

## A JACOBI–DAVIDSON ITERATION METHOD FOR LINEAR EIGENVALUE PROBLEMS\*

GERARD L. G. SLEIJPEN<sup>†</sup> AND HENK A. VAN DER VORST<sup>†</sup>

**Abstract.** In this paper we propose a new method for the iterative computation of a few of the extremal eigenvalues of a symmetric matrix and their associated eigenvectors. The method is based on an old and almost unknown method of Jacobi. Jacobi's approach, combined with Davidson's method, leads to a new method that has improved convergence properties and that may be used for general matrices. We also propose a variant of the new method that may be useful for the computation of nonextremal eigenvalues as well.

**Key words.** eigenvalues and eigenvectors, Davidson's method, Jacobi iterations, harmonic Ritz values

**AMS subject classifications.** 65F15, 65N25

**1. Introduction.** Suppose we want to compute one or more eigenvalues and their corresponding eigenvectors of the  $n \times n$  matrix  $A$ . Several iterative methods are available: Jacobi's diagonalization method [9], [22], the power method [9], the method of Lanczos [13], [22], Arnoldi's method [1], [25], and Davidson's method [4], [25], [3], [14], [17]. The latter method has been reported to be quite successful, most notably in connection with certain symmetric problems in computational chemistry [4], [5], [31]. The success of the method seems to depend quite heavily on the (strong) diagonal dominance of  $A$ .

The method of Davidson is commonly seen as an extension to Lanczos's method, but as Saad [25] points out, from the implementation point of view it is more related to Arnoldi's method. In spite of these relations, the success of the method is not well understood [25]. Some recent convergence results and improvements, as well as numerical experiments, are reported in [3], [14], [15], [17], [16], [18], [27].

Jacobi [12] proposed a method for eigenvalue approximation that essentially was a combination of (1) Jacobi rotations, (2) Gauss–Jacobi iterations, and (3) an almost forgotten method that we will refer to as Jacobi's orthogonal component correction (JOCC). Reinvestigation of Jacobi's ideas leads to another view on the method of Davidson, and this not only helps us explain the behavior of the method, it also leads to a new and robust method with superior convergence properties for nondiagonally dominant (unsymmetric) matrices as well. Special variants of this method are already known; see [18], [27] and our discussion in §4.1.

The outline of this paper is as follows. In §2 we briefly describe the methods of Davidson and Jacobi, and we show that the original Davidson's method may be viewed as an accelerated Gauss–Jacobi iteration method. Likewise, more recent approaches which include other preconditioners  $M$  can be interpreted as accelerated standard iteration methods associated with the splitting  $A = M - N$ .

In §3 we propose the new approach, which is essentially a combination of the JOCC approach and the method of Davidson for creating more general subspaces. The difference between this approach and Davidson's method may seem very subtle but it is fundamental. Whereas in Davidson's method accurate preconditioners  $M$  (accurate

\* Received by the editors June 22, 1994; accepted for publication (in revised form) by A. Greenbaum May 22, 1995.

<sup>†</sup> Mathematical Institute, University of Utrecht, Budapestlaan 6, P.O. Box 80.010, Utrecht 3508 TA, The Netherlands (sleijpen@math.ruu.nl, vorst@math.ruu.nl).

in the sense that they approximate the inverse of the given operator very well) may lead to stagnation or very slow convergence, the new approach takes advantage of such preconditioners, even if they are exact inverses. It should be stressed that in this approach we do not precondition the given eigensystem (neither does Davidson), but we precondition an auxiliary system for the corrections for the eigen approximations. The behavior of the method is further discussed in §4. There we see that for a specific choice the speed of convergence of the approximated eigenvalue is quadratic (and for symmetric problems even cubic). In practice, this requires the exact solution of a correction equation, but as we will demonstrate by simple examples (§6), this may be relaxed. We suggest using approximate solutions for the correction equations. This idea may be further exploited for the construction of efficient inner–outer iteration schemes, or by using preconditioners similar to those suggested for the Davidson method.

In §5 we discuss the harmonic Ritz values, and we show how these can be used in combination with our new algorithm for the determination of “interior” eigenvalues. We conclude with some simple but illustrative numerical examples in §6. The new method has already found its way into more complicated applications in chemistry and plasma physics modeling.

**2. The methods of Davidson and Jacobi.** Jacobi and Davidson originally presented their methods for symmetric matrices, but as is well known and as we will do in our presentation, both methods can easily be formulated for nonsymmetric matrices.

**2.1. Davidson’s method.** The main idea behind Davidson’s method is the following one. Suppose we have some subspace  $K$  of dimension  $k$ , over which the projected matrix  $A$  has a Ritz value  $\theta_k$  (e.g.,  $\theta_k$  is the largest Ritz value) and a corresponding Ritz vector  $u_k$ . Let us assume that an orthogonal basis for  $K$  is given by the vectors  $v_1, v_2, \dots, v_k$ .

