



# A real triple dqds algorithm for the nonsymmetric tridiagonal eigenvalue problem

Carla Ferreira<sup>1</sup> · Beresford Parlett<sup>2</sup>

Received: 18 December 2010 / Revised: 14 May 2021 / Accepted: 25 October 2021 /  
Published online: 18 January 2022

© The Author(s), under exclusive licence to Springer-Verlag GmbH Germany, part of Springer Nature 2021

## Abstract

The paper discusses the following topics: attractions of the real tridiagonal case, relative eigenvalue condition number for matrices in factored form, *dqds*, *triple dqds*, error analysis, new criteria for splitting and deflation, eigenvectors of the balanced form, twisted factorizations and generalized Rayleigh quotient iteration. We present our fast real arithmetic algorithm and compare it with alternative published approaches.

**Mathematics Subject Classification** 65F15

## 1 Introduction

The *dqds* algorithm was introduced in [9] as a fast and extremely accurate way to compute all the singular values of a bidiagonal matrix  $B$ . This algorithm implicitly performs the Cholesky LR iteration on the tridiagonal matrix  $B^T B$  and it is used in LAPACK. However the *dqds* algorithm can also be regarded as executing, implicitly, the LR algorithm applied to any tridiagonal matrix with 1's on the superdiagonal. Our interest here is in real unsymmetric matrices which may, of course, have some complex eigenvalues. In contrast to the QR algorithm, the LR algorithm preserves tridiagonal form and this feature makes *dqds* attractive. It is natural to try to retain real arithmetic and yet permit complex conjugate pairs of shifts. Our analogue of the *double*

---

The research of the first author was partially financed by Portuguese Funds through FCT (Fundação para a Ciência e a Tecnologia) within the Projects UIDB/00013/2020 and UIDP/00013/2020.

---

✉ Carla Ferreira  
caferrei@math.uminho.pt  
Beresford Parlett  
parlett@math.berkeley.edu

<sup>1</sup> Centro de Matemática, Universidade do Minho, 4710-057 Braga, Portugal

<sup>2</sup> Division of the Electrical Engineering and Computer Science, Department of Mathematics and the Computer Science, University of California, Berkeley, CA 94720, USA

*shift* QR algorithm of Francis [14] is the *triple step dqds* algorithm. We explain why three steps are needed. However the main goal of this paper is to derive our implicit implementation (*3dqds*) of this 3-steps process which relies on the *implicit L* analogue of the *implicit Q* theorem. See Sect. 3.4.1 and Theorem 3.1.

In order to focus on our *3dqds* algorithm we assume an extensive background for our reader. The unsymmetric eigenvalue problem can be almost ill-posed and such cases are not easily apparent. A tridiagonal matrix requires so little storage that it seems feasible to compute approximate eigenvalues together with an indication of the number of digits that are robust in the presence of computer arithmetic. We decided to provide relative condition numbers (for factored forms) for each computed eigenvalue even when the user does not request it. The extra cost, in storage and arithmetic operations is surprisingly low,  $2n$  storage and  $\mathcal{O}(n)$  computing. See Sect. 7.2 for details. We omit any history of the contributions to the field, even the seminal work of H. Rutishauser who invented the *qd* algorithm and the LR algorithm [26–28]. We must however mention that he also discovered the so-called differential form of *qd* but did not appreciate its accuracy and never published it. That understanding came much later in the computation of singular values of bidiagonal matrices. See [9]. We do describe the double LR algorithm for complex conjugate shifts because of its relation to our *triple dqds* algorithm. We say nothing about the need for an eigensolver devoted to tridiagonal matrices because that issue is covered admirably by Bini et al. [1]. We do give pseudocode for a complete program but hope it will not produce distractions from our main concern, the *3dqds* algorithm. We provide error analyses for both *dqds* and *3dqds*.

A novel feature of our approach is the usefulness of keeping matrices in factored form. We also acknowledge the preliminary work on this problem by Wu [36].

We do not follow Householder conventions except that we reserve capital Roman letters for matrices. Section 2 describes other relevant methods, Sect. 3 presents standard, but needed, material on LR, *dqds*, single and double shifts and the implicit L theorem. Section 4 develops our *3dqds* algorithm, Sect. 5 is our error analysis, Sect. 6 our splitting, deflation and shift strategy.

Section 7 analyzes applications of factored forms - the computation of eigenvectors using twisted factorizations of the balanced form, relative condition numbers and the generalized Rayleigh quotient iteration. Finally, Sect. 8 presents our numerical tests using MATLAB and Sect. 9 gives our conclusions.

## 2 Other methods relevant to *3dqds*

### 2.1 2 steps of LR = 1 step of QR

For a symmetric positive definite tridiagonal matrix 2 steps of the LR (Cholesky) algorithm produces the same matrix as 1 step of the QR algorithm. Less well known is the article by Xu [37] which extends this result when the symmetric matrix is not positive definite. The catch here is that the LR transform, if it exists, does not preserve symmetry. The remedy is to regard similarities by diagonal matrices as “trivial”, always available, operations. Indeed, diagonal similarities cannot introduce zeros into

a matrix. So, when successful, 2 steps of LR are diagonally similar to one step of QR. Even less well known is a short paper by Slemons [31] showing that for a tridiagonal matrix, not necessarily symmetric, 2 steps of LR are diagonally equivalent to 1 step of HR, see [2]. Note that when symmetry disappears then QR is out of the running because it does not preserve the tridiagonal property.

The point of listing these results is to emphasize that 2 steps of LR gives twice as many shift opportunities as 1 step of QR or HR. Thus convergence can be more rapid with LR (or *dqds*) than with QR or HR. This is one of the reasons that *dqds* is faster than QR for computing singular values of bidiagonals. This extra speed is an additional bonus to the fundamental advantage that *dqds* delivers high relative accuracy in all the singular values. The one drawback to *dqds*, for bidiagonals, is that the singular values must be computed in monotone increasing order; QR allows the singular values to be found in any order.

In our case, failure is always possible and so there is no constraint on the order in which eigenvalues are found. The feature of having more opportunities to shift leads us to favor *dqds* over QR and HR. See the list of other methods which follows. We take up the methods in historical order and consider only those that preserve tridiagonal form.

## 2.2 Cullum's complex QR algorithm

As part of a program that used the Lanczos algorithm to reduce a given matrix to tridiagonal form in [4], Jane Cullum used the fact that an unsymmetric tridiagonal matrix may always be balanced by a diagonal similarity transformation [18]. She then observed that another diagonal similarity with 1 or  $i$  produces a symmetric, but complex, tridiagonal matrix to which the (complex) tridiagonal QR algorithm may be applied. The process is not backward stable because the relation

$$\cos^2 \tau + \sin^2 \tau = 1$$

is not a constraint on  $|\cos \tau|$  and  $|\sin \tau|$  when  $\tau$  is complex. Despite the possibility of breakdown the method proved satisfactory in practice. We have not used it in our comparisons because we are persuaded by 2.1 that it is outperformed by the complex *dqds* algorithm, described below.

## 2.3 Liu's HR algorithm

In [16], Liu found a variation on the HR algorithm of Angelika Bunse–Gerstner that, in exact arithmetic, is guaranteed not to breakdown—but the price is a temporary increase in bandwidth. This procedure has only been implemented in MAPLE and we do not include it in our comparison.

## 2.4 Complex dqds

In his thesis Day [5] developed a Lanczos-style algorithm to reduce a general matrix to tridiagonal form and, as with Jane Cullum, needed a suitable algorithm to compute

its eigenvalues. He knew of the effectiveness of  $dqds$  in the symmetric positive definite case and realized that  $dqds$  extends formally to any tridiagonal that admits triangular factorization. The code uses complex arithmetic because of the possible presence of complex conjugate pairs of eigenvalues.

We compare our real  $3dqds$  algorithm with its explicit version—the three steps of  $dqds$  are computed explicitly in complex arithmetic—in a more sophisticated version of David Day's complex  $dqds$  code.

## 2.5 Ehrlich–Aberth algorithm

This very careful and accurate procedure was presented by Bini et al. [1]. It finds the zeros of the characteristic polynomial  $p(\cdot)$  and exploits the tridiagonal form to evaluate  $p'(z)/p(z)$  for any  $z$ . The polynomial solver improves a full set of approximate zeros at each step. Initial approximations are found using a divide-and-conquer procedure that delivers the eigenvalues of the top and bottom halves of the matrix  $T$ . The quantity  $p'(z)/p(z)$  is evaluated indirectly as  $[\text{trace}(zI - T)^{-1}]$  using a QR factorization of  $zI - T$ . Since  $T$  is not altered there is no deflation to assist efficiency. Very careful tests exhibit the method's accuracy - but it is very slow compared to  $dqds$ -type algorithms.

## 3 LR and $dqds$

The reader is expected to have had some exposure to the QR and/or LR algorithms so we will be brief.

### 3.1 LU factorization

Any  $n \times n$  matrix  $A$  permits unique triangular factorization  $A = LD\tilde{U}$  where  $L$  is unit lower triangular,  $D$  is diagonal,  $\tilde{U}$  is unit upper triangular, if and only if the leading principal submatrices of orders  $1, \dots, n - 1$  are nonsingular.

In this paper we follow common practice and write  $U = D\tilde{U}$  so that the “pivots” (entries of  $D$ ) lie on  $U$ 's diagonal. Throughout this paper any matrix  $L$  is unit lower triangular while  $U$  is simply upper triangular.

### 3.2 LR transform with shift

Note that  $U$  is “right” triangular and  $L$  is “left” triangular and this explains the standard name LR. For any shift  $\sigma$  let

$$A - \sigma I = LU, \quad (3.1)$$

$$\hat{A} = UL + \sigma I. \quad (3.2)$$

Then  $\hat{A}$  is the  $\text{LR}(\sigma)$  transform of  $A$ . Note that

$$\hat{A} = L^{-1}(A - \sigma I)L + \sigma I = L^{-1}AL.$$

We say that the shift is restored (in contrast to *dqds*—see below). The LR algorithm consists of repeated LR transforms with shifts chosen to enhance convergence to upper triangular form. For the theory see [28,29,33,34].

In contrast to the well known QR algorithm, the LR algorithm can breakdown and can suffer from element growth,  $\|L\| \gg \|A\|$ ,  $\|U\| \gg \|A\|$ . However LR preserves the banded form of  $A$  while QR does not (except for the Hessenberg form).

When a matrix  $A$  is represented by its entries then the shift operation  $A \rightarrow A - \sigma I$  is trivial. When a matrix is given in factored form the shift operation is not trivial.

### 3.3 The dqds algorithm

From now on we focus on tridiagonal matrices in  $J$ -form which means that entries  $(i, i + 1)$  are all 1,  $i = 1, \dots, n - 1$ . Any tridiagonal matrix  $C = \text{tridiag}(b, a, c)$  that does not split (unreduced),  $b_i c_i \neq 0$ , is diagonally similar to a  $J$ -form. Entries  $(i + 1, i)$  equal  $b_i c_i$ . Throughout this paper all  $J$  matrices have this form. See [11, Section 2.2] on representations of tridiagonals.

If  $J - \sigma I$  permits triangular factorization

$$J - \sigma I = LU$$

then  $L$  and  $U$  must have the following form

$$L = \begin{bmatrix} 1 & & & & & \\ l_1 & 1 & & & & \\ & & \ddots & \ddots & & \\ & & & l_{n-2} & 1 & \\ & & & & l_{n-1} & 1 \end{bmatrix}, \quad U = \begin{bmatrix} u_1 & 1 & & & & \\ & u_2 & 1 & & & \\ & & \ddots & \ddots & & \\ & & & u_{n-1} & 1 & \\ & & & & & u_n \end{bmatrix}. \tag{3.3}$$

It is an attractive feature of LR that

$$UL = \widehat{J}$$

is also of  $J$ -form. Thus the parameters  $l_i, i = 1, \dots, n - 1$ , and  $u_j, j = 1, \dots, n$ , determine the matrices  $L$  and  $U$  above and implicitly define two tridiagonal matrices  $LU$  and  $UL$ .

The *qds* algorithm is equivalent to the LR algorithm but only the factors  $L, U$  are formed, not the  $J$  matrices. The *progressive* transformation is from  $L, U$  to  $\widehat{L}, \widehat{U}$ ,

$$\widehat{L}\widehat{U} = UL - \sigma I. \tag{3.4}$$

Notice that the shift is not restored and so  $\widehat{U}\widehat{L}$  is not similar to  $UL$ ,

$$\widehat{U}\widehat{L} = \widehat{L}^{-1}(UL)\widehat{L} - \sigma I. \tag{3.5}$$

Equating entries in each side of equation (3.4) gives

$$\begin{aligned}
 qds(\sigma) : & \widehat{u}_1 = u_1 + l_1 - \sigma; \\
 & \mathbf{for} \ i = 1, \dots, n - 1 \\
 & \quad \widehat{l}_i = l_i u_{i+1} / \widehat{u}_i \\
 & \quad \widehat{u}_{i+1} = u_{i+1} + l_{i+1} - \sigma - \widehat{l}_i \\
 & \mathbf{end \ for.}
 \end{aligned}$$

The algorithm  $qds$  fails when  $\widehat{u}_i = 0$  for some  $i < n$ . When  $\sigma = 0$  we write simply  $qd$ , not  $qds$ .

In 1994 a better way was found to implement  $qds(\sigma)$ . These are called *differential qd* algorithms. See [21] for more history. This form uses an extra variable  $d$  but never forms matrix products.

$$\begin{aligned}
 dqds(\sigma) : & d_1 = u_1 - \sigma \\
 & \mathbf{for} \ i = 1, \dots, n - 1 \\
 & \quad \widehat{u}_i = d_i + l_i \\
 & \quad \widehat{l}_i = l_i (u_{i+1} / \widehat{u}_i) \\
 & \quad d_{i+1} = d_i (u_{i+1} / \widehat{u}_i) - \sigma \\
 & \mathbf{end \ for} \\
 & \widehat{u}_n = d_n.
 \end{aligned}$$

By definition,  $dqd = dqds(0)$ . In practice we choose to compute, and store,  $\widehat{u}_i$  and  $\widehat{l}_i$  separately from  $u_i$  and  $l_i$ . This allows us to reject a transform, choose a new shift, and proceed smoothly to another step. Only when the transform is accepted will we write  $\widehat{u}_i$  and  $\widehat{l}_i$  over  $u_i$  and  $l_i$ .

A word on terminology. In Rutishauser's original work  $q_i = u_i$ ,  $e_i = l_i$ ; and the  $q_i$ 's were certain *quotients* and the  $e_i$ 's were called *modified differences*. In fact the  $qd$  algorithm led to the LR algorithm, not vice-versa. The reader can find more information concerning  $dqds$  in [21,23]

One virtue of the  $dqds$  and QR transforms is that they work on the whole matrix so that large eigenvalues are converging near the top, albeit slowly, while the small ones are being picked off at the bottom.

We summarize some advantages and disadvantages of the factored form.

### 3.3.1 Advantages of the factored form

1.  $L, U$  determines the entries of  $J$  to greater than working-precision accuracy because the addition and multiplication of  $l$ 's and  $u$ 's is implicit. Thus, for instance, the  $(i, i)$  entry of  $J$  is given by  $l_{i-1} + u_i$  implicitly but  $fl(l_{i-1} + u_i)$  explicitly.
2. Singularity of  $J$  is detectable by inspection when  $L$  and  $U$  are given, but only by calculation from  $J$ . So,  $LU$  reveals singularity,  $J$  does not.
3.  $LU$  defines the eigenvalues better than  $J$  does (usually). There is more on this in [7].
4. Solution of  $Jx = b$  takes half the time when  $L$  and  $U$  are available.

### 3.3.2 Disadvantages of the factored form

The mapping  $J, \sigma \mapsto L, U$  is not everywhere defined for all pairs  $J, \sigma$  and can suffer from element growth. This defect is not as serious as it was when the new transforms were written over the old ones. For tridiagonals we can afford to double the storage and map  $L, U$  into different arrays  $\widehat{L}, \widehat{U}$ . Then we can decide whether or not to accept  $\widehat{L}, \widehat{U}$  and only then would  $L$  and  $U$  be overwritten. So the difficulty of excessive element growth has been changed from disaster to the non-trivial but less intimidating one of, after rejecting a transform, choosing a new shift that will not spoil convergence and will not cause another rejection.

Now we turn to our main question of  $dqds(\sigma)$ : how can complex shifts be used without having to use complex arithmetic? This question has a beautiful answer for QR and LR iterations.

### 3.4 Implicit shifted LR for $J$ matrices

#### 3.4.1 Double shift LR algorithm

We use the  $J, L$  and  $U$  notation from the previous section. Consider two steps of the LR algorithm with shifts  $\sigma_1$  and  $\sigma_2$ ,

$$\begin{aligned} J_1 - \sigma_1 I &= L_1 U_1 \\ J_2 &= U_1 L_1 + \sigma_1 I \\ J_2 - \sigma_2 I &= L_2 U_2 \\ J_3 &= U_2 L_2 + \sigma_2 I. \end{aligned} \tag{3.6}$$

Then, with matrices  $\mathcal{L} = L_1 L_2$  and  $\mathcal{U} = U_2 U_1$ , we have

$$J_3 = \mathcal{L}^{-1} J_1 \mathcal{L} \tag{3.7}$$

and the triangular factorization

$$\begin{aligned} \mathcal{L}\mathcal{U} &= L_1 (J_2 - \sigma_2 I) U_1 = L_1 (U_1 L_1 + \sigma_1 I - \sigma_2 I) U_1 \\ &= L_1 U_1 [L_1 U_1 + \sigma_1 I - \sigma_2 I] = (J_1 - \sigma_1 I) (J_1 - \sigma_2 I) \\ &= J_1^2 - (\sigma_1 + \sigma_2) J_1 + \sigma_1 \sigma_2 I =: M. \end{aligned} \tag{3.8}$$

An important observation from (3.8) is that column 1 of  $M$  is proportional to column 1 of  $\mathcal{L}$ ,

$$M e_1 = \mathcal{L} \mathcal{U} e_1 = \mathcal{L} e_1 u_{11}, \quad u_{11} = m_{11}.$$

According to the following theorem, matrix  $\mathcal{L}$  is determined by its first column and we can compute  $J_3$  from  $J_1$  without using  $J_2$ .

**Theorem 3.1** [IMPLICIT L THEOREM] *If  $H_1$  and  $H_2$  are unreduced upper Hessenberg matrices and  $H_2 = L^{-1}H_1L$ , where  $L$  is unit lower triangular, then  $H_2$  and  $L$  are completely determined by  $H_1$  and column 1 of  $L$ .*

We omit the proof and refer to [11, pp. 66–68].

