

IMPROVING THE ACCURACY OF COMPUTED EIGENVALUES AND EIGENVECTORS*

J. J. DONGARRA,† C. B. MOLER‡ AND J. H. WILKINSON§

Abstract. This paper describes and analyzes several variants of a computational method for improving the numerical accuracy of, and for obtaining numerical bounds on, matrix eigenvalues and eigenvectors. The method, which is essentially a numerically stable implementation of Newton's method, may be used to "fine tune" the results obtained from standard subroutines such as those in EISPACK [Lecture Notes in Computer Science 6, 51, Springer-Verlag, Berlin, 1976, 1977]. Extended precision arithmetic is required in the computation of certain residuals.

Introduction. The calculation of an eigenvalue λ and the corresponding eigenvector x (here after referred to as an eigenpair) of a matrix A involves the solution of the nonlinear system of equations

$$(A - \lambda I)x = 0.$$

Starting from an approximation $\bar{\lambda}$ and \bar{x} , a sequence of iterates may be determined using Newton's method or variants of it. The conditions on $\bar{\lambda}$ and \bar{x} guaranteeing convergence have been treated extensively in the literature. For a particularly lucid account the reader is referred to the book by Rall [3]. In a recent paper Wilkinson [7] describes an algorithm for determining error bounds for a computed eigenpair based on these mathematical concepts. Considerations of numerical stability were an essential feature of that paper and indeed were its main *raison d'être*. In general this algorithm provides an improved eigenpair and error bounds for it; unless the eigenpair is very ill conditioned the improved eigenpair is usually correct to the precision of the computation used in the main body of the algorithm.

In this paper we present several extensions of that algorithm which greatly increase its range of application. These extensions cover the efficient determination of the complex conjugate eigenpairs of a real matrix and the determination of appropriate invariant subspaces when individual eigenvectors are very ill conditioned, and finally give more rapid convergence when the initial eigenpair is of low accuracy. It should perhaps be emphasized that the main relevance of these algorithms in the case when the approximate eigenpairs are derived from a well designed eigenvalue package such as EISPACK [4], [2] is to provide error bounds. As in the earlier paper, the emphasis in each of the algorithms is on the problems of numerical stability.

1. The basic algorithm. We begin with a brief description of the basic algorithm described by Wilkinson. If λ, x is an approximate eigenpair, and $\lambda + \mu, x + \tilde{y}$ is a neighboring eigenpair, then

$$(1.1) \quad A(x + \tilde{y}) = (\lambda + \mu)(x + \tilde{y}),$$

* Received by the editors August 24, 1981, and in revised form May 10, 1982.

† Argonne National Laboratory, Argonne, Illinois 60439. The work of this author was supported in part by the Applied Mathematical Sciences Research Program (KC-04-02) of the Office of Energy Research of the U.S. Department of Energy under contract W-31-109-Eng-38.

‡ Department of Computer Science, University of New Mexico, Albuquerque, New Mexico 87131. The work of this author was supported in part by the Computer Science Section of the National Science Foundation under contract MCS 7603297.

§ Department of Computer Science, Stanford University, Stanford, California 94305. The work of this author was supported in part by the U.S. Air Force under contract AFOSR-79-0094 and by the National Science Foundation under contract MCS 7811985.

this relation being exact. We assume that x is normalized so that $\|x\|_\infty = 1 = x_s$, and we remove the degree of arbitrariness in \tilde{y} by requiring that $\tilde{y}_s = 0$. From (1.1)

$$(1.2) \quad (A - \lambda I)\tilde{y} - \mu x = \lambda x - Ax + \mu \tilde{y},$$

where the last term on the right will be of second order in the errors of λ, x . Equation (1.2) may be simplified by the introduction of a vector y defined by

$$(1.3) \quad y^T = (y_1, y_2, \dots, y_{s-1}, \mu, y_{s+1}, \dots, y_n),$$

so that y gives the full information on both μ and \tilde{y} . Equation (1.2) then becomes

$$(1.4) \quad By = r + y_s \tilde{y},$$

where $r = \lambda x - Ax$ is the residual vector corresponding to λ, x and B is the matrix $A - \lambda I$ with column s replaced by $-x$. For use in the analysis (but not in the computation) we may rewrite (1.4) as

$$(1.5) \quad y = \varepsilon b + y_s X \tilde{y},$$

where

$$(1.6) \quad X = B^{-1}, \quad \varepsilon b = Xr, \quad \|b\|_\infty = 1.$$

The factor ε is introduced to emphasize that $\|Xr\|_\infty$ must be at least moderately small, if the algorithm is to be satisfactory.

An essential element in all the algorithms we discuss is the solution of any linear system having as its matrix of coefficients either B as defined above or some generalization of it. Now $n - 1$ columns of B are drawn from $(A - \lambda I)$ and the remaining column is a normalized vector x . If the elements of A are very large or very small compared with unity then B is a badly scaled matrix. The ∞ -condition number $\|B\|_\infty \|B^{-1}\|_\infty$ will be very large when B is badly scaled, independent of whether or not the equations are difficult to solve accurately.

This point is perhaps best illustrated by considering a trivial example. Consider the two systems

$$Px = \begin{bmatrix} .9142 & .9825 \\ -.9475 & .9123 \end{bmatrix}, \quad x = \begin{bmatrix} .4123 \\ .4037 \end{bmatrix},$$

$$Qx = \begin{bmatrix} .9142 & 10^{10}(.9825) \\ -.9475 & 10^{10}(.9123) \end{bmatrix}, \quad y = \begin{bmatrix} .4123 \\ .4037 \end{bmatrix}.$$

The exact solutions are such that

$$x_1 = y_1, \quad x_2 = 10^{10} y_2.$$

If the systems are solved on a decimal computer, these relations are also satisfied exactly by the computed solution. The rounding errors in the mantissa are identical; only the exponents are different. The matrix P is very well conditioned (from any standpoint) and $\|P\|_\infty \|P^{-1}\|_\infty$ is of order unity. On the other hand $\|Q\|_\infty \|Q^{-1}\|_\infty$ is of order 10^{10} . Clearly the ∞ -condition number of B will not really reflect the "difficulty of solving the system" if it is badly scaled. It is the ∞ -condition number of B when A is scaled so that $\|A\|_\infty = O(1)$ that is really relevant to the behavior of our algorithm, and results are quoted in terms that are significant only when A is so scaled. It is clear that there is no need to scale A in this way in a practical algorithm any more than we need to scale Q in our trivial example. However, we must remember that in y defined by (1.3) the s th component give a corrections to λ while the remaining

$n - 1$ components give corrections to components of a normalized vector. In general, errors in the s th component will be acceptable at a level which is $\|A\|_\infty$ times as large as errors in the remaining components. It might be said that the infinity norm is not the appropriate norm when B is badly scaled, and that we should be using some *biased* norm. Although this is true, it amounts to no more than describing a simple problem in more complicated terms.

Relation (1.5) leads quite naturally to consideration of the iterative procedure

$$(1.7) \quad y^{(p+1)} = \epsilon b + y_s^{(p)} X \tilde{y}^{(p)}.$$

In practice (1.7) could be used in two rather different ways.

(i) The initial approximation may already be of quite high accuracy, and one may wish merely to use an analysis of the iterative procedure to demonstrate that the $y^{(p)}$ defined by it tends to a limit y corresponding to an exact eigenpair and to obtain a bound for $\|y - y^{(0)}\|$.

(ii) The iteration may actually be used to compute a succession of the $y^{(p)}$. The analysis of the convergence behaviour would then be employed to obtain a bound for $\|y - y^{(q)}\|$, where $y^{(q)}$ is an iterate which is considered to be of acceptable accuracy.

However, a certain volume of computation is required merely to establish that the conditions are satisfied for the iteration to converge. Once we have made this computational effort, $y^{(1)}$ is available with little additional work. Hence, even when the initial estimate has been derived using a very stable algorithm, one will normally determine $y^{(1)}$ and then obtain a bound for $\|y - y^{(1)}\|$ rather than $\|y - y^{(0)}\|$. Dongarra [1], Wilkinson [5] and Yamamoto [8], [9] have both used the iteration defined by (1.7); we present their results here, modified slightly for convenience.

THEOREM 1. *If $\kappa = \|X\|$ and $\epsilon\kappa < \frac{1}{4}$ then*

$$(1.8) \quad \|y^{(p)}\|_\infty \leq \frac{1 - (1 - 4\epsilon\kappa)^{1/2}}{2\kappa}$$

and $y^{(p)} \rightarrow y$, the solution of (1.4). The convergence is geometric in that

$$(1.9) \quad \|y^{(p+1)} - y^{(p)}\|_\infty \leq \gamma \|y^{(p)} - y^{(p-1)}\|_\infty,$$

where

$$(1.10) \quad \gamma < \frac{2^{3/2} \kappa \epsilon}{[(1 - 2\epsilon\kappa) + (1 - 4\epsilon\kappa)^{1/2}]^{1/2}} < 4\kappa\epsilon < 1.$$

In order to give greater numerical stability in the practical realization, the iteration (1.7) is first recast in the equivalent form:

$$(1.11) \quad \begin{aligned} B\delta^{(0)} &= r, & y^{(1)} &= \delta^{(0)}, \\ B\delta^{(1)} &= y_s^{(1)} \tilde{y}^{(1)}, & y^{(2)} &= y^{(1)} + \delta^{(1)}, \\ B\delta^{(2)} &= y_s^{(2)} \tilde{y}^{(2)} - y_s^{(1)} \tilde{y}^{(1)} = y_s^{(1)} \tilde{\delta}^{(1)} + \delta_s^{(1)} \tilde{y}^{(2)}, & y^{(3)} &= y^{(2)} + \delta^{(2)}, \\ &\dots & &\dots \\ B\delta^{(p)} &= y_s^{(p)} \tilde{y}^{(p)} - y_s^{(p-1)} \tilde{y}^{(p-1)} = y_s^{(p-1)} \tilde{\delta}^{(p-1)} + \delta_s^{(p-1)} \tilde{y}^{(p)}, & y^{(p+1)} &= y^{(p)} + \delta^{(p)}. \end{aligned}$$

Here each correction to y is derived by solving a linear system with the matrix B . To diminish the errors made in solving each of these systems, we include one step of iterative refinement of the solution of the k th system in the solution of the $(k + 1)$ st system. Thus we obtain as the typical equation

$$(1.12) \quad B\delta^{(p)} = [r^{(p-1)} - B\delta^{(p-1)}] + y_s^{(p-1)} \tilde{\delta}^{(p-1)} + \delta_s^{(p-1)} \tilde{y}^{(p)} \equiv r^{(p)},$$

where double precision accumulation of inner product is used in the computation of the residual, $r \equiv r^{(0)}$ and of all the $r^{(p-1)} - B\delta^{(p-1)}$. Rather surprisingly, this technique is just as effective as continuing iterative refinement to its conclusion in each individual step.

