

A Divide and Conquer Method for the Symmetric Tridiagonal Eigenproblem

J.J.M. Cuppen

Universiteit van Amsterdam, Instituut voor Toepassingen van de Wiskunde,
Roetersstraat 15, 1018 WB Amsterdam, The Netherlands

Summary. A method is given for calculating the eigenvalues of a symmetric tridiagonal matrix. The method is shown to be stable and for a large class of matrices it is, asymptotically, faster by an order of magnitude than the *QR* method.

Subject Classifications: AMS (MOS): 65F15, 68C25; CR: 5.14, 5.25.

1. Introduction

Let T be a symmetric tridiagonal matrix of order $n \geq 2$. T can be written as

$$T = \left(\begin{array}{c|c} T_1 & 0 \\ \hline 0 & T_2 \end{array} \right) + \alpha \left(\begin{array}{c|c} 1 & \\ \hline & 1 \end{array} \right) = \left(\begin{array}{c|c} T_1 & 0 \\ \hline 0 & T_2 \end{array} \right) + \alpha b b^T \tag{1.1}$$

where T_1 and T_2 are of order $n_1 \geq 1$ and $n_2 \geq 1$ with $n_1 + n_2 = n$. α is the n_1 -th off-diagonal element of T and the vector b is given by

$$\begin{aligned} b_i &= 1 && \text{if } i = n_1 \quad \text{or} \quad i = n_1 + 1, \\ b_i &= 0 && \text{otherwise.} \end{aligned} \tag{1.2}$$

Suppose that the solutions of the eigenvalue problems of T_1 and T_2 are given by

$$\begin{aligned} T_1 &= Q_1 D_1 Q_1^T \\ T_2 &= Q_2 D_2 Q_2^T \end{aligned} \tag{1.3}$$

where Q_1 and Q_2 are orthogonal matrices and D_1 and D_2 are diagonal matrices. If we denote the first row of Q_2 by f_2^T and the last row of Q_1 by l_1^T we have

$$\begin{aligned} T &= \left(\begin{array}{c|c} T_1 & \\ \hline & T_2 \end{array} \right) + \alpha b b^T = \left(\begin{array}{c|c} Q_1 & \\ \hline & Q_2 \end{array} \right) \left(\left(\begin{array}{c|c} D_1 & \\ \hline & D_2 \end{array} \right) + 2\alpha z z^T \right) \left(\begin{array}{c|c} Q_1 & \\ \hline & Q_2 \end{array} \right)^T, \\ z &= \frac{1}{\sqrt{2}} \left(\begin{array}{c|c} Q_1^T & \\ \hline & Q_2^T \end{array} \right) b = \frac{1}{\sqrt{2}} \begin{pmatrix} l_1 \\ f_2 \end{pmatrix}, \quad \|z\|_2 = 1. \end{aligned} \tag{1.4}$$

ACC/EFF PLOT FOR RANDOM (-1,1) TRIDIAGONAL MATRICES

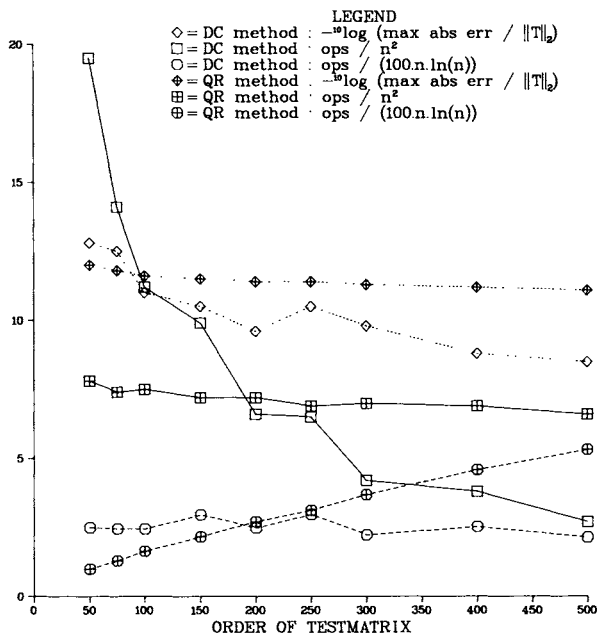


Fig. 1

The eigenvalues of T are therefore equal to those of $\left[\left(\begin{array}{c|c} D_1 & \\ \hline & D_2 \end{array} \right) + 2\alpha z z^T \right]$. These can be calculated with an algorithm for the rank-one modification of the symmetric eigenproblem which was proposed by Bunch, Nielsen and Sorensen [1], based on Golub [3].

A modified version of their algorithm, with an alternative derivation, will be given below.

A recursive application of the strategy described above leads to a Divide and Conquer method (D.C. method) for the eigenvalue problem of T . At the deepest level of recursion we can either carry on till we arrive at trivial 1×1 or 2×2 eigenvalue problems, or use a QR method to calculate the eigenvalues and the first and last rows of the eigenvector matrix of small $n_0 \times n_0$ blocks.

To be competitive with the QR method, any method should have an operation count of $6n^2$ operations (multiplications and divisions) or less. An order k rank-one modification costs however $\pm 9k^2$ operations. This seems to add up to $\pm 18n^2$ operations for the whole problem. The experiments discussed in Sect. 5 show that for many matrices (though not for all) the picture is drastically changed by deflation effects. There the method has a $c.n \log n$ behaviour, where the constant c depends on the type of matrix. For instance $c \approx 140$ for Wilkinson's matrices W_n^+ and W_n^- , $c \approx 220$ for random tridiagonal matrices and $c \approx 90$ for random matrices where the off-diagonal elements are on the average 10 times smaller than the diagonal elements.

Figure 1 gives the results obtained for random tridiagonal matrices (see also Sect. 5, Experiments).

Note: The method can be easily generalised to symmetric bandmatrices, requiring not more than m rank-one modifications per recursion step for a bandmatrix with bandwidth $2m + 1$. The first and last m rows of the transformation matrix will then have to be updated at each rank-one modification.

2. Rank One Modification of the Symmetric Eigenproblem

Let D_1 and D_2 be diagonal matrices of order n_1 and n_2 , let ρ be a scalar and f_1, l_2 and z be vectors of order n_1, n_2 and $n = n_1 + n_2$, respectively. We want to calculate a diagonal matrix Λ and vectors f and l such that there exists an orthogonal matrix P with

$$\left(\begin{array}{c|c} D_1 & \\ \hline & D_2 \end{array} \right) + \rho z z^T = P \Lambda P^T, \tag{2.1}$$

$$f = P^T \begin{pmatrix} f_1 \\ 0 \end{pmatrix}, \quad l = P^T \begin{pmatrix} 0 \\ l_2 \end{pmatrix}.$$

Let

$$D = \left(\begin{array}{c|c} D_1 & \\ \hline & D_2 \end{array} \right), \quad f' = \begin{pmatrix} f_1 \\ 0 \end{pmatrix}, \quad l' = \begin{pmatrix} 0 \\ l_2 \end{pmatrix}.$$

Following Bunch, Nielsen and Sorensen [1], we shall first show that deflation is possible if some of the elements of D are equal, or any of the elements of z is zero.