So the application to  $J_1$  and  $J_3$ , using (3.7), is to choose column 1 of  $\mathcal{L}$  (which has only three nonzero entries since  $J_1$  is tridiagonal and  $J_1^2$  is pentadiagonal) to be

$$\mathcal{L}_1 = I + \mathbf{m}_1 \mathbf{e}_1^T \tag{3.9}$$

where  $\mathbf{m}_1 = [0 \ m_{21}/m_{11} \ m_{31}/m_{11} \ 0 \ \dots \ 0]^T$  and perform a first explicit similarity transform on  $J_1$ ,

$$\mathcal{L}_1^{-1} J_1 \mathcal{L}_1 =: K.$$

Observe that  $K$  is not tridiagonal. In the  $6 \times 6$  case

$$K = \begin{bmatrix} x & 1 & & & & \\ x & x & 1 & & & \\ + & x & x & 1 & & \\ + & & x & x & 1 & \\ & & & x & x & 1 \\ & & & & x & x \end{bmatrix}. \tag{3.10}$$

Next we apply a sequence of elementary similarity transformations such that each transformation pushes the  $2 \times 1$  bulge one row down and one column to the right. Finally the bulge is chased off the bottom to restore the  $J$ -form. In exact arithmetic, the implicit L theorem ensures that this technique of *bulge chasing* gives

$$J_3 = (\mathcal{L}_1 \dots \mathcal{L}_{n-1})^{-1} J_1 (\mathcal{L}_1 \dots \mathcal{L}_{n-1}) \quad \text{and} \quad \mathcal{L} = \mathcal{L}_1 \dots \mathcal{L}_{n-1}.$$

In Sect. 4.1 below we will see the details on  $\mathcal{L}_j$ ,  $j = 2, \dots, n - 1$ .

If matrix  $J_1$  and shifts  $\sigma_1$  and  $\sigma_2$  are real then factors  $L_1, U_1, L_2, U_2$  and matrices  $J_2, J_3$  will all be real. Now suppose that  $J_1$  is real and  $\sigma_1$  is complex. Then, by (3.8),  $J_3$  will be real if, and only if,  $\sigma_2 = \bar{\sigma}_1$ . The reason is that  $M$  is real,

$$M = J_1^2 - 2(\Re \sigma_1) J_1 + |\sigma_1|^2 I,$$

so that  $\mathcal{L}$  and  $\mathcal{U}$  are real and  $J_3$  is the product of real matrices. Note however that  $J_2$ , and factors  $L_2, U_2$ , will be complex, given that  $L_1, U_1$  and shifts  $\sigma_1, \sigma_2$  are all complex. As we have described above, it is possible to skip this complex matrix  $J_2$  and go straight from real  $J_1$  to real  $J_3$ . So, in the case of complex eigenvalues (which for real matrices occur in complex conjugate pairs) we will be able to apply a complex conjugate pair of shifts implicitly and avoid complex arithmetic. Recall that we are seeking an algorithm that uses only real arithmetic and converges to real Schur form.



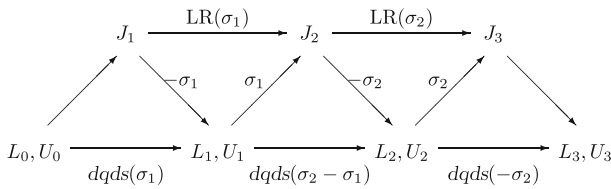


Fig. 1 Implicit two steps of LR and three steps of dqds

### 3.4.2 Connection to dqds algorithm

In Fig. 1 we examine the two steps of the LR transform derived in the previous section but with a significant difference. Instead of  $J_1$  being an arbitrary matrix in  $J$ -form, we assume that it is given to us in the form  $U_0L_0$ . A different way of introducing this factorization is saying that our initial matrix is  $J_0$  (not  $J_1$ ) and we always consider an additional unshifted LR step for constructing real factored forms

$$J_0 = L_0U_0 \text{ and } J_1 = U_0L_0$$

so that *dqds* starts with factors  $L_0, U_0$ . The *dqds* algorithm can not start with  $J_1$  unless its  $UL$  factorization is available.

The crucial observation in Fig. 1 is that the implicit LR algorithm forms only the  $J$  matrices while *dqds*, on the bottom line, works only on the factors  $L, U$ . So with LR the only  $J$  matrix which is skipped by an implicit double step is  $J_2$  and we go from  $J_1 = U_0L_0$  to  $J_3 = L_3U_3$ . The *dqds* algorithm cannot stop with  $L_2, U_2$  because it is only when we take the product  $U_2L_2$  and add back the shift  $\sigma_2$  that we get the matrix  $J_3$ ; it requires a third step to obtain  $L_3, U_3$  which define  $J_3$ . The *triple dqds* algorithm will skip the factors  $L_1, U_1, L_2, U_2$  and will go from  $L_0, U_0$  to  $L_3, U_3$  performing implicitly three *dqds* steps.

Here is another way to describe the diagonal arrows in Fig. 1 for the relation between double shift LR and *triple dqds*:

double shift LR	triple dqds	
$J_1 = U_0L_0$		(3.11)
$J_1 - \sigma_1 I = L_1U_1$	$L_1U_1 = U_0L_0 - \sigma_1 I$	
$J_2 = U_1L_1 + \sigma_1 I$		
$J_2 - \sigma_2 I = L_2U_2$	$L_2U_2 = U_1L_1 - (\sigma_2 - \sigma_1)I$	
$J_3 = U_2L_2 + \sigma_2 I$	$L_3U_3 = U_2L_2 - (-\sigma_2)I$	

So the similarity (3.7) corresponds to

$$L_3U_3 = \mathcal{L}^{-1}(U_0L_0)\mathcal{L} \tag{3.12}$$

and, in contrast to a single *dqds* step, a *triple dqds* step (implicit) restores the shifts.

Observe that in the *triple dqds* step (3.11) we find factors  $L_3, U_3$  such that  $J_3 = L_3U_3$  and these factors (different factors) would only occur in LR in the following step with a new shift  $\sigma_3$ ,

$$\left. \begin{aligned} J_3 &= U_2L_2 + \sigma_2I \\ J_3 - \sigma_3I &= L_3U_3 \end{aligned} \right\} \quad L_3U_3 = U_2L_2 - (\sigma_3 - \sigma_2)I \quad (3.13)$$

$$J_4 = U_3L_3 + \sigma_3I.$$

So to make explicit *dqds* equivalent to LR with shifts  $\sigma_i$  and  $\sigma_{i+1}$  it is necessary to use the differences  $(\sigma_{i+1} - \sigma_i)$  with *dqds*. In other words, successive shifts  $\sigma_i$  and  $\sigma_{i+1}$  in LR lead to the *dqds* step

$$L_{i+1}U_{i+1} = U_iL_i - (\sigma_{i+1} - \sigma_i)I.$$

### 3.4.3 Single shift LR and double *dqds*

Analogously to a double shift, a single shift LR step is equivalent to two steps of *dqds* when we consider the implicit implementation of these shifted algorithms.

<b>single shift LR</b>	<i>double dqds</i>	
$\left. \begin{aligned} J_1 &= U_0L_0 \\ J_1 - \sigma_1I &= L_1U_1 \end{aligned} \right\}$	$L_1U_1 = U_0L_0 - \sigma_1I$	(3.14)
$J_2 = U_1L_1 + \sigma_1I$	$L_2U_2 = U_1L_1 - (-\sigma_1)I$	

Similar to (3.7) and (3.12),

$$J_2 = \mathcal{L}^{-1}J_1\mathcal{L} \quad \text{and} \quad L_2U_2 = \mathcal{L}^{-1}(U_0L_0)\mathcal{L} \quad (3.15)$$

with  $\mathcal{L} = L_1$ . Here matrix  $M$  is tridiagonal,

$$\mathcal{L}\mathcal{U} = L_1U_1 = J_1 - \sigma_1 =: M$$

and matrix  $K$  corresponding to (3.10) has only one bulge in entry (3, 1) (instead of a  $2 \times 1$  bulge).

Recall from Sect. 3.4.1 that the implicit double LR algorithm uses the technique of *bulge chasing*. This technique is also applied for implicit single shifts.

The next section develops a form of bulge chasing for the *triple dqds* algorithm (*3dqds*). We did not develop this technique for the *double dqds* algorithm (*2dqds*) because in our shift strategy we will always use double shifts (only initially we use single *dqds* with zero shifts). See Sect. 6.3 for details on the shift strategy in our complete algorithm.

## 4 Triple *dqds* algorithm

We use the term *3dqds* as a shorthand for our *triple dqds* algorithm which, using bulge chasing, implements implicitly the three *dqds* steps (3.11) equivalent to an implicit double shift LR step. Although the *3dqds* algorithm has been primarily developed to avoid complex arithmetic in the case of consecutive complex shifts  $\sigma_1$  and  $\sigma_2 = \bar{\sigma}_1$

in the presence of complex eigenvalues, it can be applied to the case of two real shifts  $\sigma_1$  and  $\sigma_2$ . To cover both cases, all we need is the sum and the product of the pair of shifts,  $\text{sum} = \sigma_1 + \sigma_2$  and  $\text{prod} = \sigma_1\sigma_2$ , to form matrix  $M$  in (3.8),

$$M = (U_0L_0)^2 - \text{sum}(U_0L_0) + \text{prod}I. \tag{4.1}$$

Using (3.12) the idea is to transform  $U_0$  into  $L_3$  and  $L_0$  into  $U_3$  by bulge chasing in each matrix,

$$L_3U_3 = \underbrace{\mathcal{L}^{-1}U_0}_{\mathcal{L}_1} \underbrace{L_0\mathcal{L}}_{\mathcal{L}_2}.$$

Notice that we need to transform an upper bidiagonal into a lower bidiagonal and vice-versa. From the uniqueness of the  $LU$  factorization, when it exists, it follows, see [20], that there is a unique hidden matrix  $X$  such that

$$L_3 = \mathcal{L}^{-1}U_0X^{-1}, \quad XL_0\mathcal{L} = U_3.$$

The matrix  $\mathcal{L}$  is given, from Sect. 3.4.1, as a product

$$\mathcal{L} = \mathcal{L}_1 \dots \mathcal{L}_{n-1}\mathcal{L}_n$$

( $\mathcal{L}_n = I$ ) and we will gradually construct the matrix  $X$  in corresponding factored form  $X_n \dots X_2X_1$ . In fact we will write each  $X_i$  as a product

$$X_i = Y_iZ_i.$$

Matrices  $\mathcal{L}_i$  and  $Y_i$  are elementary matrices,  $\mathcal{L}_i = I + \mathbf{m}_i\mathbf{e}_i^T$  and  $Y_i = I + \mathbf{w}_i\mathbf{e}_i^T$ , but  $Z_i$  is not. The details are quite complicated and will be shown in the following sections.

### 4.1 Chasing the bulges

Starting with the factors  $L_0, U_0$  and the shifts  $\sigma_1, \sigma_2$ , we normalize column 1 of  $M$  in (4.1) to form  $\mathcal{L}_1$ , spoil the bidiagonal form with

$$\underbrace{\mathcal{L}_1^{-1}U_0}_{\mathcal{L}_1} \underbrace{L_0\mathcal{L}_1}_{\mathcal{L}_2}$$

and at each *minor* step  $i, i = 1, \dots, n$ , matrices  $Z_i, \mathcal{L}_i$  and  $Y_i$  are chosen to chase the bulges. After  $n$  minor steps, we obtain  $L_3$  and  $U_3$ ,

$$\begin{aligned} L_3U_3 &= \underbrace{\mathcal{L}_n^{-1} \dots \mathcal{L}_1^{-1}U_0Z_1^{-1}Y_1^{-1} \dots Z_n^{-1}Y_n^{-1}}_{\mathcal{L}_1} \underbrace{Y_nZ_n \dots Y_1Z_1L_0\mathcal{L}_1 \dots \mathcal{L}_n}_{\mathcal{L}_2} \\ &= \underbrace{\mathcal{L}_n^{-1} \dots \mathcal{L}_1^{-1}U_0X_1^{-1} \dots X_n^{-1}}_{\mathcal{L}_1} \underbrace{X_n \dots X_1L_0\mathcal{L}_1 \dots \mathcal{L}_n}_{\mathcal{L}_2} \\ &= \underbrace{\mathcal{L}^{-1}U_0X^{-1}}_{\mathcal{L}_1} \underbrace{XL_0\mathcal{L}}_{\mathcal{L}_2} \end{aligned}$$

All the work of bulge chasing will be confined to two matrices  $F$  and  $G$ . Initially,

$$F = U_0, \quad G = L_0$$

and, finally,

$$F = L_3, \quad G = U_3.$$

For a pair of shifts  $\sigma_1$  and  $\sigma_2$  (real or a complex conjugate pair), the *triple dqds* algorithm has the following matrix formulation:

```

3dqds( $\sigma_1, \sigma_2$ ) :
% step 1
 $F = U_0; G = L_0$ 
 $F = FZ_1^{-1}; G = Z_1G$ 
 $F = \mathcal{L}_1^{-1}F; G = G\mathcal{L}_1$  [form  $\mathcal{L}_1$  using (3.9) and (4.1)]
 $F = FY_1^{-1}; G = Y_1G$ 

% steps 2 to n-3
for  $i = 2, \dots, n-3$ 
     $F = FZ_i^{-1}; G = Z_iG$ 
     $F = \mathcal{L}_i^{-1}F; G = G\mathcal{L}_i$ 
     $F = FY_i^{-1}; G = Y_iG$  [Zi with one, Li with two and Yi with three
end for [nonzero off-diagonal entries]

% step n-2
 $F = FZ_{n-2}^{-1}; G = Z_{n-2}G$ 
 $F = \mathcal{L}_{n-2}^{-1}F; G = G\mathcal{L}_{n-2}$ 
 $F = FY_{n-2}^{-1}; G = Y_{n-2}G$  [Yn-2 with two nonzero off-diagonal entries]

% step n-1
 $F = FZ_{n-1}^{-1}; G = Z_{n-1}G$ 
 $F = \mathcal{L}_{n-1}^{-1}F; G = G\mathcal{L}_{n-1}$ 
 $F = FY_{n-1}^{-1}; G = Y_{n-1}G$  [Yn-1 and Ln-1 with one nonzero off-diagonal entry]

% step n
 $\mathcal{L}_n = I; Y_n = I$ 
 $F = FZ_n^{-1}; G = Z_nG$  [Zn diagonal]
 $L_3 = F; U_3 = G$ 

```



- (b)  $F \leftarrow \mathcal{L}_i^{-1}F$  puts 0 into  $F_{i+1,i-1}$  and  $F_{i+2,i-1}$ , and moves the bulge to column  $i$   $G \leftarrow G\mathcal{L}_i$  creates  $\widehat{u}_i$  in  $G_{i,i}$  and changes  $G_{i+1,i}$ ,  $G_{i+2,i}$  and  $G_{i+3,i}$  below it

$$\begin{aligned}
 \begin{matrix} x_l \leftarrow -x_l/\widehat{l}_{i-1}, \\ y_l \leftarrow -y_l/\widehat{l}_{i-1}, \end{matrix} \quad \mathcal{L}_i^{-1} = \begin{bmatrix} \ddots & & & & & \\ & 1 & & & & \\ & x_l & 1 & & & \\ & y_l & & 1 & & \\ & & & & \ddots & \\ & & & & & 1 \end{bmatrix}, \quad \mathcal{L}_i = \begin{bmatrix} \ddots & & & & & \\ & 1 & & & & \\ & -x_l & 1 & & & \\ & -y_l & & 1 & & \\ & & & & \ddots & \\ & & & & & 1 \end{bmatrix} \\
 \\
 \mathcal{L}_i^{-1}F = \begin{bmatrix} \ddots & & & & & \\ & 1 & & & & \\ & \widehat{l}_{i-1} & \mathbf{1} & \mathbf{0} & & \\ & \mathbf{0} & x_l & u_{i+1} & 1 & \\ & \mathbf{0} & y_l & & u_{i+2} & \ddots \\ & & & & & \ddots \\ & & & & & & 1 \end{bmatrix}, \quad G\mathcal{L}_i = \begin{bmatrix} \ddots & & & & & \\ & \widehat{u}_{i-1} & 1 & & & \\ & \widehat{u}_i & \mathbf{1} & & & \\ & x_r & 1 & & & \\ & y_r & l_{i+1} & 1 & & \\ & z_r & & l_{i+2} & 1 & \\ & & & & & \ddots \\ & & & & & & 1 \end{bmatrix} \\
 \\
 \begin{matrix} \widehat{u}_i \leftarrow x_r - x_l \\ x_r \leftarrow y_r - x_l \\ y_r \leftarrow z_r - y_l - x_l * l_{i+1} \\ z_r \leftarrow -y_l * l_{i+2} \end{matrix}
 \end{aligned}$$

- (c)  $G \leftarrow Y_i G$  puts 0 into  $G_{i+1,i}$ ,  $G_{i+2,i}$  and  $G_{i+3,i}$ , and moves the bulge to column  $i+1$   $F \leftarrow FY_i^{-1}$  creates  $\widehat{l}_i$  in  $F_{i+1,i}$  and changes  $F_{i+2,i}$  and  $F_{i+3,i}$  (bulge in  $F$ ) below it

$$\begin{aligned}
 \begin{matrix} x_r \leftarrow x_r/\widehat{u}_i \\ y_r \leftarrow y_r/\widehat{u}_i, \\ z_r \leftarrow z_r/\widehat{u}_i \end{matrix} \quad Y_i^{-1} = \begin{bmatrix} \ddots & & & & & \\ & 1 & & & & \\ & x_r & 1 & & & \\ & y_r & & 1 & & \\ & z_r & & & 1 & \\ & & & & & \ddots \\ & & & & & & 1 \end{bmatrix}, \quad Y_i = \begin{bmatrix} \ddots & & & & & \\ & 1 & & & & \\ & -x_r & 1 & & & \\ & -y_r & & 1 & & \\ & -z_r & & & 1 & \\ & & & & & \ddots \\ & & & & & & 1 \end{bmatrix} \\
 \\
 FY_i^{-1} = \begin{bmatrix} \ddots & & & & & \\ & 1 & & & & \\ & \widehat{l}_{i-1} & \mathbf{1} & \mathbf{0} & & \\ & \mathbf{0} & \widehat{l}_i & u_{i+1} & 1 & \\ & \mathbf{0} & x_l & & u_{i+2} & 1 \\ & & y_l & & & u_{i+3} & \ddots \\ & & & & & & \ddots \\ & & & & & & & 1 \end{bmatrix}, \quad Y_i G = \begin{bmatrix} \ddots & & & & & \\ & \widehat{u}_{i-1} & 1 & & & \\ & \widehat{u}_i & \mathbf{1} & & & \\ & \mathbf{0} & x_r & & & \\ & \mathbf{0} & y_r & 1 & & \\ & \mathbf{0} & z_r & l_{i+2} & 1 & \\ & & & & & \ddots \\ & & & & & & 1 \end{bmatrix}
 \end{aligned}$$

$$\begin{array}{ll}
 \widehat{l}_i \leftarrow x_l + y_r + x_r * u_{i+1} & x_r \leftarrow 1 - x_r \\
 x_l \leftarrow y_l + z_r + y_r * u_{i+2} & y_r \leftarrow l_{i+1} - y_r \\
 y_l \leftarrow z_r * u_{i+3} & z_r \leftarrow -z_r
 \end{array}$$

The result of this minor step is that the active windows of  $F$  and  $G$  shown in (4.2) have been moved down and to the right by one place.

Naturally steps 1,  $n - 2$ ,  $n - 1$ ,  $n$  are slightly different and their practical implementation may be found in [11, pp.147–157].

### 4.3 Comparison of dqds and 3dqds

In this section we present a detailed version of the inner loop of 3dqds and compare one step of 3dqds with three steps of simple dqds in terms of arithmetic effort.

Here is the inner loop of 3dqds. See “Appendix A” for the whole 3dqds algorithm.