2. Extensions of the basic algorithm. The basic algorithm can be extended and/or improved in several directions. To this end we make the following observations.

(i) The convergence rate of the basic iteration and the error bounds depend on the condition of B with respect to inversion. If λ is an approximation to a multiple eigenvalue or to one of a number of close eigenvalues, then B is ill conditioned. (B is singular when λ, x is exact and λ is a multiple eigenvalue.) Both the performance of the algorithm and the error bounds suffer from this ill conditioning, though multiple or pathologically close eigenvalues may be quite well conditioned. (If A is normal, all eigenvalues are well conditioned.) Although in the case of close eigenvalues the individual eigenvectors are ill conditioned, the invariant subspace associated with a cluster is well determined if the cluster is well separated from the remaining eigenvalues. This suggests that an algorithm for finding generators of such invariant subspaces is advisable.

(ii) When λ is one of a set of r ill conditioned eigenvalues (including possibly some defective eigenvectors), one should still be able to determine accurately an $n \times r$ matrix X and an $r \times r$ matrix M such that

$$(2.1) \quad AX = XM,$$

where the columns of X accurately define the relevant invariant subspace [7].

(iii) Although B^{-1} need not be computed explicitly in the basic algorithm, each step requires the solution of a linear system with the matrix B . This requires some stable factorization of B . Thus, if A is a full dense matrix, $O(n^3)$ multiplications and additions are required; and if p approximate eigenpairs are to be improved, $O(pn^3)$ operations are needed. When the approximate eigenpairs have been found by a reduction of A by similarity transformations, the reduced form can be used to achieve a more economical algorithm.

(iv) The basic algorithm uses the same matrix B throughout. It is natural to think in terms of updating λ and x in B at each stage, thereby greatly improving the rate of convergence. This procedure, however, would require a complete refactorization of B at each stage. (If the initial λ, x is an accurate eigenpair, refactorization may not be important because one or at most two iterations may suffice.) Success on the lines discussed in (iii) could make modifications in B less formidable.

(v) When A is real but λ is one of a complex conjugate pair, one would hope that the improvement of λ, x would require only twice as much work as the improvement of a real λ and x . (The factor two is reasonable because two eigenvalues are effectively being determined simultaneously.) Straightforward execution of the algorithm, however, requires four times as much work and storage of an $n \times n$ complex matrix.

In this paper we discuss modifications designed to cover the above weaknesses. It should be appreciated that some of the modifications can be coupled together; to cover them all effectively would require a substantial number of programs.

3. Invariant subspaces (linear elementary divisors). We begin by finding generators that give a good determination of an invariant subspace in the case where the eigenvalues are well conditioned (i.e., A is not close to being defective). For simplicity we restrict ourselves initially to two approximate eigenpairs λ_1, x_1 and λ_2, x_2 , where

$|\lambda_1 - \lambda_2|/\|A\|$ is small and x_1 and x_2 are substantially different, so that the two-space in which they lie is numerically well determined. Here it is not necessary that x_1 and x_2 should be orthogonal but only that the angle between them be substantial; λ_1 and λ_2 may correspond to a double eigenvalue.

Although x_1 and x_2 may have substantial errors, they should belong reasonably accurately in the appropriate two-space. Hence we have

$$(3.1) \quad \begin{aligned} A(x_1 + \tilde{y}_1) &= (\lambda_1 + \mu_{11})(x_1 + \tilde{y}_1) + \mu_{21}(x_2 + \tilde{y}_2), \\ A(x_2 + \tilde{y}_2) &= \mu_{21}(x_1 + \tilde{y}_1) + (\lambda_2 + \mu_{22})(x_2 + \tilde{y}_2) \end{aligned}$$

where \tilde{y}_1 , \tilde{y}_2 and μ_{ij} are expected to be small. Because (3.1) implies that

$$(3.2) \quad A[x_1 + \tilde{y}_1 | x_2 + \tilde{y}_2] = [x_1 + \tilde{y}_1 | x_2 + \tilde{y}_2] \begin{pmatrix} \lambda_1 + \mu_{11} & \mu_{12} \\ \mu_{21} & \lambda_2 + \mu_{22} \end{pmatrix},$$

the vectors $x_1 + \tilde{y}_1$, $x_2 + \tilde{y}_2$ are exact generators of an invariant two-space, the corresponding eigenvalues being those of the 2×2 matrix on the right. We assume that $\|x_1\|_\infty = \|x_2\|_\infty = 1$. To select specific vectors in the subspace, we must prescribe some form of "normalization" of $x_1 + \tilde{y}_1$ and $x_2 + \tilde{y}_2$ analogous to our requirement that $\tilde{y}_s = 0$ in the basic algorithm. We shall require that

$$(3.3) \quad \tilde{y}_{1p} = \tilde{y}_{2p} = \tilde{y}_{1q} = \tilde{y}_{2q} = 0,$$

where p and q are such that

$$(3.4) \quad |x_{1p}| = \max |x_{1i}|$$

and

$$(3.5) \quad |x_{1p}x_{2q} - x_{1q}x_{2p}| = \max |x_{1p}x_{2i} - x_{1i}x_{2p}|.$$

Notice that this selection ensures that $q \neq p$. Further, if the maximum in (3.5) is small, x_1 and x_2 are nearly parallel, and small perturbations in them could not provide us with two vectors giving a good numerical determination of the invariant subspace.

From equation (3.2), we obtain

$$(3.6) \quad \begin{aligned} (A - \lambda_1 I)\tilde{y}_1 - \mu_{11}x_1 - \mu_{21}x_2 &= r_1 + \mu_{11}\tilde{y}_1 + \mu_{21}\tilde{y}_2, \\ (A - \lambda_2 I)\tilde{y}_2 - \mu_{12}x_1 - \mu_{22}x_2 &= r_2 + \mu_{12}\tilde{y}_1 + \mu_{22}\tilde{y}_2, \end{aligned}$$

where

$$(3.7) \quad r_i = \lambda_i x_i - A x_i.$$

Because components p and q of both \tilde{y}_1 and \tilde{y}_2 are zero, it is convenient to define vectors y_1 and y_2 in which two zeros are replaced by μ_{11} and μ_{21} and μ_{12} and μ_{22} , respectively. The equations then become

$$(3.8) \quad B_i y_i = r_i + y_{ip}\tilde{y}_1 + y_{iq}\tilde{y}_2 \quad (i = 1, 2),$$

where B_i is $A - \lambda_i I$ with columns p and q replaced with $-x_1$ and $-x_2$. If the given eigenpairs are reasonably accurate, equation (3.8) is a coupled pair of mildly nonlinear equations and may be solved by the iterative procedure

$$(3.9) \quad B_i y_i^{(s+1)} = r_i + y_{ip}^{(s)} \tilde{y}_1^{(s)} + y_{iq}^{(s)} \tilde{y}_2^{(s)} \quad (i = 1, 2)$$

with $y_i^{(0)} = 0$. (Here we have changed the upper suffix from p , used in (1.7), to s to avoid conflict with p and q above.) These equations may be expressed in the form

$$(3.10) \quad y_i^{(s+1)} = \varepsilon_i b_i + X_i (y_{ip}^{(s)} \tilde{y}_1^{(s)} + y_{iq}^{(s)} \tilde{y}_2^{(s)}) \quad (i = 1, 2),$$

precisely as in (1.5), (1.6), and (1.7).

The analysis of this algorithm is similar to that of the basic algorithm. We have the following theorem:

THEOREM 2. *Let $\kappa_i = \|X_i\|$, $\kappa = \max(\kappa_1, \kappa_2)$, $\varepsilon = \max(\varepsilon_1, \varepsilon_2)$, and*

$$\alpha = 4 / [(1 - 4\varepsilon\kappa) + (1 - 8\varepsilon\kappa)^{1/2}].$$

Then if $\varepsilon\kappa < \frac{1}{8}$, we have

$$\|y_i^{(s)}\|_\infty \leq \varepsilon_i + \alpha\kappa\varepsilon^2 \quad (i = 1, 2)$$

for all s and $y_i^{(s)} \rightarrow y_i$. The convergence is geometric in that

$$\|\delta_1^{(s+1)}\| + \|\delta_2^{(s+1)}\| \leq \gamma [\|\delta_1^{(s)}\| + \|\delta_2^{(s)}\|]$$

where $\gamma = 4\kappa\varepsilon(1 + \alpha\kappa\varepsilon) = 4\kappa\varepsilon(\alpha/2)^{1/2} < 1$.

For numerical purposes equation (3.10) may be recast as in (1.11) and (1.12). The s th equation corresponding to the set (1.11) then becomes

$$(3.11) \quad B_i \delta_i^{(s)} = y_{ip}^{(s-1)} \tilde{\delta}_1^{(s-1)} + \delta_{ip}^{(s-1)} \tilde{y}_1^{(s)} + y_{iq}^{(s-1)} \tilde{\delta}_2^{(s-1)} + \delta_{iq}^{(s-1)} \tilde{y}_2^{(s)} \quad (i = 1, 2),$$

and the s th equation corresponding to the set (1.12) becomes

$$(3.12) \quad \begin{aligned} B_i \delta_i^{(s)} &= [r_i^{(s-1)} - B_i \delta_i^{(s-1)}] + y_{ip}^{(s-1)} \tilde{\delta}_1^{(s-1)} + \delta_{ip}^{(s-1)} \tilde{y}_1^{(s)} + y_{iq}^{(s-1)} \tilde{\delta}_2^{(s-1)} + \delta_{iq}^{(s-1)} \tilde{y}_2^{(s)} \\ &\equiv r_i^{(s)}, \end{aligned}$$

where we have incorporated one stage of iterative refinement in the equations for deriving $\delta_i^{(s-1)}$ in the equation determining $\delta_i^{(s)}$.

When λ_1 and λ_2 are well separated from the other eigenvalues and are well conditioned, the matrices B_1 and B_2 will be well conditioned even though the matrices B arising in the use of the basic algorithm for improving λ_1, x_1 and λ_2, x_2 independently are very ill conditioned. The two generators $x_1 + \tilde{y}_1$ and $x_2 + \tilde{y}_2$ will be accurately determined. Because the final accepted values accurately satisfy the relation (3.2), the eigenvalues of the 2×2 matrix should give very accurate approximations to the eigenvalues. In a KDF9 program in which all the right-hand sides of equations (3.12) are derived using double precision accumulation inner product and the corrections $\delta_i^{(s)}$ are added to the $y_i^{(s)}$ in double precision, the final accuracy is appropriate to double precision computation throughout, even though no multiplication of double precision numbers is used and a very high percentage of all computation is therefore in single precision. When A really does have a double root γ (say) corresponding to linear elementary divisors, we must have

$$(3.13) \quad \begin{pmatrix} \lambda_1 + \mu_{11} & \mu_{12} \\ \mu_{21} & \lambda_2 + \mu_{22} \end{pmatrix} = \begin{pmatrix} \gamma & 0 \\ 0 & \gamma \end{pmatrix}$$

with only double precision errors. Hence the μ_{12} and μ_{21} should be of order β^{-2t} for computation with precision β^{-t} . If the exact λ_1 and λ_2 are such that $|\lambda_1 - \lambda_2| / \|A\| \approx \beta^{-p-t}$, then μ_{21} and μ_{12} should reflect this closeness and be of order $\|A\| \beta^{-p-t}$. Improved individual eigenvectors are found in the KDF9 program from the eigenvalues and eigenvectors of the 2×2 matrix in (3.2).

