47) \quad \|X_{k+1}^{-1}\| \leq \|(I + D_k)^{-1}\| \|X_k^{-1}\| \leq \frac{1}{1 - \|D_k\|} \|X_k^{-1}\| \leq \frac{1}{1 - cr^k} \|X_k^{-1}\|.$$

The relation  $1/(1 - cr^k) \leq 1 + 2cr^k$  for  $k$  sufficiently large along with (47) shows that the sequence  $X_k^{-1}$  is uniformly bounded. Thus,  $X_\infty$  is invertible and  $X_k^{-1}$  approaches  $X_\infty^{-1}$ ; moreover,  $\Lambda_{k+1} = \text{diag } X_k^{-1} A X_k$  approaches  $X_\infty^{-1} A X_\infty = \Lambda_\infty$ . This completes the proof.  $\square$

During the proof of Theorem 2, it was shown that no two diagonal elements of  $\Lambda_\infty$  are equal. Hence, Theorem 1 implies that the convergence in Theorem 2 is locally quadratic.

**5. Numerical experiments.** To illustrate Theorems 1 and 2, we consider the matrix  $A$  defined by  $A_{ij} = 3^{-|i-j|}$  for  $i \neq j$ , and  $A_{ii} = i$ . We apply the iteration (35), stopping when the following inequality holds:

$$\|\text{off } A_k\|_\infty = \|\text{off } X_k^{-1} A X_k\|_\infty \leq 10^{-6}.$$

Here  $\|\cdot\|_\infty$  denotes the matrix  $\infty$ -norm (maximum absolute row sum). By Gerschgorin's theorem the diagonal elements of  $A_k$  approximate the eigenvalues of  $A$  with error at most  $\|\text{off } A_k\|_\infty$ . The number of iterations needed for convergence appears in Table 1 for various values of the dimension  $n$  of  $A$ . In Table 2 we show how the error depends on the iteration number for  $n = 10$ . Observe that the convergence appears to be at least quadratic as the theory predicts.

In Table 3 we compare the execution times (in seconds on a Sun Sparc 20 workstation, compiled with f77 -O) for the block diagonalization algorithm with the corresponding times of both Dongarra's SICEDR [6] (an implementation of the Newton algorithm contained in ACM algorithm 589) and the QR algorithm as implemented in EISPACK [9], [14]. Recall that the flop count for one iteration of the block algorithm is  $\frac{10}{3}n^3$  while the Newton approach requires about  $\frac{22}{3}n^3$  flops. As a rough estimate, a diagonalization computed by the QR method requires  $9n^3$  flops—about  $5n^3$  flops to reduce the matrix to Hessenberg form and to form the orthogonal matrix used in the reduction,  $2n^3$  flops to update the orthogonal matrix during the reduction to Schur form,  $n^3$  flops to reduce the Hessenberg matrix to triangular form, and  $n^3$  flops

TABLE 3

$n$	Block	QR	SICEDR
10	.01	.04	.05
40	.10	.12	.81
160	6.00	4.04	57.97
640	329.92	284.88	4561.26

TABLE 4

$\varepsilon$	Block	QR
.05	6	151
.01	3	148
.001	2	140
.0001	2	138

to diagonalize the triangular matrix and multiply by the orthogonal factor. Hence, three iterations of the block algorithm are comparable to the QR algorithm, which agrees with the observed computing times in Table 3. Observe though that the Newton approach was somewhat slower than was predicted by the flop counts. Moreover, some of the Newton iterates converged to the same eigenvalue; hence, one did not actually achieve a diagonalization of the starting matrix but rather achieved a refined approximation to most of the eigenpairs.

As both the flop counts and the execution times in Table 3 indicate, one or two iterations of the block algorithm execute quicker than the QR algorithm. Hence, on a serial computer, the block algorithm is only faster than the QR algorithm when we have a good starting guess for the block diagonalization, for example, when a matrix is modified incrementally and re-diagonalized after each change. On the other hand, the block algorithm is much better suited to parallel computing than the QR method since matrix multiplication and matrix inversion are readily implemented in a parallel computing environment. In comparing code size the block algorithm was implemented in 120 lines of Fortran code, SICEDR contains 1000 statements, and the QR code contains 1370 statements (excluding comment and matrix generation statements).

