

Accurate singular values and differential qd algorithms

K. Vince Fernando^{1,*,**}, Beresford N. Parlett^{2,**}

¹ NAG Ltd, Wilkinson House, Jordan Hill, Oxford OX2 8DR, UK

² Department of Mathematics, University of California, Berkeley, CA 94720, USA

Received August 8, 1993/Revised version received May 26, 1993

This work is dedicated to the memory of Heinz Rutishauser

Summary. We have discovered a new implementation of the qd algorithm that has a far wider domain of stability than Rutishauser's version. Our algorithm was developed from an examination of the Cholesky LR transformation and can be adapted to parallel computation in stark contrast to traditional qd. Our algorithm also yields useful a posteriori upper and lower bounds on the smallest singular value of a bidiagonal matrix.

The zero-shift bidiagonal QR of Demmel and Kahan computes the smallest singular values to maximal relative accuracy and the others to maximal absolute accuracy with little or no degradation in efficiency when compared with the LINPACK code. Our algorithm obtains maximal relative accuracy for all the singular values and runs at least four times faster than the LINPACK code.

Mathematics Subject Classification (1991): 65F15

1. Introduction

In September 1991 J. W. Demmel and W. M. Kahan were awarded the second SIAM prize in numerical linear algebra for their paper 'Accurate Singular Values of Bidiagonal Matrices' [11], referred to as DK hereafter. Among several valuable results was the observation that the standard bidiagonal QR algorithm used in LINPACK [12], and in many other SVD programs, can be simplified when the shift is zero and, of greater importance, no subtractions occur. The last feature permits very small singular values to be found with (almost) all the accuracy permitted by the data and at no extra cost.

In this paper we show that the DK zero shift algorithm can be further simplified and this simplicity has several benefits. One is that a new algorithm can be implemented in either parallel or pipelined format and each iteration nominally takes $O(\log_2 n)$ operations. Another benefit is that some or all singular vectors may be calculated accurately in a second phase after the singular values are known. See Sect. 13.

* Supported by NSF, under grant ASC-9005933

** Supported by ONR, contract N000014-90-J-1372

Correspondence to: B.N. Parlett

Our investigations began with the modest goal of showing that it was preferable to replace the DK zero-shift QR transform by two steps of zero-shift LR implemented in a qd (quotient-difference) format. Root-free algorithms run considerably faster than standard ones. The surprise here is that to keep the high relative accuracy property it is necessary to use a little known variant of qd (the differential form of the progressive qd algorithm or dqd [35], [34]). The standard qd will not suffice as we show in Sect. 4. There are no subtractions in dqd. We suspect that Rutishauser discovered dqd in 1968, just two years before his death, and we say more about its history in Sects. 4 and 11.

What we want to stress here is that, for reasons we may never know, Rutishauser did not consider the shifted version of dqd. Instead he reserved dqd for rectifying unsatisfactory behaviour of his qd in certain circumstances, see Sect. A4.2 of [35], [34]. Incidentally this differential qd is not to be confused with the continuous analogue of qd (see [31]) and more recent work on QR flows. The trouble with the shifted version of the ordinary qd algorithm is that it cannot recover from a shift that is too large. Consequently qd algorithms have been shackled with very conservative shift strategies, such as Newton's method, and earned the reputation of being slow compared to the QR algorithm. Had Rutishauser considered shifts with differential qd (dqds hereafter) he would have realized, as we soon did, that the transformation may be split into two parts. The parts depend on whether the machine is of sequential or parallel type but, in each case, a shift that is too big reveals itself before the old matrix is overwritten and so need not be invoked. An unused shift is not wasted because it gives an improved upper bound on the smallest singular value and the inertia count at a cost less than one qd transformation as well as contributing to an improved shift.

Our approach frees the algorithm to exploit powerful shift strategies while preserving high relative accuracy all the time. In contrast the QR algorithm delivers high relative accuracy only with a zero shift.

Even though our algorithms must find the singular values in order we can use shift strategies that are at least quadratically convergent. This is better than fourth order convergence for QR. When only the smallest few singular values are needed this ordering constraint is a great advantage. Another rather subtle feature is that it is not necessary to make an extra $O(n)$ check for splitting of the matrix into a direct sum. The necessary information is provided by the auxiliary quantities.

In June 1992 we discovered that our dqds algorithm enjoys mixed high relative stability for all shifts provided that they avoid underflow, overflow or divide by zero. Consequently it can be used in a variety of applications (eigenvalues of symmetric or unsymmetric tridiagonals, zeros of polynomials, poles and zeros of transfer functions and many applications involving continued fractions) where Rutishauser's qd has been abandoned because of its instability in the general case.

Our error bounds for singular values are significantly smaller than those in DK and our approach is quite transparent. It was this analysis, in Sect. 7, that showed us the possibility of violating positivity while still maintaining maximal relative accuracy for all singular values, not just the small ones.

It gradually dawned on us as we developed the algorithm that we were breaking away from the *orthogonal paradigm* that has dominated the field of matrix computations (often called numerical linear algebra by highbrows) since the 1960's. It seems to be sacrilegious to be achieving greater accuracy and on average, a four fold speed-up¹ by simply abandoning QR for something equivalent to LR. See Sect. 9.3

¹ All our computations are performed on a DECstation 5000/120 using double precision arithmetic (53-bit mantissa)

for details. High accuracy comes from the fact that `dqds` spends most of its time transforming lower triangular 2×2 s into upper triangular 2×2 s by premultiplication.

Rutishauser gave no direct explanation for the way shifts are introduced into `qd`. We have supplied one in terms of matrix factorizations in Sect. 5.1 and we go on to list the possible choices for a shift in Sects. 6 and 9. We do not offer a preferred shift strategy here because that aspect of the code is still evolving.

Section 3 presents the unifying general result which shows that it is possible to implement the Cholesky LR algorithm of Rutishauser [32], [36] using orthogonal transformations only. Since the term Cholesky LR over describes the algorithm we simply refer to it as the Cholesky Algorithm. Our orthogonal Cholesky algorithm is applicable to dense matrices; this more general case is studied elsewhere [17].

We want to point out the unusual historical lineage of this algorithm. The `qd` algorithm beget the LR algorithm which then gave rise to the QR algorithm of Francis. This in turn led to the Golub-Kahan and Golub-Reinsch algorithms for singular values of bidiagonal matrices which lead to the DK zero-shift variant. This inspired our orthogonal algorithm of which differential `qd` is the root-free version. We are back to `qd` again but with a new implementation.

As a service to busy readers we have included a brief account of the origins of `qd` and a summary of the DK paper. When reading [35] we regretted that the link between continued fractions and our matrices was not made explicit. We provide the connection in the final section.

2. Notation and normalization

This paper does not involve vectors very much and so we do not follow Householder conventions. However capital roman letters denote matrices while lower case Roman and Greek letters denote scalars. On the rare occasions when a vector is needed it is denoted by a lower case roman letter in boldface.

As usual the singular values of an $n \times n$ matrix C are arranged in monotone decreasing order and denoted by $\sigma_1, \sigma_2, \dots, \sigma_n$, their union is $\sigma[C]$.

- We make reference to the QR factorization of a matrix. This is the matrix form of the Gram-Schmidt orthonormalizing process applied to the columns of the matrix in natural order. By convention the diagonal entries of the upper triangular factor R are taken nonnegative. See Golub and Van Loan [19] for details.
- We make reference to the Cholesky factorization of a positive definite matrix into the product of a lower triangular matrix and its conjugate transpose. The factors are unique.
- We make references to the LR and QR algorithms. These are defined in the appropriate places.

We shall be concerned mainly with bidiagonal matrices which we call B and take them to be *upper* bidiagonal. To save space we write the bidiagonal matrix

$$B = \begin{bmatrix} a_1 & b_1 & & & & & \\ & a_2 & b_2 & & & & \\ & & & \ddots & & & \\ & & & & \ddots & & \\ & & & & & a_{n-1} & b_{n-1} \\ & & & & & & a_n \end{bmatrix}$$

as

$$B = \text{bidiag} \left\{ \begin{array}{cccccccc} & b_1 & b_2 & \cdot & \cdot & b_{n-2} & b_{n-1} & \\ a_1 & & a_2 & & & & a_{n-1} & a_n \end{array} \right\}.$$

2.1. Normalization

Consider now the effect of a zero value among the parameters of a $n \times n$ bidiagonal B .

2.1.1 Superdiagonal. Suppose that $b_k = 0$, $k < n$. Then B may be written as a direct (or diagonal) sum of two bidiagonals B_1 and B_2 . Moreover

$$\sigma[B] = \sigma[B_1] \cup \sigma[B_2].$$

This case makes the calculation of singular values easier. Even more important is the fact that our algorithms do not suffer from the failure to detect such a split when it occurs. However, the transition from a linearly convergent shift to a quadratic shift will not occur if the split lies undetected for too long.

2.1.2 Diagonal. Let $a_k = 0$, $k < n$. Since $|\det B| = \prod_{i=1}^n |a_i| = \prod_{i=1}^n \sigma_i$ it follows that $\sigma_n = 0$. However some work is needed in order to *deflate* this value, i.e. to find a new B of order $n - 1$ yielding the remaining singular values of B . In exact arithmetic one iteration of any of the unshifted algorithms given later is guaranteed to produce the desired B and so this case does not need special treatment. The zero diagonal entry may be driven to the closest end of the matrix.

If $a_k = 0$, $k < n$, at one step of our algorithm and if $a_n = 0$ at the next step then b_{k-1} will also vanish and so produce a split into two bidiagonals.

2.1.3 Signs. If the matrix is real, then using pre and post multiplication by matrices of the form $\text{diag}\{\pm 1\}$ any sign pattern may be imposed on the entries of B without changing the singular values. If the matrix is complex, then it could be transformed to a real matrix by pre and post multiplication by matrices of the form $\text{diag}\{\exp(i\omega)\}$ where $i^2 = -1$ and ω is real.

There is little loss of generality in assuming, when necessary, that B is of real positive type; all its parameters exceed 0. However in Sect. 5.3 we address the practical question of when to relax the requirement of positivity.

3. Orthogonal form of the Cholesky LR algorithm

For the next few paragraphs we consider full complex matrices. Recall that the Cholesky factorization of a positive definite Hermitian matrix $A (= A^*)$ may be written as $A = LL^*$ where L is lower triangular.

Definition. The Cholesky LR *transform* of a symmetric positive definite matrix $A = LL^*$ is

$$\hat{A} := L^*L$$

The Cholesky LR *algorithm*, consisting of successive applications of the Cholesky transformation, is a special case of Rutishauser's LR algorithm.

The following lemma is elementary and has been known for a long time. It has become the cornerstone of square root filtering in control and signal processing (see

for example [24], [6], [16]), a trend which was pioneered by Potter (see Battin [3]). In Sect.4 we show that for the bidiagonal matrix case, the ‘‘square root’’ approach can be bypassed with a huge gain in accuracy and efficiency.

Lemma 1. *If any two invertible matrices M_1 and M_2 satisfy $M_1^* M_1 = M_2^* M_2$ then $M_2 = Q M_1$ for some orthogonal matrix Q .*

The result given in the theorem below is implicit in proofs that one step of the QR algorithm is equal to two steps of the Cholesky algorithm. It was explicitly proved by Faddeev, Kublanovskaya and Faddeeva [15] and we thank Ilse Ipsen for bringing this relatively unknown work to our attention.

Theorem 1. *Let $\hat{A} = \hat{L}\hat{L}^*$ be the Cholesky factorization of the Cholesky transform of positive definite $A = LL^*$. Then*

$$L = Q\hat{L}^*$$

is the QR factorization of L .

Proof. Of course, this theorem is a direct corollary of Lemma 1 but we include a proof because it is both constructive and instructive.

Since A is positive definite all factors mentioned below are unique. By definition of \hat{L}

$$L^* L = \hat{L}\hat{L}^*.$$

We seek invertible F such that

(1)
$$L = F\hat{L}^*,$$

(2)
$$L^* = \hat{L}F^{-1}.$$

Transpose and conjugate (1) and use invertibility of \hat{L} in (2) to find

$$F^* = \hat{L}^{-1}L^* = F^{-1}.$$

So F is unitary and since \hat{L}^* is upper triangular with positive diagonal Eq.(1) above gives the QR factorization of L , as claimed. \square

The theorem shows that \hat{L} may be obtained from L by orthogonal transformations without forming \hat{A} . Moreover just as QR may be performed with column pivoting so can we obtain the Cholesky factor of a permutation of \hat{A} .

In the general dense case people have used the method of transforming triangular matrices from upper to lower form and back again using appropriate orthogonal transformations without realizing that this is equivalent to the Cholesky algorithm. See [41] and [7]. Conversely, if authors of [42] had appreciated Theorem 1 they could have used, with advantage, more preconditioning LR steps before invoking the one-sided Jacobi algorithm. We plan to pursue this general dense case in another paper [17].