We first permute the elements of D in such a way that

$$d_1 \leq d_2 \leq \dots \leq d_n,$$

meanwhile permuting z, f' and l' correspondingly.

Now if k of the elements of D are equal:

$$d_{i+1} = d_{i+2} = \dots = d_{i+k}$$

the following identity holds for any orthogonal $k \times k$ matrix P' :

$$\left(\begin{array}{c|c|c} I_i & & \\ \hline & P' & \\ \hline & & I_{n-i-k} \end{array} \right)^T D \left(\begin{array}{c|c|c} I_i & & \\ \hline & P' & \\ \hline & & I_{n-i-k} \end{array} \right) = D.$$

Therefore we can choose P' to be a Householder transformation transforming $(z_{i+1}, z_{i+2}, \dots, z_{i+k})$ to $(*, 0, \dots, 0)$ while leaving D invariant. Consequently we may assume that if $d_i = d_j$ for some i and j then $z_i = 0$ or $z_j = 0$.

Let us now consider the coordinate directions i with $z_i = 0$. It is trivial that if $z_i = 0$ then d_i is an eigenvalue of $D + \rho z z^T$ with eigenvector e_i :

$$(D + \rho z z^T) e_i = d_i e_i + \rho z z_i = d_i e_i. \tag{2.2}$$

Therefore we can choose p^i , the i -th column of P , to be e_i and we have $\lambda_i = d_i, f_i = f'_i$ and $l_i = l'_i$ in (2.1) and we may ignore all coordinate directions with $z_i = 0$.

Hence, we can reduce the problem to a smaller problem when D contains mutually equal elements and/or some components of z are zero. So, without loss of generality we may assume that all elements of D are different and that all components of z are nonzero. Moreover we may assume that $\rho > 0$ since the case $\rho = 0$ is trivial and a system with $\rho < 0$ is easily transformed into one with $\rho > 0$ by considering $-D$ instead of D . The following theorem now applies.

Theorem 2.1. *If D is a diagonal matrix, $D = \text{diag}(d_1, \dots, d_n)$, $n \geq 2$ with $d_1 < d_2 < \dots < d_n$, $z \in \mathbb{R}^n$ is a vector with $z_i \neq 0$ for $i = 1, \dots, n$ and $\rho > 0$ a scalar, then the eigenvalues of the matrix $D + \rho z z^T$ are equal to the n roots $\lambda_1 < \dots < \lambda_n$ of the rational function*

$$w(\lambda) = 1 + \rho z^T (D - \lambda I)^{-1} z \tag{2.3}$$

$$= 1 + \rho \sum_{j=1}^n \frac{z_j^2}{d_j - \lambda}.$$

The corresponding eigenvectors p^1, \dots, p^n of $D + \rho z z^T$ are given by

$$p^i = (D - \lambda_i I)^{-1} z / \|(D - \lambda_i I)^{-1} z\|_2 \tag{2.4}$$

and the d_i strictly separate the eigenvalues λ_i as follows:

$$d_1 < \lambda_1 < d_2 < \lambda_2 < \dots < d_n < \lambda_n < d_n + \rho z^T z. \tag{2.5}$$

This theorem merely restates results of Golub [3] and results of Bunch, Nielsen and Sorensen [1] who added the explicit formula for the eigenvectors. The theorem is given here with a simple proof because this proof may provide some insight in the essence of the rank-one modification method (others may prefer a derivation via the characteristic equation).

Proof. An eigenvalue-eigenvector pair (λ, p) of $D + \rho z z^T$ satisfies

$$(D + \rho z z^T) p = \lambda p,$$

so

$$(D - \lambda I) p = -\rho z^T p z.$$

We now show that $D - \lambda I$ is nonsingular. Indeed, assuming $D - \lambda I$ singular, we have $\lambda = d_i$ for some i , so $0 = ((D - \lambda I) p)_i = -\rho z^T p z_i$, so $z^T p = 0$, so $(D - \lambda I) p = -\rho z^T p z = 0$, so $(d_j - \lambda) p_j = 0$ for all j , so $p_j = 0$ for $j \neq i$, so $0 = z^T p = z_i p_i$, so $z_i = 0$ which contradicts the assumptions. Therefore we have

$$z^T p \neq 0$$

$$p = -\rho z^T p (D - \lambda I)^{-1} z. \tag{2.6}$$

Consequently, p satisfies (2.4); multiplying both sides of (2.6) by z^T yields that λ is a root of (2.3).

(2.5) easily follows from the behaviour of $w(\lambda)$. \square

As an example Fig. 2 gives a plot of $w(\lambda)$ in case $n = 6$, $(d_j)_{j=1, \dots, 6} = (0, 1, 2, 2.7, 3.4, 5.4)$, $\rho = 2$ and $(z_j^2)_{j=1, \dots, 6} = (0.1, 0.02, 0.4, 0.4, 0.03, 0.05)$. Note that for small z_i^2 either λ_i or λ_{i+1} is close to d_i .

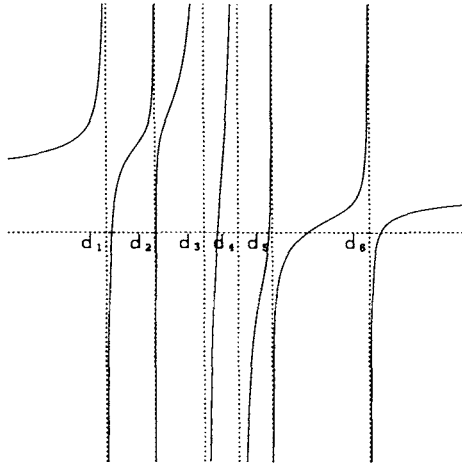


Fig. 2

With Eq. (2.3) and (2.4) we have an algorithm for the calculation of the matrix A and the vectors f and l in (2.1) which are needed in each recursion step of the D.C. method. Its stability in the presence of round-off errors and the use of a nonzero tolerance for neglecting elements of z and comparing close elements of D is far from trivial. This is readily seen from (2.4) as a serious loss of significant digits may occur in subtracting λ_i from d_i or d_{i+1} . The next section deals with this problem.

3. Accuracy

The D.C. method consists of a recursive application of rank-one modification steps (ROM steps), each of them preceded by deflation on mutually close intermediate eigenvalues and on small elements of the modification vector z . Since each time the matrix $D + \rho z z^T$ under consideration is derived from a part of the original matrix T by an orthogonal transformation we can apply backward error analysis on errors due to finite precision arithmetic and perturbations caused by deflation on close but unequal elements of D and small but nonzero components of z .

According to the theorem of Wielandt-Hoffman the errors introduced in the calculated eigenvalues of the original matrix are not larger than those introduced in $D + \rho z z^T$ itself (Wilkinson [6], p. 104).

This section considers the accuracy of one ROM-step, deflation inclusive. It extends the results of Bunch, Nielsen and Sorensen in two ways. Firstly it analyses the influence of round-off errors and secondly it bounds the disorthogonality of the eigenvectors with a factor $\min|\lambda_i - \lambda_k|$ instead of $\min|d_j - \lambda_i|$ in the denominator ($|d_j - \lambda_i|$ may become arbitrarily small, compare with [1], Theorem 6).