Quite naturally the problem of how to expand the subspace in order to find a successful update for  $u_k$  arises. To that end we compute the defect  $r = Au_k - \theta_k u_k$ . Then Davidson, in his original paper [4], suggests computing  $t$  from  $(D_A - \theta_k I)t = r$ , where  $D_A$  is the diagonal of the matrix  $A$ . The vector  $t$  is made orthogonal to the basis vectors  $v_1, \dots, v_k$ , and the resulting vector is chosen as the new  $v_{k+1}$ , by which  $K$  is expanded.

It has been reported that this method can be quite successful in finding dominant eigenvalues of (strongly) diagonally dominant matrices. The matrix  $(D_A - \theta_k I)^{-1}$  can be viewed as a preconditioner for the vector  $r$ . Davidson [6] suggests that his algorithm (more precisely, the Davidson–Liu variant of it) may be interpreted as a Newton–Raphson scheme, and this has been used as an argument to explain its fast convergence. It is tempting to also see the preconditioner as an approximation for  $(A - \theta_k I)^{-1}$ , and, indeed, this approach has been followed for the construction of more complicated preconditioners (see, e.g., [16], [3], [14], [17]). However, note that  $(A - \theta_k I)^{-1}$  would map  $r$  onto  $u_k$ , and hence it would not lead to an expansion of our search space. Clearly, this is a wrong interpretation for the preconditioner.

**2.2. The methods of Jacobi.** In his paper of 1846 [12], Jacobi introduced a combination of two iterative methods for the computation of approximations of eigenvalues of a symmetric matrix.<sup>1</sup> He proposed the combination as an entity, but

<sup>1</sup> This came to our attention by reading A. den Boer’s Master’s thesis [7].

at present the methods are only used separately. The first method is well known and is referred to as the Jacobi method (e.g., §8.4 in [9]). It is based on Jacobi plane rotations, which are used to force the matrix  $A$  to diagonal dominance. We will refer to this method as Jacobi's diagonalisation method. The second method is much less well known and is related to the Davidson method. For ease of discussion we will call this second method the JOCC. It turns out that Davidson's method can be interpreted as an accelerated JOCC method, just as Arnoldi's method can be seen as an accelerated power method.

**2.2.1. The JOCC method.** Jacobi considered an eigenvalue problem as a system of linear equations for which his iterative linear solver [11], the Jacobi or Gauss-Jacobi iteration (e.g., §10.1 in [9]), might be applicable.

Suppose we have a diagonally dominant matrix  $A$ , of which  $a_{1,1} = \alpha$  is the largest diagonal element. Then  $\alpha$  is an approximation for the largest eigenvalue  $\lambda$ , and  $e_1$  is an approximation for the corresponding eigenvector  $u$ . In modern matrix notation (which was unknown in Jacobi's time), his approach can be presented as follows.

Consider the eigenvalue problem

$$(1) \quad A \begin{bmatrix} 1 \\ z \end{bmatrix} = \begin{bmatrix} \alpha & c^T \\ b & F \end{bmatrix} \begin{bmatrix} 1 \\ z \end{bmatrix} = \lambda \begin{bmatrix} 1 \\ z \end{bmatrix},$$

where  $F$  is a square matrix,  $\alpha$  is a scalar, and  $b$ ,  $c$  and  $z$  are vectors of appropriate size. We are interested in the eigenvalue  $\lambda$  that is close in some sense to  $\alpha$ , and in the corresponding eigenvector  $u = (1, z^T)^T$ , with component  $z$  orthogonal to  $e_1$ . Problem (1) is equivalent with

$$(2) \quad \lambda = \alpha + c^T z,$$

$$(3) \quad (F - \lambda I)z = -b.$$

Jacobi proposed to solve (3) iteratively by his Jacobi iteration, with  $z_1 = 0$ , and an updated approximation for  $\lambda$ , using (2):<sup>2</sup>

$$(4) \quad \begin{cases} \theta_k &= \alpha + c^T z_k, \\ (D - \theta_k I)z_{k+1} &= (D - F)z_k - b, \end{cases}$$

where  $D$  is the diagonal of  $F$  (although  $\theta_k$  is not a Ritz value, we have used it to characterize it as an eigenvalue approximation).

Jacobi solved, as is also customary now, the update  $y_k$  for  $z_k$  from the diagonal system rather than solving  $z_{k+1}$  directly. Therefore, a better representation of the JOCC method would be

$$(5) \quad \begin{cases} \mu_k &= c^T y_k, \\ \theta_{k+1} &= \theta_k + \mu_k, \\ z_{k+1} &= z_k + y_k, \\ (D - \theta_{k+1} I)y_{k+1} &= (D - F)y_k + \mu_k z_{k+1}, \end{cases}$$

with  $z_1 = 0$ ,  $\theta_1 = \alpha$ , and  $y_1 = -(D - \alpha I)^{-1}b$ . However, in connection with Davidson's method, representation (4) simplifies the discussion.

<sup>2</sup> Actually, Jacobi updated the approximation of  $\lambda$  only at every even step. There is no explanation for why he did not update in the odd steps as well.