The basic equation $\hat{L}\hat{L}^T = L^T L$ guarantees that the Cholesky algorithm preserves bandwidth. In particular, bidiagonal B gives rise to tridiagonal $A = B^T B$ and a bidiagonal \hat{B} . In order to study how \hat{B} is derived from B , let

$$B = \text{bidiag} \left\{ \begin{matrix} & b_1 & b_2 & \dots & b_{n-2} & b_{n-1} & \\ a_1 & & a_2 & & & a_{n-1} & a_n \end{matrix} \right\}.$$

$$\hat{B} = \text{bidiag} \left\{ \begin{array}{cccccccc} & \hat{b}_1 & \hat{b}_2 & \cdot & \cdot & \hat{b}_{n-2} & \hat{b}_{n-1} & \\ \hat{a}_1 & & & & & & & \\ & \hat{a}_2 & & & & & & \\ & & \cdot & \cdot & & & & \\ & & & & & \hat{a}_{n-1} & & \\ & & & & & & \hat{a}_n & \end{array} \right\}.$$

where $\hat{B}^T \hat{B} = BB^T$. By Theorem 1

$$B^T = Q\hat{B}.$$

The matrix Q may be written as a product of $(n-1)$ plane rotation matrices [19],

$$Q = G_1 G_2 \dots G_{n-1}.$$

Before the annihilation of the subdiagonal element b_k , the active part of the matrix is of the form,

$$(3) \quad \begin{array}{cccccc} 0 & \hat{a}_{k-1} & \hat{b}_{k-1} & & & \\ & 0 & \tilde{a}_k & 0 & & \\ & & b_k & a_{k+1} & 0 & \\ & & & b_{k+1} & a_{k+2} & \end{array}$$

and after the plane rotation G_k^T , the matrix becomes

$$(4) \quad \begin{array}{cccccc} 0 & \hat{a}_{k-1} & \hat{b}_{k-1} & & & \\ & 0 & \hat{a}_k & \hat{b}_k & & \\ & & 0 & \tilde{a}_{k+1} & 0 & \\ & & & b_{k+1} & a_{k+2} & \end{array}$$

Formally we may set $B^{(0)} = B^T$ and, for $k = 1, \dots, n-1$

$$(5) \quad B^{(k)} = G_k^T B^{(k-1)}.$$

Finally $\hat{B} = B^{(n-1)}$ and, from (3) and (4), with $\tilde{a}_1 = a_1$ and $c_k^2 + s_k^2 = 1$,

$$(6) \quad \hat{a}_k = \sqrt{\tilde{a}_k^2 + b_k^2} = \tilde{a}_k / c_k$$

$$(7) \quad s_k = b_k / \hat{a}_k$$

$$(7) \quad c_k = \tilde{a}_k / \hat{a}_k$$

$$(8) \quad \hat{b}_k = s_k a_{k+1} = b_k a_{k+1} / \hat{a}_k$$

$$\tilde{a}_{k+1} = c_k a_{k+1} = \tilde{a}_k a_{k+1} / \hat{a}_k.$$

There is some redundancy in the equations given above but their most important property is the absence of subtractions. This ensures high relative accuracy in the new entries \hat{a}_i and \hat{b}_i . Observe that neither s_k nor c_k is needed explicitly to compute the new entries. To the best of our knowledge the algorithm given below is new. For reasons that appear in the next section we call it the *Orthogonal qd-Algorithm* or *oqd*. It is convenient to use

$$(9) \quad \text{cabs}(x, y) = \sqrt{x^2 + y^2}$$

whose name stands for the complex absolute value of $x + iy$. In numerical computing (e.g. Eispack), an alternative name for *cabs* is *pythag*.

Algorithm 1 (oqd)

```

 $\tilde{a} := a_1$ 
for  $k = 1, n - 1$ 
   $\hat{a}_k := \text{cabs}(\tilde{a}, b_k)$ 
   $\hat{b}_k := b_k * (a_{k+1}/\hat{a}_k)$ 
   $\tilde{a} := \tilde{a} * (a_{k+1}/\hat{a}_k)$ 
end for
 $\hat{a}_n := \tilde{a}$ 

```

This algorithm will undergo several transformations in the following pages before we are ready to implement it. Nevertheless, even at this stage, two applications of it are slightly better (fewer multiplications) than the DK Zero Shift QR algorithm [11] described briefly in our Sect. 10. This result was the initial impetus for our study of qd.

The inner loop comparisons given in Table 1 are based on one QR step which is equal to two LR steps. We have taken into account the common sub-expression a_{k+1}/\hat{a}_k in the estimation of the complexity of oqd (Algorithm 1).

Table 1. Complexity of Demmel-Kahan and oqd

	DK	oqd
Cabs	2	1*2
Divisions	2	1*2
Multiplications	6	2*2
Conditionals	1	0
Assignments	7	3*2
Auxiliary variables	6	1

DK uses six auxiliary variables while oqd needs only one. The memory traffic is essentially determined by the number of variables, arithmetic operations and assignment statements. In most advanced architectures, memory access is more expensive than floating-point operations and in such machines the oqd will be very advantageous because of fewer read and write operations.

4. The quotient difference algorithm

It is easy to avoid taking the square roots that appear in oqd (Algorithm 1). Define $b_n := 0$ and $q_k = a_k^2$, $e_k = b_k^2$, $k = 1, 2, \dots, n$. By simply squaring each assignment in oqd (Algorithm 1) one obtains an algorithm that turns out to be a little known variant of the quotient difference algorithm. Rutishauser developed his qd algorithm in several papers from 1953 or 1954 (e.g. [30]) until his early death in 1970 but this variant appeared in English only in 1990 in [35] which is the translation of the German original [34] published in 1976. The full list of the papers on qd by Rutishauser can be found in the above mentioned books which were published posthumously.

In the notes at the end of [30] and at the end of volume 2 of [34] this variant is called the *differential form* of the progressive qd algorithm or dqd. These notes were based on unfinished manuscripts of Rutishauser.

Algorithm 2 (dqd)

```

d := q1
for k := 1, n - 1
  q̂k := d + ek
  êk := ek * (qk+1/q̂k)
  d := d * (qk+1/q̂k)
end for
q̂n := d

```

The implementation of dqd (Algorithm 2) requires only 1 division, 2 multiplies, and 1 addition in the inner loop. No subtractions occur.

The intermediate variable d may be removed. At step k , $d = d_k$ and the trick is to write it as a difference.

$$d_{k+1} = c_k^2 q_{k+1} = q_{k+1} - s_k^2 q_{k+1} = q_{k+1} - \hat{e}_k.$$

The resulting algorithm is Rutishauser's qd algorithm.

Algorithm 3 (qd)

```

ê0 = 0
for k := 1, n - 1
  q̂k := (qk - êk-1) + ek
  êk := ek * qk+1/q̂k
end for
q̂n := qn - ên-1

```

Table 2 compares the complexity of orthogonal, differential and standard qd algorithms.

Table 2. Complexity of oqd, dqd and qd

	oqd	dqd	qd
Cabs	1	0	0
Civisions	2	1	1
Multiplications	4	2	1
Additions	1	1	1
Subtractions	0	0	1
Assignments	3	3	2
Auxiliary variables	1	1	0

We hasten to add that Rutishauser did not derive the qd algorithm from our Theorem 1 but from ideas described in Sect. 11.

For positive B , dqd and qd are stable in the sense that all intermediate quantities are bounded by $\|B\|^2$. Singular value errors provoked by finite precision arithmetic will be tiny compared to σ_1^2 . This is satisfactory for many purposes and it was not generally appreciated until the DK paper appeared that bidiagonal matrices do determine all their singular values, however small, to the same relative precision enjoyed by the matrix entries. Since such accuracy can be achieved for little extra cost it seems only right to do so. These considerations lead us to abandon qd and concentrate on dqd and oqd.

Example 1. Here is a bidiagonal Toeplitz matrix with $a_i = 1$, $b_i = 256$ ($q_i = 1$, $e_i = 65536$) for all i . The results of our dqd algorithm are given in Table 3. Note that $\sqrt{q_{64}} = 1.9093060930437717 \times 10^{-152} \approx 2^{-504}$ gives σ_{64} correct to full machine precision.

The results for qd were identical to dqd except that the crucial element q_{64} became zero in both steps. Hence qd is not suitable for computation of small singular values with high relative accuracy.

Table 3. Numerical results for Example 1

	After the first pass	After the second pass
q_1	6.553700000000000D+04	6.5537999984741444D+04
q_2	6.5536000015258556D+04	6.5536000061032595D+04
q_3	6.5536000000000233D+04	6.553600000001397D+04
q_4 to q_{63}	6.553600000000000D+04	6.553600000000000D+04
q_{64}	3.6455053829317361D-304	3.6454497569340717D-304
e_1	9.9998474144376459D-01	9.9995422572819948D-01
e_2	9.999999976717291D-01	9.999999883589297D-01
e_3	9.99999999999645D-01	9.99999999997513D-01
e_4 to e_{62}	1.000000000000000D+00	1.000000000000000D+00
e_{63}	1.000000000000000D+00	5.5625997664363648D-309

Example 2. We have rerun Example 1 but with a smaller value of $n(= 5)$ and the results are given in Table 4. For this example, $\sigma_5 = \sqrt{q_5} = 2.3282709094019085 \times 10^{-10}$ which is correct to full machine precision. For comparison, the answer given by the LINPACK SVD routine dsdvc (which is based on the Golub-Reinsch algorithm) is $2.3282704794711363 \times 10^{-10}$ which gets 7 of the 15 digits correct.

Using qd we got almost identical results except that q_5 is zero in both sweeps. Thus, σ_5 is zero according to the qd algorithm. Thus, qd does not deliver as much accuracy as Golub-Reinsch; in fact it can be shown that qd sometimes delivers zero for singular values as large as $\sqrt{macheps} * \|B\|$. \square

Table 4. Numerical results for Example 2

	After the first pass	After the second pass
q_1	6.553700000000000D+04	6.5537999984741449D+04
q_2	6.5536000015258551D+04	6.5536000061032593D+04
q_3	6.5536000000000238D+04	6.553600000001395D+04
q_4	6.553600000000000D+04	6.553600000000000D+04
q_5	5.4209281443662679D-20	5.4208454275671899D-20
e_1	9.9998474144376457D-01	9.9995422572819948D-01
e_2	9.999999976717293D-01	9.999999883589292D-01
e_3	9.99999999999642D-01	9.99999999997509D-01
e_4	1.000000000000000D+00	8.2716799077854419D-25

Some people do not like root free algorithms (e.g. dqd) because they limit the domain of the matrices to which they can be applied. In IEEE conforming computers the exponent range is approximately 2^{-1022} to 2^{1023} in double precision arithmetic. On such machines dqd can diagonalize bidiagonals with condition numbers up to 2^{1022} whereas oqd can deal with condition numbers up to 2^{2045} . For most applications $2^{1022} \approx 10^{308}$ is more than adequate.

On machines such as DEC Vax in D-floating mode the effective limit on the condition number is approximately $2^{127} \approx 10^{38}$ and if this is restrictive then it is only necessary to use oqd to find the smallest singular values and then switch to dqd when the effective condition number (easily approximated) is within range.

We conclude this section by pointing out that qd (Algorithm 3), the standard qd algorithm, consists of the so-called *rhombus rules* arranged in computational form and these rules are a direct consequence of the defining equation

$$BB^T = \hat{B}^T \hat{B}.$$

Equate the (k, k) entry on each side to obtain

$$(10) \quad a_k^2 + b_k^2 = \hat{b}_{k-1}^2 + \hat{a}_k^2, \quad q_k + e_k = \hat{e}_{k-1} + \hat{q}_k.$$

and equate the $(k, k+1)$ entry on each side to obtain

$$(11) \quad b_k a_{k+1} = \hat{a}_k \hat{b}_k, \quad e_k q_{k+1} = \hat{q}_k \hat{e}_k.$$

The rhombus rules can be also derived from $B^T = Q\hat{B}$ by noting that orthogonal transformation Q changes neither the norms nor the inner products of the columns. The reason for the name rhombus rule is indicated in Fig. 4 of Sect. 11.

5. Incorporation of shifts

Rutishauser introduced shifts into the qd almost from the beginning and we could simply quote him. Unfortunately he does not give any explanation of how he derived the appropriate modification of qd as given in Sect. 4. So we provide one at the end of Sect. 5.1. In fact the use of shifts in root finding algorithms can be traced back to Schröder's classic work [37], [38] which was published in 1870.

5.1. Shifted qd algorithms

In eigenvalue calculations, shifts are natural and can be easily incorporated since

$$\lambda(A - \tau^2 I) = \lambda(A) - \tau^2$$

where τ^2 is the shift and $\lambda(A)$ indicates an eigenvalue of A . Thus, by subtracting τ^2 from the diagonals of the matrix, we can introduce origin shifts into the Cholesky algorithm.

A shift τ can be introduced into oqd (Algorithm 1, Sect. 3) by modifying statements involving \hat{a} and \tilde{a} .

Algorithm 4 (oqds)

```

 $\tilde{a} := a_1$ 
for  $k = 1, n - 1$ 
   $\hat{a}_k := \sqrt{\tilde{a}^2 + b_k^2 - \tau^2}$ 
   $\hat{b}_k := b_k * (a_{k+1} / \hat{a}_k)$ 
   $\tilde{a} := \sqrt{\tilde{a}^2 - \tau^2} * (a_{k+1} / \hat{a}_k)$ 
end for
 $\hat{a}_n := \sqrt{\tilde{a}^2 - \tau^2}$ 

```

It may be verified that $\hat{B}^T \hat{B} = BB^T - \tau^2 I$. To keep \hat{B} real the shift must satisfy

$$(12) \quad \tau \leq \sigma_n[B]$$

but this constraint is not formally necessary for dqd (Algorithm 2) which uses

$$\hat{q}_k := d_k + e_k - \tau^2.$$

By defining $d = d_k$ as $\tilde{a}_k^2 - \tau^2$ an addition can be saved.