The ROM step we shall consider calculates Λ, f and l from the data D, ρ, z, f' and l' (order n , cf. (2.1)).

We may assume that $\|z\|_2 = 1, \rho > 0$ and that the elements of D are given in ascending order.

Let the constant L be defined by

$$L = 2 \|T\|_\infty \tag{3.1}$$

where T is the original tridiagonal matrix. Since the d_j and λ_i are eigenvalues of matrices ultimately derived from T in the same way as T_1 and T_2 are in (1.1), while $\|T_1\|_\infty \leq \|T\|_\infty$ and $\|T_2\|_\infty \leq \|T\|_\infty$, we have

$$L \geq \max \{|d_j - \lambda_i| : 1 \leq i \leq n, 1 \leq j \leq n\} \tag{3.2}$$

and

$$L \geq \rho.$$

It is clear that a lower bound for the magnitude of the derivative of $w(\lambda)$ at the zero's λ_i of $w(\lambda)$:

$$w'(\lambda_i) = \rho \sum_{j=1}^n \frac{z_j^2}{(d_j - \lambda_i)^2} \tag{3.3}$$

determines the accuracy to which the eigenvalues λ_i can be calculated. It shall be shown that errorbounds for the eigenvectors p^i depend on an upperbound for $w'(\lambda_i)$. This is not unnatural since the errors in p^i are caused by cancellation in the quantities which appear in denominators of $w'(\lambda_i)$ (cf. (2.4)).

Theorem 3.1. *Let $\delta > 0$ and $\zeta > 0$ such that for all i*

$$\begin{aligned} |d_{i+1} - d_i| &> L\delta \\ |z_i| &> \zeta, \end{aligned} \tag{3.4}$$

then the following bounds hold for all i :

$$\rho w'(\lambda_i) > 1, \tag{3.5}$$

$$\rho w'(\lambda_i) < \frac{8}{\delta^2 \zeta^2} \tag{3.6}$$

where the dot above the relation $<$ in (3.6) denotes a first order approximation for small δ and ζ .

Proof. i) Lowerbound: Using the inequality of Cauchy-Schwartz we have

$$\rho w'(\lambda_i) = \rho^2 \sum_{j=1}^n \frac{z_j^2}{(d_j - \lambda_i)^2} \sum_{j=1}^n z_j^2 \geq \left(\rho \sum_{j=1}^n \frac{z_j}{(d_j - \lambda_i)} z_j \right)^2.$$

The expression between brackets equals $w(\lambda_i) - 1$ so (3.5) follows from $w(\lambda_i) = 0$.

ii) Upperbound: Rearranging the terms of $w(\lambda)$ [cf. (2.3)] we get that λ_i is the unique root in the interval (d_i, d_{i+1}) of the equation

$$\frac{z_i^2}{d_i - \lambda} + \frac{z_{i+1}^2}{d_{i+1} - \lambda} = \varphi, \tag{3.7}$$

where φ is given by

$$\varphi = -\frac{1}{\rho} - \sum_{j \neq i, i+1} \frac{z_j^2}{d_j - \lambda_i}.$$

From (3.4) it follows that

$$-\frac{1}{\rho} - \frac{1}{L\delta} < \varphi < \frac{1}{L\delta}. \tag{3.8}$$

If we solve (3.7) for $\lambda \in (d_i, d_{i+1})$ we get, writing $(d_{i+1} - d_i) = L\delta_i$ and $\frac{1}{2}(z_i^2 + z_{i+1}^2) = \psi$ so $\delta_i > \delta$ and $\psi \leq \frac{1}{2}$,

$$(\lambda_i - d_i) = z_i^2 L \delta_i / (-\frac{1}{2} \varphi L \delta_i + \psi + ((\frac{1}{2} \varphi L \delta_i - \psi)^2 + \varphi z_i^2 L \delta_i)^{\frac{1}{2}}). \tag{3.9}$$

Since the left-handside of (3.7) is strictly ascending between the poles d_i and d_{i+1} we get a lower bound for $\lambda_i - d_i$ if we take φ equal to $-\frac{1}{\rho} - \frac{1}{L\delta}$ in (3.9). This yields

$$\lambda_i - d_i > z_i^2 L \delta_i / (2\psi - \varphi L \delta_i),$$

so

$$\begin{aligned} \frac{\rho^2 z_i^2}{(\lambda_i - d_i)^2} &< \frac{\rho^2}{z_i^2} \left(\frac{2\psi - \varphi L \delta_i}{L \delta_i} \right)^2 < \frac{\rho^2}{z_i^2} \left(\frac{2\psi}{L \delta_i} + \frac{1}{\rho} + \frac{1}{L\delta} \right)^2 \\ &\leq \frac{1}{\zeta^2} \left(\frac{1}{\delta_i} + 1 + \frac{1}{\delta} \right)^2 < \frac{4}{\zeta^2 \delta^2}. \end{aligned} \tag{3.10}$$

This also yields

$$\frac{1}{\rho} (\lambda_i - d_i) > \frac{1}{2} \zeta^2 \delta. \tag{3.11}$$

Analogously it follows that

$$\frac{\rho^2 z_{i+1}^2}{(d_{i+1} - \lambda_i)^2} < \frac{4}{\zeta^2 \delta^2} \tag{3.12}$$

and

$$\frac{1}{\rho} (d_{i+1} - \lambda_i) > \frac{1}{2} \zeta^2 \delta. \tag{3.13}$$

With (3.10) and (3.12) we derive (3.6) as follows

$$\begin{aligned} \rho w'(\lambda_i) &= \rho^2 \sum_{j=1}^n \frac{z_j^2}{(d_j - \lambda_i)^2} \leq \frac{\rho^2}{(L\delta)^2} + \frac{\rho^2 z_i^2}{(d_i - \lambda_i)^2} \\ &+ \frac{\rho^2 z_{i+1}^2}{(d_{i+1} - \lambda_i)^2} < \frac{1}{\delta^2} + \frac{4}{\zeta^2 \delta^2} + \frac{4}{\zeta^2 \delta^2} = \frac{8}{\zeta^2 \delta^2}. \quad \square \end{aligned}$$

Corollary 3.1. *The eigenvalues λ_i can be calculated from $w(\lambda) = 0$ to almost full precision, relative to the largest eigenvalue of the matrix T .*

Proof. Trivial from (3.5) and $\rho \leq L$.

We shall now consider the influence of finite precision arithmetic on the calculation of the vectors f and l . Recall that they are formed as innerproducts of f' and l' with the eigenvectors p^i of $D + \rho z z^T$ (cf. (2.1)). These eigenvectors are

given by

$$\begin{aligned} p^i &= v^i / \|v^i\|_2, \\ v^i &= (D - \lambda_i I)^{-1} z. \end{aligned} \tag{3.14}$$

Although the eigenvalues λ_i can be determined accurately, it is clear that problems arise when there is a loss significant digits in the calculation of $D - \lambda_i I$, i.e. if λ_i is close to either d_i or d_{i+1} or both. As we have seen in (3.9) this happens if z_i or z_{i+1} is small. It will become clear in the sequel that we will certainly meet small z -elements for matrices where the method is of use.