**2.2.2. Short discussion on the JOCC method.** Jacobi was well aware of the fact that the Jacobi iteration converges (fast) if the matrix is (strongly) diagonally dominant. Therefore, he proposed to perform a number of steps of the Jacobi diagonalization method in order to obtain a (strongly) diagonally dominant matrix before applying the JOCC method. Since he was interested in the eigenvalue closest to  $\alpha$ , he did enough steps using the diagonalization method to obtain a diagonally dominant  $F - \alpha I$  so that  $F - \theta_k I$  was also diagonally dominant. This can be done provided that  $\lambda$  is a simple eigenvalue of  $A$ . The application of the diagonalization method can be viewed as a technique to improve the initial guess  $e_1$ , i.e., the given matrix is rotated so that  $e_1$  is closer to the (rotated) eigenvector  $u$ . These rotations were done only at the start of the iteration process, and this process was carried out with fixed  $F$  and  $D$ .

However, note that in Jacobi's approach we are looking, at all stages, for the orthogonal complement to the *initial* approximation  $u_1 = e_1$ . We do not take into account that better approximations  $u_k = (1, z_k^T)^T$  become available in the process, and that it may be more efficient to try to compute the orthogonal complement  $u - (u^T u_k)u_k$ . In the JOCC framework an improved approximation would have led to a similar situation as in (1), if we would have applied plane rotations on  $A$ , such that  $u_k$  would have been rotated to  $e_1$  by these rotations. Therefore, what Jacobi did only in the first step of the process could have been done *at each step* of the iteration process. This is an exciting idea, since in Jacobi's approach the rotations were motivated by the desire to obtain stronger diagonal dominance, whereas our discussion suggests that one might take advantage of the information in the iteration process. Of course, this would have led to a different operator  $F$  in each step and this is an important observation for the formulation of our new algorithm.

**2.3. Davidson's method as an accelerated JOCC method.** We will apply Davidson's method for the same problem as before. In particular we will assume that  $A$  is as in (1) and that  $u_1 = e_1$ .

The eigenvector approximations produced in the  $k$ th step of JOCC, as well as in the Davidson method, are denoted by  $u_k$ . We assume that  $u_k$  is scaled such that its first coordinate is 1:  $u_k = (1, z_k^T)^T$ . Let  $\theta_k$  be the associated approximation to the eigenvalue. It will be clear from the context to which process an approximation refers.

The residual is given by

$$(6) \quad r_k = (A - \theta_k I)u_k = \begin{bmatrix} \alpha - \theta_k + c^T z_k \\ (F - \theta_k I)z_k + b \end{bmatrix}.$$

Davidson proposes computing  $t_k$  from

$$(7) \quad (D_A - \theta_k I)t_k = -r_k,$$

where  $D_A$  is the diagonal of  $A$ . For the component  $\hat{y}_k := (0, y_k^T)^T$  of  $t_k$  orthogonal to  $u_1$  it follows, with  $D$  the diagonal of  $F$ , that

$$(8) \quad (D - \theta_k I)y_k = -(F - \theta_k I)z_k - b = (D - F)z_k - (D - \theta_k I)z_k - b,$$

or equivalently,

$$(9) \quad (D - \theta_k I)(z_k + y_k) = (D - F)z_k - b.$$

Comparing (9) with (4), we see that  $z_k + y_k$  is the  $z_{k+1}$  that we would have obtained with one step of JOCC starting at  $z_k$ . But after this point, Davidson's method is

an improvement over the JOCC method because instead of taking  $\widehat{u}_{k+1} := (1, (z_k + y_k)^T)^T = u_k + \widehat{y}_k$  as the next approximating eigenvector (as in JOCC), Davidson suggests computing the Ritz vector of  $A$  with respect to the subspace computed so far (that is, over the subspace spanned by the old approximations  $u_1, \dots, u_k$  and the new  $\widehat{u}_{k+1}$ ). Actually, Davidson selects the correction  $t_k$ , but

$$\text{span}(u_1, \dots, u_k, \widehat{u}_{k+1}) = \text{span}(u_1, \dots, u_k, t_k).$$

For the computation of the Ritz vector it is convenient to have an orthonormal basis, and that is precisely what is constructed in Davidson's method. This orthonormal basis  $v_1, \dots, v_{k+1}$  appears if one orthonormalizes  $u_1, \dots, u_k, t_k$  by Gram-Schmidt. The  $(k+1)$ th step in either JOCC or Davidson's method can be summarized as follows.

**JOCC.** Jacobi computes the component  $\widehat{y}_k$  of  $t_k$  orthogonal to  $u_1$  and takes  $u_{k+1} = u_k + \widehat{y}_k$ ,  $\theta_{k+1} = e_1^T A u_{k+1}$ . Unlike Davidson's approach, Jacobi only computes components that are orthogonal to  $u_1 = e_1$ . However, in view of the orthogonalization step, the components in the  $u_1$ -direction (as well as in new directions) automatically vanish in Davidson's method.