Algorithm 5 (dqds)

```

d := q1 - τ2
for k := 1, n - 1
  q̂k := d + ek
  êk := ek * (qk+1 / q̂k)
  d := d * (qk+1 / q̂k) - τ2
end for
q̂n := d

```

The constraint (12) is also unnecessary for qd.

Algorithm 6 (qds)

```

ê0 = 0
for k := 1, n - 1
  q̂k := (qk - êk-1) + ek - τ2
  êk := ek * qk+1 / q̂k
end for
q̂n := qn - ên-1 - τ2

```

All that is lacking is an analogue of the orthogonal connection (Theorem 1)

$$B^T = Q\hat{B}.$$

For that it is necessary to abandon square matrices and there are two ways of doing it. In some signal processing circles one would write

$$H \begin{bmatrix} B^T \\ \tau I \end{bmatrix} = \begin{bmatrix} \hat{B} \\ 0 \end{bmatrix}$$

with $H^T H = \text{diag}[I, -I]$ and such an H can be built up from both plane rotations and plane hyperbolic rotations. We prefer to write

$$\begin{bmatrix} B^T \\ O \end{bmatrix} = Q \begin{bmatrix} \hat{B} \\ \tau I \end{bmatrix}$$

where Q is $2n \times 2n$ and orthogonal but not unique. However its first n rows are uniquely determined by B and τ for $\tau \leq \sigma_n[B]$. Moreover Q may be built up by well chosen plane rotations. Both approaches yield oqds; they provide alternative interpretations.

It is at this point that the superiority of the qd formulation becomes clear. DK showed that the standard Golub-Reinsch bidiagonal QR algorithm may be simplified when the shift is zero; see Sect. 10 for the details. Our algorithms (1, 2, or 3) are already simpler than the DK zero shift QR and they also permit use of a non-zero shift with no impediment to pipelined or parallel implementation or high relative accuracy. These natural improvements are strong evidence that our formulation is the natural one.

5.2. The two phase implementation

At first sight the auxiliary quantities d_i , $i = 1, \dots, n$ that occur in dqd are seen as the price to be paid for securing high relative accuracy. On further consideration they may be seen as an attractive feature that permits an aggressive shift strategy that also preserves high relative accuracy in the computed singular values. Moreover, as an extra bonus, we find that the vector $d = (d_1, \dots, d_n)$ may be computed in $O(\log_2 n)$ steps in a parallel computer using the technique called parallel prefix operation in computer science writings, see [9]. The numerical stability of parallel prefix is not yet well understood, but see Mathias [26].

Consider next the implementation of dqds. The auxiliary quantities d_i may be computed prior to any modification of q and e since

$$(13) \quad \begin{aligned} d_{k+1} &= d_k q_{k+1} / \hat{q}_k - \tau^2 \\ &= d_k q_{k+1} / (d_k + e_k) - \tau^2. \end{aligned}$$

An alternative formulation is

$$(14) \quad d_{k+1} = \frac{q_{k+1}}{1 + e_k/d_k} - \tau^2$$

but a division costs more than a multiplication.

It is at this point that one sees the advantage of arithmetic units that conform to the IEEE floating point standard 754: there is no need to test at each instance of (13) or (14) to prevent division by zero. The occurrence of a k with $d_k = \infty$ does no harm since it implies that $\hat{q}_{k-1} = \hat{a}_{k-1} = 0$. This signals that

$$\sigma_n^2[B] < \tau^2$$

and the transformation of B to \hat{B} (Phase 2) should not be completed. The effort in running (13) is not wasted because it yields a new upper bound on $\sigma_n^2[B]$. Furthermore, the number of negative \hat{q}_i gives the number of singular values less than the shift τ . It can be shown that the \hat{q}_i are the pivots obtained in performing Gaussian elimination without pivoting on the symmetric tridiagonal matrix $(BB^T - \tau^2 I)$. We have also used this inertia property to create a spectrum slicing (bisection) algorithm for use on a distributed memory parallel computer system.

Using (13), $d_k = \infty$ yields $d_{k+1} = \infty/\infty = NaN$ (not a number) and then $\hat{q}_i = NaN$ for $i > k+1$. Using (14), $d_k = \infty$ yields $d_{k+1} = q_{k+1} - \tau^2$ which is a better answer.

5.3. Almost positive bidiagonals

Following Rutishauser we observe that there are some conditions in which a shift τ exceeding σ_n is permissible. We go further in showing that under certain conditions the high relative stability property is also preserved.

The standard qd algorithm is well defined for most shifts but it may not be stable in an absolute sense; i.e. the new array $\{\hat{q}, \hat{e}\}$ may be far greater than old one $\{q, e\}$. Rutishauser proved stability under the assumption of positivity and took great care in his implementation to preserve this property.

Our dqds algorithm has the advantage of maintaining relative stability in the positive case and, fortunately, even beyond. For example the requirement

$$\tau^2 < \frac{1}{2}\sigma_{n-1}^2[B_{n-1}] + e_{n-1}$$

is sufficient where B_{n-1} is the leading principal submatrix of B_n because our condition ensures that the only entries in $\{\hat{q}, \hat{e}\}$ that could go negative are \hat{e}_{n-1} and \hat{q}_n . Our goal is to choose τ (actually τ^2) to make \hat{q}_n as small as possible and hence

$$\tau^2 \approx d_n = q_n(1 - e_{n-1}/\hat{q}_{n-1}).$$

Notice how strongly d_n depends on $\text{sign}(e_{n-1})$ and $\text{sign}(q_n)$ since \hat{q}_{n-1} , though unknown, remains positive. We need to know what happens to the signs of \hat{e}_{n-1} and \hat{q}_n once we abandon positivity. It turns out that it takes at least two sweeps to return to the positive case. There are four possible sign configurations in the asymptotic regime ($\tau^2 < \frac{1}{2}d_{n-1} + e_{n-1}$) and we designate them by sign pairs: $(\text{sign}(e_{n-1}), \text{sign}(q_n))$.

A careful study of the last three assignments in dqds shows the following possible paths the sign pairs could follow. These are also shown in Fig. 1.

If $\tau < \sigma_n$

$$\begin{aligned} (+, +) &\longrightarrow (+, +) \\ (+, -) &\longrightarrow (-, +) \\ (-, +) &\longrightarrow (-, -) \\ (-, -) &\longrightarrow (+, +) \end{aligned}$$

If $\tau > \sigma_n$

$$\begin{aligned} (+, +) &\longrightarrow (+, -) \\ (+, -) &\longrightarrow (-, -) \\ (-, +) &\longrightarrow (-, +) \\ (-, -) &\longrightarrow (+, -) \end{aligned}$$

Theorem 2 (Bounds for σ_{\min} without shifts). Apply the *dqd* transformation to a positive bidiagonal B (see Algorithm 1) to produce \hat{B} and $\tilde{a}_1, \tilde{a}_2, \dots, \tilde{a}_n$. Then

1. $\sigma_n \leq \min_k \{\tilde{a}_k\}$
2. $[(BB^T)^{-1}]_{k,k} = \tilde{a}_k^{-2}$
3. $(\sum_{k=1}^n \tilde{a}_k^{-1})^{-1} \leq (\sum_{k=1}^n \tilde{a}_k^{-2})^{-1/2} \leq \sigma_n$.

Proof. Since singular values are invariant under orthogonal transformations and transposition

$$\sigma_n[B] = \sigma_n[Q_k B^T] \leq \|u_k^T Q_k B^T\| = \tilde{a}_k$$

where u_k is the k th column of the identity matrix. The k th row of $Q_k B^T$ is a singleton;

$$u_k^T Q_k B^T = \tilde{a}_k u_k^T.$$

Transposing and rearranging gives

$$\tilde{a}_k^{-1} Q_k u_k = B^{-1} u_k$$

$$\tilde{a}_k^{-2} = (B^{-T} B^{-1})_{k,k}$$

as claimed. Note that

$$\sigma_n^{-2} \leq \sum_{i=1}^n \sigma_i^{-2} = \|B^{-1}\|_F^2 = \text{trace} [(BB^T)^{-1}] = \sum_{i=1}^n \tilde{a}_i^{-2}.$$

Finally we get the required result by considering the one and two norms of the vector $(\tilde{a}_1^{-1}, \tilde{a}_2^{-1}, \dots, \tilde{a}_n^{-1})$. \square

We can compute bounds on $\sigma_{\min}[B]$ even when the algorithm is used with shifts τ provided that $\tau \leq \sigma_{\min}[B]$. Formally the reduction

$$\begin{pmatrix} B^T \\ 0 \end{pmatrix} \rightarrow \begin{pmatrix} \hat{B} \\ \tau I \end{pmatrix}$$

requires $2(n-1)$ plane rotations (not just $n-1$) because the rotation G_i in $(i, i+1)$ must be preceded by a rotation \vec{G}_i in plane $(i, n+i)$ in order to introduce τ into position $(n+i, i)$. Conceptually as B^T is transformed into \hat{B} , at each step one row is in transition and all rows above it belong to \hat{B} and all rows below it belong to B^T . However the active row, row k , is still a singleton in fact

$$(16) \quad u_k^T Q_k \begin{pmatrix} B^T \\ 0 \end{pmatrix} = \tilde{a}_k u_k^T.$$

Theorem 3 (Bounds for σ_{\min} with shifts). If the *dqds* algorithm with shift τ transforms positive bidiagonal B into positive \hat{B} with auxiliary quantities $\tilde{a}_1, \dots, \tilde{a}_n$ then

1. $\sigma_n[\hat{B}] \leq \min_k \{\tilde{a}_k\}$
2. $[(BB^T)^{-1}]_{k,k} = \tilde{a}_k^{-2} \|x_k\|^2 < \tilde{a}_k^{-2}$.

where, in (16), $u_k^T Q_k := (x_k^T, y_k^T)$, and x and y each have n entries.

Proof. Since singular values are invariant under orthogonal transformation and transposition, and are not increased by the transformation $a \leftarrow \sqrt{a^2 - \tau^2}$ we have

$$\sigma_n[\hat{B}] = \sigma_n [\bar{Q}_k B^T] \leq \|u_k^T \bar{Q}_k B^T\| = \|x_k^T B^T\| \leq \tilde{a}_k$$

where \bar{Q}_k denotes the leading $n \times n$ submatrix of Q_k . The last equality uses (16). To establish the second result transpose (16) to obtain

$$\begin{bmatrix} B & 0 \end{bmatrix} Q_k^t u_k = B x_k = u_k \tilde{a}_k.$$

Since B is invertible,

$$\tilde{a}_k^{-1} x_k = B^{-1} u_k,$$

$$\tilde{a}_k^{-2} \|x_k\|^2 = u_k^T B^{-t} B^{-1} u_k = [(B B^T)^{-1}]_{k,k} \quad \square$$

Remark. Since the \tilde{a}_k are monotone decreasing in τ a successful dqds transformation produces a better upper bound and a worse lower bound than does dqd (see Theorem 2). Fortunately it is the upper bound that plays an active role in our implementation.

6.2. The Newton shift

The shift strategy used by Bauer to accelerate the rational QR algorithm RATQR is also closely related to part 3 of the above Theorem 2. See [4], [29].

We recall that the Newton shift from 0 for the characteristic polynomial of any matrix A is related to the trace of the inverse. Let

$$\chi_A(t) = \det[tI - A] = \prod_{i=1}^n (t - \lambda_i).$$

Then, by logarithmic differentiation

$$\frac{\chi'_A(t)}{\chi_A(t)} = \sum_{i=1}^n \frac{1}{t - \lambda_i}.$$

In particular

$$-\frac{\chi'_A(0)}{\chi_A(0)} = \sum_{i=1}^n \lambda_i^{-1} = \text{trace } A^{-1}$$

because the spectrum of A^{-1} is $\{\lambda_i^{-1}\}_1^n$.

In our case by Theorem 2(2), $\sum \tilde{a}_k^{-2} = \sum d_k^{-1} = \text{trace } A^{-1}$ is the Newton correction from 0 towards σ_n^2 .