Fortunately, however, we do not need the p^i to be accurate individually; it is sufficient that the matrix P of calculated eigenvectors satisfies the conditions of the following lemma for some small η .

Lemma 3.1. *If $0 < \eta < 1$ and the matrix P satisfies*

$$\|P^T P - I\|_F < \eta \tag{3.15}$$

and

$$\|(D + \rho z z^T)P - P\Lambda\|_F < \eta L,$$

then an orthogonal matrix Q exists such that

$$\|P - Q\|_F < \eta$$

and

$$\|(D + \rho z z^T)Q - Q\Lambda\|_F < 2\eta L. \tag{3.16}$$

N.B. This implies that such a P cannot introduce more than a total error $3\eta L$ in the eigenvalues of T .

Proof. Consider the singular value decomposition of P :

$$P = U\Sigma V^T.$$

Now

$$P^T P - I = V\Sigma^2 V^T - I = V(\Sigma^2 - I)V^T,$$

so

$$\|\Sigma - I\|_F \leq \|\Sigma^2 - I\|_F \|(\Sigma + I)^{-1}\|_2 \leq \|\Sigma^2 - I\|_F < \eta.$$

If we choose

$$Q = UV^T$$

then Q is orthogonal and

$$\|P - Q\|_F = \|\Sigma - I\|_F < \eta,$$

$$\|(D + \rho z z^T)Q - Q\Lambda\|_F \leq \|(D + \rho z z^T)P - P\Lambda\|_F + \|(P - Q)\Lambda\|_F$$

$$+ \|(D + \rho z z^T)(P - Q)\|_F \leq \eta L + \|P - Q\|_F (\|\Lambda\|_2 + \|D + \rho z z^T\|_2) \leq 2\eta L$$

since

$$\|\Lambda\|_2 \leq \|T\|_\infty \quad \text{and} \quad \|D + \rho z z^T\|_2 \leq \|T\|_\infty. \quad \square$$

We now proceed to bound the residual $(D + \rho z z^T)P - P\Lambda$ and the disorthogonality $P^T P - I$ of P when P is calculated in finite precision from (3.14). Let ε be the relative precision of the arithmetic involved.

Lemma 3.2. *The columns p^i of P , as calculated from (3.14), satisfy*

$$|p^{i^T} p^k| \leq \frac{2L\varepsilon}{|\lambda_i - \lambda_k|} (2 + \sqrt{\rho \max(w'(\lambda_i), w'(\lambda_k))}), \tag{3.17}$$

where λ_i approximates the exact root $\hat{\lambda}_i$ of $w(\lambda)$ in the interval (d_i, d_{i+1}) .

Proof. The elements of a computed vector v^i (cf. (3.14)) are given by

$$v_j^i = \frac{z_j}{d_j - \lambda_i - \mu_{ij}} \quad \text{for some } \mu_{ij} \text{ satisfying } |\mu_{ij}| < L\varepsilon. \tag{3.18}$$

When λ_i approximating $\hat{\lambda}_i$ is calculated as precise as possible, we have

$$\frac{1}{\rho} + \sum_{j=1}^n \frac{z_j^2}{d_j - \lambda_i - \theta_{ij}} = 0 \quad \text{with } |\theta_{ij}| < L\varepsilon.$$

Since all terms appearing in this sum are strictly increasing functions of θ_{ij} between the poles we can substitute one θ_i for all θ_{ij} and have

$$\begin{aligned} \lambda_i &= \hat{\lambda}_i - \theta_i, |\theta_i| < L\varepsilon \\ \frac{1}{\rho} w(\lambda_i + \theta_i) &= \frac{1}{\rho} + \sum_{j=1}^n \frac{z_j^2}{d_j - \lambda_i - \theta_i} = 0, \end{aligned} \tag{3.19}$$

in accordance with corollary (3.1). We now have

$$\begin{aligned} v^{i^T} v^k &= \sum_j \frac{z_j^2}{(d_j - \lambda_i - \mu_{ij})(d_j - \lambda_k - \mu_{ik})} \\ &= \frac{1}{\lambda_i - \lambda_k} \left\{ \sum_j \frac{z_j^2}{d_j - \lambda_i - \mu_{ij}} - \sum_j \frac{z_j^2}{d_j - \lambda_k - \mu_{kj}} \right. \\ &\quad \left. - \sum_j z_j^2 \frac{\mu_{ij} - \mu_{kj}}{(d_j - \lambda_i - \mu_{ij})(d_j - \lambda_k - \mu_{kj})} \right\} \\ &= \frac{1}{\lambda_i - \lambda_k} \left(\frac{1}{\rho} w(\lambda_i + \mu'_i) - \frac{1}{\rho} w(\lambda_k + \mu'_k) - \mu_{ik}' \|v^i\|_2 \|v^k\|_2 \right), \end{aligned}$$

where

$$|\mu'_i| < L\varepsilon, |\mu'_k| < L\varepsilon \quad \text{and} \quad |\mu_{ik}'| < 2L\varepsilon.$$

This is in first order approximation,

$$\begin{aligned} v^{i^T} v^k &\doteq \frac{1}{\lambda_i - \lambda_k} \left(\frac{1}{\rho} ((\mu'_i - \theta_i) w'(\lambda_i) - (\mu'_k - \theta_k) w'(\lambda_k)) \right. \\ &\quad \left. - \mu_{ik}' \|v^i\|_2 \|v^k\|_2 \right), \end{aligned}$$

$$|v^{i^T} v^k| \leq \frac{2L\varepsilon}{|\lambda_i - \lambda_k|} \left(\frac{1}{\rho} (w'(\lambda_i) + w'(\lambda_k)) + \|v^i\|_2 \|v^k\|_2 \right).$$

Also,

$$\|v^i\|^2 = \sum_j \frac{z_j^2}{(d_j - \lambda_i - \mu_{ij})^2} \doteq \frac{1}{\rho} w'(\lambda_i), \tag{3.20}$$

so

$$\begin{aligned} |p^{iT} p^k| &\doteq |v^{iT} v^k| / (\|v^i\|_2 \|v^k\|_2) \\ &\leq \frac{2L\varepsilon}{|\lambda_i - \lambda_k|} \left(\left(\frac{w'(\lambda_i)}{w'(\lambda_k)} \right)^{\frac{1}{2}} + \left(\frac{w'(\lambda_k)}{w'(\lambda_i)} \right)^{\frac{1}{2}} + 1 \right). \end{aligned} \tag{3.21}$$

Since one of the square roots in the right hand side of (3.21) is smaller than 1 (3.5) now immediately yields (3.17). \square

Lemma 3.3. *The calculated eigenvectors p^i satisfy*

$$\|(D + \rho z z^T - \lambda_i I) p^i\|_2 \leq L\varepsilon \sqrt{2 + 8\rho w'(\lambda_i)}. \tag{3.22}$$