**Davidson's method.** Davidson computes the component  $v_{k+1}$  of  $t_k$  orthogonal to  $u_1, \dots, u_k$  and takes for  $u_{k+1}$  and  $\theta_{k+1}$  the Ritz vector, respectively, Ritz value, of  $A$  with respect to the space spanned by  $v_1, \dots, v_{k+1}$ . Davidson exploits the complete subspace constructed so far, while Jacobi takes only a simple linear combination of the last vector  $z_k$  and the last correction  $y_k$  (or, taking the  $u_1$ -component into account, of  $u_1$ , the last vector  $u_k$ , and the last correction  $t_k$ ). Although  $\text{span}(u_1, \dots, u_k, \widehat{u}_{k+1}) = \text{span}(u_1, \dots, u_{k+1})$ , the Davidson approximations do not span the same subspace as the Jacobi approximations since the eigenvalue approximations  $\theta_k$  of Jacobi and of Davidson are different. Consequently, the methods have different corrections  $t_k$  and also different components orthogonal to the  $u_j$ .

Note that our formulation makes clear that Davidson's method also attempts to find the orthogonal update  $(0, z^T)^T$  for the initial guess  $u_1 = e_1$ , and it does so by a clever orthogonalization procedure. However, just as in the JOCC approach, the process works with fixed operators (in particular  $D_A$ ; other variants use different approximations for  $A$ ) and not with operators associated with the orthogonal complement of the current approximation  $u_k$  (see also §2.2.2). This characterizes the difference between our algorithm (of §3) and Davidson's method.

**3. The new Jacobi-Davidson iteration method.** From now on we will allow the matrix  $A$  to be complex, and in order to express this we use the notation  $v^*$  for the complex conjugate of a vector (if complex), or the transpose (if real), and likewise for matrices.

As we have stressed in the previous section, the JOCC method and Davidson's method can be seen as methods that attempt to find the correction to some initially given eigenvector approximation. In fact, what we want is to find the orthogonal complement for our current approximation  $u_k$  with respect to the desired eigenvector  $u$  of  $A$ . Therefore, we are interested in seeing explicitly what happens in the subspace  $u_k^\perp$ .

The orthogonal projection of  $A$  onto that space is given by  $B = (I - u_k u_k^*) A (I - u_k u_k^*)$  (we assume that  $u_k$  has been normalized). Note that for  $u_k = e_1$ , we have that  $F$  (cf. (1)) is the restriction of  $B$  with respect to  $e_1^\perp$ . It follows that

$$(10) \quad A = B + A u_k u_k^* + u_k u_k^* A - \theta_k u_k u_k^*.$$

When we are in search of an eigenvalue  $\lambda$  of  $A$  close to  $\theta_k$ , then we want to have the correction  $v \perp u_k$  to  $u_k$  such that

$$A(u_k + v) = \lambda(u_k + v),$$

or, after inserting (10) and realizing that  $Bu_k = 0$ ,

$$(11) \quad (B - \lambda I)v = -r + (\lambda - \theta_k - u_k^* A v)u_k.$$

Since the left-hand side and  $r$  have no component in  $u_k$ , it follows that the factor for  $u_k$  must vanish, and hence  $v$  should satisfy

$$(12) \quad (B - \lambda I)v = -r.$$

We replace  $\lambda$  by the current approximation  $\theta_k$  just as in JOCC and Davidson's method, but, unlike both methods, we propose to work with approximations for an operator ( $B$ ) that varies from step to step. The resulting algorithm may be viewed as a combination of Jacobi's approach to look for the orthogonal complement of a given eigenvector approximation and the Davidson algorithm for expanding the subspace in which the eigenvector approximations are constructed. This explains our choice of the name Jacobi–Davidson for the new method. Note that we are free to use any method for the (approximate) solution of (12) and that it is not necessary to require diagonal dominance of  $B$  (or  $A$ ).

Before we present our complete algorithm, we briefly review some different approaches.

1. If we approximate  $v$  simply by  $r$ , then we formally obtain the same results as with the Arnoldi method.
2. If we approximate  $v$  by  $(D_A - \theta_k I)^{-1}r$ , then we obtain the original Davidson method.
3. More recent suggestions made in [3], [14], [15], [17], [6] come down to better approximations for the inverse of  $A - \theta_k I$ , e.g., incomplete decompositions for this operator. However, as is well known, this is a risky approach (see [25], [3]) since the exact inverse of this operator leads to failure of the method.<sup>3</sup> Therefore the approximation should not be too accurate [25].
4. In our approach, we replace  $B - \lambda I$  by  $B - \theta_k I$ , and one has to select suitable approximations  $\tilde{t} \perp u_k$  for the solution of

$$(13) \quad (B - \theta_k I)t = -r \quad \text{and} \quad t \perp u_k.$$

This will lead to quadratical convergence if we take the exact solution  $t$  (see observation 3 in §4), which is in sharp contrast with what happens if one takes the exact solution of  $(A - \theta_k I)t = -r$ .

5. Another viewpoint on our modified Davidson algorithm is that it can be seen as an accelerated inexact shift and invert method (that is, the “invert” part may be inexact).

We do not know, for general systems, how to approximate a solution of (13) sufficiently well with modest computational effort. In most of our numerical experiments, we have constructed approximations by carrying out only a few steps of an iterative method (for instance, generalized minimum residual: GMRES [26]) in order

<sup>3</sup> Any progress in this case may be attributed to the effects of rounding errors.