3dqds( $\sigma_1, \sigma_2$ ) :

```

for  $i = 2, \dots, n - 3$ 
     $x_r = x_r * u_i + y_r$ 
     $x_l = -x_l / \widehat{l}_{i-1}; \quad y_l = -y_l / \widehat{l}_{i-1};$ 
     $\widehat{u}_i = x_r - x_l;$ 
     $x_r = y_r - x_l; \quad y_r = z_r - y_l - x_l * l_{i+1}; \quad z_r = -y_l * l_{i+2}$ 
     $x_r = x_r / \widehat{u}_i; \quad y_r = y_r / \widehat{u}_i; \quad z_r = z_r / \widehat{u}_i$ 
     $\widehat{l}_i = x_l + y_r + x_r * u_{i+1}$ 
     $x_l = y_l + z_r + y_r * u_{i+2}; \quad y_l = z_r * u_{i+3}$ 
     $x_r = 1 - x_r; \quad y_r = l_{i+1} - y_r; \quad z_r = -z_r$ 
end for
    
```

In contrast,

dqds( $\sigma$ ) :

$$\begin{array}{l}
 d_1 = u_1 - \sigma \\
 \mathbf{for} \ i = 1, \dots, n - 1 \\
 \quad \widehat{u}_i = d_i + l_i \\
 \quad t = u_{i+1} / \widehat{u}_i \\
 \quad d_{i+1} = d_i t - \sigma \\
 \quad \widehat{l}_i = l_i t \\
 \mathbf{end\ for} \\
 \widehat{u}_n = d_n.
 \end{array} \tag{4.3}$$

In practice, each  $d_{i+1}$  may be written over its predecessor in a single variable  $d$  and, using and auxiliary variable  $t$ , only one division is needed.

Table 1 below shows that the number of floating point operations required by one step of 3dqds is comparable to three steps of dqds (table expresses only the number of floating point operations in the inner loops).

**Table 1** Operation count of  $3dqds$  and  $3dqds$  steps

	$3dqds$	$3dqds$ steps
Divisions	5	3
Multiplications	6	6
Additions	5	3
Subtractions	6	3
Assignments	16	12
Auxiliary variables	5	2

However to accomplish a complex conjugate pair of shifts these 3  $dqds$  steps will be complex in contrast to our  $3dqds$  which uses only real arithmetic. Thus 3 steps of complex  $dqds$  take more time than one step of  $3dqds$  (complex arithmetic raises the cost by a factor of about 4 [6, p.163]).

## 5 Error analysis

We turn to the effect of finite precision arithmetic on our algorithms. First consider the  $dqds$  algorithm.

### 5.1 $dqds$

It is well known that even in exact arithmetic the  $dqds$  iteration, applied to the factors  $L, U$  of a  $J$  matrix can break down due to a zero pivot in the new factors. The  $dqds$  transform, just like the LR transform, is unstable. In the early days when the new was written over the old immediately breakdown was a disaster. Today all users can afford to store the new factors separately from the old and simply reject a transform with unacceptable element growth, choose a new shift, and continue the iteration. A rejection is a nuisance, not a disaster.

One of the attractions of  $dqds$  is that it has *high mixed relative stability*, to be explained below. One of us proved this in [9] in the context of singular values of bidiagonals and eigenvalues of real symmetric tridiagonals. Since this desired property is independent of symmetry, we take this opportunity to present the result again in the context of  $J$  matrices.

To set up notation for the proof we consider real bidiagonal factors  $L$  and  $U$  of a real  $J$  matrix together with a real shift  $\sigma$  and use the  $dqds$  transform to obtain output  $\widehat{L}, \widehat{U}$  satisfying

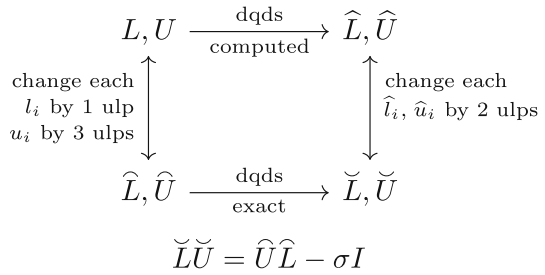
$$\widehat{L}\widehat{U} = UL - \sigma I. \quad (5.1)$$

We assume no anomalies occur, i.e. no divisions by zero, no overflow/underflow.

**Theorem 5.1** ([9], Theorem 4) *Let  $dqds(\sigma)$  map  $L, U$  into computed  $\widehat{L}, \widehat{U}$  with no anomalies. Then well chosen small relative changes in the entries of both input and*



**Fig. 2** Effects of roundoff for dqds



output matrices, of at most 3 ulps each, produces new matrices, one pair mapped into the other, in exact arithmetic, by dqds( $\sigma$ ).

Our analysis consists to write down the exact relations satisfied by the computed quantities  $\widehat{L}, \widehat{U}$  and then to work out among them an exact dqds transform whose input is close to  $L, U$  and output is close to  $\widehat{L}, \widehat{U}$ . The diagram in Fig. 2 is an excellent summary.

The model of arithmetic we assume is that the floating point result of a basic arithmetic operation  $\odot$  (one of the four binary operations  $+, -, *$  and  $/$ ) satisfies

$$fl(x \odot y) = (x \odot y)(1 + \varepsilon) = (x \odot y)/(1 + \eta) \tag{5.2}$$

where  $\varepsilon$  and  $\eta$  depend on  $x, y$ , and the operation  $\odot$ , and satisfy

$$|\varepsilon| < \epsilon, \quad |\eta| < \epsilon. \tag{5.3}$$

The quantity  $\epsilon$  is called variously *roundoff unit*, *machine precision* or *macheps*. We will choose freely the form ( $\varepsilon$  or  $\eta$ ) that suits the analysis. We will also use the acronym *ulp* which stands for *units in the last place held* and it is the natural way to refer to relative differences between numbers.

Our result is possible because of the simple form of the recurrence for the auxiliary variable  $\{d_i\}_{i=1}^n$ . In exact arithmetic

$$d_1 = u_1 - \sigma, \quad d_{i+1} = \frac{d_i u_{i+1}}{d_i + l_i} - \sigma, \quad i = 1, \dots, n - 1.$$

**Proof** We consider the floating point implementation of dqds in (4.3). The computed quantities  $\widehat{L}, \widehat{U}$  are governed by the following exact relations.

$$\begin{aligned} \widehat{u}_i &= fl(d_i + l_i) = (d_i + l_i)/(1 + \varepsilon_+) \\ t &= fl(u_{i+1}/\widehat{u}_i) = \frac{u_{i+1}(1 + \varepsilon_)}{\widehat{u}_i} = \frac{u_{i+1}(1 + \varepsilon_)(1 + \varepsilon_)}{d_i + l_i} \\ d_{i+1} &= fl(fl(d_i * t) - \sigma) = \frac{d_i t(1 + \varepsilon_*) - \sigma}{1 + \varepsilon_{i+1}} \\ \widehat{l}_i &= fl(l_i * t) = l_i t/(1 + \varepsilon_{**}) \end{aligned}$$

The symbol  $\varepsilon_{**}$  represents the rounding error in the second multiplication  $l_i * t_i$ . All the  $\varepsilon$ 's obey (5.3) and depend on  $i$  but we suppress this dependence on  $i$  except for the subtraction of the shift  $\sigma$ . Here  $\varepsilon_{i+1}$  accounts for the error in subtracting the real shift  $\sigma$ . To be consistent we must also use  $d_i(1 + \varepsilon_i)$ , where  $\varepsilon_i$  is defined in minor step  $i - 1$ , and  $(1 + \varepsilon_1)d_1 = u_1 - \sigma$ . Here  $t$  is just an auxiliary variable for the analysis.

Now we can write an exact  $dqds$  transform using  $[\cdot]$  to surround our chosen variables.

$$\begin{aligned} [\widehat{u}_i(1 + \varepsilon_+)(1 + \varepsilon_i)] &= [d_i(1 + \varepsilon_i)] + [l_i(1 + \varepsilon_i)] \\ [d_{i+1}(1 + \varepsilon_{i+1})] &= \frac{[d_i(1 + \varepsilon_i)][u_{i+1}(1 + \varepsilon_+)(1 + \varepsilon_+)(1 + \varepsilon_*)]}{[d_i(1 + \varepsilon_i)] + [l_i(1 + \varepsilon_i)]} - \sigma \\ [\widehat{l}_i(1 + \varepsilon_*)(1 + \varepsilon_{**})] &= \frac{[l_i(1 + \varepsilon_i)][u_{i+1}(1 + \varepsilon_+)(1 + \varepsilon_+)(1 + \varepsilon_*)]}{[d_i(1 + \varepsilon_i)] + [l_i(1 + \varepsilon_i)]} \end{aligned}$$

We can read off the perturbations, defining  $\widehat{l}_i$ ,  $\widehat{u}_{i+1}$  and  $\check{l}_i$ ,  $\check{u}_i$  on the way to an exact transform:

$$\begin{aligned} l_i &\longrightarrow l_i(1 + \varepsilon_i) =: \widehat{l}_i & \widehat{l}_i &\longrightarrow \widehat{l}_i(1 + \varepsilon_*)(1 + \varepsilon_{**}) =: \check{l}_i \\ u_{i+1} &\longrightarrow u_{i+1}(1 + \varepsilon_+)(1 + \varepsilon_+)(1 + \varepsilon_*) =: \widehat{u}_{i+1} & \widehat{u}_i &\longrightarrow \widehat{u}_i(1 + \varepsilon_+)(1 + \varepsilon_i) =: \check{u}_i \end{aligned}$$

The perturbations are as claimed in the theorem: 3 ulps for  $u_i$  and 1 ulp for  $l_i$ , and 2 ulps each for  $\widehat{l}_i$  and  $\widehat{u}_i$  as shown in Fig. 2. Notice that our choices of  $\widehat{L}$ ,  $\widehat{U}$  and  $\check{L}$ ,  $\check{U}$  are not in general machine representable.

When  $\sigma = 0$  the  $(1 + \varepsilon_i)$  factors are omitted but still 3 ulps are needed for  $u_{i+1}$ . □

The remarkable feature here is that element growth does not impair the result. However,

Theorem 5.1 does not guarantee that  $dqds$  returns accurate eigenvalues. For that, an extra requirement is needed such as positivity of all the parameters  $u_j$ ,  $l_j$  in the computation, as is the case for the eigenvalues of  $B^T B$  where  $B$  is upper bidiagonal.

We mention that Yao Yang considered the roundoff in  $dqds$  in his dissertation at UC, Berkeley, in 1994 [38]. He had two results. He gave an  $n = 4$  example to show that even  $dqd$  (no shift in  $dqds$ ) is not backward stable. He also produced an *a posteriori* (computable) bound on the error in the exact product  $\widehat{L}\widehat{U}$  of the output matrices. Unfortunately, his dissertation has not been published but his results are stated and proved in [21].

### 5.2 3dqds

Each minor step in the algorithm consists of 3 simple transformations on work matrices  $F$  and  $G$ . All three parts arise from similarities that chase the bulges in the transformation from  $U_0L_0$  to  $L_3U_3$ . See Sect. 4. Two of these transformations are elementary transformations of the form  $I + \mathbf{v}e_j^T$ , with inverse  $I - \mathbf{v}e_j^T$ , and  $\mathbf{v}$  has at most 3 nonzero

entries. We examine the condition number of these 3 transforms. Consult Sect. 4.2 to follow the details.

- The active part of  $Z_i$  is

$$\begin{bmatrix} u_i & 1 \\ 0 & 1 \end{bmatrix} \quad \text{and} \quad \text{cond}(Z_i) \simeq \max \left\{ |u_i|, |u_i|^{-1} \right\}.$$

- The active part of  $\mathcal{L}_i$  is

$$\begin{bmatrix} 1 & & \\ -x_l/\widehat{l}_{i-1} & 1 & \\ -y_l/\widehat{l}_{i-1} & 0 & 1 \end{bmatrix} \quad \text{and} \quad \text{cond}(\mathcal{L}_i) \simeq 1 + \left( \frac{x_l}{\widehat{l}_{i-1}} \right)^2 + \left( \frac{y_l}{\widehat{l}_{i-1}} \right)^2.$$

- The active part of  $Y_i$  is

$$\begin{bmatrix} 1 & & & \\ -x_r/\widehat{u}_i & 1 & & \\ -y_r/\widehat{u}_i & 0 & 1 & \\ -z_r/\widehat{u}_i & 0 & 0 & 1 \end{bmatrix} \quad \text{and} \quad \text{cond}(Y_i) \simeq 1 + \left( \frac{x_r}{\widehat{u}_i} \right)^2 + \left( \frac{y_r}{\widehat{u}_i} \right)^2 + \left( \frac{z_r}{\widehat{u}_i} \right)^2.$$

The variables  $x_l, y_l, x_r, y_r, z_r$  are formed from additions and multiplications of previous quantities. Note that  $u_i$  is part of the input and so is assumed to be of acceptable size. We see that it is tiny values of  $\widehat{l}_{i-1}$  and  $\widehat{u}_i$  that lead to an ill-conditioned transform at minor step  $i$ . In the simple dqds algorithm a small value of  $\widehat{u}_i$  (relative to  $u_{i+1}$ ) leads to a large value of  $\widehat{l}_i$  and  $d_{i+1}$ . In 3dqds the effect of 3 consecutive transforms is more complicated. The message is the same: reject any transform that has more than modest element growth. In practice,  $|\widehat{u}_i|$  and  $|\widehat{l}_{i-1}|$  are monitored and a transform is rejected if either quantity is too big (bigger than  $1/\sqrt{\epsilon}$ ). The computed eigenvalues are used as input for Rayleigh quotient correction in the original balanced matrix (see [19]).

In order to understand the intricate arguments below we have found it essential to absorb the contents of Sects. 4.1 and 4.2, in particular the division of the typical inner loop of 3dqds in three parts (a), (b) and (c). The three dqds similarities have morphed into a sequence of  $n - 1$  similarities of  $FG$  (implicit) each of which in its turn is composed of three transformations of  $F$  and  $G$  by matrices  $Z_i, \mathcal{L}_i, Y_i$  (with exact inverses) at minor step, or loop,  $i$  where we concentrate our attention.

**Minor step  $i$ .**

Recall that at the start the bulges  $x_l, y_l$  are in column  $i - 1$  of  $F$  while  $x_r, y_r, z_r$  are in column  $i$  of  $G$ . See (4.2). The values in these bulges change and they move one column right and one row down. In analysis, not practice, as the bulges both change value and position, new variables are created and denoted by augmentation of the subscripts (See Table 2). By the end of minor step  $i$  new values are given to all the bulge variables to be ready for the next step. The most active is  $x_r$ , the entry on the diagonal of  $G$ . The loop  $i$  updates  $x_r$  four times so we find

$$x_r, x_{r1}, x_{r2}, x_{r3}, x_{r4}$$

**Table 2** Input and output variables of  $3dqds$  algorithm

Part	Input	Output
(a)	$x_r, y_r, u_i$ $\widehat{l}_{i-1}, l_{i+1}, l_{i+2}$	$x_{r_1}$ $\widehat{u}_i$
(b)	$x_l, y_l$ $x_{r_1}, y_r, z_r$ $l_{i+1}, u_{i+2}, u_{i+3}, \widehat{u}_i$	$x_{l_1}, y_{l_1}$ $x_{r_2}, y_{r_1}, z_{r_1}$ $\widehat{l}_i$
(c)	$x_{l_1}, y_{l_1}$ $x_{r_2}, y_{r_1}, z_{r_1}$	$x_{l_2}, y_{l_2}$ $x_{r_3}, y_{r_2}, z_{r_2}$ $x_{r_4}, y_{r_3}, z_{r_3}$

and the last value  $x_{r_4}$  becomes  $x_r$  at the next loop  $i + 1$ . Its position changes from  $G_{i,i}$  to  $G_{i+1,i+1}$ . This change in position occurs at operation 14 of the 16 arithmetic operations in  $3dqds$ .

To follow the analysis below the reader should have reference to Sect. 4.2. At minor step  $i$  the inner loop transforms columns  $i - 1, i$  of  $F$  and  $i, i + 1$  of  $G$ .

To anticipate our result we are going to show that very small relative, but well chosen, perturbations in the input and output variables of each part, (a), (b), and (c), separately, of loop  $i$  yield exact, albeit implicit, transformations of  $F$  and  $G$ . Of course, the input and output variables are different for each part.

Note that the output variables for Part (b) may be perturbed (again) as input variables of Part (c). We will point out the two (of 16) operations at which our perturbations fail to give an exact implementation of the whole of loop  $i$ . That would be a result as strong as the one for real  $dqds$ .

As said above, if  $e_j$  denotes column  $j$  of  $I$  and  $v$  is a vector satisfying  $e_j^T v = 0$  then the exact inverse of  $I - ve_j^T$  is  $I + ve_j^T$  since  $(I - ve_j^T)(I + ve_j^T) = I - ve_j^T ve_j^T = I - (e_j^T v)ve_j^T = I$ . Hence the attraction of using elementary matrices for Parts (b) and (c). The matrix  $Z_i$ , whose active part is  $\begin{bmatrix} u_i & 1 \\ 0 & 1 \end{bmatrix}$ , is not elementary but the action of its inverse is implicit in creating 0 and 1 in  $F$  and it is only  $Z_i$  that acts on  $G$ . Thus  $1/u_i$  is never used explicitly and, again, there is no error in the implicit use of  $Z_i^{-1}$  on  $F$ . These observations help to explain the welcome accuracy of  $3dqds$  in practice.

In the analysis in each statement we use a subscript on  $\varepsilon$  as an indicator of the operation. For example,

$$fl(a + b + c * d) = fl(fl(a + b) + fl(c * d)) = [(a + b)(1 + \varepsilon_+) + c \cdot d(1 + \varepsilon_*)](1 + \varepsilon_{++}).$$

We find it simpler to not name the perturbed variables but to indicate them by judicious use of parentheses and square brackets. We use either a dot or juxtaposition to represent an exact multiplication.