6.3. The $(1, \infty)$ bound

The DK paper also provides lower bounds on σ_n . Two recurrences (see Sect. 10 for details) produce

$$\min_j \lambda_j = \|B^{-1}\|_{\infty}^{-1}$$

and

$$\min_j \mu_j = \|B^{-1}\|_1^{-1}.$$

Then

$$\sigma_n^{-1} = \|B^{-1}\| \leq \min\{\|B^{-1}\|_{\infty}, \|B^{-1}\|_1\}.$$

Since $\|C\| \leq \sqrt{\|C\|_1 \|C\|_{\infty}}$ for any square matrix C , we can improve the DK bound to give,

$$\sigma_n^{-1} = \|B^{-1}\| \leq \sqrt{\|B^{-1}\|_1 \|B^{-1}\|_{\infty}} \leq \min\{\|B^{-1}\|_{\infty}, \|B^{-1}\|_1\}.$$

6.4. The Johnson bound

For a general complex matrix C , a Gersgorin-type bound for σ_{\min} is given by Johnson (see [23]),

$$\sigma_{\min} \geq \max\{0, \theta\}$$

where

$$\theta = \min_i \left\{ |c_{i,i}| - \frac{1}{2} \sum_{k \neq i} |c_{k,i}| + |c_{i,k}| \right\}.$$

For a positive bidiagonal B , this simplifies to

$$\theta = \min_i \left\{ a_i - \frac{1}{2}(b_i + b_{i-1}) \right\}$$

and near convergence $a_n \rightarrow 0$, $b_{n-1} \rightarrow 0$ so that

$$\theta = a_n - \frac{1}{2}b_{n-1}.$$

7. Effects of finite precision

7.1. Error analysis – overview

One of the benefits of the simplicity of our algorithms oqd and dqd is that their analysis is relatively easy. The DK zero shift QR transformation, though simpler than the Golub/Reinsch transformation, is complicated enough to defy anything but a forward error analysis. After heroic struggles with innumerable details DK establish the error bound quoted in Sect. 10.4.

When discussing this result and our own analyses it is convenient to use the acronym *ulp* which stands for *units in the last place held*. It is the natural way to refer to *relative* differences between numbers. When a result is correctly rounded the

error is not more than half an *ulp*. In this section we usually omit the ubiquitous phrase ‘at most’ qualifying errors and modifications.

Our algorithms still do not admit a pure backward error analysis, the computed output \hat{B} is not the exact output from a matrix very close to B . Nevertheless we can use a hybrid interpretation involving both backward and forward interpretation.

Whereas DK’s zero shift guarantees that each computed singular value is in error by no more than $69n^2$ *ulps* our dqds algorithm causes no more than $4n$ *ulps* change using any properly chosen shift. However the main point is that our analysis is easy to grasp.

The next subsection establishes this strong property of dqds. A similar result holds for oqds but the square roots and squaring provoke a slightly larger bound.

The trick of the proof is to define \vec{B} (see Fig. 2) so that the computed auxiliary quantities $\{d_i\}$ are exact outputs of dqds. The difference between \vec{B} and \hat{B} is the forward error.

At the beginning of the paper we made much of the fact that algorithms oqd and dqd required no subtractions. Yet, in the interest of efficiency, we have introduced shifts and quietly brought back subtraction. The miracle is that the subtraction is in the d ’s and does not impair the high relative accuracy property. However Rutishauser’s qd does not guarantee high relative accuracy so long as q ’s are dominated by neighbouring e ’s.

Since no intermediate quantities exceed σ_1 , it is assumed that the initial data are scaled so that σ_1 (or σ_1^2 for dqds) is close to the overflow threshold. Underflow, though possible, is then a rare event.

Finally we remind the reader that the symbol \circ carries its normal mathematical meaning.

7.2. High relative (mixed) stability in the presence of shifts

We refer the reader to Sect. 5.3 where almost positive bidiagonals are introduced. Rutishauser merges the q ’s and e ’s into a single array Z ;

$$Z := \{q_1, e_1, q_2, e_2, \dots, e_{n-1}, q_n\}$$

and this is a convenient notation for the analysis which follows.

Before stating our claim we need more notation because the difficulty in the analysis is one of interpretation. Given Z the dqds algorithm in finite precision arithmetic produces representable output \hat{Z} . We introduce two ideal arrays \vec{Z} and \check{Z} such that \check{Z} is the output of dqds with shift τ acting on \vec{Z} in *exact arithmetic*. Moreover \vec{Z} is a small relative perturbation of Z and \hat{Z} is a small relative perturbation of \check{Z} . See Fig. 2. This property is called mixed stability in [10] but note that the perturbations are relative ones.

Our model of arithmetic is that the floating point result of a basic arithmetic operation \circ satisfies

$$(17) \quad fl(x \circ y) = (x \circ y)(1 + \eta) = (x \circ y)/(1 + \delta)$$

where η and δ depend on x , y , and \circ , and the arithmetic unit but satisfy

$$|\eta| < \epsilon, \quad |\delta| < \epsilon$$

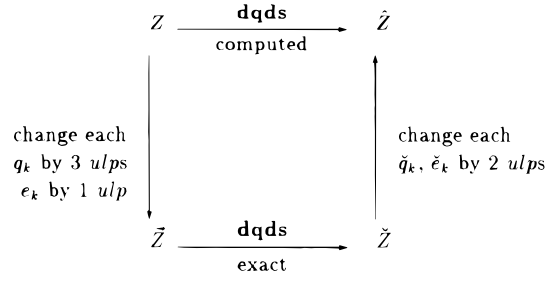


Fig. 2. Effects of roundoff. $Z \equiv B$

for a given ϵ that depends only on the arithmetic unit. We shall choose freely the form (η or δ) that suits the analysis.

A fairly simple result is possible because the only truly sequential part of dqds is the sequence $\{d_i\}_1^n$. Note that, in exact arithmetic

$$d_{k+1} = \frac{d_k q_{k+1}}{d_k + e_k} - \tau^2$$

The trick is to write down the relations governing the computed quantities and then to *discern* among them an exact dqds transform whose input is close to Z and whose output is close to \hat{Z} .

Theorem 4. *In the absence of underflow or overflow, the Z diagram given above commutes and \vec{q}_k (\vec{e}_k) differs from q_k (e_k) by 3 (1) ulps, \hat{q}_k (\hat{e}_k) differs from \check{q}_k (\check{e}_k) by 2 (2) ulps.*

Proof. We write down the exact relations satisfied by the computed quantities \hat{Z} .

$$(18) \quad \hat{q}_k = (d_k + e_k)/(1 + \epsilon_+)$$

$$(19) \quad t_k = q_{k+1}(1 + \epsilon_+)/\hat{q}_k = \frac{q_{k+1}(1 + \epsilon_+)(1 + \epsilon_+)}{d_k + e_k}$$

$$(20) \quad \hat{e}_k = e_k t_k (1 + \epsilon_*) = \frac{e_k q_{k+1} (1 + \epsilon_+)(1 + \epsilon_+)(1 + \epsilon_*)}{d_k + e_k}$$

$$d_{k+1} = \frac{\{d_k t_k (1 + \epsilon_*) - \tau^2\}}{1 + \epsilon_{k+1}}$$

Note the difference between $*$ and $*$. Of course all the ϵ 's obey (17) and depend on k but we have chosen to single out the one that accounts for the subtraction because it is the only one where the k dependence must be made explicit. In more detail the last relation is

$$\begin{aligned}
(1 + \epsilon_{k+1})d_{k+1} &= \frac{d_k q_{k+1}}{d_k + e_k} (1 + \epsilon_+)(1 + \epsilon_+)(1 + \epsilon_*) - \tau^2 \\
(21) \quad &= \frac{(1 + \epsilon_k)d_k q_{k+1} (1 + \epsilon_+)(1 + \epsilon_+)(1 + \epsilon_*)}{(1 + \epsilon_k)d_k + (1 + \epsilon_k)e_k} - \tau^2
\end{aligned}$$

This tells us how to define \vec{Z} . Note that ϵ_k arose in the previous step. Moreover

$$(22) \quad (1 + \epsilon_1)d_1 = q_1 - \tau^2$$

Our choice of \vec{Z} , in general, is not a machine representable array.

For $k \geq 1$,

$$(23) \quad \begin{aligned} \vec{d}_k &:= (1 + \epsilon_k)d_k \\ \vec{e}_k &:= (1 + \epsilon_k)e_k \end{aligned}$$

$$(24) \quad \vec{q}_{k+1} := q_{k+1}(1 + \epsilon_/) (1 + \epsilon_+) (1 + \epsilon_*) \quad , \quad (\vec{q}_1 = q_1)$$

and by (21),

$$(25) \quad \vec{d}_{k+1} = \frac{\vec{d}_k \vec{q}_{k+1}}{\vec{d}_k + \vec{e}_k} - \tau^2$$

Then, for exact dqds, we must define

$$\begin{aligned} \check{q}_k &:= \vec{d}_k + \vec{e}_k = (1 + \epsilon_k)(d_k + e_k) \\ \check{e}_k &:= \frac{\vec{e}_k \vec{q}_{k+1}}{\vec{d}_k + \vec{e}_k} . \end{aligned}$$

Finally \hat{q}_k and \hat{e}_k must be recast in terms of \check{Z} ;

$$(26) \quad \hat{q}_k = \check{q}_k / (1 + \epsilon_k)(1 + \epsilon_+) \quad , \quad \text{from (18)}$$

$$(27) \quad \begin{aligned} \hat{e}_k &= \frac{e_k q_{k+1}}{d_k + e_k} (1 + \epsilon_/) (1 + \epsilon_+) (1 + \epsilon_*) \quad , \quad \text{from (20)} \\ &= \frac{\vec{e}_k \vec{q}_{k+1}}{(\vec{d}_k + \vec{e}_k)} \frac{(1 + \epsilon_*)}{(1 + \epsilon_*)} . \end{aligned}$$

It is (23) that yields $\vec{e}_k / (\vec{d}_k + \vec{e}_k) = e_k / (d_k + e_k)$. Equations (23) and (24) give the change from \check{Z} to \hat{Z} , and equation (25) fixes the exact dqds transform. \square

Recall that, in exact arithmetic, algorithm dqds diminishes all eigenvalues (of LR) by the shift. For finite precision execution we have the following.

Corollary 1. *Algorithm dqds preserves high relative stability. When Z and $(a_i = \sqrt{q_i}$, $b_i = \sqrt{e_i}$, etc.), together with the associated ideal bidiagonals \vec{B} and \check{B} , satisfy*

$$\begin{aligned} \sigma_i[\vec{B}] &= \sigma_i[B] \exp \{2(n-1)\epsilon_1^{(i)}\}, \\ \sigma_i^2[\check{B}] &= \sigma_i^2[\vec{B}] - \tau^2, \\ \sigma_i[\hat{B}] &= \sigma_i[\check{B}] \exp \{(2n-1)\epsilon_2^{(i)}\}, \end{aligned}$$

for $i = 1, 2, \dots, n$, and $\epsilon_1^{(i)} \leq \epsilon$, $\epsilon_2^{(i)} \leq \epsilon$.

Proof. For $i = 1, \dots, n - 1$

$$\begin{aligned}\vec{a}_{i+1} &= \sqrt{\vec{q}_{i+1}} = a_{i+1} \sqrt{(1 + \epsilon_{\prime})(1 + \epsilon_{+})(1 + \epsilon_{*})} \\ \vec{b}_i &= \sqrt{\vec{e}_i} = b_i \sqrt{(1 + \epsilon_i)}.\end{aligned}$$

By Theorem 2 in DK, the relative change in any singular value in going from B to \vec{B} is the product of all the relative changes, namely

$$\prod_{i=1}^{n-1} [(1 + \epsilon_{\prime})(1 + \epsilon_{+})(1 + \epsilon_{*})(1 + \epsilon_i)]^{\frac{1}{2}} \leq \exp 2(n - 1)\epsilon.$$

Similarly

$$\begin{aligned}\hat{a}_i &= \sqrt{\hat{q}_i} = \check{a}_i / \sqrt{(1 + \epsilon_i)(1 + \epsilon_{+})}, \quad i < n \\ \hat{b}_i &= \sqrt{\hat{e}_i} = \check{b}_i \sqrt{(1 + \epsilon_{*})(1 + \epsilon_i)}, \quad i < n \\ \hat{a}_n &= \sqrt{\hat{d}_n} = \check{a}_n / \sqrt{(1 + \epsilon_n)}.\end{aligned}$$

The relative change in any singular value in the transformation from \check{B} to \hat{B} is bounded by

$$\sqrt{1 + \epsilon_n} \prod_{i=1}^{n-1} [(1 + \epsilon_{*}) / (1 + \epsilon_i)(1 + \epsilon_{+})]^{\frac{1}{2}} \leq \exp(4n - 3)\epsilon/2.$$

Since the passage from \vec{B} to \check{B} is exact the singular values diminish by τ^2 . \square

Remark. It can be shown by similar means that one dqd transformation cannot alter any singular value by more than $3(n - 1)$ ulps.

Theorem 4 is much stronger than the familiar error analysis based on norms because:

1. The perturbed matrices considered here inherit the bidiagonal structure
2. The bounds are very much smaller than those from DK (see Sect. 10) or the Golub-Reinsch algorithm (see Chapter 8 of [19]).

For multiple sweeps of dqds, our results can be stated more simply in terms of the positive sequence $\{\check{Z}_l\}$ where l denotes the sweep with $\check{Z}_1 = Z_1 = Z$. See Fig. 3 for the corresponding commutative diagram. Then by combining $\check{Z}_{l+1} \rightarrow Z_{l+1} \rightarrow \vec{Z}_{l+1} \rightarrow \check{Z}_{l+2}$ one obtains that \check{Z}_{l+2} is the exact dqds transformation of a perturbed (in relative sense) \check{Z}_{l+1} . Thus backward stability is present for $\{\check{Z}_{l+2}\}$.

Similarly it can be shown that the sequence $\{\vec{Z}_l\}$ is forward stable (in relative sense) with $\vec{Z}_f = Z_f$ where Z_f is the final computed result. On exact application of dqds on \vec{Z}_l we get \check{Z}_{l+1} instead of \vec{Z}_{l+1} (see Fig. 3) and the error between \vec{Z}_{l+1} and \check{Z}_{l+1} is small.