ALGORITHM 1. *The Jacobi-Davidson method.*

1. **Start:** Choose an initial nontrivial vector  $v$ .
  - Compute  $v_1 = v / \|v\|_2$ ,  $w_1 = Av_1$ ,  $h_{11} = v_1^* w_1$ ,  
 set  $V_1 = [v_1]$ ,  $W_1 = [w_1]$ ,  $H_1 = [h_{11}]$ ,  
 $u = v_1$ ,  $\theta = h_{11}$ , compute  $r = w_1 - \theta u$ .
2. **Iterate:** Until convergence do:
3. **Inner Loop:** For  $k = 1, \dots, m - 1$  do:
  - Solve (approximately)  $t \perp u$ ,
 
$$(I - uu^*)(A - \theta I)(I - uu^*)t = -r.$$
    - Orthogonalize  $t$  against  $V_k$  via modified Gram-Schmidt,  
 and expand  $V_k$  with this vector to  $V_{k+1}$ .
    - Compute  $w_{k+1} := Av_{k+1}$   
 and expand  $W_k$  with this vector to  $W_{k+1}$ .
    - Compute  $V_{k+1}^* w_{k+1}$ , the last column of  $H_{k+1} := V_{k+1}^* A V_{k+1}$ ,  
 and  $v_{k+1}^* W_k$ , the last row of  $H_{k+1}$  (only if  $A \neq A^*$ ).
    - Compute the largest eigenpair  $(\theta, s)$  of  $H_{k+1}$  (with  $\|s\|_2 = 1$ ).
    - Compute the Ritz vector  $u := V_{k+1}s$ ,  
 compute  $\hat{u} := Au$  ( $= W_{k+1}s$ ), and  
 the associated residual vector  $r := \hat{u} - \theta u$ .
    - Test for convergence. Stop if satisfied.
4. **Restart:** Set  $V_1 = [u]$ ,  $W_1 = [\hat{u}]$ ,  $H_1 = [\theta]$ , and goto 3.

to illustrate how powerful our approach is even when we solve the system only in very modest precision, but this is certainly not an optimal choice in many practical situations. However, our experiments illustrate that the better we approximate the solution  $\tilde{v}$  of (13), the faster convergence we obtain (and no stagnation as in Davidson's method).

The algorithm for the improved Davidson method then becomes as in Algorithm 1 (in the style of [25]). We have skipped indices for variables that overwrite old values in an iteration step, e.g.,  $u$  instead of  $u_k$ . We do not discuss implementation issues, but we note that the computational costs for Algorithm 1 are about the same as for Davidson's method (provided that the same amount of computational work is spent to approximate the solutions of the involved linear systems).

**4. Further discussion.** In this section we will discuss a convenient way of incorporating preconditioning in the Jacobi-Davidson method. We will also discuss relations with other methods, e.g., shift and invert techniques, and we will try to get some insight into the behavior of the new method in situations where Davidson's method or shift and invert methods work well. This will make the differences between these methods clearer.

In the Jacobi-Davidson method we must solve (13), or equivalently

$$(14) \quad (I - u_k u_k^*)(A - \theta_k I)(I - u_k u_k^*)t = -r \quad \text{and} \quad t \perp u_k$$

(see Algorithm 1, in which we have skipped the index  $k$ ).

Equation (14) can be solved approximately by selecting some more easily invertible approximation for the operator  $(I - u_k u_k^*)(A - \theta_k I)(I - u_k u_k^*)$ , or by some (preconditioned) iterative method. If any approximation (preconditioner) is available, then this will most often be an approximation for  $A - \theta_k I$ .

However, the formulation in (14) is not very suitable for incorporating available approximations for  $A - \theta_k I$ . We will first discuss how to construct approximate solutions orthogonal to  $u_k$  straight from a given approximation for  $A - \theta_k I$  (1-step approximation: §4.1). Then we will propose how to compute such an approximated solution efficiently by a preconditioned iterative scheme (iterative approximation: §4.2).

**4.1. 1-step approximations.** A more convenient formulation for (14) is obtained as follows. We are interested in determining  $t \perp u_k$ , and for this  $t$  we have that

$$(I - u_k u_k^*)t = t,$$

and then it follows from (14) that

$$(15) \quad (A - \theta_k I)t - \varepsilon u_k = -r$$

or

$$(A - \theta_k I)t = \varepsilon u_k - r.$$

When we have a suitable preconditioner  $M \approx A - \theta_k I$  available, then we can compute an approximation  $\tilde{t}$  for  $t$ :

$$(16) \quad \tilde{t} = \varepsilon M^{-1} u_k - M^{-1} r.$$

The value of  $\varepsilon$  is determined by the requirement that  $\tilde{t}$  should be orthogonal to  $u_k$ :

$$(17) \quad \varepsilon = \frac{u_k^* M^{-1} r}{u_k^* M^{-1} u_k}.$$

Equation (16) leads to several interesting observations.