**Loop  $i$ ,**  $1 < i < n - 2$ , in Section 4.3.

**Part (a)**

$F \leftarrow FZ_i^{-1}$  puts  $F_{i,i+1} = 0$  and  $F_{i,i} = 1$ . No errors.  
 $G \leftarrow Z_iG$  turns  $G_{i,i+1} = 1$ , updates  $x_r$  in  $G_{i,i}$ .

$$1 \quad x_{r1} = fl(fl(x_r * u_i) + y_r) = [x_r \cdot u_i(1 + \varepsilon_*) + y_r](1 + \varepsilon_+) \\ x_{r1} = [x_r(1 + \varepsilon_+)] [u_i(1 + \varepsilon_*)] + [y_r(1 + \varepsilon_+)]$$

**Part (b)**

$F \leftarrow \mathcal{L}_i^{-1}F$  puts 0 into  $F_{i+1,i-1}$  and  $F_{i+2,i-1}$ , and moves the bulge to column  $i$ .  
 $G \leftarrow G\mathcal{L}_i$  creates  $\widehat{u}_i$  in  $G_{i,i}$  (an  $LU$  output) and creates  $x_{r2}, y_{r1}, z_{r1}$  (bulge in  $G$ ) below it.

$$2 \quad x_{l1} = -fl(x_l/\widehat{l}_{i-1}) = -x_l/\widehat{l}_{i-1}(1 + \varepsilon_j) \quad x_{l1} = -[x_l(1 + \varepsilon_j)]/\widehat{l}_{i-1} \\ 3 \quad y_{l1} = -fl(y_l/\widehat{l}_{i-1}) = -y_l/\widehat{l}_{i-1}(1 + \varepsilon_j) \quad y_{l1} = -[y_l(1 + \varepsilon_j)]/\widehat{l}_{i-1}$$

Note that  $\widehat{l}_{i-1}$  is created in loop  $i - 1$  and  $(1 + \varepsilon_j)$  differs from  $(1 + \varepsilon_j)$  in Op. 3.

$$4 \quad \widehat{u}_i = fl(x_{r1} - x_{l1}) = (x_{r1} - x_{l1})/(1 + \varepsilon_-) \quad [(1 + \varepsilon_-)\widehat{u}_i] = x_{r1} - x_{l1} \\ 5 \quad x_{r2} = fl(y_r - x_{l1}) = (y_r - x_{l1})/(1 + \varepsilon_-) \quad [(1 + \varepsilon_-)x_{r2}] = y_r - x_{l1} \\ 6 \quad y_{r1} = fl(fl(z_r - y_{l1}) - fl(x_{l1} * l_{i+1})) = [(z_r - y_{l1})/(1 + \varepsilon_-) - x_{l1} \cdot l_{i+1}(1 + \varepsilon_*)]/(1 + \varepsilon_{--}) \\ [(1 + \varepsilon_-)(1 + \varepsilon_{--})y_{r1}] = z_r - y_{l1} - x_{l1}[l_{i+1}(1 + \varepsilon_*)(1 + \varepsilon_-)] \\ 7 \quad z_{r1} = -fl(y_{l1} * l_{i+2}) = y_{l1}l_{i+2}(1 + \varepsilon_*) \quad z_{r1} = -y_{l1}[l_{i+2}(1 + \varepsilon_*)]$$

**Part (c)**

$F \leftarrow FY_i^{-1}$  creates  $\widehat{l}_i$  in  $F_{i+1,i}$  (an  $LU$  output) and creates  $x_{l2}, y_{l2}$  (bulge in  $F$ ) below it.  
 $G \leftarrow Y_iG$  puts 0 into  $G_{i+1,i}, G_{i+2,i}$  and  $G_{i+3,i}$ , and moves the bulge to column  $i + 1$ .

$$8 \quad x_{r3} = -fl(x_{r2}/\widehat{u}_i) = -x_{r2}/\widehat{u}_i(1 + \varepsilon_j) \quad x_{r3} = -[x_{r2}(1 + \varepsilon_j)]/\widehat{u}_i \\ 9 \quad y_{r2} = -fl(y_{r1}/\widehat{u}_i) = -y_{r1}/\widehat{u}_i(1 + \varepsilon_j) \quad y_{r2} = -[y_{r1}(1 + \varepsilon_j)]/\widehat{u}_i \\ 10 \quad z_{r2} = -fl(z_{r1}/\widehat{u}_i) = -z_{r1}/\widehat{u}_i(1 + \varepsilon_j) \quad z_{r2} = -[z_{r1}(1 + \varepsilon_j)]/\widehat{u}_i \\ 11 \quad \widehat{l}_i = fl(fl(x_{l1} + y_{r2}) + fl(x_{r3} * u_{i+1})) = fl((x_{l1} + y_{r2})/(1 + \varepsilon_+) + x_{r3} \cdot u_{i+1}(1 + \varepsilon_*)) \\ = [(x_{l1} + y_{r2})/(1 + \varepsilon_+) + x_{r3} \cdot u_{i+1}(1 + \varepsilon_*)]/(1 + \varepsilon_{++}) \\ [(1 + \varepsilon_+)(1 + \varepsilon_{++})\widehat{l}_i] = x_{l1} + y_{r2} + x_{r3}[u_{i+1}(1 + \varepsilon_*)(1 + \varepsilon_+)]$$

**Part (c)**(cont.)

$$12 \quad x_{l_2} = fl(fl(y_{l_1} + z_{r_2}) + fl(y_{r_2} * u_{i+2})) = fl((y_{l_1} + z_{r_2})/(1 + \varepsilon_+) + y_{r_2} \cdot u_{i+2}(1 + \varepsilon_*))$$

$$= [(y_{l_1} + z_{r_2})/(1 + \varepsilon_+) + y_{r_2} \cdot u_{i+2}(1 + \varepsilon_*)]/(1 + \varepsilon_{++})$$

$$[(1 + \varepsilon_+)(1 + \varepsilon_{++})x_{l_2}] = y_{l_1} + z_{r_2} + y_{r_2}[u_{i+2}(1 + \varepsilon_*)(1 + \varepsilon_+)]$$

$$13 \quad y_{l_2} = fl(z_{r_2} * u_{i+3}) = z_{r_2} \cdot u_{i+3}(1 + \varepsilon_*) \qquad y_{l_2} = z_{r_2}[u_{i+3}(1 + \varepsilon_*)]$$

$$14 \quad x_{r_4} = fl(1 - x_{r_3}) = (1 - x_{r_3})/(1 + \varepsilon_-) \qquad (1 + \varepsilon_-)x_{r_4} = 1 - x_{r_3}$$

$$15 \quad y_{r_3} = fl(l_{i+1} - y_{r_2}) = (l_{i+1} - y_{r_2})/(1 + \varepsilon_-) \qquad (1 + \varepsilon_-)y_{r_3} = l_{i+1} - y_{r_2}$$

$$16 \quad z_{r_3} = fl(-z_{r_2}) = -z_{r_2} \qquad z_{r_3} = -z_{r_2}$$

No error in negation.

**end loop**

Note that  $\widehat{l}_{i-1}$  is created in loop  $i - 1$  and  $(1 + \varepsilon_j)$  differs from  $(1 + \varepsilon_l)$  in Op. 3.

Some comments. In Op. 4, for example, when cancellation occurs ( $x_{r_1}$  and  $x_{l_1}$  have same exponent) there is no error in subtraction but  $\widehat{u}_i$ 's uncertainty increases (Table 2).

We perturb  $x_{r_2}$  in Op. 5, as an output in Part (b), and also in Op. 8, as an input in Part (c). Similarly, we perturb  $y_{r_1}$  in Op. 6, in Part (b), and also in Op. 9 in Part (c). We use plain  $z_{r_1}$  in Op. 7, as an output in Part (b), and perturb  $z_{r_1}$  in Op. 10, as an input in Part (c).

We did not need to perturb  $x_{l_1}$  nor  $y_{l_1}$  in Ops. 2 and 3 in Part (b) and used  $x_{l_1}$  and  $y_{l_1}$  in Ops. 5 and 6, still Part (b), as well as Op. 11 and 12, in Part (c). So  $x_{l_1}$  and  $y_{l_1}$  did preserve their identities for the whole of loop  $i$ .

In Ops. 5 and 6 the perturbations we heaped on  $x_{r_2}$  and  $y_{r_1}$  were to avoid perturbing  $x_{l_1}$  and  $y_{l_1}$ . It seemed just too messy to try and carry the perturb  $x_{r_2}$  and  $y_{r_1}$  through the later operations in Part (c) that use them, such as Ops. 8 and 9.

Minor step 1 has a slightly different analysis but we omit the details which may be derived using a similar analysis.

In summary,

**Theorem 5.2** *If 3dqds is executed in standard floating point IEEE standard arithmetic with no invalid operations then suitable small perturbations (2 ulps maximum) of Parts (a), (b), and (c) produce an exact instance of each part in every minor step.*

## 6 Implementation details

### 6.1 Deflation ( $n \leftarrow n - 1$ )

Some of our criteria for deflating come from [23], others are new. Consider both matrices  $UL$  and  $LU$  and the trailing  $2 \times 2$  blocks,

$$\begin{bmatrix} l_{n-1} + u_{n-1} & 1 \\ l_{n-1}u_n & u_n \end{bmatrix}, \quad \begin{bmatrix} l_{n-2} + u_{n-1} & 1 \\ l_{n-1}u_{n-1} & l_{n-1} + u_n \end{bmatrix}.$$

Deflation ( $n \leftarrow n - 1$ ) removes  $l_{n-1}$  as well as  $u_n$ . Looking at entry  $(n - 1, n - 1)$  of  $UL$  shows that a necessary condition is that  $l_{n-1}$  be negligible compared to  $u_{n-1}$ ,

$$|l_{n-1}| < tol \cdot |u_{n-1}|, \tag{6.1}$$

for a certain tolerance  $tol$  close to roundoff unit  $\varepsilon$ .

The  $(n, n)$  entries of  $UL$  and  $LU$  suggest either  $u_n + acshift$  or  $l_{n-1} + u_n + acshift$  as eigenvalues where  $acshift$  is the accumulated shift. Recall that simple  $dqds$  is a non-restoring transform (see (3.5)). To make these consistent we require that

$$|l_{n-1}| < tol \cdot |u_n + acshift|. \tag{6.2}$$

Finally we must consider the change  $\delta\lambda$  in the eigenvalue  $\lambda$  caused by setting  $l_{n-1} = 0$ . We estimate  $\delta\lambda$  by starting from  $UL$  with  $l_{n-1} = 0$  and then allowing  $l_{n-1}$  to grow. To this end let  $J$  be  $UL$  with  $l_{n-1} = 0$  and  $(u_n, \mathbf{y}^T, \mathbf{x})$  be the eigentriple for  $J$ . Clearly  $\mathbf{y}^T = \mathbf{e}_n^T$ . Now we consider perturbation theory with parameter  $l_{n-1}$ . The perturbing matrix  $\delta J$ , as  $l_{n-1}$  grows, is

$$l_{n-1}(\mathbf{e}_{n-1} + \mathbf{e}_n u_n) \mathbf{e}_{n-1}^T.$$

By first order perturbation analysis

$$|\delta\lambda| = \frac{|\mathbf{y}^T \delta J \mathbf{x}|}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}$$

and  $\|\mathbf{y}\|_2 = 1$  in our case. So,

$$|\delta\lambda| = \frac{|l_{n-1} \mathbf{e}_n^T (\mathbf{e}_{n-1} + \mathbf{e}_n u_n) \mathbf{e}_{n-1}^T \mathbf{x}|}{\|\mathbf{x}\|_2} = \frac{|l_{n-1} u_n| |x_{n-1}|}{\|\mathbf{x}\|_2}$$

and we use the crude bound  $\frac{|x_{n-1}|}{\|\mathbf{x}\|_2} < 1$ . So, we let  $l_{n-1}$  grow until the change

$$|\delta\lambda| < |l_{n-1} u_n|$$

in eigenvalue  $\lambda = u_n$  is no longer acceptable. Our condition for deflation is then

$$|l_{n-1}u_n| < tol \cdot |acshift + u_n|. \tag{6.3}$$

A similar first order perturbation analysis for  $LU$  with  $l_{n-1} = 0$  will give our last condition for deflation. For the eigentriple  $(u_n, \mathbf{y}^T, \mathbf{x})$  we also have  $\mathbf{y}^T = \mathbf{e}_n^T$ . The perturbing matrix is now

$$l_{n-1}\mathbf{e}_n \left( \mathbf{e}_{n-1}^T u_{n-1} + \mathbf{e}_n^T \right)$$

and

$$|\delta\lambda| = \frac{|l_{n-1}\mathbf{e}_n^T \mathbf{e}_n (\mathbf{e}_{n-1}^T u_{n-1} + \mathbf{e}_n^T) \mathbf{x}|}{\|\mathbf{x}\|_2} = |l_{n-1}| \frac{|u_{n-1}x_{n-1} + x_n|}{\|\mathbf{x}\|_2} < |l_{n-1}| (|u_{n-1}| + 1).$$

Finally we require

$$|l_{n-1}| (|u_{n-1}| + 1) < tol \cdot |acshift + u_n|. \tag{6.4}$$

### 6.2 Splitting and deflation ( $n \leftarrow n - 2$ )

Recall that the implicit L theorem was invoked to justify the *3dqds* algorithm. This result fails if any  $l_k, k < n - 1$  vanishes. Consequently, checking for negligible values among the  $l_k$  is a necessity, not a luxury for increased efficiency. Consider  $J = UL$  in block form

$$\begin{bmatrix} J_1 & & & \\ & & & 1 \\ \hline & & \mu & \\ & & & J_2 \end{bmatrix}$$

where  $\mu = u_{k+1}l_k, k < n - 1$ . We can replace  $\mu$  by 0 when

$$spectrum(J_1) \cup spectrum(J_2) = spectrum(J), \text{ to working accuracy.}$$

However we are not going to estimate the eigenvalues of  $J_1$  and  $J_2$ . Instead we create a local criterion for splitting at  $(k + 1, k)$  as follows. Focus on the principal  $4 \times 4$  window of  $J$  given by

$$\left[ \begin{array}{cc|cc} u_{k-1} + l_{k-1} & 1 & & \\ u_k l_{k-1} & u_k + l_k & & 1 \\ \hline & u_{k+1} l_k & u_{k+1} + l_{k+1} & 1 \\ & & u_{k+2} l_{k+1} & u_{k+2} + l_{k+2} \end{array} \right].$$



Now  $J_1$  and  $J_2$  are both  $2 \times 2$  and our local criterion is

$$\det(J_1) \cdot \det(J_2) = \det(J), \quad \text{to working accuracy.} \tag{6.5}$$

Let us see what this yields. Perform block factorization on  $J$  and note that the Schur complement of  $J_1$  in  $J$  is

$$J'_2 = J_2 - \begin{bmatrix} 0 & \mu \\ 0 & 0 \end{bmatrix} J_1^{-1} \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}$$

with

$$J_1^{-1} = \frac{1}{\det_1} \begin{bmatrix} u_k + l_k & -1 \\ -u_k l_{k-1} & u_{k-1} + l_{k-1} \end{bmatrix}$$

where

$$\det_1 = \det(J_1) = u_{k-1}(u_k + l_k) + l_{k-1}l_k.$$

Thus

$$J'_2 = \begin{bmatrix} u_{k+1} + l_{k+1} & 1 \\ u_{k+2}l_{k+1} & u_{k+2} + l_{k+2} \end{bmatrix} - \begin{bmatrix} \mu(u_{k-1} + l_{k-1})/\det_1 & 0 \\ 0 & 0 \end{bmatrix}.$$

Since  $\det$  is linear by rows and the second rows are equal

$$\det(J_2) - \det(J'_2) = \mu(u_{k-1} + l_{k-1})(u_{k+2} + l_{k+2})/\det_1.$$

Our criterion reduces to splitting only when

$$\det(J'_2) = \det(J_2), \quad \text{to working accuracy.}$$

Thus we require

$$|l_k u_{k+1}(u_{k+2} + l_{k+2})(u_{k-1} + l_{k-1})/\det_1| < tol \cdot |\det(J_2)|.$$

Since

$$\det_2 = \det(J_2) = u_{k+1}(u_{k+2} + l_{k+2}) + l_{k+1}l_{k+2},$$

the criterion for splitting  $J$  at  $(k + 1, k)$  is then

$$|l_k u_{k+1}(u_{k+2} + l_{k+2})(u_{k-1} + l_{k-1})| < tol \cdot |\det_1 \det_2|. \tag{6.6}$$

Finally, to remove  $l_k$  we also need  $l_k$  to be negligible compared to  $u_k$ ,

$$|l_k| < tol \cdot |u_k|. \tag{6.7}$$

**Deflation ( $n \leftarrow n - 2$ )**

We use the same criterion for deflation ( $n \leftarrow n - 2$ ), but because  $l_{k+2} = l_n = 0$  there is a common factor  $det_2$  on each side of (6.6). Deflate the trailing  $2 \times 2$  submatrix when

$$|l_{n-2}| < tol \cdot |u_{n-2}| \tag{6.8}$$

and

$$|l_{n-2}(u_{n-3} + l_{n-3})| < tol \cdot |u_{n-3}(u_{n-2} + l_{n-2}) + l_{n-3}l_{n-2}|. \tag{6.9}$$

We omit the role of *acshift* here because it makes the situation more complicated. We have to recall that *3dqds* uses restoring shifts and *acshift* is always real. So, for complex shifts,  $det_2$  is not going to zero. In fact

$$|det_2| \geq |\Im(\lambda)|^2$$

where  $\lambda$  is an eigenvalue of  $J_2$ .

When  $n = 3$  these criteria simplify a lot. Both reduce to

$$|l_1| < tol \cdot |u_1|.$$

**6.3 Shift strategy**

As with LR, the *dqds* algorithm with no shift gradually forces large entries to the top and brings small entries towards the bottom. We want to use a shift as soon as the trailing  $2 \times 2$  principal submatrix appears to be converging. We use the size of the last two entries of  $L$  to make the judgement. The code executes a *dqds* transform with a zero shift if

$$l_{n-1} > 10^{-2} \quad \text{and} \quad l_{n-2} > 10^{-2}.$$

Otherwise, a *3dqds* transform is executed with

$$\text{sum} = l_{n-1} + (u_{n-1} + u_n), \quad \text{prod} = u_{n-1}u_n,$$

the trace and the determinant of the trailing  $2 \times 2$  submatrix of  $UL$ . This will let us converge to either two real eigenvalues in the bottom  $2 \times 2$  or a single  $2 \times 2$  block with a complex conjugate pair of eigenvalues.

An unexpected reward for having both transforms available is to cope with a rejected transform. Our strategy is simply to use the other transform with the current shift. More precisely, given `sum` and `prod`, if *3dqds*(`sum`, `prod`) is rejected we try *dqds*( $u_n$ ); if *dqds*(0) is rejected, we try *3dqds*( $\delta$ ,  $\delta$ ) with  $\delta = \sqrt{\epsilon}$ . We have to admit the possibility of a succession of rejections and in this case we don't want to move away from the previous shift too much, just a small amount so that the transformation does not

breakdown. See Algorithm 4 in “AppendixB” for details. The number of rejections is recorded and added to the total number of iterations.

## 7 Factored forms

### 7.1 Eigenvectors from twisted factorizations of the balanced form $\Delta T$

A salient property of an unreduced real tridiagonal matrix  $C = \text{tridiag}(\mathbf{b}, \mathbf{a}, \mathbf{c})$  (no off-diagonal entry vanishes) is that it can be balanced by a diagonal similarity easily and, once the matrix is balanced, it can be made real symmetric by changing the signs of certain rows. However, changing the signs is not a similarity transformation and would not preserve the eigenvalues. It is accomplished by premultiplying by a so-called *signature matrix*  $\Delta = \text{diag}(\delta_1, \dots, \delta_n)$ ,  $\delta_i = \pm 1$ . So we can write

$$\Delta T = SCS^{-1} \tag{7.1}$$

where  $T$  is real symmetric and  $S$  is diagonal positive definite,  $S = \text{diag}(s_1, \dots, s_n)$  with  $s_1 = 1$ ,  $s_i = (|c_1c_2 \cdots c_{i-1}|/|b_1b_2 \cdots b_{i-1}|)^{1/2}$ ,  $i = 2, \dots, n$ . See [11, Section 2.2.3]

Let  $\lambda$  be a simple eigenvalue of  $\Delta T$  with eigenvector equations

$$\Delta T \mathbf{x} = \mathbf{x}\lambda, \quad \mathbf{y}^* \Delta T = \lambda \mathbf{y}^*. \tag{7.2}$$

Recall that  $\mathbf{x}$ ,  $\mathbf{y}$  and  $\lambda$  may be complex and  $\mathbf{y}^* \mathbf{x} \neq 0$ , since  $\lambda$  is simple. An attraction of the  $\Delta T$  representation is that the row eigenvector  $\mathbf{y}^*$  is determined by the right (or column) eigenvector  $\mathbf{x}$ . Transpose  $\Delta T \mathbf{x} = \mathbf{x}\lambda$  and insert  $I = \Delta^2$  to find

$$(\mathbf{x}^\top \Delta) \Delta T = \lambda (\mathbf{x}^\top \Delta). \tag{7.3}$$

Compare with  $\mathbf{y}^* \Delta T = \lambda \mathbf{y}^*$  to see that  $\mathbf{y}^* = \mathbf{x}^\top \Delta$ . See [13,30].