For simplicity of notation, we have exposed the case of two simple eigenvalues. The algorithm extends immediately to a set of k close eigenvalues. We have then

$$(3.14) \quad A[x_1 + \tilde{y}_1, \dots, x_k + \tilde{y}_k] = [x_1 + \tilde{y}_1, \dots, x_k + \tilde{y}_k][\text{diag}(\lambda_i) + M],$$

where $m_{ij} = \mu_{ij}$ and it is expected that both \tilde{y}_i and μ_{ij} will be small. We now obtain a set of k loosely coupled nonlinear equations, the matrix B_i associated with the i th set being $A - \lambda_i I$ with k of its columns replaced by $-x_1, -x_2, \dots, -x_k$. The only new complication is how to determine the k elements of \tilde{y}_i that are to be zero. The choice can be made as follows. Let X be the $k \times n$ matrix with rows x_i . Let this be reduced to upper-trapezoidal form using Gaussian elimination with column pivoting, rather than the row pivoting involved in the standard partial pivoting algorithm. If the relevant pivotal elements are in columns p_1, p_2, \dots, p_k , respectively, then these elements are to be zero in the \tilde{y}_i . (The last p_k is chosen to be the maximum element in the final reduced row although no further reduction is to be done at this point. It will readily be verified that when $k = 2$, this gives the choice which we have described.)

At each step in the iterative solution of the nonlinear equations, we have to solve k linear systems of order n with matrices B_i ($i = 1, \dots, k$); in addition, we have the initial factorization of the B_i . The total amount of work is only marginally greater than that in the separate improvement of each of the λ_i using the basic algorithm. Indeed, if some of the approximate λ_i 's are equal, then corresponding to these we have only one B_i to factorize. Furthermore, if m of the λ_i are almost equal, we can start by replacing each of these by the mean of the m values. The invariant subspace algorithm can therefore be substantially more economical than the basic algorithm.

4. Invariant subspaces (almost defective matrices). When A is defective (or almost defective), the computed eigenvectors corresponding to the relevant eigenvalues will be almost linear dependent. For a set of vectors x_1, x_2, \dots, x_k this near linear dependence will become apparent when Gaussian elimination with column pivoting is performed. Small corrections to such a set of x_i will serve little purpose. Instead, to achieve a good determination of the invariant subspace, we proceed as follows.

Again for simplicity we concentrate on just two eigenvalues. Clearly it is essential to start with two well separated generators x_1, x_2 of the invariant subspace. We should start then with approximate x_1, x_2 and a 2×2 matrix M such that

$$(4.1) \quad A[x_1 x_2] \approx [x_1 x_2]M \approx [x_1 x_2] \begin{pmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{pmatrix},$$

and attempt to determine \tilde{y}_1, \tilde{y}_2 and μ_{ij} ($i, j = 1, 2$) such that

$$(4.2) \quad A[x_1 + \tilde{y}_1, x_2 + \tilde{y}_2] = [x_1 + \tilde{y}_1, x_2 + \tilde{y}_2] \begin{pmatrix} m_{11} + \mu_{11} & m_{12} + \mu_{12} \\ m_{21} + \mu_{21} & m_{22} + \mu_{22} \end{pmatrix}.$$

If we attempt to derive an algorithm based on this relation, however, we find that the pair of equations is no longer loosely coupled. A much more effective algorithm can be produced if $m_{21} = 0$ in (4.1), Equation (4.2) then gives

$$(4.3) \quad \begin{aligned} (A - m_{11}I)\tilde{y}_1 - \mu_{11}x_1 - \mu_{21}x_2 &= r_1 + \mu_{11}\tilde{y}_1 + \mu_{21}\tilde{y}_2, \\ m_{12}\tilde{y}_1 + (A - m_{22}I)\tilde{y}_2 - \mu_{12}x_1 - \mu_{22}x_2 &= r_2 + \mu_{12}\tilde{y}_1 + \mu_{22}\tilde{y}_2, \end{aligned}$$

where

$$(4.4) \quad r_1 = m_{11}x_1 - Ax_1, \quad r_2 = m_{12}x_1 + m_{22}x_2 - Ax_2$$

and both r_1 and r_2 are expected to be small. These equations can be expressed in simpler form as in (3.8) using B_1, B_2, y_1 , and y_2 . Thus we have obtained the iterative procedure

$$(4.5) \quad B_1 y_1^{(s+1)} = r_1 + y_{1p}^{(s)} \tilde{y}_1^{(s)} + y_{2q}^{(s)} \tilde{y}_2^{(s)},$$

$$(4.6) \quad B_2 y_2^{(s+1)} = r_2 - m_{12} \tilde{y}_1^{(s+1)} + y_{2p}^{(s)} \tilde{y}_1^{(s)} + y_{2q}^{(s)} \tilde{y}_2^{(s)}$$

in which $y_1^{(s+1)}$ is derived from (4.5) and then $y_2^{(s+1)}$ is derived from (4.6) using the $y_1^{(s+1)}$ just computed. Again we have only two matrices to factor and two linear systems of order n to solve in each step. The reformulation of the equations in the interest of numerical stability proceeds on the lines already established. The algorithm gives very accurate generators of the invariant subspace, but although the 2×2 matrix on the right of (4.2) is accurate, this accuracy is not reflected in the two eigenvalues because its eigenvalues are necessarily sensitive to small perturbations in its elements.

There remains the task of determining generators x_1 and x_2 which correspond to a zero value of m_{21} . This can be done when the eigenvalue problem of A has been solved by the orthogonal triangularization of A , as in Francis's double QR algorithm. We assume that

$$(4.7) \quad Q^T A Q = T,$$

and that the two approximate eigenvalues λ_1 and λ_2 are the diagonal elements $t_{l,l}$ and $t_{m,m}$, respectively, with $l < m$. The two eigenvectors of T may be found immediately in the form

$$(4.8) \quad z_{1,1}, z_{1,2}, \dots, z_{1,l-1}, 1, 0, \dots, 0$$

and

$$(4.9) \quad z_{2,1}, z_{2,2}, \dots, z_{2,m-1}, 1, 0, \dots, 0$$

by solving the two relevant triangular systems $(T - \lambda_1 I)z_1 = 0$, $(T - \lambda_2 I)z_2 = 0$, though in our case, of course, they will be almost parallel. In fact, we find z_1 in this way but determine z_2 as the solution of

$$(4.10) \quad (T - \lambda_2 I)z_2 = k z_1,$$

where k is chosen so that $z_{2,l} = 0$. Clearly k is given by

$$(4.11) \quad t_{l,l+1} z_{2,l+1} + t_{l,l+2} z_{2,l+2} + \dots + t_{l,m} z_{2,m} = k$$

and is determined when we reach element $z_{2,l}$ of z_2 . Then $x_1 = Qz_1$, $x_2 = Qz_2$, satisfying our relation with

$$(4.12) \quad m_{11} = \lambda_1, \quad m_{12} = k, \quad m_{22} = \lambda_2.$$

The technique extends immediately to a set of k ill conditioned eigenvalues. From the triangular matrix T we can determine z_1, z_2, \dots, z_k such that

$$(4.13) \quad (T - \lambda_1 I)z_1 = 0, \quad (T - \lambda_2 I)z_2 = m_{12} z_1, \quad (T - \lambda_3 I)z_3 = m_{13} z_1 + m_{23} z_2,$$

where each m_{ij} is chosen so that the appropriate element of z_i is zero. In general the z_i determined in this way will be highly linearly independent, and we may take $x_i = Qz_i$.

We can now define an iterative process in each step in which we solve k sets of linear equations with matrices of order n , the left-hand side being derived by

$$(4.14) \quad \begin{aligned} & B_1 y_1^{(s+1)}, \\ & B_2 y_2^{(s+1)} + m_{12} \tilde{y}_1^{(s+1)}, \\ & B_s y_3^{(s+1)} + m_{13} \tilde{y}_1^{(s+1)} + m_{23} \tilde{y}_2^{(s+1)}, \\ & \dots \end{aligned}$$

so that $y_1^{(s+1)}, y_2^{(s+1)}, \dots$ may be found in succession. The work involved is only marginally greater than in determining k individual eigenpairs by the basic algorithm. If any of the λ_i are equal, the same is true of the corresponding B_i , and the volume of work is correspondingly reduced.

5. Use of orthogonal triangularization of real A . We now turn to the case where the original approximate eigenpairs have been found by Francis' double QR algorithm. We have an orthogonal matrix Q such that

$$(5.1) \quad A = QTQ^T,$$

where T is quasi upper triangular, that is, triangular apart from possible 2×2 diagonal blocks corresponding to complex conjugate eigenvalues. For simplicity we assume that T is truly upper triangular. In presenting the basic algorithm we solved a succession of linear systems of the form

$$(5.2) \quad Bw = g$$

with various right-hand sides g . It will be convenient to introduce the generic notation $Z - \lambda I = Z_\lambda$ and $(Z - \lambda I)e_s = Z_\lambda e_s = z_{\lambda s}$. Equation (5.2) may then be written in the form

$$(5.3) \quad [A_\lambda - (x + a_{\lambda s})e_s^T]w = (A_\lambda + ce_s^T)w = g,$$

where

$$(5.4) \quad c = -x - a_{\lambda s}.$$

From (5.1) we have

$$(5.5) \quad Q[T_\lambda + Q^T ce_s^T Q]Q^T w = g,$$

giving

$$(5.6) \quad (T_\lambda + df^T)Q^T w = Q^T g,$$

where

$$(5.7) \quad d = Q^T c, \quad f^T = e_s^T Q.$$

The matrix df^T in (5.6) is a rank one modification of the triangular matrix T_λ . To solve this system, we need to re-triangularize $T_\lambda + df^T$. Accordingly, we premultiply the system by two orthogonal matrices Q_1 and Q_2 , giving

$$(5.8) \quad Q_2 Q_1 (T_\lambda + df^T) Q^T w = Q_2 Q_1 Q^T g,$$

where Q_1 and Q_2 are products of elementary plane rotations determined as follows.