In the next sequence of experiments we take a  $100 \times 100$  matrix  $A$  whose elements are randomly distributed between zero and one, and we use the QR method to obtain the real Schur decomposition  $A = QUQ^T$ , where  $Q$  is orthogonal and  $U$  is quasi-upper triangular. We then perturb  $A$  by a matrix  $E$  whose elements are randomly distributed on the interval  $[0, \varepsilon]$ . In Table 4 we compare the number of iterations for the block diagonalization algorithm with the number of iterations for the QR method applied to the matrix  $Q^T(A + E)Q = U + Q^TEQ$ . The starting  $X_0$  in the block algorithm corresponds to the invariant spaces of  $A$ , where a pair of real vectors are used to span each complex conjugate pair of eigenvectors. With this choice for  $X_0$ , the block algorithm could be implemented with real arithmetic. Since the number of QR iterations needed to obtain the original Schur decomposition was 132, it appears that the number of iterations for the QR method applied to a small perturbation of an upper Hessenberg matrix is essentially the same as the number of iterations that were required for the original random matrix. On the other hand, with the block diagonalization algorithm, the number of iterations decreases as the perturbation decreases.

## REFERENCES

- [1] T. M. APOSTOL, *Mathematical Analysis*, 2nd ed., Addison-Wesley, Reading, MA, 1974.
- [2] P. M. ANSELONE AND L. B. RALL, *The solution of characteristic value-vector problems by Newton's method*, Numer. Math., 11 (1968), pp. 38–45.
- [3] R. H. BARTELS AND G. W. STEWART, *Solution of the equation  $AX + XB = C$* , Comm. ACM, 15 (1972), pp. 820–826.
- [4] F. CHATELIN, *Simultaneous Newton's iteration for the eigenproblem*, Comput. Suppl., 5 (1984), pp. 67–74.
- [5] J. W. DEMMEL, *Three methods for refining estimates of invariant subspaces*, Computing, 38 (1987), pp. 43–57.
- [6] J. J. DONGARRA, *Algorithm 589 SICEDR: A FORTRAN subroutine for improving the accuracy of computed matrix eigenvalues*, ACM Trans. Math. Software, 8 (1982), pp. 371–375.
- [7] J. J. DONGARRA, C. B. MOLER, AND J. H. WILKINSON, *Improving the accuracy of computed eigenvalues and eigenvectors*, SIAM J. Numer. Anal., 20 (1983), pp. 23–45.
- [8] L. FLATTO, *Advanced Calculus*, Waverly Press, Baltimore, MD, 1976.
- [9] B. S. GARBOW, J. M. BOYLE, J. J. DONGARRA, AND C. B. MOLER, *Matrix Eigensystem Routines - EISPACK Guide Extension*, Springer-Verlag, Berlin, 1977.
- [10] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 2nd ed., The Johns Hopkins University Press, Baltimore, MD, 1989.
- [11] W. W. HAGER, *Bidiagonalization and diagonalization*, Comput. Math. Appl., 14 (1987), pp. 561–572.
- [12] W. W. HAGER, *Applied Numerical Linear Algebra*, Prentice-Hall, Englewood Cliffs, NJ, 1988.
- [13] P. LANCASTER AND M. TISMENETSKY, *The Theory of Matrices*, Academic Press, Orlando, 1985.
- [14] B. T. SMITH, J. M. BOYLE, J. J. DONGARRA, B. S. GARBOW, Y. IKEBE, V. C. KLEMA, AND C. B. MOLER, *Matrix Eigensystem Routines - EISPACK Guide*, Springer-Verlag, Berlin, 1976.
- [15] G. W. STEWART, *Error and perturbation bounds for subspaces associated with certain eigenvalue problems*, SIAM Rev., 15 (1973), pp. 727–764.