Example 3. The following experiment shows vividly the difference between an algorithm that obtains high relative accuracy (dqds) and one that does not (LINPACK

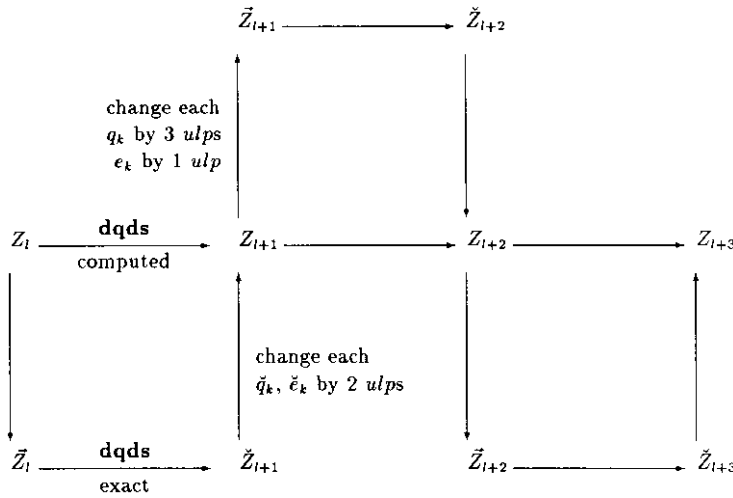


Fig. 3. Effects of roundoff for multiple sweeps

dsvdc based on the Golub-Reinsch algorithm) but which delivers excellent absolute accuracy. We took the graded matrix B_+ ,

$$a_{i-1} = \beta a_i \quad , \quad b_i = a_i$$

with $a_n = 1$, $n = 8$ and $\beta = 60$. We applied both algorithms to B_+ and its reversal B_- ,

$$a_i \rightarrow a_{n+1-i} \quad , \quad b_i \rightarrow b_{n-i}.$$

We did not allow any flipping of the matrix within the dqds algorithm although such flipping improves convergence. See next section.

In Tables 5 and 6, the third column shows,

$$abs_i := (\sigma_i[B_-] - \sigma_i[B_+]) / \sigma_1[B_+]$$

the differences between outputs scaled by the 2-norm of the nicer matrix. Recall that $macheps \approx 2^{-53} \approx 1.1 \times 10^{-16}$. For dsvdc (see Table 5), it can be seen that the absolute error is even smaller than absolute stability guarantees.

In Tables 5 and 6, the fourth column shows, rel_i ,

$$rel_i := (\sigma_i[B_-] - \sigma_i[B_+]) / \sigma_i[B_+]$$

the relative differences in the outputs. For dqds the largest magnitude is less than two $macheps$ (see Table 5) while for dsvdc (see Table 6), it is very much larger and shows that dsvdc does not give relative accuracy.

8. Convergence

8.1. Linear convergence

Convergence of the Cholesky algorithm and of the standard qd algorithm for tridiagonal matrices (in the positive case) have been given by Rutishauser and others (see

Table 5. Numerical results using dqds for Example 3

i	$\sigma_i[B_+]$	$\sigma_i[B_-]$	abs_i	rel_i
1	3.9590303657774160D+12	3.9590303657774155D+12	-1.2D-16	-1.2D- 16
2	5.7143240472800255D+10	5.7143240472800247D+10	-1.9D-18	-1.3D-1 6
3	8.9790986853271568D+08	8.9790986853271568D+08	0.0D+00	0.0D+0 0
4	1.4489876544914651D+07	1.4489876544914651D+07	0.0D+00	0.0D+0 0
5	2.3661793507020348D+05	2.3661793507020348D+05	0.0D+00	0.0D+0 0
6	3.8884661685208386D+03	3.8884661685208386D+03	-1.1D-25	-1.2D-1 6
7	6.4142972113704085D+01	6.4142972113704109D+01	3.6D-27	2.2D-1 6
8	3.5351579203702068D-01	3.5351579203702068D-01	0.0D+00	0.0D+0 0

Table 6. Numerical results using dsdvc for Example 3

i	$\sigma_i[B_+]$	$\sigma_i[B_-]$	abs_i	rel_i
1	3.9590303657774146D+12	3.9590303657774160D+12	3.7D-16	3.7D-1 6
2	5.7143240472800224D+10	5.7143240472800278D+10	1.3D-17	9.3D-1 6
3	8.9790986853271544D+08	8.9790986853271413D+08	-3.3D-19	-1.5D-1 5
4	1.4489876544914653D+07	1.4489876544914989D+07	8.5D-20	2.3D-1 4
5	2.3661793507020345D+05	2.3661793507022545D+05	5.6D-21	9.3D-1 4
6	3.8884661685208386D+03	3.8884661685173243D+03	-8.9D-22	-9.0D-1 3
7	6.4142972113704073D+01	6.4142972113772929D+01	1.7D-23	1.1D-1 2
8	3.5351579203702068D-01	3.5351579205582154D-01	4.7D-24	5.3D-1 1

[32], [39]). Nevertheless we present here a direct convergence proof for oqd because it is a new algorithm and the proof is both short and illuminating. We give a brief discussion of the effects of finite precision on the results of the theorem at the end of the section.

In finite precision a computed cosine c_k may equal unity even though the corresponding sine s_k may merely be small.

Our proof makes use of the fact that a symmetric tridiagonal matrix with nonzero superdiagonal entries (e.g. $B^T B$) cannot have multiple eigenvalues [28], [39]. In other words, the singular values of a positive bidiagonal B are distinct

$$(28) \quad \sigma_1 > \sigma_2 > \dots > \sigma_n.$$

One must bear in mind that distinct σ_i may be equal to working precision.

Theorem 5 (Convergence of oqd). *From a positive bidiagonal B_1 the unshifted oqd algorithm in exact arithmetic produces a sequences $\{B_l\}_1^\infty$ of orthogonally equivalent bidiagonals. As $l \rightarrow \infty$,*

$$B_l \rightarrow \Sigma = \text{diag}(\sigma_1, \dots, \sigma_n).$$

Furthermore, if $a_n \leq a_1$, then the sequence $\{\prod_{i=1}^{n-1} b_i^{(l)}\}$ is monotone decreasing in l from the beginning. Each $\{b_i^{(l)}\}_{l=1}^\infty$ converges linearly to 0 with convergence factor σ_{i+1}/σ_i .

Proof. Consider a typical step of oqd (Algorithm 1). Since there are no subtractions each B_l is positive since B_1 is positive.

Equation (7) can be written in the form,

$$(29) \quad \hat{a}_k = c_{k-1} a_k / c_k \quad \text{for } k = 1, \dots, n$$

provided that we set $c_0 = c_n = 1$. Take the product of the first k instances of (29) to find

$$(30) \quad \prod_{i=1}^k \hat{a}_i = \left(\prod_{i=1}^k a_i \right) / c_k \geq \prod_{i=1}^k a_i.$$

The sequence $\{\prod_{i=1}^k a_i^{(l)}\}$ is bounded above, by $\|B\|^k$, and is monotone increasing by (30) and thus convergent. The limit may be written $\prod_{i=1}^k \mu_i$ to reveal that, as $l \rightarrow \infty$,

$$(31) \quad a_i^{(l)} \rightarrow \mu_i$$

$$(32) \quad c_i^{(l)} \rightarrow 1$$

$$(33) \quad s_i^{(l)} \rightarrow 0.$$

By (8) and (33), as $l \rightarrow \infty$,

$$(34) \quad b_i^{(l+1)} = s_i^{(l)} a_{i+1}^{(l)} \rightarrow 0, \quad i = 1, \dots, n-1$$

Thus B_l converges to diagonal form and each μ_i is a singular value. To identify μ_i use the product rhombus rule to find

$$(35) \quad b_k^{(l)} / b_k^{(l-1)} = a_{k+1}^{(l-1)} / a_k^{(l)} \rightarrow \mu_{k+1} / \mu_k.$$

Since $\{b_k^{(l)}\}$ is bounded (by $\|B_1\|$) the limit in (35) cannot exceed 1. By (28) the limit is not 1 and thus $\mu_{k+1} < \mu_k$, $k = 1, \dots, n-1$, and we may identify μ_k as σ_k . Thus, (35) proves that $b_k^{(l)} \rightarrow 0$ linearly with convergence factor σ_{k+1} / σ_k , as claimed, and $B_l \rightarrow \Sigma$.

Finally consider $\prod_{i=1}^{n-1} b_i^{(l)}$. Apply (8) and (6) in turn to find

$$\begin{aligned} \hat{b}_1 \dots \hat{b}_{n-1} &= b_1 \dots b_{n-1} \frac{a_2 \dots a_n}{\hat{a}_1 \dots \hat{a}_{n-1}} \\ &= b_1 \dots b_{n-1} \frac{c_1}{a_1} \frac{c_2}{c_1} \frac{c_3}{c_2} \dots \frac{c_{n-1}}{c_{n-2}} a_n \\ &= b_1 \dots b_{n-1} \left(\frac{a_n}{a_1} \right) c_{n-1} \\ &< b_1 \dots b_{n-1}, \end{aligned}$$

if $a_n \leq a_1$. Since B may be flipped about its antidiagonal without altering the singular values there is no loss of generality in assuming that $a_n \leq a_1$. In this case $\prod_{i=1}^{n-1} b_i^{(l)}$ is monotone decreasing with l from the start. \square

It is worth mentioning that Rutishauser's proof of convergence for qd (Algorithm 3), is based on the observation that $\|\hat{B}_k\|_F > \|B_k\|_F$ and $\|B_{-k}\|_F < \|\hat{B}_{-k}\|_F$, $k = 1, 2, \dots, n$. Here $M_{\{\pm k\}}$ is the submatrix of M containing the first (last) k columns.

By taking the product of the final $n-k$ instances of (29) one finds that $\{\prod_{i=k}^n a_i^{(l)}\}$ is monotone decreasing in l for $k = n-1, n-2, \dots, 2, 1$. In particular $a_1^{(l)}$ increases and $a_n^{(l)}$ decreases so that a flipping of B is needed at most once.

The assertions of Theorem 5 bear up quite well in finite precision arithmetic. The computed sequence $\{\prod_{i=1}^k a_i^{(l)}\}$ is monotone nondecreasing and so $c_i^{(l)} \rightarrow 1$ as $l \rightarrow \infty$. Thus in considering a sequence of computed trigonometric values we do not wish to infer $s_k^{(l)} \rightarrow 0$ from $c_k^{(l)} \rightarrow 1$. So the first casualty is the conclusion that $s_i^{(l)} \rightarrow 0$.

Instead we find that $b_i^{(l)}$ becomes negligible relative to $a_{i+1}^{(l)}$ and $a_i^{(l)}$. Even so, in the absence of underflow, the diagonal entries eventually rearrange themselves in (almost) monotone nonincreasing order. Though distinct, by (28), some singular values may be equal to working accuracy and diagonal monotonicity may actually fail by one or two *ulps* (units in the last place held) because the ratio, though exceeding 1 might be too small to cause the neighbouring b -value to grow at all. All in all the practical oqd performs as closely to exact oqd as it is reasonable to expect.

8.2. Quadratic convergence

Consider the last few steps in dqds with shift τ :

$$\begin{aligned}\hat{q}_{n-1} &= d_{n-1} + e_{n-1} \\ \hat{e}_{n-1} &= e_{n-1}q_n / \hat{q}_{n-1} \\ d_n &= d_{n-1}q_n / \hat{q}_{n-1} - \tau^2 \\ \hat{q}_n &= d_n.\end{aligned}$$

Hence

$$\begin{aligned}\hat{e}_{n-1}\hat{q}_n &= \frac{e_{n-1}q_n}{\hat{q}_{n-1}} \left[q_n \left(1 - \frac{e_{n-1}}{\hat{q}_{n-1}} \right) - \tau^2 \right], \\ (36) \quad &= \frac{e_{n-1}q_n}{\hat{q}_{n-1}} \left[q_n - \tau^2 - \frac{q_n e_{n-1}}{\hat{q}_{n-1}} \right]\end{aligned}$$

In exact arithmetic, as $\tau \rightarrow \sigma_n[B]$ we have $q_n \rightarrow 0$, $e_{n-1} \rightarrow 0$, $q_{n-1} \rightarrow \sigma_{n-1}^2[B] - \sigma_n^2[B] := \text{gap} > 0$ because the singular values are distinct if the initial B is of positive type. Thus convergence will be quadratic with respect to this *gap*.

Expression (36) shows that if

$$(37) \quad 0 \leq q_n - \tau^2 \leq 2 \frac{q_n e_{n-1}}{\hat{q}_{n-1}}$$

then

$$\left| q_n - \tau^2 - \frac{q_n e_{n-1}}{\hat{q}_{n-1}} \right| \leq \frac{q_n e_{n-1}}{\hat{q}_{n-1}}$$

and so, by (36),

$$\frac{|\hat{e}_{n-1}\hat{q}_n|}{(e_{n-1}q_n)^2} \leq \frac{1}{\hat{q}_{n-1}^2} \rightarrow \frac{1}{\text{gap}^2} > 0,$$

as $\tau \rightarrow \sigma_n[B]$. Thus (37) shows a (theoretical) interval for those τ that deliver quadratic convergence. Next we seek a computable expression that will ultimately lie in that interval.