The matrix Q_1 is such that

$$(5.9) \quad Q_1 d = (P_2 P_3 \cdots P_n) d = \gamma e_1 \quad \text{where } \gamma = \|d\|_2$$

and P_i is a rotation in the $(i-1, i)$ plane designed to annihilate the i th component of $P_{i+1}P_{i+2} \cdots P_n d$. We have

$$(5.10) \quad Q_1(T_\lambda + df^T) = Q_1 T_\lambda + \gamma e_1 f^T.$$

$Q_1 T_\lambda$ is upper Hessenberg while $\gamma e_1 f^T$ is null except in the first row. Hence the right-hand side of (5.10) is also an upper Hessenberg matrix H . H may now be reduced to upper triangular form, \bar{T}_λ , by premultiplication with Q_2 defined by

$$(5.11) \quad Q_2 = P'_n \cdots P'_3 P'_2,$$

where the premultiplication by P'_i annihilates the element $(i, i-1)$ of the current matrix by a rotation in the $(i-1, i)$ plane. Hence we have left to solve the triangular system

$$(5.12) \quad \bar{T}_\lambda Q^T w = Q_2 Q_1 Q^T g.$$

Solution of a system with the matrix B may thus be solved in $O(n^2)$ operations.

6. Updating of the matrix B . When the orthogonal triangularization of A is used, it becomes practical to update the matrix B at each stage of the iteration using the current approximation to the eigenpair. Accordingly, we treat the $(p+1)$ st step of the iteration as though it were the first step in the basic iteration starting with values $\lambda^{(p)}$ and $x^{(p)}$. The algorithm then becomes

$$(6.1) \quad (A - \lambda^{(p)} I) \bar{\delta}^{(p)} - \delta_s^{(p)} x^{(p)} = r^{(p)} = (\lambda^{(p)} I - A) x^{(p)}$$

where

$$(6.2) \quad x^{(p+1)} = x^{(p)} + \bar{\delta}^{(p)}, \quad \lambda^{(p+1)} = \lambda^{(p)} + \delta_s^{(p)}.$$

We may rewrite this as

$$(6.3) \quad B^{(p)} \bar{\delta}^{(p)} = r^{(p)}$$

where

$$(6.4) \quad B^{(p)} = A - \lambda^{(p)} I + c^{(p)} e_s^T, \quad c^{(p)} = -x - a_{\lambda s}^{(p)}.$$

We now write

$$(6.5) \quad B^{(p)} = Q(T_\lambda^{(p)} + Q^T c^{(p)} e_s^T Q) Q^T = Q(T_\lambda^{(p)} + d^{(p)} f^T) Q^T$$

and solve (6.3) by the triangular system

$$(6.6) \quad \bar{T}_\lambda^{(p)} Q^T \bar{\delta}^{(p)} = Q_2^{(p)} Q_1^{(p)} Q^T r^{(p)}.$$

Note that Q and f will be independent of p if s is not changing from one iteration to the next. The rotations involved in $Q_1^{(p)}$ and $Q_2^{(p)}$, on the other hand, differ from one iteration to the next, but because the number of operations in each re-triangularization is $O(n^2)$ and some $n^2/2$ multiplication and additions are necessarily involved in the solution of a triangular system, this is quite acceptable.

7. Convergence of the updated iteration process. We can now consider the convergence of the updated iteration process. This process is precisely Newton's method applied to the system

$$(7.1) \quad (A - \lambda I)x = 0, \quad e_s^T x = 0.$$

(In fact, the basic algorithm itself is merely a recasting of the simplified Newton method in which the Jacobian matrix is not updated.) Convergence could therefore

be analyzed using the Newton–Kantorovich theorem [3]. Such an analysis would lead to a result of the following form:

If the initial error is small enough, the iteration will converge and the convergence is quadratic.

However, in the rounding error analysis given later we shall need most of the intermediate results which are derived in proving the Kantorovich theorem. It is therefore more convenient to analyze convergence directly from first principles. In doing so we can take advantage of the simple form of the equations (7.1).

First we introduce the notation used in analyzing both the exact iterates and the computed iterates. We shall assume that all approximate eigenpairs x, λ —whether exact or computed—are such that $x_s = 1$. Hence the difference between any two x 's has a zero in position s , and we can replace this with the difference between the λ 's. All the information can therefore be provided in a single error vector ξ . We may write

$$(7.2) \quad (x_2, \lambda_2) - (x_1, \lambda_1) = \xi$$

where $\tilde{\xi} = x_2 - x_1$, $\xi_s = \lambda_2 - \lambda_1$, but for brevity we shall often write

$$(7.3) \quad x_2 - x_1 = \xi.$$

Corresponding to any (x, λ) is a B matrix defined by

$$(7.4) \quad B = A - \lambda I - (a_{\lambda s} + x)e_s^T.$$

The B matrices corresponding to x_1, λ_1 and $\bar{x}, \bar{\lambda}$ will be denoted by B_1 and \bar{B} , respectively. Note that

$$(7.5) \quad By = (A - \lambda I)\tilde{y} - y_s x \quad \text{for all } y.$$

We begin the analysis with the following lemmas about a single step of the iteration.

LEMMA 1. *If x_1, λ_1 and x_2, λ_2 are any two approximate eigenpairs, the residuals r_1 and r_2 are defined by*

$$(7.6) \quad r_i = \lambda_i x_i - Ax_i \quad (i = 1, 2)$$

and ξ is the difference between their approximations. Then

$$(7.7) \quad r_2 = r_1 - B_2(x_2 - x_1) - \xi_s \tilde{\xi},$$

where B_2 is the B matrix corresponding to λ_2 and x_2 .

Proof.

$$(7.8) \quad \begin{aligned} r_2 &= \lambda_2 x_2 - Ax_2 = \lambda_2(x_1 + \tilde{\xi}) - A(x_1 + \tilde{\xi}) \\ &= (\lambda_2 I - A)\tilde{\xi} + \lambda_2 x_1 - Ax_1 \\ &= (\lambda_2 I - A)\tilde{\xi} + \lambda_1 x_1 - Ax_1 + \xi_s x_1 \\ &= (\lambda_2 I - A)\tilde{\xi} + (\lambda_1 x_1 - Ax_1) + \xi_s x_2 - \xi_s \tilde{\xi} \\ &= r_1 - B_2 \xi - \xi_s \tilde{\xi} = r_1 - B_2(x_2 - x_1) - \xi_s \tilde{\xi} \end{aligned}$$

because $B_2 \xi = (A - \lambda_2 I)\tilde{\xi} - \xi_s x_2$. \square

LEMMA 2. *With the same notation as in Lemma 1, if x_1 and x_2 are used as initial values in one step of the iteration process to and give corrections of δ_1 and δ_2 , respectively, then*

$$(7.9) \quad \delta_2 = \delta_1 - \xi + (B_2^{-1} - B_1^{-1})r_1 - B_2^{-1} \xi_s \tilde{\xi}.$$

Proof. By definition,

$$(7.10) \quad B_1 \delta_1 = r_1, \quad B_2 \delta_2 = r_2.$$

From Lemma 1 we obtain the second term

$$(7.11) \quad B_2 \delta_2 = r_1 - B_2 \xi - \xi_s \tilde{\xi},$$

$$(7.12) \quad \delta_2 = B_2^{-1} r_1 - \xi - B_2^{-1} \xi_s \tilde{\xi} = \delta_1 - \xi + (B_2^{-1} - B_1^{-1}) r_1 - B_2^{-1} \xi_s \tilde{\xi}.$$

Note that the term $-\xi$ on the right of (7.12) cancels the initial difference between (x_1, λ_1) and (x_2, λ_2) . This difference is replaced by

$$(7.13) \quad (B_2^{-1} - B_1^{-1}) r_1 - B_2^{-1} \xi_s \tilde{\xi}.$$

If the two approximations are close, ξ is small and $\|\xi_s \tilde{\xi}\| \leq \|\xi\|^2$. Hence, provided we have a satisfactory bound on $\|B_2^{-1}\|$ the second term in (7.13) is promising. For the first term we have the following:

LEMMA 3. *If B_1^{-1} exists and $2\|B_1^{-1}\|\|\xi\| < 1$ then*

$$(7.14) \quad \|B_2 - B_1\| \leq 2\|\xi\|$$

and

$$(7.15) \quad \|(B_2^{-1} - B_1^{-1})u\| \leq \frac{2\|B_1^{-1}\|\|\xi\|\|B_1^{-1}u\|}{1 - 2\|B_1^{-1}\|\|\xi\|}.$$

Proof.

$$(7.16) \quad (B_2 - B_1) = (\lambda_1 - \lambda_2)I - (x_2 - x_1)e_s^T = -\xi_s I - \tilde{\xi}_s^T,$$

$$(7.17) \quad \|B_2 - B_1\| \leq |\xi_s| + \|\tilde{\xi}\| \leq 2\|\xi\|.$$

For the second inequality we have

$$(7.18) \quad (B_2^{-1} - B_1^{-1})u = (B_2^{-1}B_1 - I)B_1^{-1}u,$$

$$(7.19) \quad B_2^{-1}B_1 = (B_1 + E)^{-1}B_1 = I + F \quad \text{where } E = B_2 - B_1$$

where

$$(7.20) \quad \|F\| \leq \|B_1^{-1}\|\|E\|/(1 - \|B_1^{-1}\|\|E\|) \leq 2\|B_1^{-1}\|\|\xi\|/(1 - 2\|B_1^{-1}\|\|\xi\|).$$

Equation (7.15) follows immediately. Using Lemma 3, we have

$$(7.21) \quad \begin{aligned} \|(B_2^{-1} - B_1^{-1})r_1\| &\leq 2\|B_1^{-1}\|\|\xi\|\|B_1^{-1}r_1\|/(1 - 2\|B_1^{-1}\|\|\xi\|) \\ &= 2\|B_1^{-1}\|\|\xi\|\|\delta_1\|/(1 - 2\|B_1^{-1}\|\|\xi\|). \end{aligned}$$

LEMMA 4. *If $\bar{x}, \bar{\lambda}$ is an approximation to an exact eigenpair (x, λ) , the difference being ξ , then*

$$(7.22) \quad \bar{r} = -\bar{B}\xi - \xi_s \tilde{\xi},$$

and the correction provided by one step of iteration is

$$(7.23) \quad -\xi - \bar{B}^{-1}\xi_s \tilde{\xi},$$

giving an approximation with error $-\bar{B}^{-1}\xi_s \tilde{\xi}$.