1. If we choose  $\varepsilon = 0$ , then we obtain the Davidson method (with preconditioner  $M$ ). In this case  $\tilde{t}$  will not be orthogonal to  $u_k$ .
2. If we choose  $\varepsilon$  as in (17), then we have an instance of the Jacobi–Davidson method. This approach has already been suggested in [18]. In that paper the method is obtained from a first-order correction approach for the eigensystem. Further experiments with this approach are reported in [27]. Note that this method requires two operations with the preconditioning matrix per iteration.
3. If  $M = A - \theta_k I$ , then (16) reduces to

$$t = \varepsilon (A - \theta_k I)^{-1} u_k - u_k.$$

Since  $t$  is made orthogonal to  $u_k$  afterwards, this choice is equivalent with  $t = (A - \theta_k I)^{-1} u_k$ . In this case the method is just mathematically equivalent to (accelerated) shift and invert iteration (with optimal shift). For symmetric  $A$  this is the (accelerated) inverse Rayleigh quotient method, which converges cubically [21]. In the unsymmetric case we have quadratical convergence [21], [19]. In view of the speed of convergence of shift and invert methods, it



may hardly be worthwhile to accelerate them in a “Davidson” manner: the overhead is significant and the gains may only be minor. Moreover, in finite precision arithmetic the vector  $(A - \theta_k I)^{-1} u_k$  may make a very small angle with  $u_k$ , so that it will be impossible then to compute a significant orthogonal search direction.

4. If  $M \neq A - \theta_k I$ , then with  $\tilde{t} = M^{-1} u_k$  we obtain an inexact shift and invert method with “Davidson” subspace acceleration. This method may be an attractive alternative for the previous one if the invert part cannot be carried out exactly. Also in this variant we have no orthogonality between  $\tilde{t}$  and  $u_k$ . If  $M$  is a good approximation for  $A - \theta_k I$  then  $M^{-1} u_k$  may also make a very small angle with  $u_k$ , so that effective subspace expansion will be impossible (as in 3).

The methods suggested in the first and the third observation are well known, and the question arises whether we may gain anything by the less well known second alternative (or the fourth one).

To get some insight into this matter, we consider a situation for which Davidson’s method converges rapidly, namely, when  $A$  is strongly diagonally dominant. We write

$$A = D_A + E,$$

where  $D_A$  denotes the diagonal of  $A$ , and we assume that  $\|E\|_\infty$  is small compared with  $\|D_A\|_\infty$  and that  $a_{11}$  is the largest diagonal element in absolute value (note that this also includes situations where only the largest diagonal element has relatively small off-diagonal elements in the same row and column).

We write  $u_k = e_1 + f$  and assume that  $\|f\|_\infty \ll \|e_1\| = 1$  (which is a natural assumption in this case).

Then for the coordinates  $(r)_i$  of  $r$  it follows that

$$(r)_1 = (a_{11} - \theta_k) + (Eu_k)_1 + (a_{11} - \theta_k)(f)_1$$

and

$$(r)_i = (Eu_k)_i + (a_{ii} - \theta_k)(f)_i.$$

Since  $\theta_k \approx a_{11}$ , this means that the coordinates  $(r)_i$  are not small relative to  $(r)_1$ . In the case that  $f = 0$  we even have that  $r = Eu_k$ , and  $(r)_1 = 0$  (since  $r \perp u_k = e_1$ ).

With Davidson’s method we obtain

$$\tilde{t} = (D_A - \theta_k I)^{-1} r = u_k + (D_A - \theta_k I)^{-1} Eu_k,$$

and the part  $(D_A - \theta_k I)^{-1} Eu_k$  of this vector will not be very small compared to  $u_k$  (for  $f = 0$  the component  $u_k$  even vanishes). This means that we may expect to recover this part in a large number of significant digits after orthogonalizing  $\tilde{t}$  with respect to  $u_k$ , and this makes Davidson’s method work well for diagonally dominant problems (since we expand the search space by a well-determined vector).

We have seen the effect of the component  $(D_A - \theta_k I)^{-1} r$ , and now we consider what happens with the component  $(D_A - \theta_k I)^{-1} u_k$  in the Jacobi–Davidson method for this situation. To this end we compute  $\tilde{t}$  as in observation 4 above:

$$\tilde{t} = (D_A - \theta_k I)^{-1} u_k.$$

For the coordinates of  $\tilde{t}$  we have that

$$(\tilde{t})_i = \frac{(u)_i}{a_{ii} - \theta_k},$$

and we see that  $\tilde{t}$  will make a very small angle with  $u_k \approx e_1$  (since  $a_{11} \approx \theta_k$ ). This implies that  $\tilde{t} - (\tilde{t}^* u_k) u_k$  may have only a few significant digits, and then it may be a waste of effort to expand the subspace with this (insignificant) vector. However, in the Jacobi–Davidson method we compute the new vector as

$$\tilde{t} = \varepsilon(D_A - \theta_k I)^{-1} u_k - (D_A - \theta_k I)^{-1} r_k.$$

The factor  $\varepsilon$  is well determined, and note that for our strongly diagonally dominant model problem we have in good approximation that

$$\begin{aligned} \|\varepsilon(D_A - \theta_k I)^{-1} u_k\|_2 &\lesssim \frac{\|u_k\|_2 \|(D_A - \theta_k I)^{-1} r\|_2}{\|u_k\|_2 \|(D_A - \theta_k I)^{-1} u_k\|_2} \|(D_A - \theta_k I)^{-1} u_k\|_2 \\ &= \|(D_A - \theta_k I)^{-1} r\|_2. \end{aligned}$$

Furthermore, since  $u_k \perp r$ , we have that  $\varepsilon(D_A - \theta_k I)^{-1} u_k$  and  $(D_A - \theta_k I)^{-1} r$  are not in the same direction, and therefore there will be hardly any cancellation in the computation of  $\tilde{t}$ . This means that  $\tilde{t}$  is well determined in finite precision and  $\tilde{t} \perp u_k$ .