The so-called twisted factorizations generalize the lower and upper bidiagonal factorizations. These factorizations gained new popularity as they were used for the purpose of computing eigenvectors of symmetric tridiagonal matrices [10,22]. The idea is to begin both a top-to-bottom and a bottom-to-top factorization until they meet at, say, the  $k$ -th row, where they will have to be glued together. The index  $k$  is called the *twist index* or the *twist position*.

Observe that the eigenvector equations  $\Delta T \mathbf{x} = \mathbf{x}\lambda$  and  $(\mathbf{x}^\top \Delta) \Delta T = \lambda (\mathbf{x}^\top \Delta)$  are equivalent to

$$(T - \lambda \Delta) \mathbf{x} = \mathbf{0} \quad \text{and} \quad \mathbf{x}^\top (T - \lambda \Delta) = \mathbf{0}$$

where  $T - \lambda \Delta$  is symmetric. Now suppose that we have  $\tilde{\lambda}$  as an approximation to an eigenvalue  $\lambda$  of  $\Delta T$  and that  $T - \tilde{\lambda} \Delta$  admits both lower and upper bidiagonal factorizations, starting the Gaussian elimination at the first row and at the last row, respectively,



The above is just inverse iteration to obtain  $z$  (and  $z^T \Delta$ ) as an approximation to  $\lambda$ 's eigenvector  $x$  (and  $x^T \Delta$ ) with residual norm

$$\frac{\|(T - \tilde{\lambda}\Delta)z\|}{\|z\|} = \frac{|\gamma_k|}{\|z\|}.$$

Therefore a natural choice for the twist index would be  $k$  such that

$$|\gamma_k| = \min_{i=1, \dots, n} |\gamma_i|.$$

This strategy to choose an initial guess for the eigenvector provides, as a by-product, the diagonal entries of  $(T - \tilde{\lambda}\Delta)^{-1}$  since  $[(T - \tilde{\lambda}\Delta)^{-1}]_{k,k} = \gamma_k^{-1}$ . See [35, Lemma 2.3].

If  $\Delta$  is definite, one important result presented in [7,8] is that we can always find a twist index  $k$  such that

$$|\gamma_k| \leq \sqrt{n}|\tilde{\lambda} - \lambda|.$$

Since (7.9) uses only multiplications, the computed vector will be very good provided that  $\tilde{\lambda}$  is accurate enough. In the general case, to judge the accuracy of the eigenvectors, we compute column (and row) residual norm relative to the eigenvalue,

$$\frac{\|\Delta Tz - \tilde{\lambda}z\|}{|\tilde{\lambda}|\|z\|} = \frac{\|(T - \tilde{\lambda}\Delta)z\|}{|\tilde{\lambda}|\|z\|} = \frac{\|z^T(T - \tilde{\lambda}\Delta)\|}{|\tilde{\lambda}|\|z^T \Delta\|}. \tag{7.10}$$

This is a stricter measure than the usual  $\frac{\|\Delta Tz - \tilde{\lambda}z\|}{\|z\|\|\Delta T\|}$ .

In [24] we show that unique tridiagonal “backward error” matrices can be designated for an approximate pair of complex eigenvectors (column and row) or two approximate real eigenvectors.

### 7.2 Relative eigenvalue condition numbers

The condition number of every eigenvalue of a real symmetric matrix is 1, but only in the absolute sense. The relative condition number can vary. In the unsymmetric case even the absolute condition numbers can rise to  $\infty$  and little is known about relative errors. In [13] several *relative condition numbers* with respect to eigenvalues were derived. Some of them use bidiagonal factorizations instead of the matrix entries and so they shed light on when eigenvalues are less sensitive to perturbations of factored forms than to perturbations of the matrix entries. These condition numbers are measures of *relative sensitivity*, i.e., measures of the relative variation of an eigenvalue with respect to the largest relative perturbation of each of the nonzero entries of the representation of the matrix. So the perturbations we consider are of the form  $|\delta p_i| \leq \eta|p_i|, 0 < \eta \ll 1$ . In this section we present the relative condition number for the entries of the matrix  $C$  and for the  $LU$  factorization of the  $J$ -form.

Assume that  $\lambda \neq 0$  is a simple eigenvalue of real tridiagonal matrix  $C = \text{tridiag}(\mathbf{b}, \mathbf{a}, \mathbf{c})$ . Let  $\Delta T = SC S^{-1}$  be the balanced form (7.1) of  $C$  and  $(\lambda, \mathbf{x}, \mathbf{x}^\top \Delta)$  be an eigen triple of  $\Delta T$ ,

$$\Delta T \mathbf{x} = \mathbf{x} \lambda, \quad (\mathbf{x}^\top \Delta) \Delta T = \lambda (\mathbf{x}^\top \Delta), \quad \lambda \neq 0, \tag{7.11}$$

and recall that  $\Delta T$  and  $C$  eigenvectors are simply related by

$$C(S^{-1} \mathbf{x}) = (S^{-1} \mathbf{x}) \lambda, \quad (\mathbf{x}^\top \Delta S) C = \lambda (\mathbf{x}^\top \Delta S). \tag{7.12}$$

The relative condition number with respect to  $\lambda$  for the entries of  $C$  is

$$\text{relcond}(\lambda; C) = \frac{|\mathbf{x}^\top \Delta S| |C| |S^{-1} \mathbf{x}|}{|\lambda| |(\mathbf{x}^\top \Delta S)(S^{-1} \mathbf{x})|},$$

where  $|M|_{ij} = |M_{ij}|$ , for any matrix  $M$ . Since  $S$  is diagonal it follows that

$$\text{relcond}(\lambda; C) = \frac{|\mathbf{x}^\top \Delta| |S| |C| |S^{-1}| |\mathbf{x}|}{|\lambda| |\mathbf{x}^\top \Delta \mathbf{x}|} = \frac{|\mathbf{x}^\top \Delta| |\Delta T| |\mathbf{x}|}{|\lambda| |\mathbf{x}^\top \Delta \mathbf{x}|} = \text{relcond}(\lambda; \Delta T). \tag{7.13}$$

We have just shown that, in general, for any scaling matrix  $X$  invertible and diagonal, the expression for  $\text{relcond}(\lambda; C)$  yields  $\text{relcond}(\lambda; X C X^{-1}) = \text{relcond}(\lambda; C)$ . See [13, Lemma 6.2].

When  $C$  is unreduced it is also diagonally similar to a  $J$ -form,

$$J = D C D^{-1} = \text{tridiag}(\mathbf{b}, \mathbf{a}, \mathbf{1})$$

where  $D = \text{diag}(1, c_1, c_1 c_2, \dots, c_1 c_2 \dots c_{n-1})$  and  $\mathbf{b} = \text{diag}(b_1 c_1, b_2 c_2, \dots, b_{n-1} c_{n-1})$ . Now assume that  $J$  permits bidiagonal factorization  $J = LU$  and write

$$\Delta T = F J F^{-1}, \quad F = S D^{-1}, \tag{7.14}$$

to obtain

$$LU (F^{-1} \mathbf{x}) = (F^{-1} \mathbf{x}) \lambda, \quad (\mathbf{x}^\top \Delta F) \mathcal{L} U = \lambda (\mathbf{x}^\top \Delta F), \quad \lambda \neq 0.$$

Recall that  $L = I + \mathring{L}$  and  $U = \text{diag}(u_1, \dots, u_n) (I + \mathring{U})$  with

$$\mathring{L} = \begin{bmatrix} 0 & & & & & \\ l_1 & 0 & & & & \\ & \ddots & \ddots & & & \\ & & l_{n-2} & 0 & & \\ & & & l_{n-1} & 0 & \end{bmatrix} \quad \text{and} \quad \mathring{U} = \begin{bmatrix} 0 & u_1^{-1} & & & & \\ & 0 & u_2^{-1} & & & \\ & & \ddots & \ddots & & \\ & & & 0 & u_{n-1}^{-1} & \\ & & & & & 0 \end{bmatrix}.$$

For the cost of solving two bidiagonal linear systems,

$$v^T (I + \mathcal{U}) = (x^T \Delta F) \text{ for } v^T \quad \text{and} \quad \mathcal{L}w = \mathring{L} (F^{-1}x) \text{ for } w,$$

we obtain the following expression of the relative condition number for the entries of  $L$  and  $U$ ,

$$relcond(\lambda; L, U) := \frac{|v^T|F^{-1}x| + |x^T \Delta F||w|}{|x^T \Delta x|}. \tag{7.15}$$

See [13, Section 6.3]. Next we deal with a case of a simple zero eigenvalue. Although the right hand side of (7.15) is a nonzero finite number for a simple eigenvalue  $\lambda = 0$ , observe that the perturbations we consider for  $U$ , that is,  $|\delta u_i| \leq \eta|u_i|$ , produce  $u_n + \delta u_n = 0$  whenever  $u_n = 0$ . This means that singularity is preserved or, equivalently, that the zero eigenvalue is preserved. Therefore, it seems appropriate to set  $relcond(0; L, U) = 0$ .

We use eigenvalue approximation  $\tilde{\lambda}$  from *Rayleigh Quotient Iteration* (RQI) and eigenvector approximations  $z$  and  $z^T \Delta$  obtained from (7.8) to compute the relative condition numbers (7.13) and (7.15). The residual norms for  $\Delta T$  are given by (7.10) but for  $J = LU$ , with eigenvector approximations  $F^{-1}z$  and  $z^T \Delta F$ , by

$$\frac{\|F^{-1}(T - \tilde{\lambda}\Delta)z\|}{|\tilde{\lambda}|\|F^{-1}z\|} \quad \text{and} \quad \frac{\|z^T(T - \tilde{\lambda}\Delta)F\|}{|\tilde{\lambda}|\|z^T \Delta F\|}. \tag{7.16}$$

### 7.3 Generalized Rayleigh quotient iteration

In addition to computing both column and row eigenvector approximations from twisted factorizations of  $\Delta T$ , the algorithm described in Sect. 7.1 can also be used to improve the accuracy of the eigenvalue approximation  $\tilde{\lambda}$  by performing a Rayleigh Quotient Iteration. So, our code will return an eigenpair approximation  $(S^{-1}z, z^T \Delta S)$  for  $C$  together with an improved eigenvalue estimate, the generalized Rayleigh quotient,

$$\frac{(z^T \Delta S)C(S^{-1}z)}{(z^T \Delta S)(S^{-1}z)} = \frac{(z^T \Delta)\Delta Tz}{z^T \Delta z} = \tilde{\lambda} + \frac{(z^T \Delta)(\Delta T - \tilde{\lambda}I)z}{z^T \Delta z}. \tag{7.17}$$

Given the twisted factorization in (7.4) and (7.7), the *Rayleigh quotient correction* is given by

$$\rho := \frac{(z^T \Delta)(\Delta N_k G_k N_k^T)z}{z^T \Delta z} = \frac{z^T \gamma_k e_k}{z^T \Delta z} = \frac{\gamma_k}{z^T \Delta z},$$

since  $z_k = 1$ , where  $\gamma_k$  is given in (7.6) and (7.5).

Recall that for  $x \in \mathbb{C}$ ,  $\Re(x)$  denotes the real part of  $x$ . The following lemma extends Lemma 12 in [8, pg. 886] to the unsymmetric case.

**Lemma 7.1** *Let  $T - \tilde{\lambda}\Delta = N_k G_k N_k^T$  and  $N_k G_k N_k^T z = \gamma_k e_k$ ,  $z_k = 1$ . Then the Rayleigh quotient  $\rho$  with respect to  $\Delta T - \tilde{\lambda}I$  is*

$$\rho = \frac{\gamma_k}{z^T \Delta z}$$

and

$$\frac{\|(\Delta T - (\tilde{\lambda} + \rho)I)z\|}{\|z\|} = \frac{|\gamma_k|}{\|z\|} \left( \frac{|z^T \Delta z|^2 - \omega_k}{|z^T \Delta z|^2} \right)^{1/2} \tag{7.18}$$

where  $\omega_k = 2\delta_k \Re(z^T \Delta z) - \|z\|^2$ .

**Proof** Let  $\Delta = \text{diag}(\delta_1, \dots, \delta_n)$ ,  $\delta_i = \pm 1$ . Since

$$(\Delta T - (\tilde{\lambda} + \rho)I)z = (\Delta T - \tilde{\lambda}I)z - \rho z = \Delta N_k G_k N_k^T z - \rho z = \delta_k \gamma_k e_k - \rho z,$$

then

$$\begin{aligned} \|(\Delta T - (\tilde{\lambda} + \rho)I)z\|^2 &= \|\delta_k \gamma_k e_k - \rho z\|^2 = (\delta_k \overline{\gamma_k} e_k - \overline{\rho z})^T \cdot (\delta_k \gamma_k e_k - \rho z) \\ &= |\gamma_k|^2 + |\rho|^2 \|z\|^2 - 2\delta_k \Re(\overline{\gamma_k} \rho) \\ &= |\gamma_k|^2 + \frac{|\gamma_k|^2}{|z^T \Delta z|^2} \|z\|^2 - \frac{2\delta_k |\gamma_k|^2}{|z^T \Delta z|^2} \Re(z^T \Delta z) \\ &= \frac{|\gamma_k|^2}{|z^T \Delta z|^2} \left( |z^T \Delta z|^2 + \|z\|^2 - 2\delta_k \Re(z^T \Delta z) \right). \end{aligned}$$

Thus,

$$\frac{\|(\Delta T - (\tilde{\lambda} + \rho)I)z\|}{\|z\|} = \frac{|\gamma_k|}{\|z\|} \left( \frac{|z^T \Delta z|^2 - (2\delta_k \Re(z^T \Delta z) - \|z\|^2)}{|z^T \Delta z|^2} \right)^{1/2}.$$

Observe that, since  $z_k = 1$ ,

$$|z^T \Delta z|^2 + \|z\|^2 - 2\delta_k \Re(z^T \Delta z) = \|z\|^2 + |z^T \Delta z - \delta_k|^2 - 1 > 0.$$

If the easily checked condition

$$\omega_k := 2\delta_k \Re(z^T \Delta z) - \|z\|^2 > 0 \tag{7.19}$$

is satisfied, we obtain a decrease in the residual norm by using the Rayleigh quotient; the pair  $(\tilde{\lambda} + \rho, z)$  is a better approximate eigenpair than  $(\tilde{\lambda}, z)$ .  $\square$

When  $\Delta = I$  (symmetrizable case) the condition (7.19) reduces to  $2\|z\|^2 - \|z\|^2 = \|z\|^2 > 0$  and the Rayleigh quotient correction always gives an improvement. In this case (7.18) simplifies to



$$\frac{\|(T - (\tilde{\lambda} + \rho)I)z\|}{\|z\|} = \frac{|\gamma_k|}{\|z\|} \left( \frac{\|z\|^4 - \|z\|^2}{\|z\|^4} \right)^{1/2} = \frac{|\gamma_k|}{\|z\|} \left( \frac{\|z\|^2 - 1}{\|z\|^2} \right)^{1/2}. \tag{7.20}$$

Given an approximation  $\tilde{\lambda}$  to an eigenvalue  $\lambda$  of  $\Delta T$  we compute the twisted factorization of  $T - \tilde{\lambda}\Delta$  and use inverse iteration (7.8) to obtain  $\lambda$ 's column and row eigenvector approximations,  $z$  and  $z^T \Delta$ . The Rayleigh quotient correction (7.17) gives a new approximation  $\tilde{\lambda} + \rho$  for  $\lambda$ . We may repeat this process until there is no improvement in the residual (7.18). Although RQI can misbehave for non-normal matrices, we can use (7.19) to judge improvement (see [15,34]). Our code *3dqds* examines  $\omega_k$  and whenever it is greater than zero we apply RQI, otherwise not.

### 8 Numerical examples

The need for a tridiagonal eigensolver is admirably covered in Bini, Gemignani and Tisseur [1], many parts of which have been of great help to us. We refer to the Ehrlich-Aberth algorithm (see Sect. 2.5) as *BGT* and to our code simply as *3dqds*, although we combine *3dqds* with real *dqds* as described in Sect. 6.3.

Here the *exact* eigenvalue  $\lambda_i$  is computed in quadruple precision, using MATLAB Symbolic Math Toolbox with variable-precision arithmetic, and  $\tilde{\lambda}_i$  denotes the computed eigenvalue in double precision (unit roundoff  $2.2 \cdot 10^{-16}$ ). We compare our *3dqds* algorithm with its explicit version, referred as *ex3dqds* (the three steps of *dqds* are computed explicitly in complex arithmetic, see Fig. 1), with a MATLAB implementation of *BGT* and with the *QR* algorithm on an upper Hessenberg matrix (MATLAB function *eig*).

All the experiments were performed in MATLAB (R2020b) on a LAPTOP-KVSVAAU8 with an Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz and 8 GB RAM, under Windows 10 Home. No parallel operations were used. We acknowledge that the MATLAB tests do not reflect Fortran performance, but even in MATLAB environment the ratio of elapsed times is an important feature.

#### 8.1 Bessel matrix

Bessel matrices, associated with generalized Bessel polynomials [25], are nonsymmetric tridiagonal matrices defined by  $B_n^{(a,b)} = \text{tridiag}(\beta, \alpha, \gamma)$  with

$$\begin{aligned} \alpha_1 &= -\frac{b}{a}, & \gamma_1 &= -\alpha_1, & \beta_1 &= \frac{\alpha_1}{a+1}, \\ \alpha_j &:= -b \frac{a-2}{(2j+a-2)(2j+a-4)}, & j &= 2, \dots, n, \\ \gamma_j &:= b \frac{j+a-2}{(2j+a-2)(2j+a-3)}, \\ \beta_j &:= -b \frac{j}{(2j+a-1)(2j+a-2)}, & j &= 2, \dots, n-1. \end{aligned}$$

**Table 3** Relative errors for computed eigenvalues from  $B_n^{(-8.5,2)}$ ,  $B_n^{(12,2)}$ , and  $(B_n^{(-4.5,2)})$  with one RQI

$(a, b); n$	eig		BGT		3dqds	
	$rel_{min}$	$rel_{max}$	$rel_{min}$	$rel_{max}$	$rel_{min}$	$rel_{max}$
$(-8.5, 2); 18$	$1.6 \cdot 10^{-6}$	$2.7 \cdot 10^{-1}$	$7.1 \cdot 10^{-7}$	$2.3 \cdot 10^{-1}$	$5.9 \cdot 10^{-7}$	$2.3 \cdot 10^{-1}$
$(-8.5, 2); 25$	$2.5 \cdot 10^{-1}$	$1.9 \cdot 10^0$	$1.3 \cdot 10^{-1}$	$1.9 \cdot 10^0$	$2.4 \cdot 10^{-1}$	$1.8 \cdot 10^0$
$(-4.5, 2); 20$	$4.1 \cdot 10^{-7}$	$3.2 \cdot 10^{-1}$	$1.0 \cdot 10^{-7}$	$2.1 \cdot 10^{-1}$	$1.5 \cdot 10^{-8}$	$1.2 \cdot 10^{-1}$
$(-4.5, 2); 25$	$2.0 \cdot 10^{-1}$	$1.3 \cdot 10^0$	$5.7 \cdot 10^{-2}$	$1.2 \cdot 10^0$	$2.0 \cdot 10^{-1}$	$7.3 \cdot 10^{-1}$
$(12, 2); 40$	$3.3 \cdot 10^{-15}$	$1.3 \cdot 10^{-1}$	$1.1 \cdot 10^{-15}$	$1.9 \cdot 10^{-1}$	$2.1 \cdot 10^{-15}$	$1.7 \cdot 10^{-1}$
$(12, 2); 50$	$7.0 \cdot 10^{-15}$	$3.5 \cdot 10^{-1}$	$8.5 \cdot 10^{-16}$	$4.3 \cdot 10^{-1}$	$6.5 \cdot 10^{-15}$	$3.4 \cdot 10^{-1}$

Parameter  $b$  is a scaling factor and most authors take  $b = 2$  and so do we. The case  $a \in \mathbb{R}$  is the most investigated in literature. The eigenvalues of  $B_n^{(a,b)}$ , well separated complex eigenvalues, suffer from ill-conditioning that increases with  $n$  - close to a defective matrix. In Pasquini [25] it is mentioned that the ill-conditioning seems to reach its maximum when  $a$  ranges from  $-8.5$  to  $-4.5$ . We pay a lot of attention to these matrices because they are an interesting family for our purposes. Each picture teaches us a lot about the behavior of eigenvalues.