Proof. Using Lemma 1 with (x, λ) and $(\bar{x}, \bar{\lambda})$ as the two approximations, we have

$$(7.24) \quad \bar{r} = r - \bar{B}\xi - \xi_s \tilde{\xi} = -\bar{B}\xi - \xi_s \tilde{\xi} \quad \text{since } r = 0.$$

One step of iteration gives

$$(7.25) \quad \overline{B}\delta = \bar{r} = -\bar{B}\xi - \xi_s \bar{\xi}, \quad \bar{\delta} = -\xi - \bar{B}^{-1} \xi_s \bar{\xi}.$$

Now the update iteration gives a sequence of $x^{(p)}$, $\lambda^{(p)}$ defined by

$$(7.26) \quad \begin{aligned} B^{(0)}\delta^{(0)} &= r^{(0)}, & x^{(1)} &= x^{(0)} + \tilde{\delta}^{(0)}, & \lambda^{(1)} &= \lambda^{(0)} + \delta_s^{(0)}, \\ B^{(1)}\delta^{(1)} &= r^{(1)}, & x^{(2)} &= x^{(1)} + \tilde{\delta}^{(1)}, & \lambda^{(2)} &= \lambda^{(1)} + \delta_s^{(1)}, \\ & \dots & & \dots & & \\ B^{(p)}\delta^{(p)} &= r^{(p)}, & x^{(p+1)} &= x^{(p)} + \tilde{\delta}^{(p)}, & \lambda^{(p+1)} &= \lambda^{(p)} + \delta_s^{(p)}. \end{aligned}$$

We have the following lemma:

LEMMA 5. *If the $x^{(i)}$, $\lambda^{(i)}$ are iterates provided by the algorithm starting from $x^{(0)}$, $\lambda^{(0)}$, then*

$$(7.27) \quad r^{(p+1)} = \delta_s^{(p)} \tilde{\delta}^{(p)}.$$

Proof. Applying Lemma 1 with x_1 , $\lambda_1 = x^{(p+1)}$, $\lambda^{(p+1)}$ and x_2 , $\lambda_2 = x^{(p)}$, $\lambda^{(p)}$ so that $\xi = -\delta^{(p)}$ we have

$$(7.28) \quad r^{(p)} = r^{(p+1)} - B^{(p)}(-\delta^{(p)}) - \delta_s^{(p)} \tilde{\delta}^{(p)}$$

giving

$$(7.29) \quad r^{(p+1)} = r^{(p)} - B^{(p)}\delta^{(p)} + \delta_s^{(p)} \tilde{\delta}^{(p)} = \delta_s^{(p)} \tilde{\delta}^{(p)}$$

because by definition $B^{(p)}\delta^{(p)} = r^{(p)}$. Hence, the exact algorithm is defined by

$$(7.30) \quad B^{(0)}\delta^{(0)} = r^{(0)}, \quad B^{(p)}\delta^{(p)} = \delta_s^{(p-1)} \tilde{\delta}^{(p-1)} \quad (p = 1, 2, \dots).$$

We are now in a position to combine these results for a single step to establish the quadratic convergence of the overall iterative process. We naturally assume that $B^{(0)}$ is nonsingular. For the process to be defined, all $B^{(p)}$ must be nonsingular. From Lemma 3 we know

$$(7.31) \quad \|B^{(p)} - B^{(0)}\| \leq 2\|\delta^{(0)} + \delta^{(1)} + \dots + \delta^{(p-1)}\|,$$

$$(7.32) \quad \|B^{(p)} - B^{(p-1)}\| \leq 2\|\delta^{(p-1)}\|.$$

We note that $B^{(p)}$ will certainly be nonsingular if

$$(7.33) \quad 2\|(B^{(p)})^{-1}\| \|\delta^{(0)} + \delta^{(1)} + \dots + \delta^{(p-1)}\| < 1$$

and hence certainly if

$$(7.34) \quad 2\|(B^{(0)})^{-1}\| (\|\delta^{(0)}\| + \|\delta^{(1)}\| + \dots + \|\delta^{(p-1)}\|) < 1.$$

We write $\|(B^{(0)})^{-1}\| = \kappa^{(0)} = \kappa$ and $\|\delta^{(0)}\| = \varepsilon = \eta^{(0)}$. We see immediately that unless $2\kappa\varepsilon < 1$, even $B^{(1)}$ may be singular. It is intuitively obvious that for convergence we require $\kappa\varepsilon$ to be small enough.

We introduce the quantities $\eta^{(p)}$ and $\kappa^{(p)}$ defined by

$$(7.35) \quad \kappa^{(p)} = \kappa^{(p-1)} / (1 - 2\kappa^{(p-1)}\eta^{(p-1)}), \quad \eta^{(p)} = \kappa^{(p)}(\eta^{(p-1)})^2 \quad (p = 1, 2, \dots).$$

Provided all $\kappa^{(p)}$ are positive, it is evident from (7.30) and (7.32) that the $\kappa^{(p)}$ and $\eta^{(p)}$ majorize $\|(B^{(p)})^{-1}\|$ and $\|\delta^{(p)}\|$ respectively, and ensure the existence of the former.

Quadratic convergence of the iteration can therefore be established by proving that the $\kappa^{(p)}$ remain and that the $\eta^{(p)}$ approach zero quadratically. To this end, we introduce

$$(7.36) \quad \beta_p = \kappa^{(p)}\eta^{(p)},$$

so that

$$(7.37) \quad \kappa^{(p)} = \kappa^{(p-1)} / (1 - 2\beta_{p-1}).$$

It is certainly inadvisable that the β_i should increase because this endangers the nonsingularity of $B^{(p)}$. If $\beta_0 = \beta_1$ we have

$$(7.38) \quad \beta_0 = \beta_1 = \kappa^{(1)} \eta^{(1)} = (\kappa^{(1)} \eta^{(0)})^2 = (\kappa^{(0)} \eta^{(0)} / (1 - 2\beta_0))^2 = (\beta_0 / (1 - 2\beta_0))^2,$$

giving $\beta_0 = \frac{1}{4}$. (The solutions $\beta_0 = 0$, $\beta_1 = -1$ are of no interest.) Because in general

$$(7.39) \quad \beta_p = \kappa^{(p)} \eta^{(p)} = (\kappa^{(p)} \eta^{(p-1)})^2 = (\kappa^{(p-1)} \eta^{(p-1)} / (1 - 2\beta_{p-1}))^2 = (\beta_{p-1} / (1 - 2\beta_{p-1}))^2,$$

we have $\beta_p = \frac{1}{4}$ ($p = 1, 2, \dots$). On the other hand,

$$(7.40) \quad \eta^{(p)} = \kappa^{(p)} (\eta^{(p-1)})^2 = [\kappa^{(p-1)} \eta^{(p-1)} / (1 - 2\kappa^{(p-1)} \eta^{(p-1)})] \eta^{(p-1)} = \frac{1}{2} \eta^{(p-1)}.$$

Hence,

$$(7.41) \quad \eta^{(0)} + \eta^{(1)} + \dots + \eta^{(p-1)} = \eta^{(0)} [1 + \frac{1}{2} + \dots + \frac{1}{2}^{p-1}] = \eta^{(0)} [2 - \frac{1}{2}^{p-1}],$$

and $\sum_{i=0}^{\infty} \eta^{(i)}$ converges to $2\eta^{(0)}$. For the $\kappa^{(p)}$ we have

$$(7.42) \quad \kappa^{(p)} = \frac{1}{4} \eta^{(p)} = 2^{p-2} / \eta^{(0)},$$

showing that although all $B^{(p)}$ are certainly nonsingular, they may be tending to singularity. The value $\beta_0 = \frac{1}{4}$ is therefore a borderline case.

When $\beta_0 < \frac{1}{4}$, (7.38) shows that $\beta_1 < \frac{1}{4}$, and then (7.39) shows by induction that $\beta_p < \frac{1}{4}$. Similarly $\eta^{(p)} < \frac{1}{2} \eta^{(p-1)}$, $\sum_{i=0}^{\infty} \eta^{(i)}$ converges, and

$$(7.43) \quad \sum_{i=0}^{\infty} \eta^{(i)} < 2\eta^{(0)} = 2\varepsilon.$$

This convergence, however, conceals the essential difference between the condition $\beta_0 = \frac{1}{4}$ and $\beta_0 < \frac{1}{4}$. From (7.39) we have

$$(7.44) \quad \beta_p < 4\beta_{p-1}^2 < 4^3\beta_{p-2}^4 < \dots < \frac{1}{4}(4\beta_0)^{2^p}$$

showing that β_p is converging quadratically to zero. Because

$$(7.45) \quad \beta_p = (1 - 2\beta_{p-1})^{-2} \beta_{p-1}^2 \approx \beta_{p-1}^2,$$

we see that even (7.43) severely underestimates the convergence rate.

For the $\kappa^{(p)}$, repeated use of (7.35) gives:

$$(7.46) \quad \begin{aligned} \kappa^{(p)} &= \kappa^{(p-1)} / [1 - 2\kappa^{(p-1)} \eta^{(p-1)}] = \kappa^{(p-2)} / [1 - 2\kappa^{(p-2)} (\eta^{(p-1)} + \eta^{(p-2)})] \\ &= \kappa^{(0)} / [1 - 2\kappa^{(0)} [\eta^{(0)} + \eta^{(1)} + \dots + \eta^{(p-1)}]]. \end{aligned}$$

Because $\sum \eta^{(i)}$ converges and its sum is less than 2ε , we see that $\kappa^{(p)} \rightarrow \kappa^{(\infty)}$, which is finite. However,

$$(7.47) \quad \kappa^{(p)} \eta^{(p)} = \beta_p < (4\beta_0)^{2^p} / 4,$$

and hence $\eta^{(p)}$ also tends quadratically to zero. When $\kappa\varepsilon$ is significantly less than $\frac{1}{4}$, we would expect $\sum \eta^{(i)}$ to be significantly less than 2ε , and this is indeed true. To establish a simple explicit expression for $\sum \eta^{(i)}$, let

$$(7.48) \quad f_p = [1 - (1 - 4\beta_p)^{1/2}] / 2\kappa_p.$$

We establish a simple relation between successive f_i . From (7.39),

$$(7.49) \quad 1 - 4\beta_p = (1 - 4\beta_{p-1}) / (1 - 2\beta_{p-1})^2,$$

$$(7.50) \quad 1 - (1 - 4\beta_p)^{1/2} = 1 - \frac{(1 - 4\beta_{p-1})^{1/2}}{(1 - 2\beta_{p-1})} = \frac{1 - 2\beta_{p-1} - (1 - 4\beta_{p-1})^{1/2}}{1 - 2\beta_{p-1}},$$

$$(7.51) \quad f_p = \frac{1 - (1 - 4\beta_p)^{1/2}}{2\kappa^{(p)}} = \frac{1 - 2\beta_{p-1} - (1 - 4\beta_{p-1})^{1/2}}{2\kappa^{(p)}(1 - 2\beta_{p-1})} = \frac{1 - 2\beta_{p-1} - (1 - 4\beta_{p-1})^{1/2}}{2\kappa^{(p-1)}}.$$