Proof. Using (3.18-20) we get

$$\begin{aligned} &((D + \rho z z^T - \lambda_i I) p^i)_j \\ &= \left((d_j - \lambda_i) \frac{z_j}{d_j - \lambda_i - \mu_{ij}} + \rho z_j \sum_k \frac{z_k^2}{d_k - \lambda_i - \mu_{ik}} \right) / \|v^i\| \\ &\doteq \left(z_j + z_j \frac{\mu_{ij}}{d_j - \lambda_i - \mu_{ij}} + z_j (w(\lambda_i + \mu_i) - 1) \right) / \sqrt{\frac{1}{\rho} w'(\lambda_i)} \\ &\doteq z_j \left(\frac{\mu_{ij}}{d_j - \lambda_i} + (\mu_i - \theta_i) w'(\lambda_i) \right) \sqrt{\rho} / \sqrt{w'(\lambda_i)}. \end{aligned}$$

Therefore

$$\begin{aligned} \|(D + \rho z z^T - \lambda_i I) p^i\|_2^2 &\doteq \frac{\rho}{w'(\lambda_i)} \sum_j \left(\mu_{ij} \frac{z_j}{d_j - \lambda_i} + (\mu_i - \theta_i) w'(\lambda_i) z_j \right)^2 \\ &\leq \frac{2\rho}{w'(\lambda_i)} \sum_j \left((L\varepsilon)^2 \frac{z_j^2}{(d_j - \lambda_i)^2} + (2L\varepsilon w'(\lambda_i))^2 z_j^2 \right) \\ &\doteq 2(L\varepsilon)^2 + 8\rho(L\varepsilon)^2 w'(\lambda_i), \end{aligned}$$

which completes the proof.

The factor $\frac{1}{|\lambda_i - \lambda_k|}$ in the bound (3.17) on the inner product of two calculated eigenvectors p^i and p^k can of course not be avoided. From (3.11) and (3.13) we have that

$$\frac{1}{|\lambda_i - \lambda_k|} \leq \frac{1}{\rho \zeta^2 \delta},$$

but this upperbound only helps us if we use very crude tolerances in the deflation and for neglecting a small ρ . Therefore I prefer to perform a (modified Gram-Schmidt) orthogonalisation on eigenvectors $p^i, p^{i+1}, \dots, p^{i+k}$ if the eigenvalues λ_{j-1} and λ_j are close to each other for $j = i + 1, \dots, i + k$. This orthogonalisation

sation disturbs, however, the bound (3.22) we derived for the residuals of the p^i as shall be considered in the following lemma.

Lemma 3.4. *Let the resulting eigenvectors be given by*

$$\begin{aligned}\tilde{p}^j &= \left(p^j - \sum_{l=1}^{j-1} \alpha_{jl} \tilde{p}^l \right) / \beta_j, \quad j = i+1, \dots, i+k, \\ \alpha_{jl} &= p^{jT} \tilde{p}^l, \quad l = i, \dots, j-1, \\ \beta_j &= \left(1 - \sum_{l=1}^{j-1} (\alpha_{jl})^2 \right)^{\frac{1}{2}}.\end{aligned}$$

Let

$$\omega = \max_{1 \leq l \leq n} \sqrt{2 + 8\rho w'(\lambda_l)} \quad (3.23)$$

and

$$\begin{aligned}\gamma_j &= \left(1 + \sum_{l=i}^{j-1} |\alpha_{jl}| (\gamma_l + |\lambda_j - \lambda_l| / L\varepsilon \omega) \right) / \beta_j, \quad j = i+1, \dots, i+k, \\ \gamma_i &= 1,\end{aligned}$$

then

$$\| (D + \rho z z^T - \lambda_j I) \tilde{p}^j \|_2 \leq \gamma_j L\varepsilon \omega, \quad j = i+1, \dots, i+k. \quad (3.24)$$

Proof. Note that $2\omega > 2 + \sqrt{\rho w'(\lambda_l)}$ for all l . Now (3.24) follows from

$$\begin{aligned}(D + \rho z z^T - \lambda_j I) \tilde{p}^j \\ = \left((D + \rho z z^T - \lambda_j I) p^j - \sum_{l=i}^{j-1} \alpha_{jl} ((D + \rho z z^T - \lambda_l) \tilde{p}^l + (\lambda_j - \lambda_l) \tilde{p}^l) \right) / \beta_j. \quad \square\end{aligned}$$

No satisfactory theoretical answer was found to the question how large the factors γ_j can become. Therefore we define

$$\gamma = \max_i \gamma_i$$

and monitor on γ in practice. In my experiments γ never was larger than 1.5.

So if reorthogonalisation is performed between any two eigenvectors p^i and p^k with $|\lambda_i - \lambda_k| < Lr$ we have using (3.17), for any i and k

$$|p^{iT} p^k| < \frac{2\varepsilon}{r} (2 + \sqrt{\rho \max_j w'(\lambda_j)}). \quad (3.25)$$

We now seek to bound, in a backward manner, the errors introduced and propagated in one recursion step, e.g. one ROM step, deflation included. We may assume that we are given D_1 , D_2 , \hat{f}_1 , \hat{f}_2 , \tilde{l}_1 and \tilde{l}_2 such that there exist orthogonal matrices Q_1 and Q_2 with

$$T_i = Q_i D_i Q_i^T + E_i, \quad \tilde{f}_i = f_i + \Delta f_i, \quad \tilde{l}_i = l_i + \Delta l_i, \quad i = 1, 2, \quad (3.26)$$

where f_i^T is the first row and l_i^T is the last row of Q_i and the E_i , Δf_i and Δl_i have a small norm (cf. (1.4)).

Let

$$T = \left(\begin{array}{c|c} T_1 & \rho/2 \\ \hline \rho/2 & T_2 \end{array} \right) = \left(\begin{array}{c|c} Q_1 & \\ \hline & Q_2 \end{array} \right) \left(\left(\begin{array}{c|c} D_1 & \\ \hline & D_2 \end{array} \right) + \rho z z^T \right) \left(\begin{array}{c|c} Q_1 & \\ \hline & Q_2 \end{array} \right)^T + \left(\begin{array}{c|c} E_1 & \\ \hline & E_2 \end{array} \right)$$

$$z = \frac{1}{\sqrt{2}} \begin{pmatrix} l_1 \\ \vdots \\ f_2 \end{pmatrix}, \quad \rho > 0. \quad (3.27)$$

A ROM step is applied to calculate A , \tilde{f} and \tilde{l} with

$$T = Q A Q^T + E, \quad \tilde{f} = f + \Delta f, \quad \tilde{l} = l + \Delta l, \quad (3.28)$$

where Q is an orthogonal matrix and f^T and l^T are its first and last row. Theorem 3.2 gives upperbounds for $\|E\|$, $\|\Delta f\|$ and $\|\Delta l\|$.