Our examples take  $B_n^{(-8.5,2)}$ ,  $B_n^{(-4.5,2)}$  and  $B_n^{(12,2)}$  for  $n = 40, 50$ . We show pictures for MATLAB (double precision), *BGT* and *3dqds* to illustrate the extreme sensitivity of some of the eigenvalues. The results of *ex3dqds* are visually identical to *3dqds*, so we don't show them. In exact arithmetic the spectrum lies on an arc in the interior of the moon-shaped region. Our pictures show this region and the eigenvalues computed in quadruple precision (labeled as *exact*).

Table 3 shows the minimum and maximum relative errors, respectively,  $rel_{min} = \min_i |\lambda_i - \tilde{\lambda}_i|/|\lambda_i|$  and  $rel_{max} = \max_i |\lambda_i - \tilde{\lambda}_i|/|\lambda_i|$ . The relative condition numbers  $relcond(\lambda; C)$  and  $relcond(\lambda; L, U)$  (see (7.13) and (7.15)) and residual norms (see (7.10) and (7.16)) are shown in Table 4. We show both condition numbers because MATLAB and *BGT* only use matrix entries and *3dqds* uses  $L, U$  factors.

The results for  $B_{18}^{(-8.5,2)}$ ,  $n = 18, 25$ , are shown in Fig. 3, without RQI (Rayleigh Quotient Iteration) and with one RQI. Observe on the real line that our approximations with one RQI (c) lie on top of the *BGT* approximations. We include the pictures (b) and (d) to show well the extreme sensitivity of the eigenvalues. Note how the eigenvalues move out of the moon-shaped inclusion region.

In Fig. 4 we show the results for  $B_n^{(-4.5,2)}$ ,  $n = 20, 25$ , and  $B_n^{(12,2)}$ ,  $n = 40, 50$ , without RQI. The reader is invited to see the large effect of changing  $n$  from 20 to 25, in (a) and (b), and from 40 to 50, in (c) and (d). Notice that our results are slightly but consistently better than those of the other two methods.

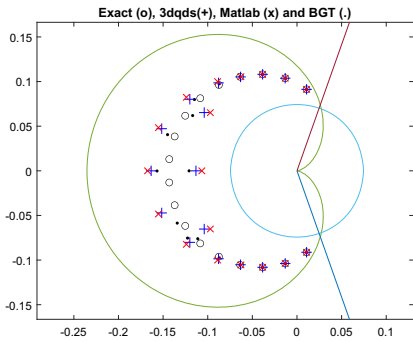
### 8.2 Clement matrix

The so-called *Clement* matrices (see [3])

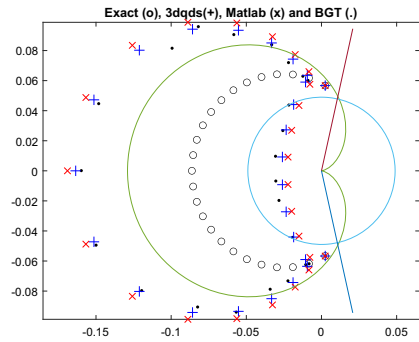
$$C = tridiag(b, \mathbf{0}, c)$$

**Table 4** Relative condition numbers and residual norms for  $B_n^{(-8.5,2)}$ ,  $B_n^{(-4.5,2)}$ , and  $B_n^{(12,2)}$

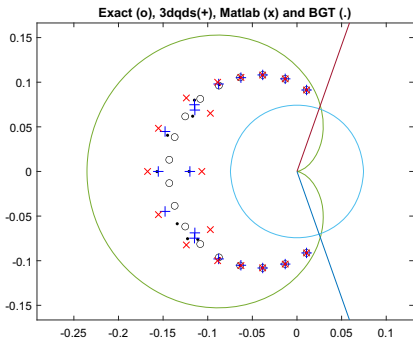
$(a, b); n$	$relcond(\lambda; L, U)$		$relcond(\lambda; C)$		$max\ residuals$	
	$min$	$max$	$min$	$max$	$J = LU$	$\Delta T$
$(-8.5, 2); 18$	$4.3 \cdot 10^8$	$6.6 \cdot 10^{13}$	$1.5 \cdot 10^{10}$	$2.2 \cdot 10^{15}$	$2.1 \cdot 10^{-14}$	$1.8 \cdot 10^{-14}$
$(-8.5, 2); 25$	$8.0 \cdot 10^{10}$	$4.3 \cdot 10^{13}$	$4.3 \cdot 10^{12}$	$2.6 \cdot 10^{15}$	$3.9 \cdot 10^{-14}$	$1.2 \cdot 10^{-14}$
$(-4.5, 2); 20$	$8.1 \cdot 10^7$	$3.5 \cdot 10^{14}$	$3.5 \cdot 10^9$	$1.7 \cdot 10^{16}$	$3.5 \cdot 10^{-14}$	$3.1 \cdot 10^{-14}$
$(-4.5, 2); 25$	$6.3 \cdot 10^8$	$2.6 \cdot 10^{14}$	$4.6 \cdot 10^{10}$	$1.9 \cdot 10^{16}$	$4.5 \cdot 10^{-14}$	$5.0 \cdot 10^{-14}$
$(12, 2); 40$	$9.3 \cdot 10^1$	$7.0 \cdot 10^{15}$	$1.4 \cdot 10^2$	$1.8 \cdot 10^{16}$	$1.3 \cdot 10^{-15}$	$4.6 \cdot 10^{-14}$
$(12, 2); 50$	$1.8 \cdot 10^2$	$1.1 \cdot 10^{16}$	$2.7 \cdot 10^2$	$3.8 \cdot 10^{16}$	$1.9 \cdot 10^{-15}$	$7.7 \cdot 10^{-15}$



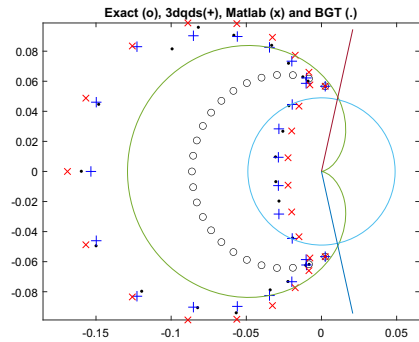
(a)  $n = 18$ ; no RQI



(b)  $n = 25$ ; no RQI



(c)  $n = 18$ ; one RQI

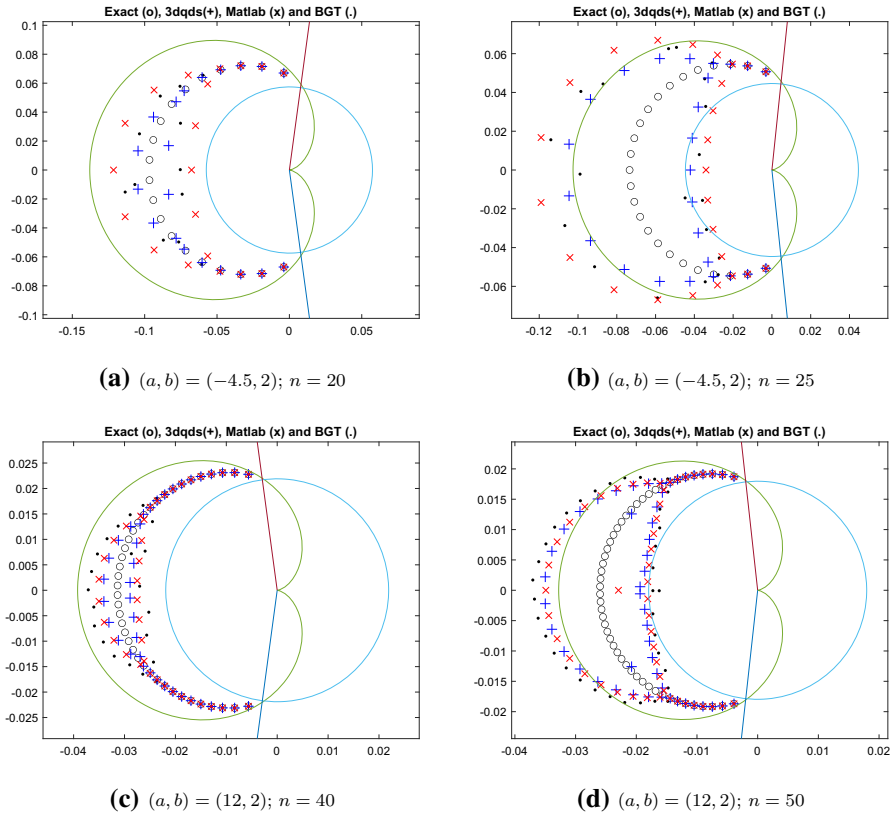


(d)  $n = 25$ ; one RQI

**Fig. 3** Eigenvalues of  $B_n^{(-8.5,2)}$ ,  $n = 18, 25$

with  $b_j = j$  and  $c_j = b_{n-j}$ ,  $j = 1, \dots, n - 1$ , have real eigenvalues,

$$\begin{aligned} &\pm n - 1, n - 3, \dots, 1, && \text{for } n \text{ even,} \\ &\pm n - 1, n - 3, \dots, 0, && \text{for } n \text{ odd.} \end{aligned}$$



**Fig. 4** Eigenvalues of  $B_n^{(-4.5,2)}$  and  $B_n^{(12,2)}$  (without RQI)

These matrices posed no serious difficulties. The initial zero diagonal obliges the *dqds* based methods to take care when finding an initial *LU* factorization.

The *3dqds* and *ex3dqds* codes use only real shifts as they should and the accuracy (approximately the same) is less than *BGT* but satisfactory. One RQI reduces errors to  $\mathcal{O}(\varepsilon)$ .

Our numerical tests have  $n = 50, 100, 200, 400, 800$ . The relative condition number  $relcond(\lambda; C)$  ranges from  $10^0$  to  $4 \cdot 10^2$  and it is smaller at the ends of the spectrum. The maximum residual norm for  $C$  is  $\mathcal{O}(10^{-11})$ . The minimum and maximum relative errors,  $rel_{min}$  and  $rel_{max}$ , are shown in Table 5. Note the poor performance of MATLAB’s *eig* (so much for backward stability).

The CPU elapsed times are presented in Table 6. We put (+) whenever a RQI is used. Since we compare MATLAB versions of all the codes we acknowledge that the elapsed times are accurate to only about 0.02 seconds. However, this is good enough to show the striking time ratios between *BGT* and the *dqds* codes.

We draw the readers attention that for  $n = 400$  our algorithm is about 200 times faster than *BGT* but when  $n$  rises to 1000 it is over 600 times faster. This is finding the

**Table 5** Relative errors for Clement matrices (without RQI)

$n$	eig		$BGT$		$3dqds$	
	$rel_{min}$	$rel_{max}$	$rel_{min}$	$rel_{max}$	$rel_{min}$	$rel_{max}$
50	$4.4 \cdot 10^{-16}$	$7.4 \cdot 10^{-11}$	0	0	$1.7 \cdot 10^{-16}$	$4.7 \cdot 10^{-15}$
100	$1.6 \cdot 10^{-15}$	$1.6 \cdot 10^{-3}$	0	$1.8 \cdot 10^{-16}$	0	$2.1 \cdot 10^{-14}$
200	$4.3 \cdot 10^{-16}$	$1.6 \cdot 10^1$	0	$1.1 \cdot 10^{-15}$	0	$9.4 \cdot 10^{-14}$
400	$5.7 \cdot 10^{-16}$	$5.5 \cdot 10^1$	0	$5.6 \cdot 10^{-16}$	0	$7.6 \cdot 10^{-13}$
800	$2.7 \cdot 10^{-15}$	$4.4 \cdot 10^2$	0	$1.2 \cdot 10^{-15}$	0	$1.8 \cdot 10^{-12}$

**Table 6** CPU time in seconds for Clement matrices

$n$	eig	$BGT$	$ex3dqds (+)$	$3dqds (+)$
100	0.011	0.83	0.014	0.009
200	0.097	2.01	0.020	0.014
400	0.28	8.33	0.036	0.025
800	0.90	35.70	0.080	0.066
1000	1.49	67.02	0.120	0.094

eigenvalues to the same accuracy, namely  $\mathcal{O}(\varepsilon)$ . In addition we provide eigenvectors and condition numbers.

An important further comment which illustrates challenges of the unsymmetric eigenvalue problem is that in these examples for  $n \geq 200$  the scaling matrices  $F$  used above (see (7.14)) are not representable. This limitation indicates why we use the  $\Delta T$  form for computing the eigenvectors. The overflow problem, which also arises for  $S^{-1}$  (see (7.12)), although not so quickly, explains why  $BGT$  confines its attention to  $n = 50$ , but we go further because of our approach.

### 8.3 Matrix with clusters

Matrix in Test 5 in [1],

$$C = D^{-1}tridiag(\mathbf{1}, \boldsymbol{\alpha}, \mathbf{1}), \quad D = diag(\boldsymbol{\beta}), \quad \boldsymbol{\alpha}, \boldsymbol{\beta} \in \mathbb{R}^n$$

$$\alpha_k = 10^{5(-1)^k} \cdot (-1)^{\lfloor k/4 \rfloor}, \quad \beta_k = (-1)^{\lfloor k/3 \rfloor}, \quad k = 1, \dots, n,$$

seems to be a challenging test matrix. It was designed to have large, tight clusters of eigenvalues around  $10^{-5}$ ,  $-10^5$  and  $10^5$ . Half the spectrum is around  $10^{-5}$  and the rest is divided unevenly between  $-10^5$  and  $10^5$ . The diagonal alternates between entries of absolute value  $10^5$  and  $10^{-5}$  and so, for  $dqds$  codes, there is a lot of rearranging to do. When  $n \geq 100$  it is not clear what is meant by accuracy.

The matrix has a repetitive structure and the diagonal entries are a good guide to the eigenvalues. For  $n = 100$  and for the large real eigenvalues near  $\pm 10^5$  the eigenvectors have spikes  $(-10^{-5}, -1, 10^{-5})$  (complex conjugate pairs have spikes  $(10^{-5}, 1, -10^{-7}, -1, 10^{-5})$ , at the appropriate places, and negligible elsewhere.

Hence the numerical supports for many eigenvectors are disjoint. The essential structure of the matrix is exhibited with  $n = 10$ ,

$$C = \begin{pmatrix} 10^{-5} & 1 & & & & & & & & & \\ & 1 & 10^5 & & & & & & & & \\ & & -1 & -10^{-5} & & & & & & & \\ & & & -1 & 10^5 & & & & & & \\ & & & & -1 & 10^{-5} & & & & & \\ & & & & & -1 & 10^5 & & & & \\ & & & & & & 1 & -10^{-5} & & & \\ & & & & & & & 1 & -10^{-5} & & \\ & & & & & & & & 1 & 10^5 & \\ & & & & & & & & & -1 & -10^{-5} & \\ & & & & & & & & & & -1 & -10^5 \end{pmatrix}$$

and it has 5 eigenvalues near 0, 3 eigenvalues near  $10^5$  and 2 near  $-10^5$ . All the eigenvalues are well-conditioned and the three codes obtain the correct number of eigenvalues in each cluster.

When  $n = 20$  there are 10 eigenvalues near  $10^{-5}$ , 6 near  $-10^5$  and 4 near  $10^5$ ;  $relcond(\lambda; C)$  and  $relcond(\lambda; L, U)$  are all less than  $1.2 \cdot 10^1$ ; the maximum residual norm for  $C$  and  $J = LU$  is  $\mathcal{O}(10^{-2})$ . For the eigenvalues with small modulus, *BGT* and *3dqds* (with 2 RQI, in average) compute approximations with relative errors of  $\mathcal{O}(\varepsilon)$ , whereas *eig* yield larger relative errors, as large as  $10^{-6}$ . See Table 7.