The perfect shift is

$$\tau^2 = q_n \left(1 - \frac{e_{n-1}}{\hat{q}_{n-1}} \right)$$

so that the natural strategy is to estimate \hat{q}_{n-1} by assuming that $\hat{q}_{n-k} = q_{n-k}$ for some modest value of k , like 6, and then evaluate the approximate d 's from $n-k$ to n .

$$\begin{aligned}\hat{q}_{n-1} &= \left(1 - \frac{e_{n-2}}{\hat{q}_{n-2}}\right)q_{n-1} - \tau^2 - e_{n-1} \\ &\approx \left(1 - \frac{e_{n-2}}{q_{n-2}}\right)q_{n-1},\end{aligned}$$

when e_{n-2} and e_{n-1} are small enough.

8.3. Cubic convergence

The assertion in (37) of Sect. 8.2 that the shift $\tau^2 = q_n$ yields quadratic convergence for qds and dqds appears to contradict the result of Rutishauser [33] that this choice yields cubic convergence. See also Rutishauser and Schwarz [36] and Chapter 8 of Wilkinson [43]. Actually, there is no anomaly because the shift strategies are not quite the same. In our terminology what Rutishauser suggests is that the qd transform with $\tau^2 = q_n$ should not be formed explicitly. The only item wanted from it is \hat{q}_n and it is assumed to be the only negative q_i . Then it is shown that $q_n + \hat{q}_n$ is a fourth order approximation to σ_n^2 from below. A qd transform of $\{q, e\}$ with shift $\tau^2 = q_n + \hat{q}_n$ will yield cubic convergence.

The point that is stressed by neither Rutishauser nor Wilkinson is that the computation of $\hat{q}_n = d_n$ costs $O(n)$ operations, very close to 1 step of qds. From another perspective Rutishauser's analysis is a disguised derivation of the cubic convergence of the tridiagonal QR algorithm with the Rayleigh quotient shift.

For our algorithm dqds Rutishauser's (late failure) shift strategy described above is more appealing. Only the first phase of our algorithm is needed to compute d_n and the cost is about $\frac{2}{3}$ of a dqds step. Moreover the positive form is preserved a little longer.

For the sake of completeness we indicate why $q_n + d_n$ is a fourth order lower bound when \hat{q}_n is second order in $q_n e_{n-1}$. The relevant tridiagonal matrix is BB^T and its leading principal $(n-1) \times (n-1)$ submatrix is called V . Recall that u_j is column j of I .

Fact 1: Provided $(V - \sigma_n^2 I)$ is invertible the singularity of $BB^T - \sigma_n^2 I$ yields

$$q_n - \sigma_n^2 = (q_n e_{n-1}) u_{n-1}^T (V - \sigma_n^2 I)^{-1} u_{n-1}$$

Fact 2: With shift $\tau^2 = q_n$,

$$d_n = \hat{q}_n = -(q_n e_{n-1}) u_{n-1}^T (V - q_n I)^{-1} u_{n-1}$$

Conclusion,

$$q_n + \hat{q}_n = \sigma_n^2 + (q_n e_{n-1}) u_{n-1}^T [(V - \sigma_n^2 I)^{-1} - (V - q_n I)^{-1}] u_{n-1}.$$

By Hilbert's first resolvent identity (see p. 90 of [8]),

$$q_n + \hat{q}_n = \sigma_n^2 + (q_n e_{n-1}) (\sigma_n^2 - q_n) u_{n-1}^T (V - q_n I)^{-1} (V - \sigma_n^2 I)^{-1} u_{n-1}.$$

Using Fact 1 again, $q_n + \hat{q}_n =$

$$(38) \quad \sigma_n^2 - (q_n e_{n-1})^2 \{ u_{n-1}^T (V - \sigma_n^2 I)^{-1} u_{n-1} u_{n-1}^T (V - q_n I)^{-1} (V - \sigma_n^2 I)^{-1} u_{n-1} \}.$$

The gap conditions ensure that $\{ \cdot \} \rightarrow O(1/gap^3)$ as $q_n \rightarrow 0$, $e_{n-1} \rightarrow 0$.

In contrast to Rutishauser and Wilkinson, our expression (38) is exact provided that the inverse matrices exist. However it is only for small enough e_{n-1} and q_n that the quantity in $\{ \cdot \}$ can be bounded away from ∞ .

9. A preliminary implementation

9.1. Choice of shifts

The standard singular value codes in LINPACK and LAPACK need about 2 QR steps per singular value, in most cases, and that provides a hard target to beat. Moreover each of our qd transformations needs only $O(n)$ flops and no square roots so we are reluctant to spend $O(n)$ flops on shift selection.

A strategy used to generate the numerical results in this paper may not be the best but it is based on the following somewhat surprising observation. The upper bound

$$\text{sup} = d_k = \min_i d_i$$

is an increasingly good estimate for $\sigma_{\min}^2[\hat{B}]$. Our code maintains bounds on $\sigma_{\min}[B]$ at all times. Following Rutishauser, we call them *sup* and *inf*. Moreover, at each step, we will know the index k from the previous step. This index points to the largest diagonal entry of $(B^T B)^{-1}$ and helps tell us whether σ_{\min}^2 has yet migrated to the bottom of $B^T B$. If $k \geq n - 1$ we expect the trailing 2×2 principal submatrix of $B B^T$ to give a good approximation to σ_{\min}^2 . When $k < n - 1$ the matrix is not yet in asymptotic form and the situation is more difficult. However we do know that $\sqrt{q_k + e_k}$ is the smallest (leftmost) center of all the Johnson discs for the matrix of equation (15). If this disc were separated from the rest of the discs then it would certainly contain σ_{\min} . Even if it is not isolated this disc may still contain σ_{\min} so we use it.

9.2. Splitting and deflation

In Sect. 2 it was noted that if $e_i = 0$ (i.e. $b_i = 0$), for $i < n$, then the bidiagonal B splits into two complementary submatrices. Consider now the case when e_i (or b_i) is small enough to permit such a splitting without making a relative change in any singular value exceeding a given tolerance η . Our situation is a little more complicated than the one studied in DK because of the non-restoring shift. Let σ^2 denote the cumulative sum of all shifts used on the given matrix in the dqds algorithm (which computes the squared singular values).

Our criterion (39) is based on Weyl's monotonicity theorem for eigenvalues of symmetric matrices. Consider the zero diagonal representation of B ,

$$T = \text{tridiag} \left\{ \begin{array}{cccccccc} & a_1 & & b_1 & & a_2 & & \cdot & & a_n & & \\ 0 & & 0 & & 0 & & 0 & & 0 & & 0 & \\ & a_1 & & b_1 & & a_2 & & \cdot & & a_n & & \end{array} \right\}.$$

The spectral norm of a matrix of the form

$$\begin{pmatrix} 0 & b_i \\ b_i & 0 \end{pmatrix}$$

is b_i and it is submatrices of this form that we may remove. Weyl's theorem states that this removal changes no eigenvalues of T by more than b_i . The output of the algorithm is numbers of the form $\sqrt{\sigma^2 + \lambda_i^2[T]}$ and hence for a relative tolerance η we must require

$$b_i < \sqrt{\sigma^2 + \lambda_i^2 [T]}$$

and this is guaranteed when

$$(39) \quad b_i < \eta \sqrt{\sigma^2 + \text{inf}}$$

where inf is our current lower bound on $\sigma_{\min}^2 [B]$.

The suppression of b_k leaves two tridiagonals which correspond to bidiagonal matrices which may be processed separately. There are less stringent checks than (39) for splitting but they require more computational effort.

9.3. Performance of a prototype implementation

We have developed and implemented dqds in FORTRAN 77 to study and exploit the theory we have developed in this paper. This prototype program is built in modular fashion.

We have run our code on a broad test bed of bidiagonals. Here we report on comparisons on three interesting classes using our dqds and LINPACK's dsydc (with reduction to bidiagonals removed). There is little difference in timing between the Demmel-Kahan code and dsydc since DK reverts to QR after tiny singular values are found.

Examples 1,2 and 3 were given in earlier sections.

Example 4 (nice matrices). We considered the graded matrix B_+ defined earlier in Example 3 with the parameter $\beta = 2$ and $n = 30$. Table 7 gives the performance on this example and other examples in this section. The speedup is the ratio of dsydc time to dqds time.

We have also tested this problem with $n = 40$ and in that case the LINPACK dsydc returned with an error flag as it could not compute σ_{40} within 30 iterations. We were prevented from comparing with larger values of n because dsydc reported errors.

Example 5 (perversely graded). To make conditions artificially difficult for dqds, we also ran the programs with the reversely graded matrix B_- as the input with $n = 30$ and $\beta = 2$. See Table 7 for details.

Usually, our dqds will flip B_- to obtain B_+ . If the user does not flip B_- before calling dsydc then the time ratio goes up to 18.9.

dsydc also failed to converge for many combinations of β and n .

Example 6. Let B_w be the Wilkinson-type bidiagonal matrix where

$$a_i = \left| i - \frac{n}{2} + 1 \right|, \quad i = 1, \dots, n$$

$$b_i = 1, \quad i = 1, \dots, n-1$$

This matrix has close singular values (twins) and our current coding does not fully exploit this structure. Hence the low performance compared with the previous results.

Example 7. Doubled Wilkinson-type matrices, B_{2w} ,

$$a_i = \left| i - \frac{n}{4} + 1 \right|, \quad i = 1, \dots, \frac{n}{2}$$

Table 7. Performance comparison

Example	Matrix	n	dqd sweeps	dsvdc sweeps	speedup
4	B_+	30	52	60	10.2
5	B_-	30	79	101	12.3
6	B_w	21	78	62	4.8
7	B_{2w}	41	230	120	4.8
8	B_t	100	374	308	11.0

$$a_{i+\frac{n}{2}} = a_i \quad , \quad i = \frac{n}{2} + 1, \dots, n$$

$$b_i = 1 \quad , \quad i = 1, \dots, n-1$$

with $n = 41$. This matrix has close singular values (quads) and some of them are exactly equal. Table 7 gives the details.

Example 8. Toeplitz matrix B_t ,

$$a_i = 1 \quad , \quad b_i = 2.$$

For $n = 100$, the matrix has a tiny singular value; others are between 1 and 3.

10. The Demmel/Kahan paper

We summarize the highlights of [11].

10.1. High relative accuracy

Corollary 2 of Theorem 2 of [11]. Suppose that $(B + \delta B)_{i,i} = \alpha_{2i-1} a_i$, $(B + \delta B)_{i,i+1} = \alpha_{2i} b_i$, $\alpha_i \neq 0$. Define

$$\bar{\alpha} := \prod_{i=1}^{2n-1} \max\{|\alpha_i|, |\alpha_i^{-1}|\}.$$

Let $\sigma'_1 \geq \sigma'_2 \geq \dots \geq \sigma'_n$ be the singular values of $B + \delta B$. Then

$$\sigma_i / \bar{\alpha} \leq \sigma'_i \leq \sigma_i \bar{\alpha} \quad , \quad i = 1, 2, \dots, n.$$

This shows that bidiagonal matrices determine their singular values to high relative accuracy.

10.2. Bounds for σ_n

It is possible to compute $\|B^{-1}\|_\infty$ and $\|B^{-1}\|_1$ using $2(n-1)$ divisions and multiplications. The algorithm is

$$\lambda_j := a_j(\lambda_{j+1}/(\lambda_{j+1} + b_j)) \quad , \quad j = n-1, n-2, \dots, 1 \quad \text{with} \quad \lambda_n := a_n$$

$$\mu_{j+1} := a_{j+1}(\mu_j/(\mu_j + b_j)) \quad , \quad j = 1, 2, \dots, n-1 \quad \text{with} \quad \mu_1 := a_1$$

$$\|B^{-1}\|_\infty = 1 / \min_j \lambda_j$$

$$\|B^{-1}\|_1 = 1 / \min_j \mu_j.$$

Finally,

$$n^{-1/2} \max\{\|B^{-1}\|_\infty^{-1}, \|B^{-1}\|_1^{-1}\} \leq \sigma_n \leq n^{1/2} \min\{\|B^{-1}\|_\infty^{-1}, \|B^{-1}\|_1^{-1}\}$$

and

$$\min\{\|B^{-1}\|_\infty^{-1}, \|B^{-1}\|_1^{-1}\} \leq \sigma_n$$

10.3. A stopping criterion

Let $\eta \ll 1$ be the desired relative accuracy of the computed singular values. Then if either

$$b_j / \lambda_{j+1} \leq \eta \quad \text{or} \quad b_j / \mu_j \leq \eta$$

set b_j to zero and the two pieces into which B splits may be processed separately. The criteria used in LINPACK [12] can sometimes deliver a zero singular value when it should not and can sometimes fail to suppress a negligible off diagonal entry b_j .

10.4. Bidiagonal QR with zero shift

The standard Golub/Reinsch algorithm [19], [18] used in LINPACK may be simplified when no shifts are used. Of more importance is the fact that in this case all round off errors arise multiplicatively. Moreover for the calculation of tiny singular values zero is a good shift and it pays to compute them first rather than letting the standard shift strategy dictate the order in which the singular values are found. The arithmetic effort in the innermost loop is

Golub/Reinsch: 2 calls to ROT + 12 multiplications + 4 additions

Demmel/Kahan: 2 calls to ROT + 4 multiplications.