Hence, because $\beta_{p-1} = \kappa^{(p-1)}\eta^{(p-1)}$, we have

$$(7.52) \quad f_p = f_{p-1} - \eta^{(p-1)},$$

giving

$$(7.53) \quad \eta^{(0)} + \eta^{(1)} + \dots + \eta^{(p-1)} = f_0 - f_p.$$

But $f_p \rightarrow 0$ (because $\beta_p \rightarrow 0$) and $\kappa^{(p)} \rightarrow \kappa^{(\infty)}$. Hence,

$$(7.54) \quad \sum_{i=0}^{\infty} \eta^{(i)} = f_0$$

$$(7.55) \quad = [1 - (1 - 4\beta_0)^{1/2}] / 2\kappa_0$$

$$(7.56) \quad = \frac{1 - (1 - 4\beta_0)^{1/2}}{2\beta_0} \varepsilon$$

$$(7.57) \quad = \frac{2}{1 + (1 - 4\beta_0)^{1/2}} \varepsilon.$$

This analysis can be summarized in the following theorem:

THEOREM 3. *When $\beta_0 = \kappa\varepsilon < \frac{1}{4}$, then $\kappa^{(p)}$ is bounded and $\eta^{(p)}$ and β_p approach zero quadratically. In fact,*

$$(7.58) \quad \sum_{i=0}^{\infty} \eta^{(i)} = \frac{2}{1 + (1 - 4\beta_0)^{1/2}} \varepsilon.$$

Note that we may regard the process as starting with any of the $x^{(p)}$. Because $\beta_p < \frac{1}{4}$, for all p when $\beta_0 < \frac{1}{4}$ we see that

$$(7.59) \quad \|x^{(q)} - x^{(p)}\| \leq f_p, \quad q \geq p$$

and

$$(7.60) \quad \|x - x^{(p)}\| \leq f_p,$$

where x is the limit, i.e. the exact solution. In fact the process will often converge when $\beta_0 > \frac{1}{4}$ and some later β_p will satisfy the requirement. At that stage we shall have a ball containing all subsequent iterations and the limit.

8. Complex eigenvalues for real matrices. When we have a real matrix with complex eigenvalues, the previously developed approach for improving the accuracy runs into a problem. While we could use the procedures as described with complex arithmetic throughout, we would end up doing four times as much computation and using twice as much storage. The various components needed in solving this problem for the most part are real; only the diagonal of $T - \lambda I$ and the vector x are complex.

In the real eigenvalue case, in order to find the improvements we need to solve a system based on the matrix

$$(8.1) \quad \begin{pmatrix} A - \lambda I & -x \\ e_s^T & 0 \end{pmatrix}.$$

This matrix can be transformed by $\begin{pmatrix} Q & 0 \\ 0 & 1 \end{pmatrix}$ and its transpose to arrive at

$$(8.2) \quad \bar{T} = \begin{pmatrix} T - \lambda I & p \\ q^T & 0 \end{pmatrix}$$

where $p = -Q^T x$ and $q = Q^T e_s$. One diagonal element of $T - \lambda I$ is zero, say at the k th position. We will replace this zero by the value 1 through a rank one change and remove the row q^T similarly. Then the resulting matrix, say \bar{T}_+ has the form

$$(8.3) \quad \bar{T}_+ = \begin{pmatrix} T_+ & p \\ 0 & 1 \end{pmatrix}$$

where $T_+ = T - \lambda I + e_k e_k^T$, and \bar{T}_+ differs from \bar{T} by a rank 2 change.

When λ and x are complex the matrix T is real and quasi-triangular, Q is real and $T - \lambda I$ is complex on the diagonal only. $T - \lambda I$ has a singular 2×2 block on the diagonal (this corresponds in the real case to a zero on the diagonal). The 2×2 block has the form

$$(8.4) \quad \begin{pmatrix} \alpha - (\lambda_r + i\lambda_i) & b \\ c & d - (\lambda_r + i\lambda_i) \end{pmatrix}.$$

To force this block to be nonsingular, a rank one change is made by adding 1 to the 1, 1 element of that block. The resulting matrix is $T_+ = T - \lambda I + e_k e_k^T$. The matrix then has the form

$$(8.5) \quad \begin{pmatrix} T_+ & p \\ q^T & 1 \end{pmatrix}.$$

The row q^T is removed by a rank one change to arrive at

$$(8.6) \quad \bar{T}_+ = \begin{pmatrix} T_+ & p \\ 0 & 1 \end{pmatrix}.$$

We wish to solve systems of the form $\bar{T}_+ z = v$. For this system the matrix \bar{T}_+ is real except for its diagonal and last column, and the right-hand side vector v is complex. Since \bar{T}_+ is almost completely real, hardly any complex arithmetic is involved. To correspond to a 1×1 diagonal block of \bar{T}_+ we have to solve

$$(8.7) \quad (t_{kk} - \lambda_r - i\lambda_i)z_k = v_k.$$

This will involve a complex division of v_k by the quantity $t_{kk} - \lambda_r - i\lambda_i$. For the 2×2 diagonal block associated with $\lambda_r + i\lambda_i$, a 2×2 complex linear system will have to be solved.

By using such a procedure the work needed in this case goes up by a factor of two over the case where there is a simple real eigenvalue. This factor of two is not surprising since the improvement process will produce an improved eigenpair and its conjugate. The total additional storage needed will be modest, only a few additional vectors.

9. Influence of rounding errors. As a result of our calculations, we have a sequence of computed $\bar{x}^{(p)}, \bar{\lambda}^{(p)}$ that are contaminated by rounding errors. One might attempt to find a bound for each $\|\delta^{(p)} - \bar{\delta}^{(p)}\|$ and then sum them to obtain a bound at any stage for $\|x^{(i)} - \bar{x}^{(i)}\|$. This procedure would be reasonable if the iterates were steadily drifting apart. In fact, however, even if some intermediate $\bar{x}^{(p)}, \bar{\lambda}^{(p)}$ is appreciably distant from $x^{(p)}, \lambda^{(p)}$, the iterates must subsequently move together again. This tendency for current differences in iterates to be largely cancelled in the next step is apparent from Lemmas 2–4. Moreover, we may think in terms of starting fresh with each $\bar{x}^{(p)}, \bar{\lambda}^{(p)}$. Thus the previous history is, in a sense, irrelevant, provided the iterates do not drift so far away that there is a chance of homing in on some different solution.

In practice the boundary value $\beta_0 = \frac{1}{4}$ is rather less important than it may seem from the above analysis for a number of reasons.

(i) The iteration often converges starting from a value of β_0 which is much greater than $\frac{1}{4}$. When this is true, one soon reaches a stage when $\beta_p < \frac{1}{4}$ is well satisfied and we can then regard the current values as the initial values.

(ii) Suppose the initial approximation has been derived by a specific algorithm. If that algorithm is used on two different computers, one of which has a mantissa with one binary digit more than that on the other, the initial values of β on the two machines will almost certainly differ by a factor of 2. This puts the relevance of the value $\beta = \frac{1}{4}$ in perspective.

(iii) Consider the behavior of the iteration with $\beta = \frac{1}{5}$. From (ii) above we may regard this as “only marginally smaller than $\frac{1}{4}$ ”. Yet the sequence of β_i ’s is now

$$(9.1) \quad \beta_0 = \frac{1}{5}, \quad \beta_1 = \frac{1}{9}, \quad \beta_2 = \frac{1}{49}, \quad \beta_3 = \frac{1}{2209}, \quad \beta_4 < 2.1 \times 10^{-7}, \quad \beta_5 < 5.0 \times 10^{-14},$$

and the $\eta_i \rightarrow 0$ in much the same way.

In the rounding error analysis there is little point in obscuring the essential simplicity of the argument by sailing too close to the wind. On the other hand, realistic bounds are essential for the rounding errors made at each step of the computed sequence.

We assume that $\|A\|, \|A - \lambda I\|, \|B\|$ are all of order unity throughout, and, thus we shall replace them by unity whenever they occur. For the computation of Pq , where P is an $n \times n$ matrix and q is an n vector, we make the following assumption for single-precision floating point computation to the base β with a t digit mantissa:

$$(9.2) \quad fl(Pq) = Pq + \xi, \quad \|\xi\| \leq n\beta^{-t}\|P\|\|q\|.$$

If, on the other hand, all inner products are accumulated in double precision and rounded to single precision on completion, we assume that

$$(9.3) \quad fl_2(Pq) = Pq + \xi, \quad \|\xi\| \leq \beta^{-t}\|Pq\| + n\beta^{-2t}\|P\|\|q\|,$$

where the second term in the bound for ξ comes from the rounding in the double precision part, and the first term comes from the final rounding to single precision. The first term is, of course, the dangerous one: Its omission removes all realism from the analysis. Finally, we assume that the computed solution of $Px = q$ satisfies exactly the relation

$$(9.4) \quad (P + E)x = q, \quad \|E\| \leq n\beta^{-t}\|P\|.$$

Of the three assumptions, the first two are strict for a computer using a standard rounding procedure. The third assumption is partly empirical in nature, but is likely to be very conservative for stable methods of solving linear systems.

To demonstrate that our results converge to working accuracy, we must show that the error vector ξ satisfies $\|\xi\| \leq \beta^{-t}$ or is only marginally larger than this. Before attempting a detailed error analysis, however, we shall show that higher precision is essential in the computation of the residual. Suppose we already have an \bar{x} , $\bar{\lambda}$ which is correct to working accuracy, and we then perform one further iteration. We have then

$$(9.5) \quad \|\bar{x} - x\| = \|\xi\| \leq \beta^{-t}.$$

If we were to perform this iteration exactly, then from Lemma 4 we would have

$$(9.6) \quad \bar{r} = -\bar{B}\xi - \xi_s \tilde{\xi},$$

$$(9.7) \quad \overline{B\delta} = \bar{r},$$

$$(9.8) \quad \bar{\delta} = -\xi - \bar{B}^{-1} \xi_s \tilde{\xi},$$

where \bar{r} , \bar{B} are exact results corresponding to the given \bar{x} , $\bar{\lambda}$. The error in the correct \bar{x} , $\bar{\lambda}$ is therefore $-\bar{B}^{-1} \xi_s \tilde{\xi}$, the error ξ having been cancelled. Clearly,

$$(9.9) \quad \|\bar{B}^{-1} \xi_s \tilde{\xi}\| \leq \|\bar{B}^{-1}\| \|\xi\|^2 \leq \beta^{-2t} \|\bar{B}^{-1}\|.$$

Naturally we require that $\|\bar{B}^{-1}\| \leq \beta^t$ or the error may actually be increased. However, something appreciably stronger than this is needed in any case when rounding errors are involved or $\bar{B} + E$ might well be singular. If \hat{r} is the computed residual using single precision arithmetic, we have from (9.2)

$$(9.10) \quad \hat{r} = \bar{r} + f, \quad \|f\| \leq n\beta^{-t},$$

and from (9.6)

$$(9.11) \quad \|\bar{r}\| \leq \|\xi\| + \|\xi\|^2 \leq \beta^{-t} + \beta^{-2t}.$$

The computed residual is unlikely to have any correct significant figures and even if we solve the linear equation exactly, we can guarantee only that

$$(9.12) \quad \|\bar{B}^{-1} f\| \leq \|\bar{B}^{-1}\| n\beta^{-t}.$$

The corrected solution is likely to be much less accurate than its predecessor. Because we cannot even conserve a correctly rounded solution, there is little chance that it would ever be attained in the first place. We shall assume, therefore, that all residuals are determined by double precision accumulation with rounding to single precision on completion. We have then from (9.3), (9.10), and (9.11)

$$(9.13) \quad \hat{r} = \bar{r} + f,$$

$$(9.14) \quad \|f\| \leq \beta^{-t} [\|\xi\| + \|\xi\|^2] + n\beta^{-2t} \leq \beta^{-t} [\beta^{-t} + \beta^{-2t}] + n\beta^{-2t} \leq (n+2)\beta^{-2t}.$$

Solving exactly, we obtain the error from f bounded by

$$(9.15) \quad (n+2) \|\bar{B}^{-1}\| \beta^{-2t}.$$

The solution may be degraded unless $\|\bar{B}^{-1}\| \beta^{-t}$ is appreciably less than unity. We have already seen that such a demand is, in any case, inevitable.