**Table 7** Relative errors for the three clusters in Test 5, with  $n = 20$

$\lambda$	<i>eig</i>		<i>BGT</i>		<i>3dqds (+ +)</i>	
	$rel_{min}$	$rel_{max}$	$rel_{min}$	$rel_{max}$	$rel_{min}$	$rel_{max}$
$\lambda \approx -10^5$	$2.0 \cdot 10^{-20}$	$1.9 \cdot 10^{-15}$	$2.0 \cdot 10^{-20}$	$7.1 \cdot 10^{-17}$	$2.0 \cdot 10^{-20}$	$8.6 \cdot 10^{-11}$
$\lambda \approx 10^5$	$5.0 \cdot 10^{-31}$	$2.2 \cdot 10^{-13}$	$5.0 \cdot 10^{-31}$	$2.2 \cdot 10^{-14}$	$5.0 \cdot 10^{-31}$	$1.0 \cdot 10^{-10}$
$ \lambda  \approx 10^{-5}$	$1.2 \cdot 10^{-8}$	$2.0 \cdot 10^{-6}$	$3.0 \cdot 10^{-17}$	$2.7 \cdot 10^{-16}$	$8.0 \cdot 10^{-17}$	$2.0 \cdot 10^{-16}$

### 8.4 Other scaled test matrices

Here we consider other test matrices from [1]. The eigenvalues of these matrices have a variety of distributions, in particular, the eigenvalues in Test 4 and Test 7 are distributed along curves. See Fig. 5. All these matrices are given in the form

$$C = D^{-1}tridiag(\mathbf{1}, \alpha, \mathbf{1}), D = \text{diag}(\beta), \alpha, \beta \in \mathbb{R}^n.$$

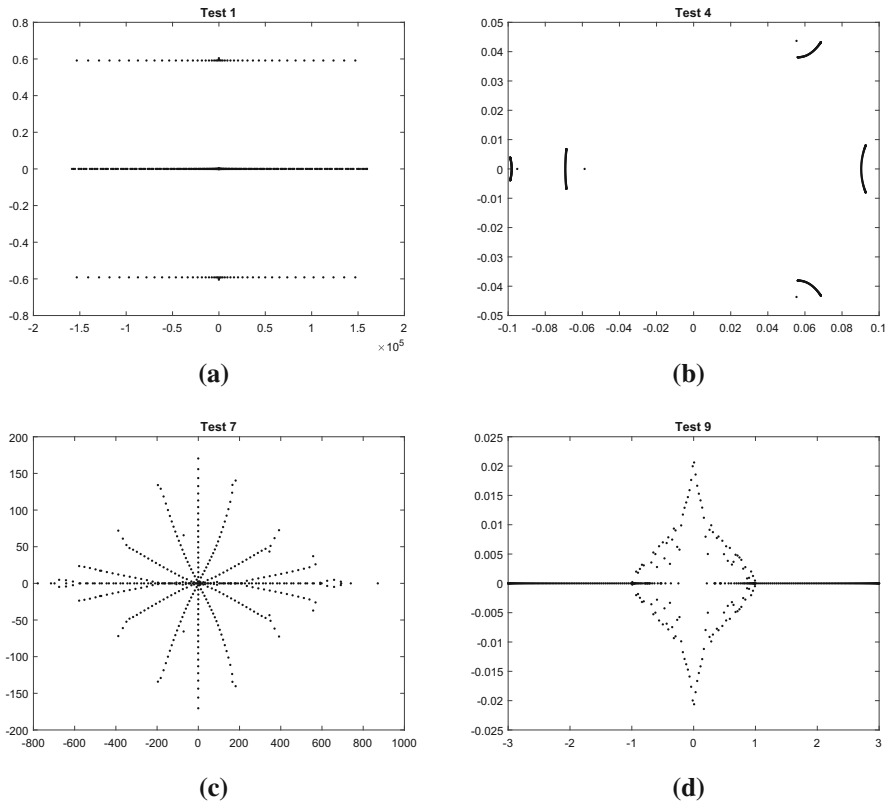


Fig. 5 Eigenvalues of matrices in Tests 1, 4, 7 and 9 for  $n = 400$

$$\begin{aligned}
 \text{Test 1} : \alpha_k &= (-1)^{\lfloor k/8 \rfloor}, \beta_k = (-1)^k/k, \quad k = 1, \dots, n. \\
 \text{Test 3} : \alpha_k &= k, \beta_k = n - k + 1, \quad k = 1, \dots, n. \\
 \text{Test 4} : \alpha_k &= (-1)^k, \beta_k = 20 \cdot (-1)^{\lfloor k/5 \rfloor}, \quad k = 1, \dots, n. \\
 \text{Test 6} : \alpha_k &= 2, \beta_k = 1, \quad k = 1, \dots, n. \\
 \text{Test 7} : \alpha_k &= \frac{1}{k} + \frac{1}{n - k + 1}, \beta_k = \frac{1}{k}(-1)^{\lfloor k/9 \rfloor}, \quad k = 1, \dots, n. \\
 \text{Test 9} : \alpha_k &= 1, \quad k = 1, \dots, n; \quad \beta_k = \begin{cases} 1 & \text{if } k < n/2 \\ -1 & \text{if } k \geq n/2 \end{cases}
 \end{aligned}
 \tag{8.1}$$

The extreme relative errors, condition numbers and residual norms for the three codes, MATLAB’s eig, BGT and 3dqds, are shown in Tables 8 and 9.

Table 10 reports the CPU time in seconds required by 3dqds versus the time required by MATLAB’s eig, BGT and ex3dqds with  $n$  ranging from 400 to 1000. Examples were chosen to represent the best, worst, and average efficiency of BGT.

**Table 8** Relative errors for matrices in (8.1) for  $n = 100$

Test	eig		BGT		3dqds		
	<i>rel<sub>min</sub></i>	<i>rel<sub>max</sub></i>	<i>rel<sub>min</sub></i>	<i>rel<sub>max</sub></i>	<i>rel<sub>min</sub></i>	<i>rel<sub>max</sub></i>	
1	$2.3 \cdot 10^{-17}$	$3.3 \cdot 10^{-13}$	$9.8 \cdot 10^{-19}$	$2.8 \cdot 10^{-16}$	$9.8 \cdot 10^{-19}$	$1.0 \cdot 10^{-15}$	(+)
3	0	$1.1 \cdot 10^{-14}$	0	$1.1 \cdot 10^{-14}$	0	$1.1 \cdot 10^{-14}$	(+)
4	$2.8 \cdot 10^{-16}$	$5.5 \cdot 10^{-15}$	$4.3 \cdot 10^{-18}$	$1.1 \cdot 10^{-16}$	$6.6 \cdot 10^{-18}$	$1.4 \cdot 10^{-16}$	(+ +)
6	$2.8 \cdot 10^{-18}$	$4.5 \cdot 10^{-13}$	$3.0 \cdot 10^{-19}$	$1.3 \cdot 10^{-13}$	$3.0 \cdot 10^{-19}$	$3.3 \cdot 10^{-14}$	(+)
7	$2.2 \cdot 10^{-17}$	$6.1 \cdot 10^{-14}$	$1.9 \cdot 10^{-18}$	$3.5 \cdot 10^{-16}$	$1.5 \cdot 10^{-18}$	$8.0 \cdot 10^{-16}$	(+)
9	$3.1 \cdot 10^{-17}$	$1.2 \cdot 10^{-14}$	$1.4 \cdot 10^{-18}$	$6.7 \cdot 10^{-16}$	$1.4 \cdot 10^{-18}$	$3.2 \cdot 10^{-15}$	(+)

**Table 9** Relative condition numbers and residual norms for matrices in (8.1) for  $n = 100$

Test	<i>relcond</i> ( $\lambda$ ; $L, U$ )		<i>relcond</i> ( $\lambda$ ; $C$ )		<i>max residuals</i>	
	<i>min</i>	<i>max</i>	<i>min</i>	<i>max</i>	$J = LU$	$\Delta T$
1	$1.0 \cdot 10^0$	$3.8 \cdot 10^2$	$1.0 \cdot 10^0$	$1.6 \cdot 10^2$	$4.8 \cdot 10^{-12}$	$1.8 \cdot 10^{-11}$
3	$1.0 \cdot 10^0$	$2.3 \cdot 10^0$	$1.0 \cdot 10^0$	$1.1 \cdot 10^1$	$4.9 \cdot 10^{-14}$	$1.3 \cdot 10^{-12}$
4	$1.3 \cdot 10^0$	$2.8 \cdot 10^0$	$1.3 \cdot 10^0$	$2.4 \cdot 10^1$	$3.5 \cdot 10^{-8}$	$1.3 \cdot 10^{-7}$
6	$1.0 \cdot 10^0$	$5.0 \cdot 10^1$	$1.0 \cdot 10^0$	$4.1 \cdot 10^3$	$1.3 \cdot 10^{-10}$	$1.3 \cdot 10^{-10}$
7	$1.5 \cdot 10^0$	$5.5 \cdot 10^2$	$1.3 \cdot 10^0$	$2.4 \cdot 10^1$	$3.0 \cdot 10^{-10}$	$1.5 \cdot 10^{-9}$
9	$1.1 \cdot 10^0$	$7.2 \cdot 10^2$	$1.0 \cdot 10^0$	$2.1 \cdot 10^2$	$3.3 \cdot 10^{-9}$	$3.3 \cdot 10^{-9}$

**Table 10** CPU time in seconds for matrices in Tests 3, 6 and 9

Test; $n$	eig	BGT	ex3dqds (+)	3dqds (+)
3; 400	0.11	3.12	0.07	0.03
6; 400	0.003	53.0	0.04	0.03
9; 400	0.39	19.5	0.52	0.32
3; 800	0.34	13.5	0.13	0.08
6; 800	0.01	360.2	0.12	0.08
9; 800	1.28	84.5	1.28	0.94
3; 1000	0.77	18.76	0.14	0.10
6; 1000	0.02	443.3	0.14	0.09
9; 1000	2.12	145.0	1.68	1.31

### 8.5 Liu matrix

Liu [16] devised an algorithm to obtain one-point spectrum unreduced tridiagonal matrices of arbitrary dimension  $n \times n$ . These matrices have only one eigenvalue, zero with multiplicity  $n$ , the Jordan form consists of one Jordan block and so the eigenvalue condition number is infinite. Our code *3dqds* computes this eigenvalue exactly (and also the generalized eigenvectors) using the following method which is part of the prologue. See [12].



The best place to start looking for eigenvalues of a tridiagonal matrix  $C = \text{tridiag}(\mathbf{b}, \mathbf{a}, \mathbf{c})$  is at the arithmetic mean which we know ( $\mu = \text{trace}(C)/n$ ). Before converting to  $J$ -form and factoring, we check whether  $\mu$  is an eigenvalue by using the 3-term recurrence to solve

$$(\mu I - C)\mathbf{x} = \mathbf{e}_n p_n(\mu) / \prod_{i=1}^{n-1} c_i.$$

Here

$$x_1 = 1, \quad x_2 = (\mu - a_2)/c_1, \quad x_{j+1} = \frac{1}{c_j} [(\mu - a_j)x_j - b_{j-1}x_{j-1}], \quad j = 2, \dots, n-1,$$

and

$$v := (\mu - a_n)x_n - b_{n-1}x_{n-1} \left( = p_n(\mu) / \prod_{i=1}^{n-1} c_i \right).$$

If, by chance,  $v$  vanishes, or is negligible compared to  $\|\mathbf{x}\|$ , then  $\mu$  is an eigenvalue (to working accuracy) and  $\mathbf{x}$  is an eigenvector. To check its multiplicity we differentiate with respect to  $\mu$  and solve

$$(\mu I - C)\mathbf{y} = \mathbf{x}$$

with  $y_1 = 0, y_2 = 1 = x'_2 (= x_1)$ . If

$$v' = p'_n(\mu) / \prod_{i=1}^{n-1} c_i := (\mu - a_n)y_n - b_{n-1}y_{n-1} + x_n$$

vanishes, or is negligible w.r.t.  $\|\mathbf{y}\|$ , then we continue the same way until the system is inconsistent or there are  $n$  generalized eigenvectors.

Usually  $v \neq 0$  and the calculation appears to have been a waste. This is not quite correct. In exact arithmetic, triangular factorization of  $\mu I - C$  or  $\mu I - J$ , where  $J = DCD^{-1}$ , will break down if, and only if,  $x_j$  vanishes for  $1 < j < n$ . So our code examines  $\min_j |x_j|$  and if it is too small w.r.t. its neighbors and w.r.t.  $\|\mathbf{x}\|$  then we do not choose  $\mu$  as our initial shift. Otherwise we do obtain initial  $L$  and  $U$  from  $J - \mu I = LU$ .

For comparison purposes, we ignored our prologue and give to our *3dqds* code the Liu matrices for  $n = 14$  and  $n = 28$ ,  $\text{tridiag}(\mathbf{1}, \boldsymbol{\alpha}^n, \boldsymbol{\gamma}^n)$  defined by

$$\begin{aligned} \boldsymbol{\alpha}^{14} &= [0, 0, 0, 0, 0, 0, -1, 1, 0, 0, 0, 0, 0, 0], \\ \boldsymbol{\gamma}^{14} &= [-1, 1, 1, -1, 1, -1, -1, -1, 1, -1, 1, 1, -1], \end{aligned}$$

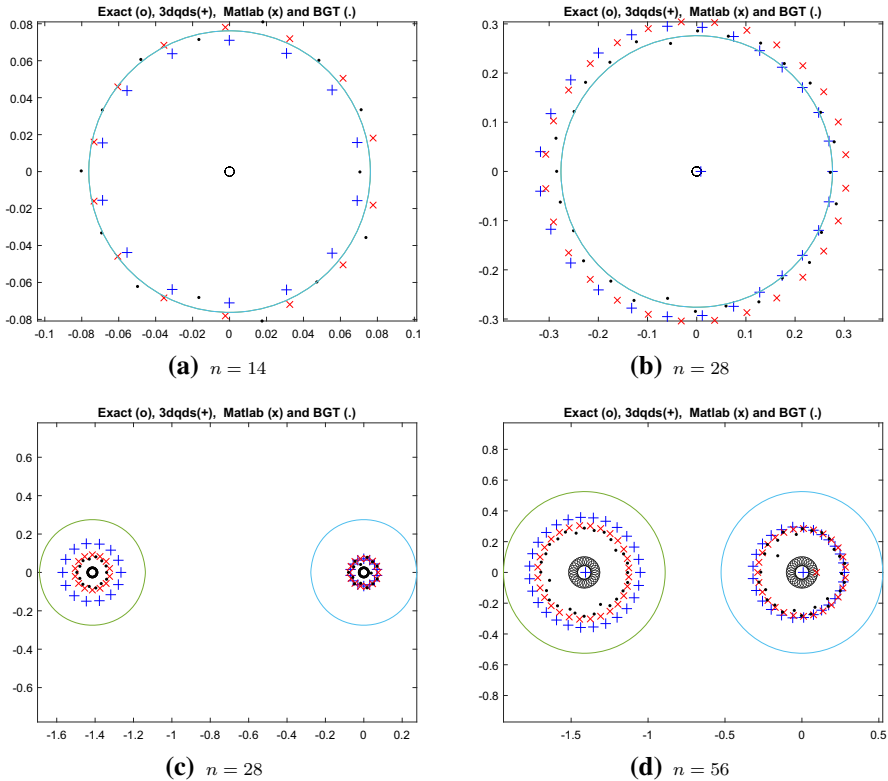


Fig. 6 Eigenvalues of Liu matrices (a, b), and glued Liu matrices (c, d)

and

$$\alpha^{28} = [0, 0, 0, 0, 0, 0, -1, 1, 0, 0, 0, 0, 0, -1, 1, 0, 0, 0, 0, 0, 0, 1, -1, 0, 0, 0, 0, 0, 0],$$

$$\gamma^{28} = \begin{bmatrix} -1, 1, 1, -1, 1, -1, -1, -1, 1, -1, 1, 1, -1, \\ -1, -1, 1, 1, -1, 1, -1, -1, -1, 1, -1, 1, 1, -1 \end{bmatrix}.$$

The accuracy of the approximations delivered by *3dqds* is as good as the accuracy of those provided by *MATLAB* and *BGT*. The absolute errors are  $\mathcal{O}(10^{-2})$  for  $n = 14$  and  $\mathcal{O}(10^{-1})$  for  $n = 28$ . The number of iterations needed for *3dqds* to converge is less than  $3n$ . See Fig. 6a, b. We show the numerical results along with the circles  $z = \sqrt[n]{\varepsilon}$ .

We also considered *glued* Liu matrices which are defined as the direct sum of two Liu matrices, shifting one of them by  $\sqrt{2}$  and letting the glue between them be  $\varepsilon$ . Roundoff will give us two clusters, one around 0, the other around  $\sqrt{2}$ . This is not a one-point spectrum matrix and all three methods give the results expected by perturbation theory. See Fig. 6c, d. This is a very unstable example, the condition numbers all exceed  $10^{10}$ .

## 9 Conclusions

Following the broad success of the HQR algorithm to compute eigenvalues of real square matrices it seems natural to use a sequence of similarity transforms to reduce an initial real matrix to eventual triangular form and also deflate eigenvalues from the bottom of the matrix as they converge. Any real (unreduced) tridiagonal matrix is easily put into  $J$ -form (all superdiagonal entries are 1) and such matrices ask for the use of the LR (not QR) algorithm since it preserves the  $J$ -form. The potential breakdown of the LR transform, from a 0 pivot, was a strong deterrent in the early days (1960s) [20] but today is a mild nuisance as explained in Sect. 5.1. A further incentive is that the whole procedure can be carried out in real arithmetic since complex conjugate pairs of eigenvalues are determined from  $2 \times 2$  submatrices that converge and may be deflated in a manner similar to real eigenvalues. The more recent success of the  $dqds$  transform in computing singular values of bidiagonal matrices encouraged us to keep out  $J$  matrices in factored form:  $J - \sigma I = LU$ ,  $\widehat{J} = UL$ , because, in exact arithmetic, the two algorithms, LR and  $dqds$ , are equivalent. In addition the  $dqds$  transform of today is numerically superior to the original, and seminal,  $qd$  transform discovered by Rutishauser [26] and which gave rise to the LR algorithm itself.

In order to hasten convergence we will need to apply complex conjugate pairs of shifts to our current  $LU = J$  matrix. It is well known how to do this entirely in real arithmetic in the context of the LR algorithm. To the best of our knowledge this has not been tried in the context of  $dqds$ . The main contribution of this paper is the solution to this challenge. We realized that three, not two, transforms are required to return to real factors  $L$  and  $U$  when complex shifts are applied consecutively. This is the nature of our explicit version, a local detour invoking complex arithmetic. We went further and produced a subprogram  $3dqds$  that accomplishes the same goal but in (exact) real arithmetic. This implicit version is more efficient than the explicit but is sensitive to roundoff error in its initial step. Experts will recall the papers on “washout of the shift” in the implicit shift HQR algorithm in the 1980s. We can not prove that our algorithm is backward stable. In fact we doubt that it is. However we do show that the three parts of the inner loop separately enjoy high mixed relative stability.

In the process of implementing our new features we were led to a novel and detailed criterion for deciding when our  $J = LU$  matrix has split into two or more unreduced submatrices. We check for splits at every iteration. Our new subprograms must only be applied to unreduced matrices. We also gave attention to the choice of a new shift when a factorization fails and when to start using the bottom  $2 \times 2$  submatrix for shifting.

We save a lot of space by confining our eigenvector calculations to the  $\Delta T$  form so that only one vector need be stored. From it we can compute the relative condition numbers that we need. Instructions are given how to generate the eigenvectors for the original and the  $J$ -form representations.

**Acknowledgements** The authors would like to thank Associate Editor Martin H. Gutknecht and the anonymous referees for forcing us to look more deeply into an error analysis of our *triple dqds* algorithm (first version) and to give a clearer presentation of its mathematical analysis and implementation details (last version).