The procedure ROT computes the sine and cosine needed for a plane rotation using 2 divisions, 3 multiplications, and 1 square root. Here is the algorithm.

```

oldcs := cs := 1
for i := 1, n - 1
  call ROT(a_i * cs, b_i, cs, sn, r)
  if (i ≠ 1) b_{i-1} := oldsn * r
  call ROT(oldcs * r, a_{i+1} * sn, oldcs, oldsn, a_i)
end for
(40) h := a_n * cs; b_{n-1} := h * oldsn; a_n := h * oldcs

```

In the absence of underflow the error bound on singular values after one zero shift bidiagonal QR transform is

$$|\sigma_i - \sigma'_i| \leq \frac{w}{1-w} \sigma_i \quad i = 1, \dots, n$$

where

$$w := 69n^2\epsilon < 1.$$

See Theorem 6, p. 906, of [11].

10.5. The overall algorithm

```

if (roundoff in  $\sigma_1$  exceeds  $tol * bound$  on  $\sigma_n$ ) then
  use zero-shift QR or QL
else
  use shifted QR or QL
end if.

```

10.6. Other improvements

The new code uses either QL or QR as appropriate according to the way B is graded.

An efficient accurate subroutine is provided to return the singular values and vectors of 2×2 bidiagonal matrices.

Deflation when a diagonal entry a_i vanishes is automatic and occurs either at the bottom or top of B .

11. Evolution of qd

Some of the available presentations of the qd-algorithm, see [32], [39], [21] show its close connection with factorization of tridiagonal matrices but some do not [22], [40]. Nevertheless its discovery had nothing to do with matrix decompositions and a knowledge of the origins helps us to understand the somewhat neglected status of the algorithm. In the next few paragraphs we sketch an earlier paper [27] which described the gradual evolution of the qd-algorithm.

The story begins with Daniel Bernoulli in 1728 when he showed that the largest and the smallest roots of an n th order polynomial can be obtained by iterating an n th order difference equation. See [5]. The work of Bernoulli was extended by Euler in 1748. See Chapter 17 of [14] (English translation [13]).

We are given a rational function of a complex variable z ,

$$f(z) = \sum_{k=0}^{\infty} h_k z^k ,$$

assumed to be regular (analytic) at both $z = 0$ and $z = \infty$. The Taylor series converges to $f(z)$ only within a circle (in \mathcal{E}) centered at $z = 0$ and extending up to the nearest pole p_1 . However, by analytic continuation, the Taylor coefficients $\{h_k\}$ actually define a unique rational function f on all of \mathcal{E} except the poles p_1, p_2, p_3, \dots . The problem is to determine the poles directly from the $\{h_k\}$ without having recourse to analytic continuation.

In 1884 König [25] showed that if p_1 is a simple pole and smaller than all the others then

$$\lim_{k \rightarrow \infty} (h_k / h_{k+1}) = p_1 .$$

Exactly one hundred years ago (i.e. in 1892), in his dissertation, Hadamard [20] showed that

$$\lim_{k \rightarrow \infty} \left(\frac{H_m^{k+1}}{H_m^k} \right) = \prod_{i=1}^m p_i$$

where

$$H_m^k = \det \begin{bmatrix} h_k & h_{k+1} & \cdots & h_{k+m-1} \\ h_{k+1} & h_{k+2} & \cdots & h_{k+m} \\ & & \cdots & \\ h_{k+m-1} & h_{k+m} & \cdots & h_{k+2m-2} \end{bmatrix}$$

The H_m^k are now called Hankel determinants but Hadamard did not give them a name. It follows that

$$p_m = \lim_{k \rightarrow \infty} \left(\frac{H_m^{k+1}/H_m^{k+1}}{H_m^k/H_m^{k-1}} \right).$$

The solution is brilliant but does not give us a practical algorithm.

During the 1920s, in Scotland, A. C. Aitken rediscovered for himself a remarkable connection among Hankel determinants that was known to Hadamard but which was not fully exploited by him;

$$(41) \quad (H_m^k)^2 + H_{m+1}^{k-1} H_{m-1}^{k+1} = H_m^k H_{m-1}^{k+1}.$$

See [1], [2].

The relation (41) permits the computation of all the H_m^k without being drowned in determinantal evaluations.

The blemish in (41) is that the H_m^k are not of direct interest. We want to compute

$$q_m^{(k)} := \frac{H_m^{k+1}/H_m^k}{H_{m-1}^{k+1}/H_{m-1}^k}.$$

Rutishauser's clever observation was that if one introduces an auxiliary quantity

$$e_m^{(k)} := \frac{H_{m-1}^{k+1} H_{m+1}^k}{H_m^k H_m^{k+1}}$$

then (41) can be interpreted as

$$q_m^{(k)} + e_m^{(k)} = q_m^{(k+1)} + e_{m-1}^{(k+1)},$$

the additive rhombus rule, while the definitions of q and e give the product rhombus rule

$$q_m^{(k+1)} e_m^{(k+1)} = q_{m+1}^{(k)} e_m^{(k)}.$$

The rhombus rules were introduced at the end of Sect. 4. The qs and es are best laid out in a tableau that is like a difference table. See Fig. 4. Rutishauser called the e 's *modified* differences and so chose the letter e rather than d .

This qd table may be built up via the rhombus rules either from column 1 or from the top diagonal. The first column, $\{q_1^{(k)}\}$ is at hand, since

$$q_1^{(k-1)} = H_1^k / H_1^{k-1} = h_k / h_{k-1}, \quad k \geq 1$$

and

$$e_1^{(k)} = q_1^{(k+1)} - q_1^{(k)},$$

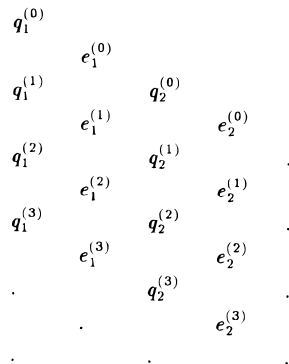


Fig. 4. qd in a (modified) difference table

a true difference.

This is far simpler than Hadamard’s solution but, in finite precision arithmetic, it is hopelessly unstable because the later e ’s are (modified) differences of converging values.

Fortunately computation along descending diagonals is stable but here the difficulty is the calculation of the top diagonal. This is not as daunting as it appears at first. If the function $f(\zeta)$ has only n poles then all q (and e) columns beyond the n th vanish. Then it suffices to build the $n \times n$ Hankel matrix \tilde{H}_n^0 (recall that H_n^0 is a number) and compute its triangular factorization

$$\tilde{H}_n^0 = L_n D_n L_n^T$$

where $D_n = \text{diag}(d_1, \dots, d_n)$ holds the successive pivots. It turns out that

$$q_k^{(0)} = H_k^0 / H_{k-1}^0 = d_k, \quad k = 1, \dots, n$$

The $e_k^{(0)}$ are found from the pivots of \tilde{H}_n^1 . This is the practical way to compute the poles from the Taylor coefficients. In fact a careful form of row interchanges (not partial pivoting) may be used to improve the accuracy of the factorization.

Next we relate the qd tableau to the computation of eigenvalues. Given a square matrix C the appropriate rational function f comes from the resolvent,

$$f(z) = x^T (I - zC)^{-1} y,$$

where x and y are arbitrary column vectors. A technical assumption is needed to guarantee that the qd tableau is well defined. In the language of control theory, see [24], the linear dynamical system $\mathcal{S}(C, y, x^T)$ must be minimal. If it is not minimal, then we might not be able to find all the poles of the system.

For this function f ,

$$h_k = x^T C^k y / x^T C^{k-1} y$$

and so the $\{h_k\}$ could be computed by the power method. However, it would be preferable to compute $\{q_1^{(0)}, e_1^{(0)}, q_2^{(0)}, e_2^{(0)}, \dots\}$ directly from C and we now know that this can be done by invoking the Lanczos algorithm on C and using the resulting tridiagonal matrix J . It turns out that the pivots that occur in computing the triangular factorization of J are the $\{q_k^{(0)}\}$ and their reciprocals are the $\{e_k^{(0)}\}$. The details are given in [27]. In other words,

$$J = \begin{bmatrix} 1 & & & & & \\ e_1^{(0)} & 1 & & & & \\ & e_2^{(0)} & 1 & & & \\ & & e_3^{(0)} & 1 & & \\ & & & \cdot & \cdot & \\ & & & & \cdot & \cdot \end{bmatrix} \begin{bmatrix} q_1^{(0)} & 1 & & & & \\ & q_2^{(0)} & 1 & & & \\ & & q_3^{(0)} & 1 & & \\ & & & q_4^{(0)} & 1 & \\ & & & & \cdot & \cdot \\ & & & & & \cdot \end{bmatrix}.$$

We see here how the LR algorithm on tridiagonals was hidden in the qd table.

12. The continued fraction connection

There is an intimate connection between our bidiagonal matrix B , the tridiagonal matrix $T = B^T B$, and a continued fraction associated with them. Properties of continued fractions influenced the qd algorithm initially and only later did the LR transformation emerge and nearly displace the continued fraction. We can not find any discussion by Rutishauser of the connection between the continued fraction and $B^T B$ so we supply it here.

Recall the notation from Sects. 2, 3, and 4.

$$B = \text{bidiag} \left\{ \begin{array}{cccccccc} & b_1 & & b_2 & & \cdot & & b_{n-2} & & b_{n-1} & & \\ a_1 & & a_2 & & \cdot & & \cdot & & a_{n-1} & & a_n & \end{array} \right\}.$$

$$q_i = a_i^2$$

$$e_i^2 = b_i^2, \quad e_0 = e_n = 0$$

$$T = \text{tridiag} \left\{ \begin{array}{ccccccccccc} & \sqrt{q_1 e_1} & & \sqrt{q_2 e_2} & & \sqrt{q_3 e_3} & & \sqrt{q_{n-1} e_{n-1}} & & & & \\ q_1 & & q_2 + e_1 & & q_3 + e_2 & & \dots & & q_n + e_{n-1} & & & \\ & \sqrt{q_1 e_1} & & \sqrt{q_2 e_2} & & \sqrt{q_3 e_3} & & \sqrt{q_{n-1} e_{n-1}} & & & & \end{array} \right\}.$$

Rutishauser associates with T the continued fraction

$$(42) \quad F(\zeta) = \frac{1}{\zeta - q_1 - \frac{e_1 q_1}{\zeta - q_2 - e_1 - \frac{e_2 q_2}{\zeta - q_3 - e_2 - \dots}}}$$

It is not obvious how $F(\zeta)$ relates to T . The answer is

$$F(\zeta) = [(\zeta I - T)^{-1}]_{1,1}$$

or more generally,

$$F(\zeta) = x^T (\zeta I - T)^{-1} y$$

with \mathbf{x} and \mathbf{y} as defined near the end of Sect. 11. The inverse is well defined for all ζ with $|\zeta|$ exceeding the spectral radius of T . The particular form of the continued fraction arises from the triangular factorization of $\zeta I - T$ from the bottom up:

$$\zeta I - T = \tilde{L}^T \tilde{D} \tilde{L}$$

where \tilde{L} is unit lower triangular and

$$\tilde{D} = \text{diag}(\tilde{d}_1, \tilde{d}_2, \dots, \tilde{d}_n)$$

$$\tilde{d}_i = \tilde{d}_i(\zeta).$$

Then

$$(\zeta I - T)^{-1} = \tilde{L}^{-1} \tilde{D}^{-1} \tilde{L}^{-T}$$

and

$$F(\zeta) = \tilde{d}_1^{-1}$$

The recurrence for the \tilde{d}_j is

$$(43) \quad \begin{aligned} \tilde{d}_n &:= \zeta - q_n - e_{n-1} \\ \tilde{d}_j &:= \zeta - q_j - e_{j-1} - q_j e_j / \tilde{d}_{j+1} \quad \text{for } j = n-1, \dots, 2, 1. \end{aligned}$$

This establishes (42).

There is a simpler continued fraction expansion for $F(\zeta)$. It corresponds to a recurrence for $\tilde{d}_j + e_{j-1}$. From (43)

$$(44) \quad \begin{aligned} \tilde{d}_j + e_{j-1} &= \zeta - q_j \left(1 + \frac{e_j}{\tilde{d}_{j+1}}\right) \\ &= \zeta - q_j / \left(\frac{\tilde{d}_{j+1} + e_j - e_j}{\tilde{d}_{j+1} + e_j}\right) \\ &= \zeta - q_j / (1 - e_j / (\tilde{d}_{j+1} + e_j)) \end{aligned}$$

Since $e_0 = 0$, (44) gives

$$F(\zeta) = 1/\tilde{d}_1 = \frac{1}{\zeta -} \frac{q_1}{1 -} \frac{e_1}{\zeta -} \frac{q_2}{1 -} \frac{e_2}{\zeta -} \dots$$

This form is remarkable for the direct connection of q_i and e_i to the $(1, 1)$ entry of $(\zeta I - T)^{-1}$.

13. Singular vectors

We sketch three methods for this task. Some parts of these methods are independent of the bidiagonal form. We plan a separate communication discussing in detail the computation of singular vectors and eigenvectors.