Theorem 3.2. *If deflation is performed to obtain (3.4) for certain $\delta > 0$ and $\zeta > 0$ and reorthogonalisation is performed between any two calculated eigenvectors p^i , p^k with $|\lambda_i - \lambda_k| < Lr$, $0 < r \leq 1$, then the orthogonal matrix Q in (3.28) can be chosen in such a way that*

$$\|E\|_F \leq \|E_1\|_F + \|E_2\|_F + L\delta\sqrt{n} + \rho(\sqrt{2}\|\Delta l_1\| + \sqrt{2}\|\Delta f_2\| + 2\zeta\sqrt{n}) + 2L\frac{n\varepsilon\omega}{r\sqrt{2}},$$

$$\|\Delta f\| \leq \|\Delta f_1\| + \frac{n\varepsilon\omega}{r\sqrt{2}}, \quad \|\Delta l\|_2 \leq \|\Delta l_2\|_2 + \frac{n\varepsilon\omega}{r\sqrt{2}}, \quad (3.29)$$

where $\omega = \max_l \sqrt{2 + 8\rho w'(\lambda_l)} < \frac{8}{\delta\zeta}$.

Proof. We have to deal with errors introduced during a) the process of deflation and b) the ROM step. We will not take into account the errors due to Householder transformations and other stable processes which are negligible compared to those due to a) and b).

Ad a): Deflation can be regarded as first perturbing D by an amount ΔD , making close elements of D exactly equal, then applying the appropriate Householder transformations and finally perturbing z by an amount $\Delta z'$, making small elements of z exactly zero. Taking into account that we do not have the exact z , but an approximation of it equal to $z + \frac{1}{\sqrt{2}} \begin{pmatrix} \Delta l_1 \\ \Delta f_2 \end{pmatrix}$ (cf. 3.26–27) we see that after deflation with the tolerances $L\delta$ and ζ (cf. 3.4) we have

$$T = Q_3(D + \rho \tilde{z} \tilde{z}^T) Q_3^T + E_3, \quad (3.30)$$

$$\|E_3\| \leq \left\| \begin{pmatrix} E_1 & \\ & E_2 \end{pmatrix} \right\|_F + \|\Delta D\|_F + \rho \|z z^T - \tilde{z} \tilde{z}^T\|_F$$

$$\leq \|E_1\|_F + \|E_2\|_F + L\delta\|I\|_F + 2\rho \|z - \tilde{z}\|_2 \quad (3.31)$$

$$\leq \|E_1\|_F + \|E_2\|_F + L\delta\sqrt{n} + \rho(\sqrt{2}(\|\Delta l_1\| + \|\Delta f_2\|) + 2\zeta\sqrt{n}).$$

Since in this stage on \tilde{f}_1 and \tilde{l}_2 only Householder transformations are applied, Δf_1 and Δl_2 have changed, but not their norms.

Ad b): (3.25) yields that for any i and k :

$$|p^{i^T} p^k| \leq \frac{2L\varepsilon\omega}{Lr\sqrt{8}} = \frac{\varepsilon\omega}{r\sqrt{2}},$$

so

$$\|P^T P - I\|_F \leq \frac{n\varepsilon\omega}{r\sqrt{2}}, \tag{3.32}$$

and with Lemmas 3.3 and 3.4 we get

$$\|(D + \rho \tilde{z} \tilde{z}^T) P - P\Lambda\|_F \leq \gamma L\varepsilon\omega \sqrt{n}, \tag{3.33}$$

where γ is an upperbound for all γ_j (cf. (3.24)). Assuming that $\gamma\sqrt{n} \leq \frac{n}{r\sqrt{2}}$ we can apply Lemma 3.1 with $\eta = \frac{n\varepsilon\omega}{r\sqrt{2}}$ which yields that there is an orthogonal matrix Q_4 such that

$$\begin{aligned} Q_4^T (D + \rho \tilde{z} \tilde{z}^T) Q_4 &= \Lambda + E_4 \\ \|E_4\| &< 2L \frac{n\varepsilon\omega}{r\sqrt{2}}, \\ \|P - Q_4\|_F &< \frac{n\varepsilon\omega}{r\sqrt{2}} \end{aligned} \tag{3.34}$$

Combining (3.34), (3.30) and (3.31) we get

$$T = Q_5^T A Q_5 + E,$$

$$\|E\|_F \leq \|E_1\|_F + \|E_2\|_F + L\delta\sqrt{n} + \rho\sqrt{2}(\|A l_1\|_2 + \|A f_2\|_2 + \zeta\sqrt{2n}) + 2L \frac{n\varepsilon\omega}{r\sqrt{2}},$$

$$\|A f\|_2 \leq \|A f_1\|_2 + \|P - Q_4\|_F \|f\|_2 \leq \|A f_1\|_2 + \frac{n\varepsilon\omega}{r\sqrt{2}},$$

$$\|A l\|_2 \leq \|A l_2\|_2 + \frac{n\varepsilon\omega}{r\sqrt{2}}.$$

Theorem 3.1 gives the upperbound for ω . \square

Corollary 3.2. *If δ and ζ are both chosen to be $2\varepsilon^{1/3}$ then the accuracy of the method is of the order of $\varepsilon^{1/3}$ since $\omega < 2\varepsilon^{1/3}$.*

It should be noted, however, that this choice of δ and ζ is based on the worst case analysis given above. It guarantees a minimum accuracy but a sharper analysis may be possible and a smaller choice of δ and ζ gives much better accuracy in practice. Therefore in the implementation of the method δ and ζ are chosen in accordance with a requested accuracy and then the value of $\sqrt{2 + 8\rho w'(\lambda_i)}$ is monitored. δ and ζ are readjusted if necessary and their final value can be used to bound the precision achieved in the results.

4. Time-Complexity of the DC Method

Our goal was to compute efficiently the eigenvalues of a symmetric tridiagonal matrix. This was to be done in a recursive way, in each step merging two eigenvalue systems by means of a rank one modification. In this section we consider the total number of multiplications and divisions required.

Assume that the zero-finding algorithm requires on the average m function evaluations per zero (m turned out to be ± 5 when using zero-inrat of Bus and Dekker [2]). In one rank one modification step of order n this contributes mn^2 operations (cf. (2.3)). The calculation of n eigenvectors costs $2n^2$ ops (cf. (3.14)) and the transformation of f and l another $2n^2$ ops (or n^2 ops plus a considerable bookkeeping overhead, cf. (2.1)). Since all other contributions are of a lower order we may estimate the cost of one rank one modification as $(m+4)n^2$ ops. For the total system this means that $(2m+4)n^2$ ops are required (no vectors and no transformation in the last step). The QR method does the job in about $2n$ sweeps of diminishing length, $6k$ ops per sweep of length k , so it needs for the total problem $\pm 6n^2$ ops. The conclusion is that the DC method is at most about 3 times slower than QR .

In the experiments that were performed, see section 5, the behaviour described above showed up for a special class of matrices (examples 7.6 and 7.10 in Gregory and Karney [4]). However, for Wilkinson's matrices W_n^\pm and also for random tridiagonal symmetric matrices the DC method was faster by an order of magnitude than QR once the order of the testmatrix became high enough ($\sim n \log n$ versus $\sim n^2$).

This remarkable phenomenon was caused by substantial deflation taking place in the ROM steps of the DC method. It turned out that usually only a constant number $k \ll n$ of the z_i in a ROM step of order n were non negligible (cf. (2.2)), which reduced the number of ops in this ROM step from $(m+2)n^2$ to $(m+2)k^2$.