Our discussion can be adapted for nondiagonally dominant matrices as well, when we restrict ourselves to situations where the approximations  $u_k$  and  $\theta_k$  are sufficiently close to their limit values and where we have good preconditioners (e.g., inner iteration methods).

We will illustrate our observations by a simple example taken from [3]: Example 5.1. In that example the matrix  $A$  of dimension 1000 has diagonal elements  $a_{j,j} = j$ . The elements on the sub- and super-diagonal ( $a_{j-1,j}$  and  $a_{j,j+1}$ ) are all equal to 0.5, as well as the elements  $a_{1,1000}$  and  $a_{1000,1}$ .

For this matrix we compute the largest eigenvalue ( $\approx 1000.225641$ ) with (a) the standard Lanczos method, (b) Davidson's method with diagonal preconditioning  $((D_A - \theta_k I)^{-1})$ , and (c) the Jacobi–Davidson method with the same diagonal preconditioning, carried out as in (16).

With the same starting vector as in [3] we obtain, of course, the same results: a slowly converging Lanczos process, a much faster Davidson process, and Jacobi–Davidson is just as fast as Davidson in this case. The reason for this is that the starting vector  $e_1 + e_{1000}$  for Davidson and  $\approx e_1 + 2000e_{1000}$  for Lanczos are quite good for these processes, and the values for  $\varepsilon$ , which describe the difference between (b) and (c), are very small in this case. Shift and invert Lanczos with shift 1001.0 takes 5 steps for full convergence, whereas Jacobi–Davidson with exact inverse for  $A - \theta_k I$  takes 3 steps.

In Table 1 we see the effects when we take a slightly different starting vector  $u_1 = (0.01, 0.01, \dots, 0.01, 1)^T$ , that is, we have taken a starting vector which still has a large component in the dominating eigenvector. This is reflected by the fact that the Ritz value in the first step of all three methods is equal to  $954.695 \dots$ . In practical situations we will often not have such good starting vectors available. The Lanczos process converges slowly again, as might be expected for this uniformly distributed spectrum. In view of our discussion in §2.3 we may interpret the new starting vector as a good starting vector for a perturbed matrix  $A$ . This implies that the diagonal preconditioner may not be expected to be a very good preconditioner. This is reflected by the very poor convergence behavior of Davidson's method. The difference with the Jacobi–Davidson method is now quite notable (see the values of  $\varepsilon$ ), and for this method we observe rather fast convergence again.

TABLE 1  
Approximation errors  $\lambda - \theta_k$ .

Iteration	Lanczos	Davidson	Jacobi-Davidson	$ \varepsilon $
0	0.45e+02	0.45e+02	0.45e+02	
1	0.56e+01	0.40e+02	0.25e+02	0.50e+02
2	0.16e+01	0.40e+02	0.74e+01	0.12e+03
3	0.71e+00	0.40e+02	0.15e+01	0.11e+02
4	0.43e+00	0.40e+02	0.14e+01	0.14e+01
5	0.32e+00	0.40e+02	0.55e-01	0.49e+00
6	0.26e+00	0.39e+02	0.13e-02	0.72e-01
7	0.24e+00	0.38e+02	0.29e-04	0.29e-02
8	0.22e+00	0.37e+02	0.33e-06	0.14e-03
9	0.21e+00	0.36e+02	0.25e-08	0.34e-05
10	0.20e+00	0.36e+02		
11	0.19e+00	0.35e+02		
12	0.19e+00	0.34e+02		
13	0.18e+00	0.33e+02		
14	0.17e+00	0.32e+02		
15	0.16e+00	0.31e+02		

Although this example may seem quite artificial, it displays quite nicely the behavior that we have seen in our experiments, as well as what we have tried to explain in our discussion.

In conclusion, Davidson's method works well in these situations where  $\tilde{t}$  does not have a strong component in the direction of  $u_k$ . Shift and invert approaches work well if the component in the direction of  $u$  in  $u_k$  is strongly increased. However, in this case this component may dominate so strongly (when we have a good preconditioner) that it prevents us from reconstructing in finite precision arithmetic a relevant orthogonal expansion for the search space. In this respect, the Jacobi-Davidson method is a compromise between the Davidson method and the accelerated (inexact) shift and invert method, since the factor  $\varepsilon$  properly controls the influence of  $u_k$  and makes sure that we construct the orthogonal expansion of the subspace correctly. In this view Jacobi-Davidson offers the best of two worlds, and this will be illustrated by our numerical examples.

**4.2. Iterative approximations.** If a preconditioned iterative method is used to solve (14), then, in each step of the linear solver, a "preconditioning equation" has to be solved.