13.1. Method 1

Let $B = U \Sigma V^T$. Sometimes only a few columns of V are wanted, those corresponding to selected singular values. Sometimes all of V is wanted and, on other occasions, both U and V are required. Consequently there will not be just one procedure for computing singular vectors.

Consider first the oqd algorithm with which we began our development of dqds. Let $\{B_j\}$ be a sequence of upper bidiagonals generated from $B_1 (= B)$ by oqd. By Theorem 1,

$$B_2 = Q_1^T B_1^T = Q_1^T V \Sigma U^T$$

and

$$B_3 = Q_2^T B_2^T = Q_2^T U \Sigma V^T Q_1.$$

Continuing this process one finds that

$$B_{2k+1} = Q_{2k}^T \cdots Q_2^T U \Sigma V^T Q_1 Q_3 \cdots Q_{2k-1}.$$

As $k \rightarrow \infty$, $B_{2k+1} \rightarrow \Sigma$ and so

$$U = \lim_{k \rightarrow \infty} Q_2 Q_4 \cdots Q_{2k},$$

$$V = \lim_{k \rightarrow \infty} Q_1 Q_3 \cdots Q_{2k-1}.$$

This suggests one method for producing U , V , or both.

Method 1 (To find U and V). Apply oqd to B and accumulate the plane rotations from the odd and even passes to build up V and U .

13.2. Method 2

Since oqd converges slowly Method 1 is not a serious contender for the preferred algorithm but see Method 3 where it plays a role.

Shifts are used to accelerate convergence and because they are non-restoring, no n-orthogonal transformations are used and one set of singular vectors is lost and the other set is preserved if alternative shifts are zero. The motivation for Method 2, given next, is simple. For a bidiagonal B with last entry $a_n = 0$ the oqd transformation gives bidiagonal \hat{B} where

$$B^T = Q \hat{B}$$

and both $\hat{a}_n = 0$ and $\hat{b}_{n-1} = 0$. So the last column of Q is a right singular vector for $0 (= \sigma_n[B])$. It is only necessary to form the product of the plane rotations that build up Q and the last column is the desired vector.

When the singular values are known we can use their differences, in turn, as shifts in oqds to obtain suitable singular matrices.

Before describing Method 2 in more detail we emphasize a subtle benefit that follows from finding the singular vectors in the order given by $\sigma_n, \sigma_{n-1}, \dots, \sigma_1$: the sines and cosines of the plane rotations used in our oqd and oqds algorithms can be recovered from their squares because they are all positive. Consequently we can work with the root free versions dqd and dqds although we still refer below to the underlying matrices B_j rather than their squared components given by qd arrays Z_j which we actually use.

Method 2 (to find U , not V)

1. Initialize: $B_n = B$; $U = I$, $j = n$, $\sigma_{n+1}^2 = 0$.
2. Apply dqds to B_j (using $e_i = b_i^2$, $q_i = a_i^2$) with shift $\sigma_j^2 - \sigma_{j+1}^2$ obtaining \hat{B}_j .
3. Apply dqd to \hat{B}_j obtaining B_{j-1} and save $\{c_i^2, ; s_i^2, i = 1, j-1\}$ where

$$c_i^2 = \frac{d_i}{d_i + e_i} \quad , \quad s_i^2 = \frac{e_i}{d_i + e_i}$$

4. Take positive square roots to obtain $\{c_i, s_i, i = 1, j-1\}$
5. Accumulate the corresponding plane rotations G_1, \dots, G_{n-1} of oqd into U according to $U \leftarrow U G_i$, $i = 1, j-1$
6. $j \leftarrow j-1$, if $j = 1$ stop else go to 2.

Since the last row and column of B_{j-1} are zero we deflate them and consider B_i as an $i \times i$ matrix, $i = n, n-1, \dots, 2, 1$.

The same algorithm may be used to compute V , not U , if one dqd transformation is applied to B initially to give \hat{B} . Then initialize by $B_n = \hat{B}$, $V = Q$, $j = n$.

Here is the justification for each of the steps in Method 2.

1. Initialization
2. The output of dqds satisfies

$$\hat{B}_j^T \hat{B}_j = B_j B_j^T - (\sigma_j^2 - \sigma_{j+1}^2) I \quad , \quad B_j \text{ is } j \times j.$$

Consequently, in exact arithmetic, $\sigma_j[\hat{B}_j] = 0$ and the right singular vectors of \hat{B}_j are the left singular vectors of B_j .

- 3.

$$B_{j-1} = Q_j^T \hat{B}_j^T \quad , \quad Q_j^T = G_{j-1}^T \cdots G_1^T.$$

So the last column of Q_j is a right singular vector of \hat{B}_j and hence a left singular vector of B_{j-1} . At this point B_{j-1} is a $j \times j$ matrix but its zero last column will be discarded at the next iteration.

4. There is an advantage to taking the square roots all together in a pipelined or vector machine. In a sequential computer Step 4 could be merged with Step 3.
- 5.

$$U \leftarrow U G_i \quad , \quad i = 1, \dots, j \quad . \quad G_i = \begin{pmatrix} c_i & -s_i \\ s_i & c_i \end{pmatrix}.$$

13.3. Method 3

This approach can be used on a distributed memory parallel computer. It employs plane rotations exclusively.

Consider the permuted form of $\begin{bmatrix} 0 & B \\ B^T & 0 \end{bmatrix}$,

$$T = \text{tridiag} \left\{ \begin{array}{cccccccc} & a_1 & b_1 & a_2 & \cdot & a_n & & \\ 0 & & 0 & 0 & 0 & & 0 & 0 \\ & a_1 & b_1 & a_2 & \cdot & a_n & & \end{array} \right\}.$$

The eigenvalues of T are $\pm\sigma_i$, $i = 1, 2, \dots, n$ and the right and the left singular vectors of B may be recovered from the eigenvectors of T .

Since the eigenvalues of T are known one may obtain the eigenvectors for σ_j by one QR factorization of $T - \sigma_j I$ and accumulating the effect of each plane rotation on a single vector. In more detail, the plane rotations build up an upper Hessenberg orthogonal matrix and one updates the last active column in this matrix at each minor step. If $s = \sin \theta_j$ and $c = \cos \theta_j$ at a step j then, for $j < n$,

$$v(i) = v(i) * (-s) \quad , \quad i = 1, \dots, j \quad , \quad v(j+1) = c$$

Initially, $v = e_1$.

If this vector v is not adequate then we solve

$$R^T(Rw) = v$$

for w and finally normalize w . Thus the triple diagonal matrix R must be saved.

Acknowledgements. KVF wishes to thank Professor James W. Demmel for inviting him to join the LAPACK group at UC Berkeley from July 1991 to July 1993. Our thanks also go to Professors Stanley C. Eisenstat and Ilse C. F. Ipsen for their careful reading and detection of errors.

References

1. Aitken, A.C., On Bernoulli's numerical solution of algebraic equations, Proc. Roy. Soc. Edinburgh, Section A, **46**: 289–305, 1926
2. Aitken, A.C., Further numerical studies in algebraic equations and matrices, Proc. Roy. Soc. Edinburgh, Section A, **51**: 80–90, 1931
3. Battin, R.H., Astronautical Guidance, McGraw-Hill, New York, 1964
4. Bauer, F.L., qd-method with Newton shift, Technical Report, 56, Computer Science Department, Stanford University, 1967
5. Bernoulli, Daniel, Observationes de seriebus quae formantur ex additione vel subtractione quacunqu terminorum su mutuo consequentium, ubi praesertim earundem insignis usus pro inveniendis radicum omnium aequationum algebraicarum ostenditur, Commentarii Academiae Scientiarum Imperialis Petropolitanae **3**: 85–100, 1732 (1728)
6. Bierman, G.J., Factorization Methods for Discrete Sequential Estimation, Academic Press, New York, 1977
7. Chandrasekaran, S., Ipsen, I.C.F., Analysis of a QR algorithm for computing singular values. Research Report 917, Department of Computer Science, Yale University, 1992
8. Chatelin, Françoise, Spectral Approximation of Linear Operators, Academic Press, New York, 1983
9. Codenotti, B., Leoncini, M., Parallel Complexity of Linear System Solution, World Scientific, Singapore, 1991
10. DeJong, L.S., Towards a formal definition of numerical stability, Numer. Math. **28**: 211–220, 1977
11. Demmel, James, Kahan, W., Accurate singular values of bidiagonal matrices, SIAM J. Sci. Sta. Comput. **11**: 873–912, 1990
12. Dongarra, J.J., Bunch, J.R., Moler, C.B., Stewart, G.W., LINPACK User' Guide, SIAM, Philadelphia, 1979
13. Euler, Leonard, Introduction to Analysis of the Infinite, volume 1. Springer-Verlag, New York, 1988
14. Eulero, Leonhardo, Introductio in Analysin Infinitorum, volume 8 of series 1. Marcum-Michaellem Bousquet, Lausannæ, 1748
15. Faddeev, D.K., Kublanovskaya, V.N., Faddeeva, V.N., Sur les systèmes Linéaires Algébriques de Matrices Rectangulaires et Mal-Conditionnées, Colloq. Int. du C.N.R.S. Besançon 1966, No. 165: 161–70, 1968
16. Fernando, K. Vince and Nicholson. H., Identification of linear systems with input and output noise: The Koopmans-Levin method, Proc. IEE Part D Control Theory and Applications **132**: 30–36. 1985
17. Fernando, K. Vince, Parlett, Beresford N., Implicit Cholesky Algorithm for singular values and vectors, Technical Report under preparation, Centre for Pure and Applied Mathematics, University of California at Berkeley, 1993
18. Golub, Gene H. and Reinsch, C., Singular value decomposition and least squares solutions, Numer. Math. **14**: 403–420, 1970
19. Golub, Gene H., Van Loan, Charles F., Matrix Computations, Johns Hopkins University Press, Baltimore, 1989
20. Hadamard, M.J., Essai sur l'étude des fonctions données par leur développement de Taylor, Gauthier-Villars, Paris, 1891
21. Henrici, P., The Quotient-Difference Algorithm, Nat. Bur. Standards Appl. Math. Series **19**: 23–46 1958
22. Householder, A. S., The Numerical Treatment of a Single Nonlinear Equation, McGraw-Hill, New York, 1970
23. Johnson, Charles. R., A Gersgorin-type lower bound for the smallest singular value, Linear Algebra and Its Applications **112**: 1–7, 1989

24. Kailath, Thomas, Linear Systems, Prentice-Hall, Englewood Cliffs, NJ, 1980
25. König, Julius, Ueber eine Eigenschaft der Potenzreihen, Math. Ann. **23**: 447–449, 1884
26. Mathias, R., The stability of parallel prefix matrix multiplication – with applications to tridiagonal matrices, Technical Report, 1993
27. Parlett, Beresford. N., The development and use of methods of LR type, SIAM Review **6**: 275–295, 1964
28. Parlett, Beresford. N., The Symmetric Eigenvalue Problem, Prentice-Hall, Englewood Cliffs, NJ, 1980
29. Reinsch, C., Bauer, F. L., Rational QR transformation with Newton shift for symmetric tridiagonal matrices, Numer. Math. **11**: 264–272, 1968
30. Rutishauser, H., Der Quotienten-Differenzen-Algorithmus, Z. Angew. Math. Phys. **5**: 233–251, 1954
31. Rutishauser, H., Ein infinitesimales Analogon zum Quotienten-Differenzen-Algorithmus, Arch. Math. **5**: 132–137, 1954
32. Rutishauser, H., Solution of eigenvalue problems with the LR-transformation, Nat. Bur. Standards Appl. Math. Series **49**: 47–81, 1958
33. Rutishauser, H., Über eine kubisch konvergente Variante der LR-Transformation, Z. Angew. Math. Mech. **11**: 49–54, 1960
34. Rutishauser, H., Vorlesungen über numerische Mathematik, Birkhäuser, Basel, 1976
35. Rutishauser, H., Lectures on Numerical Mathematics, Birkhäuser, Boston, 1990
36. Rutishauser, H., Schwarz H.R., The LR transformation method for symmetric matrices. Numer. Math. **5**: 273–289, 1963
37. Schröder, E., Über unendlich viele Algorithmen zur Auflösung der Gleichungen, Math. Ann. **2**: 317–365, 1870
38. Schröder, E., On infinitely many algorithms for solving equations (Translated by G. W. Stewart), Technical Report, UMIACS-TR-92-12, CS-TR 2990, Dept. of Computer Science and Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742, 1992
39. Schwarz, H.R., Rutishauser, H., Stiefel, E., Numerical Analysis of Symmetric Matrices, Prentice-Hall, Englewood Cliffs, NJ, 1973
40. Stewart, G.W., On a companion operator for analytic functions, Numer. Math. **18**: 26–43, 1971
41. Stewart, G. W, On an algorithm for refining a rank-revealing URV factorization and a perturbation theorem for singular values, Technical Report, UMIACS-TR-91-38, CS-TR 2626, Dept. of Computer Science and Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742, 1991
42. Veselić, K., Hari, V., A note on a one-sided Jacobi algorithm, Numer. Math. **56**: 627–633, 1989
43. Wilkinson, J.H., The Algebraic Eigenvalue Problem, Clarendon Press, Oxford, 1965