An impression of how and why this reduction takes place can be gained from the following lemma's. Let T be a symmetric tridiagonal matrix of order n with diagonal elements d_1, \dots, d_n and off-diagonal elements b_1, \dots, b_{n-1} . For convenience, let $b_0 = b_n = 0$.

Lemma 4.1. *Let p be an eigenvector of T to the eigenvalue λ , $\|p\|=1$, and for some k , $1 \leq k \leq n$*

$$|b_{j-1}| + |b_j| \leq |d_j - \lambda| \quad \text{for } j=1, \dots, k, \quad (4.1)$$

then p_1 (which is an element of the first row of the eigenvectormatrix of T) satisfies

$$|p_1| \leq \prod_{j=1}^k \frac{|b_j|}{|d_j - \lambda| - |b_{j-1}|} \leq \prod_{j=1}^k \frac{|b_{j-1}| + |b_j|}{|d_j - \lambda|}. \quad (4.2)$$

This means that we have a decay with k which is more or less exponential since each of the terms in the products in (4.2) is smaller than 1.

Proof. For each j we have, since $TP = \lambda p$, that

$$p_j = -\frac{b_{j-1} p_{j-1} + b_j p_{j+1}}{d_j - \lambda}.$$

We first prove, by induction that $|p_1| \leq |p_2| \leq \dots \leq |p_k|$.

Assume that for some $j < k$ $|p_{j-1}| \leq |p_j|$. Then

$$\begin{aligned} |p_j| &\leq \frac{|b_{j-1}|}{|d_j - \lambda|} |p_{j-1}| + \frac{|b_j|}{|d_j - \lambda|} |p_{j+1}| \\ &\leq \frac{|b_{j-1}|}{|d_j - \lambda|} |p_j| + \frac{|b_j|}{|d_j - \lambda|} |p_{j+1}| \leq \frac{|b_j|}{|d_j - \lambda| - |b_{j-1}|} |p_{j+1}|. \end{aligned} \tag{4.3}$$

The assumption (4.1) now yields that $|p_j| \leq |p_{j+1}|$ whereas the induction can be started at $j=1$ since $0 = |p_0| \leq |p_1|$.

(4.2) now easily follows from (4.3) and the fact that (4.1) implies that

$$0 \leq \frac{|b_{j-1}|}{|d_j - \lambda| - |b_j|} \leq \frac{|b_{j-1}| + |b_j|}{|d_j - \lambda|} \quad \text{for } j=1, \dots, k. \quad \square$$

Lemma 4.2. *If $\xi_2 > \xi_1 > 0$ and*

$$\begin{aligned} \frac{|b_{j-1}| + |b_j|}{|d_j|} &< \xi_1, \\ |d_j - d_i| &> \xi_2 |d_j| + \xi_1 |d_i|, \end{aligned} \tag{4.4}$$

for all $i, j=1, \dots, n$, then the elements f_i of the first row of the eigenvectormatrix of T satisfy (in some ordering of the eigenvectors)

$$|f_i| < \left(\frac{\xi_1}{\xi_2}\right)^{i-1}, \quad i=1, \dots, n \tag{4.5}$$

Proof. The Gerschgorin theorem says that each eigenvalue of T lies in one of the disks

$$\mathcal{D}_i = \{\lambda \mid |\lambda - d_i| < \xi_1 |d_i|\}.$$

The assumptions imply that these disks are disjoint so we have that each disk contains exactly one eigenvalue. If the eigenvalues are ordered in such a way that $\lambda_i \in \mathcal{D}_i$ we have

$$|\lambda_i - d_i| < \xi_1 |d_i|,$$

so

$$|d_j - \lambda_i| = |d_j - (1 + \theta) d_i|,$$

for some θ with $|\theta| < \xi_1$ and for $j=1, \dots, i-1$.

Therefore

$$\frac{|b_{j-1}| + |b_j|}{|d_j - \lambda_i|} < \frac{|b_{j-1}| + |b_j|}{\xi_2 |d_j|} < \frac{\xi_1}{\xi_2}, \quad j=1, \dots, i-1.$$

Application of Lemma 4.1 yields (4.5). \square

An analogous property holds for the last row of the eigenvectormatrix of T . We can conclude that if ξ_1/ξ_2 is not very close to 1 these rows cannot have many non-negligible elements (they decay exponentially) and since the vectors z (cf. (1.4)) are composed of these rows we get an enormous deflation for these matrices ($k \leq 6 \log(\xi_1/\xi_2)/(-\log \varepsilon)$ if we neglect $|z_i| < \sqrt[3]{\varepsilon}$. cf. Corollary 3.2).

It is clear that even for matrices which do not satisfy the global condition (4.4) of Lemma 4.2 deflation becomes important if the local condition (4.1) of Lemma 4.1 is satisfied for quite a number of indices j and eigenvalues λ . These matrices probably have an eigenvectormatrix which is close to a bandmatrix. The effect described above is easily shown to be absent in matrices with a constant diagonal and a constant subdiagonal which have an eigenvectormatrix which is not close to a bandmatrix (Gregory and Karney [4] p. 137).

5. Experiments

The method was implemented in Algol 68 and tested on a CDC-Cyber 73/173 system ($\varepsilon \approx 10^{-14}$) against a stable version of the square-root free QR method as included in the Numal program library [5] (translated by hand into Algol 68). In the tests reference values were produced using a double length version of the same QR method in order to compute the maximum errors in the eigenvalues calculated by both methods. Multiplications and divisions were counted for both methods. Note that the actual performance of the DC method depends on some parameters. These were set to values which experience indicated to be good, but no attempt was made to optimise them.

The complete program and all results of the testruns are available on a microfiche and can be requested from the author.

Figure 3 gives the results for Wilkinson's matrices W_n^- for several values of n , and Fig. 4 the results for W_n^+ (Wilkinson [6], p. 308). It is clear that for these matrices the DC method shows an $n \log n$ behaviour. Its accuracy is, for W_n^+ , less than that of QR , but much better than $\varepsilon^{1/3}$ (compare Sect. 3).

The next test was performed on examples 7.6 and 7.10 from Gregory and Karney ([4], p. 138, 140). Example 7.6 is a tridiagonal matrix with all off-diagonal elements equal to a constant b (here 0.3) and all diagonal elements equal to a constant a (here 1) except the first and the last diagonal element which are equal to $a-b$ and $a+b$, respectively. Since the eigenvalues of this matrix are given by $\lambda_k = a + 2b \cos((2k-1)\pi/2n)$ it is easily seen that the condition of Lemma 4.1 is not satisfied for any index and eigenvalue.

Example 7.10 is a tridiagonal matrix with all diagonal elements equal to zero and the off-diagonal elements b_j equal to $\sqrt{j(n-j)}$, $j=1, \dots, n-1$. The eigenvalues of this matrix are $-n, -n+2, \dots, n$. The matrix is not at all diagonally dominant.

For both these matrices the $n \log n$ behaviour did not show up as can be seen from the operation count of the DC method which is given in Table 1.