## AppendixA 3dqds algorithm

```

 $\widehat{l}, \widehat{u}$  = 3dqds( $l, u, \text{sum}, \text{prod}$ )
% sum = ( $\sigma_1 + \sigma_2$ ); prod =  $\sigma_1 \sigma_2$ 
%  $l = [l_1, l_2, \dots, l_{n-1}]$ ;  $u = [u_1, u_2, \dots, u_n]$ 
%  $\widehat{l} = [\widehat{l}_1, \widehat{l}_2, \dots, \widehat{l}_{n-1}]$ ;  $\widehat{u} = [\widehat{u}_1, \widehat{u}_2, \dots, \widehat{u}_n]$ 

% step 1
 $x_r = 1$ ;  $y_r = l_1$ ;  $z_r = 0$ 
% the effect of  $Z_1$ 
 $x_r = x_r * u_1 + y_r$ 
% the matrix  $\mathcal{L}_1^{-1}$ 
 $x_l = (u_1 + l_1)^2 + u_2 l_1 - \text{sum}(u_1 + l_1) + \text{prod}$ 
 $y_l = -u_2 l_1 u_3 l_2 / x_l$ 
 $x_l = -u_2 l_1 (u_1 + l_1 + u_2 + l_2 - \text{sum}) / x_l$ 
% the effect of  $\mathcal{L}_1$ 
 $\widehat{u}_1 = x_r - x_l$ ;
 $x_r = y_r - x_l$ ;  $y_r = z_r - y_l - x_l * l_{i+1}$ ;
 $z_r = -y_l * l_3$ 
% the matrix  $Y_1^{-1}$ 
 $x_r = x_r / \widehat{u}_1$ ;  $y_r = y_r / \widehat{u}_1$ ;  $z_r = z_r / \widehat{u}_1$ 
% the effect of  $Y_1^{-1}$ 
 $\widehat{l}_1 = x_l + y_r + x_r * u_2$ 
 $x_l = y_l + z_r + y_r * u_3$ ;  $y_l = z_r * u_4$ 
% the effect of  $Y_1$ 
 $x_r = 1 - x_r$ ;  $y_r = l_2 - y_r$ ;  $z_r = -z_r$ 

% steps 2 to n-3
for  $i = 2, \dots, n - 3$ 
% the effect of  $Z_i$ 
 $x_r = x_r * u_i + y_r$ 
% the matrix  $\mathcal{L}_i^{-1}$ 
 $x_l = -x_l / \widehat{l}_{i-1}$ ;  $y_l = -y_l / \widehat{l}_{i-1}$ ;
% the effect of  $\mathcal{L}_i$ 
 $\widehat{u}_i = x_r - x_l$ ;
 $x_r = y_r - x_l$ ;  $y_r = z_r - y_l - x_l * l_{i+1}$ ;
 $z_r = -y_l * l_{i+2}$ 
% the matrix  $Y_i^{-1}$ 
 $x_r = x_r / \widehat{u}_i$ ;  $y_r = y_r / \widehat{u}_i$ ;  $z_r = z_r / \widehat{u}_i$ 
% the effect of  $Y_i^{-1}$ 
 $\widehat{l}_i = x_l + y_r + x_r * u_{i+1}$ 
 $x_l = y_l + z_r + y_r * u_{i+2}$ ;  $y_l = z_r * u_{i+3}$ 
% the effect of  $Y_i$ 
 $x_r = 1 - x_r$ ;  $y_r = l_{i+1} - y_r$ ;  $z_r = -z_r$ 
end for

% step n-2
% the effect of  $Z_{n-2}$ 
 $x_r = x_r * u_{n-2} + y_r$ 
% the matrix  $\mathcal{L}_{n-2}^{-1}$ 
 $x_l = -x_l / \widehat{l}_{n-3}$ ;  $y_l = -y_l / \widehat{l}_{n-3}$ ;
% the effect of  $\mathcal{L}_{n-2}$ 
 $\widehat{u}_{n-2} = x_r - x_l$ ;
 $x_r = y_r - x_l$ ;  $y_r = z_r - y_l - x_l * l_{n-1}$ 
% the matrix  $Y_{n-2}^{-1}$ 
 $x_r = x_r / \widehat{u}_{n-2}$ ;  $y_r = y_r / \widehat{u}_{n-2}$ 
% the effect of  $Y_{n-2}^{-1}$ 
 $\widehat{l}_{n-2} = x_l + y_r + x_r * u_{n-1}$ 
 $x_l = y_l + y_r * u_n$ 
% the effect of  $Y_{n-2}$ 
 $x_r = 1 - x_r$ ;  $y_r = l_{n-1} - y_r$ 

% step n-1
% the effect of  $Z_{n-1}$ 
 $x_r = x_r * u_{n-1} + y_r$ 
% the matrix  $\mathcal{L}_{n-1}^{-1}$ 
 $x_l = -x_l / \widehat{l}_{n-2}$ 
% the effect of  $\mathcal{L}_{n-1}$ 
 $\widehat{u}_{n-1} = x_r - x_l$ ;
 $x_r = y_r - x_l$ 
% the matrix  $Y_{n-1}^{-1}$ 
 $x_r = x_r / \widehat{u}_{n-1}$ 
% the effect of  $Y_{n-1}^{-1}$ 
 $\widehat{l}_{n-1} = x_l + x_r * u_n$ 
% the effect of  $Y_{n-1}$ 
 $x_r = 1 - x_r$ 

% step n
% the effect of  $Z_n$ 
 $x_r = x_r * u_n$ 
% the matrix  $\mathcal{L}_n^{-1} = I$ 
% the effect of  $\mathcal{L}_n$ 
 $\widehat{u}_n = x_r$ ;
% the matrix  $Y_n^{-1} = I$ 

```

## Appendix B Pseudocode for the whole algorithm

---

### Algorithm 1 wrapper for 3dqds

---

**Input:** vectors  $a, b, c$

**Output:** eigenvalues of  $\text{tridiag}(b, a, c)$

```

top = 1
split(1) = top
indsplit = 1
nits = 0
itmax = 100n
acshift = 0

find  $l$  and  $u$  of  $J$  form [Algorithm 5]
while (top + 1 < n and its < itmax) do
  deflate as warranted [Algorithm 2]
  find splits, if any [Algorithm 3]
  if ( $l_{n-1} > 10^{-2}$  and  $l_{n-2} > 10^{-2}$ ) then
    [ $l_1, u_1, \text{fail}$ ] = dqds( $l(\text{top} : n), u(\text{top} : n), 0$ )
    if fail then
      [ $l_1, u_1, \text{shift}, \text{fail}$ ] = recover( $l(\text{top} : n), u(\text{top} : n)$ ) [Algorithm 4]
    end if
  else
    sum =  $l_{n-1} + (u_{n-1} + u_n)$ 
    prod =  $u_{n-1}u_n$ 
    [ $l_1, u_1, \text{fail}$ ] = 3dqds( $l(\text{top} : n), u(\text{top} : n), \text{sum}, \text{prod}$ )
    if fail then
      [ $l_1, u_1, \text{shift}, \text{fail}$ ] = recover( $l(\text{top} : n), u(\text{top} : n), \text{sum}, \text{prod}$ )
    end if
  end if
  if fail then
    return "too many failures, no convergence."
  end if
   $l = l_1; u = u_1$ 
  acshift = acshift + shift
  its = its + 1
end while

```

▷ code works on submatrix  $\text{top} : n$   
 ▷ vector  $\text{split}$  saves all active  $\text{top}'s$   
 ▷ index for  $\text{split}$   
 ▷ number of iterations  
 ▷ maximum number of iterations  
 ▷ accumulated shift; simple  $\text{dqds}$  is not restoring  
 ▷ vectors  $l$  and  $u$  for  $J = LU$   
 ▷ code should maintain  $\text{top} + 1 < n$   
 ▷ deflation may reduce  $n$   
 ▷ splits may increase  $\text{top}$   
 ▷ entries at the bottom are not small  
 ▷ simple  $\text{dqds}$  with zero shift  
 ▷  $\text{fail}$  is a boolean for failure  
 ▷ triple  $\text{dqds}$   
 ▷ update accumulated shift

---

**Algorithm 2** deflation body

---

▷ deflate as warranted; deflation may reduce  $n$

$\text{tol} = 10\varepsilon$  ▷ tolerance for deflation;  $\varepsilon =$  roundoff unit

**repeat**

▷ deflation  $2 \times 2$  criteria (9.8) and (9.9)

**if**  $\text{criteria}_{2 \times 2}$  **then**

▷ discriminant

$\text{ssum} = (l_{n-1} + (u_{n-1} + u_n))/2$

$\text{disc} = ((l_{n-1} + (u_{n-1} - u_n))/2)^2 + u_n l_{n-1}$  ▷ discriminant

$t = \sqrt{|\text{disc}|}$

**if**  $\text{disc} < 0$  **then** ▷ complex conjugate pair

$x_1 = \text{ssum} + it$

$x_2 = \text{ssum} - it$  ▷ no use of complex arithmetic

**else if**  $\text{ssum} == 0$  **then** ▷ real pair

$x_1 = t$

$x_2 = -t$

**else**

$x_1 = \text{sign}(\text{ssum}) * (|\text{ssum}| + t)$  ▷ no subtractions

$x_2 = u_{n-1} u_n / x_1$

**end if**

$\text{eigvals}([n-1, n]) = [x_1, x_2] + \text{acshift}$  ▷ eigvals stores the eigenvalues

$n = n - 2$

**else if**  $\text{criteria}_{1 \times 1}$  **then** ▷ deflation  $1 \times 1$  criteria (6.1)–(6.4)

$\text{eigvals}(n) = \text{acshift } u_n + \text{acshift}$

$n = n - 1$

**end if**

**until** deflation criteria not met

---

**Algorithm 3** splitting body

---

▷ find splits, if any, define  $\text{top}$  after a split

$\text{tol} = 10\varepsilon$  ▷ tolerance for splitting;  $\varepsilon =$  roundoff unit

**if**  $n > \text{top} + 2$  **then**

$k = n - 3$

**while** ( $k > \text{top}$  and  $\text{criteria}_{\text{split}}$  not met) **do** ▷ splitting criteria (9.6) and (9.7)

$k = k - 1$

**end while**

**if**  $k > \text{top}$  **then** ▷ there is a split

$\text{indsplit} = \text{indsplit} + 1$

$\text{split}(\text{indsplit}) = \text{top}$

$l_k = \text{acshift}$  ▷  $l_k$  saves accumulated shift of the previous segment

$\text{top} = k + 1$

**end if** ▷ if the condition for splitting is not met, there is nothing to do

**end if**

---

**Algorithm 4** recover**Input:** vectors  $l, u$ , real  $sum, prod$  (or vectors  $l, u$ )**Output:** vectors  $l_1, u_1$ , real  $shift$ , boolean  $fail$ 

```

 $\delta = \sqrt{\varepsilon}$ 
if  $nargin == 2$  then
   $simple = true$ 
   $sum = \delta$ 
   $prod = \delta$ 
   $shift = 0$ 
else
   $simple = false$ 
   $shift = u_n$ 
end if

 $fail = true$ 
 $nfail = 0$ 
 $maxfail = 10n$ 

while ( $fail$  and  $nfail < maxfail$ ) do
  if  $simple$  then
     $sum = sum(1 + \delta)$ 
     $prod = prod(1 + \delta)^2$ 
     $[l_1, u_1, fail] = 3dqds(l, u, sum, prod)$ 
     $simple = false$ 
    if  $fail == false$  then
       $shift = 0$ 
    end if
  else
     $shift = shift + \delta$ 
     $[l_1, u_1, fail] = dqds(l, u, shift)$ 
     $simple = true$ 
  end if
   $nfail = nfail + 1$ 
end while

```

▷ shift increment;  $\varepsilon =$  roundoff unit  
 ▷  $nargin =$  number of input arguments  
 ▷ failure in *dqds* with zero shift  
 ▷  $sum$  and  $prod$  for *3dqds*  
 ▷ in case of successive failures  
 ▷ failure in *3dqds*  
 ▷  $shift$  for simple *dqds*  
 ▷ boolean for failure  
 ▷ number of failures  
 ▷ maximum number of failures allowed  
 ▷ increase shift and reverse choice of transform  
 ▷ switch to *3dqds*  
 ▷ successful recovery with *3dqds*  
 ▷ switch to simple *dqds*  
 ▷ after a failure the opposite transform is used next

**Algorithm 5** initial *LU* factorization

---

**Input:** vectors  $a, b, c$  ▷ *LU* factorization of *tridiag*( $b, a, c$ ) in *J* form  
**Output:** vectors  $l, u$ , real shift, boolean fail

$nfail = 0$  ▷ number of failures  
 $maxfail = 10n$  ▷ maximum number of failures allowed  
 $b = b \cdot *c$  ▷ element-wise product; off-diagonal of *J*  
 $delta = \min(1/2, 2 * \min(abs(a(a \sim = 0))))$  ▷ shift increment in case of failure  
▷ one eighth of the minimum nonzero diagonal element

$shift = 0$  ▷ in case of failure take  $J - shift \cdot I = LU$   
 $[l, u, fail] = LUfact(a, b, shift)$  [Algorithm 6] ▷ *LU* factorization of *J*  
**while** (fail **and**  $nfail < maxfail$ ) **do**  
     $nfail = nfail + 1$   
     $shift = shift + delta$  ▷ after a failure the shift is increased  
     $[l, u, fail] = LUfact(a, b, shift)$   
**end while**

**if** fail **then**  
    **return** "Too many failures, no initial factorization."  
**end if**

---

**Algorithm 6** *LUfact*


---

**Input:** vectors  $a, b$ , real shift ▷ *LU* factorization of  $J = tridiag(b, a, 1)$  without pivoting  
**Output:** vectors  $l, u$ , boolean fail

$tolg = 1/\sqrt{\varepsilon}$  ▷ tolerance for element growth;  $\varepsilon$  = roundoff unit  
 $fail = false$  ▷ boolean for failure

$u(1) = a(1)$   
**for**  $i = 1 : n - 1$  **do**  
     $l(i) = b(i)/u(i)$   
     $u(i + 1) = a(i + 1) - l(i)$   
**end for**

**if** ( $any(isnan([l, u]))$  **or**  $any(abs([l, u])) > tolg$ ) **then** ▷ checking for element growth  
     $fail = true$   
**end if**

---

**References**

1. Bini, D.A., Gemignani, L., Tisseur, F.: The Ehrlich–Aberth method for the nonsymmetric tridiagonal eigenvalue problem. *SIAM J. Matrix Anal. Appl.* **27**(1), 153–175 (2005)
2. Bunse-Gerstner, A.: An analysis of the HR algorithm for computing the eigenvalues of a matrix. *Linear Algebra Appl.* **35**, 155–173 (1981)
3. Clement, P.A.: A class of triple-diagonal matrices for test purposes. In: *SIAM Review*, vol. 1 (1959)
4. Cullum, J.K.: A QL procedure for computing the eigenvalues of complex symmetric tridiagonal matrices. *SIAM J. Matrix Anal. Appl.* **17**(1), 83–109 (1996)
5. Day, D.: Semi-duality in the Two-sided Lanczos Algorithm. Ph.D Thesis, University of California, Berkeley (1993)
6. Demmel, J.W.: *Applied Numerical Linear Algebra*, Society for Industrial and Applied Mathematics (1997)
7. Dhillon, I.S., Parlett, B.N.: Multiple representations to compute orthogonal eigenvectors of symmetric tridiagonal matrices. *Linear Algebra Appl.* **387**, 1–28 (2004)



8. Dhillon, I.S., Parlett, B.N.: Orthogonal eigenvectors and relative gaps. *SIAM J. Matrix Anal. Appl.* **25**, 858–899 (2004)
9. Fernando, K.V., Parlett, B.: Accurate singular values and differential QD algorithms. *Numer. Math.* **67**, 191–229 (1994)
10. Fernando, K.V.: On computing an eigenvector of a tridiagonal matrix. Part I: basic results. *SIAM J. Matrix Anal. Appl.* **18**, 1013–1034 (1997)
11. Ferreira, C.: The Unsymmetric Tridiagonal Eigenvalue Problem. Ph.D Thesis, University of Minho (2007). <http://hdl.handle.net/1822/6761>
12. Ferreira, C., Parlett, B.: Convergence of LR algorithm for a one-point spectrum tridiagonal matrix. *Numer. Math.* **113**(3), 417–431 (2009)
13. Ferreira, C., Parlett, B., Froilán, M.D.: Sensitivity of eigenvalues of an unsymmetric tridiagonal matrix. *Numer. Math.* (2012). <https://doi.org/10.1007/s00211-012-0470-z>
14. Francis, J.G.F.: The QR transformation—a unitary analogue to the LR transformation, parts I and II. *Comput. J.* **4**, 265–272; 332–245 (1961/1962)
15. Kahan, W., Parlett, B.N., Jiang, E.: Residual bounds on approximate eigensystems of non-normal matrices. *SIAM J. Numer. Anal.* **19**, 470–484 (1982)
16. Liu, Z.A.: On the extended HR algorithm. Technical Report PAM-564, Center for Pure and Applied Mathematics, University of California, Berkeley, CA, USA (1992)
17. Parlett, B.N., Reinsch, C.: Balancing a matrix for calculation of eigenvalues and eigenvectors. *Numer. Math.* **13**, 292–304 (1969)
18. Parlett, B.N.: The Rayleigh quotient iteration and some generalizations for non-normal matrices. *Math. Comput.* **28**(127), 679–693 (1974)
19. Parlett, B.N.: The contribution of J. H. Wilkinson to numerical analysis. In: Nash, S.G. (ed.), *A History of Scientific Computing*, ACM Press, p. 25 (1990)
20. Parlett, B.N.: Reduction to tridiagonal form and minimal realizations. *SIAM J. Matrix Anal. Appl.* **13**, 567–593 (1992)
21. Parlett, B.N.: The new QD algorithms. *Acta Numer.* **4**, 459–491 (1995)
22. Parlett, B.N., Dhillon, I.S.: Fernandos solution to Wilkinsons problem: an application of double factorization. *Linear Algebra Appl.* **267**, 247–279 (1997)
23. Parlett, B.N., Marques, O.A.: An implementation of the DQDS algorithm. *Linear Algebra Appl.* **309**, 217–259 (2000)
24. Parlett, B., Dopico, F.M., Ferreira, C.: The inverse eigenvector problem for real tridiagonal matrices. *SIAM J. Matrix Anal. Appl.* **37**, 577–597 (2016)
25. Pasquini, L.: Accurate computation of the zeros of the generalized Bessel polynomials. *Numerische Mathematic* **86**, 507–538 (2000)
26. Rutishauser, H.: Der Quotienten-Differenzen-Algorithmus. *Z. Angew. Math. Physik* **5**, 233–251 (1954)
27. Rutishauser, H.: Der Quotienten-Differenzen-Algorithmus. *Mitt. Inst. Angew. Math. ETH*, vol. 7, Birkhäuser, Basel (1957)
28. Rutishauser, H.: Solution of eigenvalue problems with the LR-transformation. *Natl. Bur. Stand. Appl. Math. Ser.* **49**, 47–81 (1958)
29. Rutishauser, H., Schwarz, H.R.: The LR transformation method for symmetric matrices. *Numer. Math.* **5**, 273–289 (1963)
30. Slemons, J.: Toward the Solution of the Eigenproblem: Nonsymmetric Tridiagonal Matrices. Ph.D Thesis, University of Washington, Seattle (2008)
31. Slemons, J.: The result of two steps of the LR algorithm is diagonally similar to the result of one step of the HR algorithm. *SIAM J. Matrix Anal. Appl.* **31**(1), 68–74 (2009)
32. Trefethen, L.N., Embree, M.: Spectra and pseudospectra. In: *The Behavior of Nonnormal Matrices and Operators*, Princeton University Press (2005)
33. Watkins, D.S.: QR-like algorithms—an overview of convergence theory and practice. *Lect. Appl. Math.* **32**, 879–893 (1996)
34. Watkins, D.S., Elsner, L.: Convergence of algorithms of decomposition type for the eigenvalue problem. *Linear Algebra Appl.* **143**, 19–47 (1991)
35. Willems, P.R., Lang, B.: Twisted factorizations and QD-type transformations for the MR<sup>3</sup> algorithm—new representations and analysis. *SIAM J. Matrix Anal. Appl.* **33**(2), 523–553 (2012)
36. Wu, Z.: The Triple DQDS Algorithm for Complex Eigenvalues. Ph.D Thesis, University of California, Berkeley (1996)

37. Xu, H.: The relation between the QR and LR algorithms. *SIAM J. Matrix Anal. Appl.* **19**(2), 551–555 (1998)
38. Yao, Y.: Error Analysis of the QDs and DQDs Algorithms. Ph.D Thesis, University of California, Berkeley (1994)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.