In order to get some idea of the behavior under conditions which are almost borderline for Theorem 3 to apply, we analyze the case

$$(9.16) \quad \beta_0 = \kappa^{(0)} \eta^{(0)} = \kappa \varepsilon < \frac{1}{5}, \quad \kappa n \beta^{-t} < 0.01.$$

The second of these is reasonable since if κ were appreciably larger than this we should be computing an invariant subspace rather than a simple eigenpair. We

emphasize that relations (9.16) need not necessarily be satisfied by the first approximation, but provided there is some q for which

$$(9.17) \quad \beta_q = \kappa^{(q)} \eta^{(q)} = \kappa \varepsilon < \frac{1}{5},$$

then this q th computed \bar{x}_q could, for the purposes of the analysis, be regarded as the initial value. From the condition $\beta_0 = \frac{1}{5}$ we know that the exact process converges to a solution x, λ and

$$(9.18) \quad \|x^{(p)} - x^{(0)}\| \leq [2/(1 + (0.2)^{1/2})] \varepsilon < 1.4\varepsilon \quad \text{for all } p,$$

$$(9.19) \quad \|x - x^{(0)}\| < 1.4\varepsilon,$$

$$(9.20) \quad \kappa^{(p)} \leq \kappa^{(\infty)} = \kappa^{(0)}/(0.2)^{1/2} < 2.24\kappa \quad \text{for all } p.$$

We wish to show that the $\bar{x}^{(i)}$ never deviates far from the $x^{(i)}$. The analysis of the first step is different from that for the general step since $x^{(0)} = \bar{x}^{(0)}$ and the simple expression given in (7.27) for the residual does not apply. We have then

$$(9.21) \quad B^{(0)}\delta^{(0)} = r^{(0)}, \quad \eta^{(0)} = \|\delta^{(0)}\| = \varepsilon,$$

$$(9.22) \quad (B^{(0)} + E^{(0)})\bar{\delta}^{(0)} = r^{(0)} + f^{(0)}, \quad \|f^{(0)}\| \leq \beta^{-t}\|r^{(0)}\| + n\beta^{-2t},$$

where we assume $\|E^{(0)}\| \leq n\beta^{-t}$ covers the errors made during the solution of the linear system and also the trivial errors made in computing $A - \lambda^{(0)}I$. From (9.1) we have

$$(9.23) \quad \|r^{(0)}\| \leq \|B^{(0)}\|\|\delta^{(0)}\| = \varepsilon,$$

and from (9.22) we have

$$(9.24) \quad \bar{\delta}^{(0)} = (B^{(0)})^{-1}r^{(0)} - M^{(0)}(B^{(0)})^{-1}r^{(0)} + (B^{(0)} + E^{(0)})^{-1}f^{(0)},$$

where

$$(9.25) \quad \|M^{(0)}\| \leq \|(B^{(0)})^{-1}\|\|E^{(0)}\|/(1 - \|(B^{(0)})^{-1}\|\|E^{(0)}\|) \leq \kappa n\beta^{-t}/(1 - \kappa n\beta^{-t}),$$

$$(9.26) \quad \|(B^{(0)} + E^{(0)})^{-1}\| \leq \|(B^{(0)})^{-1}\|/(1 - \|(B^{(0)})^{-1}\|\|E^{(0)}\|) \leq \kappa/(1 - \kappa n\beta^{-t}).$$

Hence from (9.24) we get

$$(9.27) \quad \bar{\delta}^{(0)} = \delta^{(0)} - M^{(0)}\delta^{(0)} + (B^{(0)} + E^{(0)})^{-1}f^{(0)}$$

and

$$(9.28) \quad \begin{aligned} \|-M^{(0)}\delta^{(0)} + (B^{(0)} + E^{(0)})^{-1}f^{(0)}\| &\leq \frac{\kappa}{1 - \kappa n\beta^{-t}}[n\beta^{-t}\varepsilon + \|f^{(0)}\|] \\ &\leq \frac{\kappa}{1 - \kappa n\beta^{-t}}[n\beta^{-t}\varepsilon + \beta^{-t}\varepsilon + n\beta^{-2t}] \\ (9.29) \quad &\leq 1.02[\frac{1}{5}n\beta^{-t} + \frac{1}{5}\beta^{-t} + 0.01\beta^{-t}], \end{aligned}$$

where we have used condition (9.17). To simplify the analysis, we assume $n > 10$ (this merely ensures $\beta^{-t} < 0.1n\beta^{-t}$); condition (9.29) then gives

$$(9.30) \quad \|\bar{\delta}^{(0)} - \delta^{(0)}\| < 0.23n\beta^{-t}.$$

We still have to add $\bar{\delta}^{(0)}$ to $x^{(0)}$; this step gives a further error, bounded by $\beta^{-t} < 0.1n\beta^{-t}$ and

$$(9.31) \quad \|\xi^{(0)}\| = \|\bar{x}^{(1)} - x^{(1)}\| \leq 0.33n\beta^{-t}.$$

In subsequent steps we have

$$(9.32) \quad B^{(p)}\delta^{(p)} = r^{(p)} = \delta_s^{(p-1)}\tilde{\delta}^{(p-1)}, \quad \|\delta^{(p)}\| \leq \eta^{(p)}, \quad \|r^{(p)}\| \leq (\eta^{(p-1)})^2,$$

$$(9.33) \quad (\bar{B}^{(p)} + E^{(p)})\bar{\delta}^{(p)} = \text{computed}(\bar{r}^{(p)}),$$

where $\bar{B}^{(p)}$ is the exact B corresponding to $\bar{x}^{(p)}$, and $E^{(p)}$ covers the rounding errors made both in the computation of $A - \lambda^{(p)}I$ and in the solution of the linear system. We assume

$$(9.34) \quad \|E^{(p)}\| \leq n\beta^{-t}.$$

From Lemma 5 we have

$$(9.35) \quad \bar{r}^{(p)} = r^{(p)} - \bar{B}^{(p)}\xi^{(p-1)} + \xi_s^{(p-1)}\tilde{\xi}^{(p-1)}.$$

(Here $\bar{r}^{(p)}$ is the exact residual corresponding to $\bar{x}^{(p)}$.) The computed results show that

$$(9.36) \quad \text{computed}(\bar{r}^{(p)}) = \bar{r}^{(p)} + f^{(p)},$$

$$(9.37) \quad \|f^{(p)}\| \leq \beta^{-t}\|\bar{r}^{(p)}\| + n\beta^{-2t} \leq \beta^{-t}[\|r^{(p)}\| + \|\xi^{(p-1)}\| + \|\xi^{(p-1)}\|^2] + n\beta^{-2t}.$$

Hence we get

$$(9.38) \quad (\bar{B}^{(p)} + E^{(p)})\bar{\delta}^{(p)} = r^{(p)} - \bar{B}^{(p)}\xi^{(p-1)} + g^{(p)},$$

$$(9.39) \quad g^{(p)} = \xi_s^{(p-1)}\tilde{\xi}^{(p-1)} + f^{(p)},$$

$$(9.40) \quad \|g^{(p)}\| \leq \|\xi^{(p-1)}\|^2 + \beta^{-t}[(\eta^{(p-1)})^2 + \|\xi^{(p-1)}\| + \|\xi^{(p-1)}\|^2] + n\beta^{-2t}.$$

We now need to show that the matrices $\bar{B}^{(p)} + E^{(p)}$ are nonsingular and have bounds for their inverses. We know that the $B^{(p)}$ are nonsingular; for information on $\bar{B}^{(p)} + E^{(p)}$, we need bounds for $\|\xi^{(p-1)}\| = \|\bar{x}^{(p)} - x^{(p)}\|$, remembering that

$$(9.41) \quad \|\bar{B}^{(p)} - B^{(p)}\| = \|F^{(p)}\| \leq 2\|\xi^{(p-1)}\|.$$

A rough preliminary analysis indicates that $\xi^{(p)}$ will not exceed a modest multiple of $n\beta^{-t}$, and we now prove by induction that

$$(9.42) \quad \|\xi^{(p)}\| \leq 0.5n\beta^{-t} \quad \text{for all } p$$

(and indeed that the later $\xi^{(p)}$ are much smaller). We already know it is true for $p = 0$, and we assume that it is true up to $\xi^{(p-1)}$. Notice that this would merely extend the bounds in (9.18) and (9.19) to $1.4\epsilon + 0.5n\beta^{-t}$. From our inductive hypothesis we have

$$(9.43) \quad \begin{aligned} \bar{B}^{(p)} &= B^{(p)} + F^{(p)}, & \|F^{(p)}\| &\leq 2 \times 0.5n\beta^{-t} = n\beta^{-t}, \\ \bar{B}^{(p)} + E^{(p)} &= B^{(p)} + E^{(p)} + F^{(p)}, & \|E^{(p)} + F^{(p)}\| &\leq 2n\beta^{-t}. \end{aligned}$$

From (9.20) and (9.19) we therefore have

$$(9.44) \quad \|(\bar{B}^{(p)} + E^{(p)})^{-1}\| \leq \frac{\kappa^{(p)}}{1 - (2.24)\kappa(2n\beta^{-t})} < 1.05\kappa^{(p)} < 2.36\kappa.$$