Further tests were performed with random symmetric tridiagonal matrices. Here the diagonal elements were chosen one by one as $(2 \times \text{random} - 1)$ and the off-diagonal elements as $f \times (2 \times \text{random} - 1)$ where f is a factor governing the

ACC/EFF PLOT FOR WILKINSON'S MATRICES W-

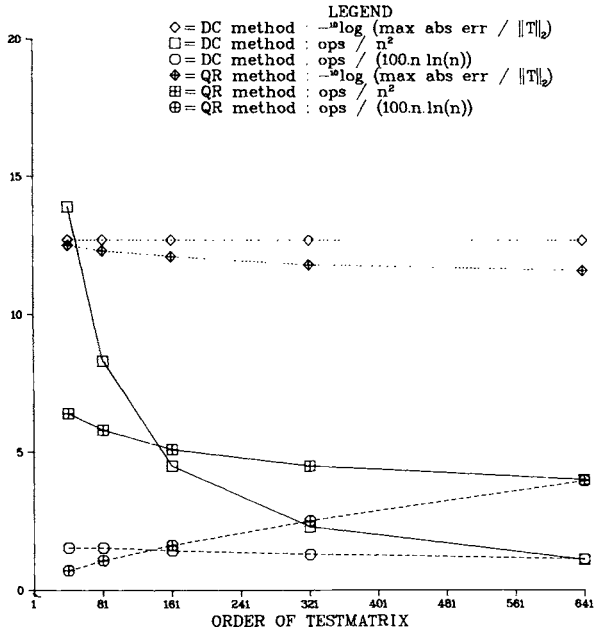


Fig. 3

ACC/EFF PLOT FOR WILKINSON'S MATRICES W+

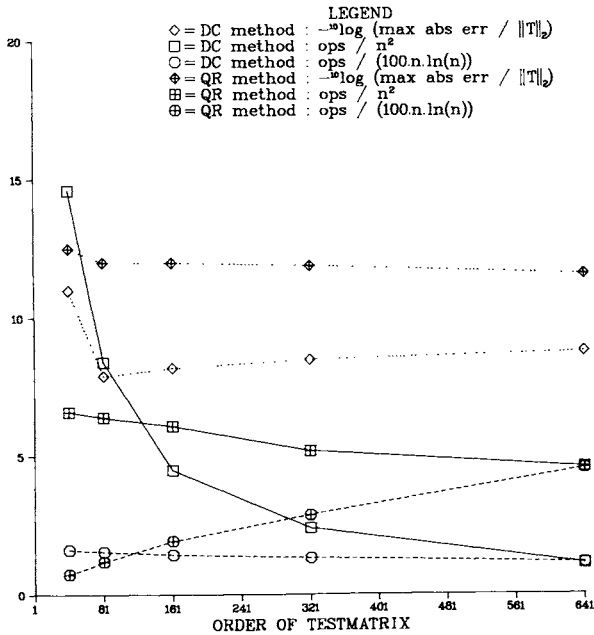


Fig. 4

ACC/EFF PLOT FOR RANDOM TYPE MATRICES WITH F=5

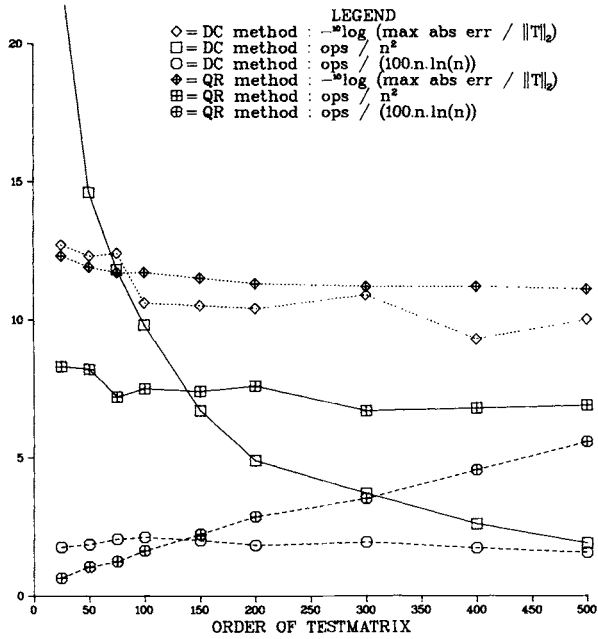


Fig. 5

ACC/EFF PLOT FOR RANDOM TYPE MATRICES WITH F=1

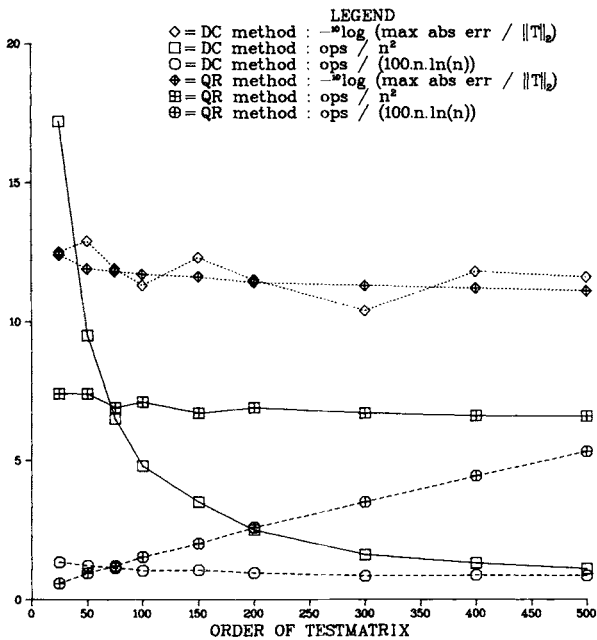


Fig. 6

Table 1

Ex. 7.6		Ex. 7.10	
n	ops/ n^2	n	ops/ n^2
300	15.23	300	11.85
500	13.37	500	11.70
700	13.03	700	11.29

“diagonal dominance” of the matrix and random is a (systems-)procedure which generates a pseudo-random sequence of numbers between 0 and 1 (uniform distribution). The results with $f=1$ were already given in Fig. 1, Fig. 5 gives the results with $f=0.5$ (weakly diagonally dominant, on the average), and Fig. 6 gives the results with $f=0.1$.

We see that in all these cases the *DC* method had a very clear $n \log n$ behaviour. Also the influence of the factor f is clear: the more “diagonally dominant” a matrix is, the more efficiently the *DC* method operates.

Acknowledgement. The author is grateful to T.J. Dekker and W. Hoffmann for helpful discussions and comments.

References

1. Bunch, J.R., Nielsen, C.P., Sorensen, D.C.: Rank one modification of the symmetric eigenproblem. *Numer. Math.* **31**, 31–48 (1978)
2. Bus, J.C., Dekker, T.J.: Two efficient algorithms with guaranteed convergence for finding a zero of a function. *TOMS* **1**, 330–345 (1975)
3. Golub, G.H.: Some modified matrix eigenvalue problems. *SIAM Rev.* **15**, 318–334 (1973)
4. Gregory, R.T., Karney, D.L.: A collection of matrices for testing computational algorithms. New York: John Wiley, 1969
5. Numal, A library of numerical procedures in Algol 60, second revision. Mathematisch Centrum Amsterdam, 1977
6. Wilkinson, J.H.: The algebraic eigenvalue problem. Oxford: Clarendon Press 1965

Received March 24, 1980