(Sometimes we shall use the former and sometimes the latter of the two bounds.) Similarly we have

$$(9.45) \quad (\bar{B}^{(p)} + E^{(p)})^{-1}u = (B^{(p)})^{-1}u - L^{(p)}(B^{(p)})^{-1}u \quad \text{for all } u$$

where

$$(9.46) \quad \|L^{(p)}\| < 1.05\kappa^{(p)}(2n\beta^{-t}) < 2.36\kappa(2n\beta^{-t})$$

and

$$(9.47) \quad (\bar{B}^{(p)} + E^{(p)})^{-1} \bar{B}^{(p)} u = u - M^{(p)} u$$

where

$$(9.48) \quad \|M^{(p)}\| < 1.05\kappa^{(p)}(n\beta^{-t}) < 2.36\kappa n\beta^{-t} u.$$

The inductive hypothesis guarantees the nonsingularity of all relevant matrices up to and including the case we are about to use in determining $\bar{\delta}^{(p)}$. Returning now to (9.38) we have first, using (9.45) and (9.46) with $u = r^{(p)}$,

$$(9.49) \quad (\bar{B}^{(p)} + E^{(p)})^{-1} r^{(p)} = \delta^{(p)} - L^{(p)} \delta^{(p)} = \delta^{(p)} + a^{(p)}$$

where

$$(9.50) \quad \|a^{(p)}\| \leq 1.05\kappa^{(p)}(2n\beta^{-t})\eta^{(p)};$$

then, using (9.47)

$$(9.51) \quad -(\bar{B}^{(p)} + E^{(p)})^{-1} \bar{B}^{(p)} \xi^{(p-1)} = -\xi^{(p-1)} + M^{(p)} \xi^{(p-1)} = -\xi^{(p-1)} + b^{(p)}$$

where

$$(9.52) \quad \|b^{(p)}\| \leq 2.36\kappa n\beta^{-t} \|\xi^{(p-1)}\|.$$

Finally, we wish to bound $(\bar{B}^{(p)} + E^{(p)})^{-1} g^{(p)}$ using (9.40). We treat the term in (9.40) involving $(\eta^{(p-1)})^2$ differently from the rest. We have

$$(9.53) \quad \|(\bar{B}^{(p)} + E^{(p)})^{-1}\| \beta^{-t} (\eta^{(p-1)})^2 \leq 1.05\kappa^{(p)} \beta^{-t} (\eta^{(p-1)})^2 = c^{(p)},$$

while from the rest of $g^{(p)}$ we have the bound

$$(9.54) \quad 2.36\kappa [\|\xi^{(p-1)}\|^2 + \beta^{-t} (\|\xi^{(p-1)}\| + \|\xi^{(p-1)}\|^2) + n\beta^{-2t}] = d^{(p)}.$$

Combining these results, we get

$$(9.55) \quad \bar{\delta}^{(p)} = \delta^{(p)} - \xi^{(p-1)} + h^{(p)},$$

where

$$(9.56) \quad \|h^{(p)}\| \leq \|a^{(p)}\| + c^{(p)} + \|b^{(p)}\| + d^{(p)}.$$

Clearly

$$\|\xi^{(p)}\| \leq \|h^{(p)}\| + \beta^{-t} < \|h^{(p)}\| + 0.1n\beta^{-t},$$

the β^{-t} coming from the addition of $\bar{\delta}^{(p)}$ to $\bar{x}^{(p)}$. Note that the term $-\xi^{(p-1)}$ in (9.55) annihilates the previous error. In (9.56) the first two terms involve $\eta^{(p)}$ and $\eta^{(p-1)}$, those “belonging” to the exact process, and tend quadratically to zero. They are of significance only in the first one or two iterations. The terms $\|b^{(p)}\|$ and $d^{(p)}$ ultimately control the behavior completely. Using the inductive hypothesis (and the assumption that $n > 10$), we have

$$\begin{aligned} \|b^{(p)}\| &\leq 2.36\kappa n\beta^{-t} (0.5n\beta^{-t}) < 0.0118n\beta^{-t}, \\ d^{(p)} &< 2.36\kappa [0.25n^2\beta^{-2t} + 0.5n\beta^{-2t} + 0.25n^2\beta^{-3t} + n\beta^{-2t}] \\ &< 2.36[.0025n\beta^{-t} + .0005n\beta^{-t} + .0000025n\beta^{-t} + .001n\beta^{-t}] \\ &< .0092n\beta^{-t}. \end{aligned}$$

The terms $\|a^{(p)}\|$ and $c^{(p)}$ present their greatest danger when $p = 1$ and we have

$$\|a^{(1)}\| \leq 1.05\kappa(2n\beta^{-t}) = \frac{1.05}{9}2n\beta^{-t} < .24n\beta^{-t},$$

$$c^{(1)} \leq 1.05\kappa^{(1)}\beta^{-t}(\eta^{(0)})^2 = (1.05)_{\frac{5}{3}}\kappa\varepsilon^2\beta^{-t} < \frac{1.05}{9} \frac{\beta^{-t}}{.01}n\beta^{-t},$$

which is negligible for any reasonable β^{-t} .

Summing the contributions, we see that

$$\|\xi^{(1)}\| \leq 0.4n\beta^{-t}.$$

Hence the result is improving, but the proof shows that after the first two iterations the $\xi^{(i)}$ are decreasing rapidly. The computation of the successive $\xi^{(i)}$ is so tedious that we have contented ourselves with showing that the $\xi^{(i)}$'s remain less than $0.5n\beta^{-t}$. However, it is evident that after the first two iterations the $\|\xi^{(i)}\|$ decreases rapidly. Although tedious, the computations are perfectly straightforward; all the error analysis is covered in the assumptions embodied in (9.2), (9.3) and (9.4). It is more in the spirit of practical numerical analysis to produce a program for computing the $\|\xi^{(i)}\|$.

The program parameters are

$$\beta_0, n, \beta^{-t} \quad \text{and} \quad \omega = \kappa^{(0)}n\beta^{-t}.$$

It is more convenient to work in terms of γ_i and α_i defined by

$$\gamma_i = \kappa_i/\kappa_0 \quad \text{and} \quad \|\xi^{(i)}\| \leq \alpha_in\beta^{-t}.$$

The first step is special and α_0 is given by

$$\alpha_0 = \left[\beta_0 + \frac{1}{n}(\beta_0 + \omega) \right] / (1 - \omega).$$

For subsequent steps we have the relations

$$\beta_i = \left(\frac{\beta_{i-1}}{1 - 2\beta_{i-1}} \right)^2, \quad \gamma_i = \frac{\gamma_{i-1}}{1 - 2\beta_{i-1}},$$

$$\bar{\gamma}_i = \frac{\gamma_i}{(1 - \vartheta_i)}, \quad \text{where } \vartheta_i = \gamma_i(1 + 2\alpha_{i-1})\omega$$

and

$$\alpha_i = \frac{(1 + 2\alpha_{i-1})\beta_i}{1 - \vartheta_i} + \frac{\beta^{-t}\beta_i}{\gamma_i\omega(1 - \vartheta_i)} + \bar{\gamma}_i\omega \left\{ \alpha_{i-1} + (1 + \beta^{-t})\alpha_{i-1}^2 + \frac{\alpha_{i-1} + 1}{n} \right\} + \frac{1}{n},$$

of which the last term comes merely from the addition of $\bar{\delta}^{(p)}$ to $x^{(p)}$. The α_i computed in this way give realistic upper bounds for the $\|\xi^{(i)}\|$. In fact, if $\xi^{(p-1)} = \alpha n\beta^{-t}$ we have

$$\|b^{(p)}\| \leq 0.0236\alpha n\beta^{-t},$$

$$\|d^{(p)}\| \leq 2.36[0.01\alpha^2 + 0.001\alpha + .0001\alpha^2]n\beta^{-t} + 0.0236\beta^{-t},$$

where the last term arises from the $n\beta^{-2t}$ and we have not made the substitution of $0.1n\beta^{-t}$ for β^{-t} . This last term and the term β^{-t} coming from the addition of $\bar{\delta}^{(p)}$ to $x^{(p)}$ finally remain when all the other terms have receded. The ultimate $\bar{x}^{(p)}$ will almost certainly be the correct rounded solution, though the term $0.0236\beta^{-t}$ poses a slight danger that it will not be the correctly rounded value and could even oscillate between two values that differ by β^{-t} .

As an example, we display these quantities as obtained from the above formulae with $\kappa\varepsilon = \frac{1}{5}$, $\kappa n\beta^{-t} = 0.01$, $n = 9$ (Table 1).

TABLE 1
Error analysis quantities

i	κ_i	γ_i	β_i	α_i
0	6.5×10^3	1.0	.2	.202
1	1.08×10^4	1.67	.111	.274
2	1.39×10^4	2.14	.0204	.151
3	1.45×10^4	2.23	4.5×10^{-4}	.116
4	1.45×10^4	2.24	2.05×10^{-7}	.115
5	1.45×10^4	2.24	4.21×10^{-14}	.115

Naturally when $\kappa\varepsilon$ is significantly smaller than $\frac{1}{5}$ and/or $\kappa n\beta^{-t}$ is significantly smaller than 0.01, all our results will be significantly stronger. If $\kappa n\beta^{-t} = \beta^{-m}$ (say), then if we iterate until convergence to working accuracy and add the final computation to $x^{(p)}$ using precision, we shall have an eigenpair with an error of approximately β^{-t-m} .

REFERENCES

- [1] J. J. DONGARRA, *Improving the accuracy of computed matrix eigenvalues*, ANL-80-84, Applied Mathematics Division, Argonne National Laboratory, Argonne, IL, August 1980.
- [2] B. S. GARBOW, et al., *Matrix Eigensystem Routines—EISPACK Guide Extension*, Lecture Notes in Computer Science, 51, Springer-Verlag, Berlin, 1977.
- [3] L. B. RALL, *Computational Solution of Nonlinear Operator Equations*, John Wiley, New York, 1969.
- [4] B. T. SMITH, et al. *Matrix Eigensystem Routines—EISPACK Guide*, Lecture Notes in Computer Science, 6, 2nd ed., Springer-Verlag, Berlin, 1976.
- [5] H. J. SYMM AND J. H. WILKINSON, *Realistic error bounds for a simple eigenvalue and its associated eigenvector*, Numer. Math., 35 (1980), pp. 113–126.
- [6] J. H. WILKINSON, *The Algebraic Eigenvalues Problem*, Oxford University Press, London, 1965.
- [7] ———, *Error bounds for computed invariant subspaces*, Proc. of Rutishauser Symposium on Numerical Analysis, Research Report 81-02, Eidgenössische Technische Hochschule, Zürich, Switzerland, February, 1981.
- [8] T. YAMAMOTO, *Error bounds for computed eigenvalues and eigenvectors*, Numer. Math., 34 (1980), pp. 189–199.
- [9] ———, private communication, 1982.