If  $M$  is some approximation of  $A - \theta_k I$  then the projected matrix

$$M_d := (I - u_k u_k^*) M (I - u_k u_k^*)$$

can be taken as an approximation of  $(I - u_k u_k^*)(A - \theta_k I)(I - u_k u_k^*)$  and, in each iterative step, we will have to solve an equation of the form  $M_d z = y$ , where  $y$  is some given vector orthogonal to  $u_k$  and  $z \perp u_k$ , has to be computed. Of course,  $z$  can be computed as (cf. (16)–(17))

$$(M_d^{-1} y =) \quad z = \alpha M^{-1} u_k - M^{-1} y \quad \text{with} \quad \alpha = \frac{u_k^* M^{-1} y}{u_k^* M^{-1} u_k}.$$

In this approach, we have to solve, except for the first iteration step, only one system involving  $M$  in each iteration step. The inner product  $u_k^* M^{-1} u_k$ , to be computed only once, can also be used in all steps of the iteration process for (14).

The use of a (preconditioned) Krylov subspace iteration method for (14) does not lead to the same result as when we apply this iterative method to the two equations in (16) separately. For instance, if  $p$  is a polynomial such that  $p(A - \theta_k I) \approx (A - \theta_k I)^{-1}$  then, with  $M^{-1} = p(A - \theta_k I)$ , (16) can be used to find an approximate solution of (14) leading to

$$(18) \quad \tilde{t} = \varepsilon p(A - \theta_k I) u_k - p(A - \theta_k I) r = (I - u_k u_k^*) p(A - \theta_k I) (I - u_k u_k^*) r,$$

while using  $p$  directly for (14) would yield

$$(19) \quad \tilde{t} = p \left( (I - u_k u_k^*) (A - \theta_k I) (I - u_k u_k^*) \right) r.$$

Clearly, these expressions are not identical. For Krylov subspace methods that automatically (and implicitly) determine such polynomials, the differences may be even more significant. Most importantly, such a method for (14) would be aiming for an approximation of the inverse of  $(I - u_k u_k^*) (A - \theta_k I) (I - u_k u_k^*)$  on the space orthogonal to  $u_k$ , rather than for an approximation of  $(A - \theta_k I)^{-1}$  as the method for (16) would do. If  $\theta_k$  is an accurate approximation of the eigenvalue  $\lambda$ ,  $A - \theta_k I$  will be almost singular, while that will not be the case for the projected matrix  $(I - u_k u_k^*) (A - \theta_k I) (I - u_k u_k^*)$  (as a map on  $u_k^\perp$ , if  $\lambda$  is simple). This means that the iterative solution of (14) may be easier than iteratively solving systems such as  $(A - \theta_k I)z = y$ .

By iteratively solving (14) we expect more stable results: by putting the intermediate approximations orthogonal to  $u_k$  (as, for instance, in (19)) we may hope to have less cancellation by rounding errors than when putting only the final approximation  $\tilde{t}$  orthogonal to  $u_k$  (as, for instance, in (18)).

We cannot answer the question of how accurately (14) should be solved in order to have convergence for the Jacobi–Davidson method. Our experiences, as well as experiences reported in [2], seem to indicate that even a modest error reduction in the solution of (14) suffices and more work spent in this (inner) iteration for (14) often leads to a reduction in the number of Jacobi–Davidson iteration steps. For some numerical evidence, see §6.

**5. Jacobi–Davidson with harmonic Ritz values.** In the previous sections we have used the Galerkin approach for the approximation of eigenvectors and eigenvalues. In this approach,  $H_k$  is the orthogonal projection of the matrix  $A$  onto  $\mathcal{V}_k = \{v_1, \dots, v_k\}$ , and its eigenvalues are called the Ritz values of  $A$  with respect to  $\mathcal{V}_k$  [22]. The Ritz values converge monotonically to the extremal eigenvalues when  $A$  is symmetric. If  $A$  is nonsymmetric, the convergence is in general not monotonical, but the convergence behavior is still often quite regular with respect to the extremal eigenvalues. Interesting observations for the nonsymmetric case have been made in [24], [10].

For the “interior” (the non extremal) eigenvalues the situation is less clear. The convergence can be very irregular, even in the symmetric situation (due to rounding errors). This behavior makes it difficult to approximate interior eigenvalues or to design algorithms that select the correct Ritz values and handle rounding errors well (see, e.g., [24]).

In [14] the author suggested using a minimum residual approach for the computation of interior eigenvalues. We follow a slightly different approach which leads to identical results for symmetric matrices. In this approach, as we will show in the next section, we use orthogonal projections onto  $A\mathcal{V}_k$ . The obtained eigenvalue

approximations differ from the standard Ritz values for  $A$ . In a recent paper [20], these approximations were called *harmonic Ritz values*, and they were identified as inverses of Ritz approximations for the inverse of  $A$ . It was also shown in [20] that, for symmetric matrices, these harmonic Ritz values exhibit a monotonic convergence behavior with respect to the eigenvalues with smallest absolute value. This further supports the observation made in [14] that, for the approximation of “interior” eigenvalues (close to some  $\mu \in \mathbb{C}$ ), more regular convergence behavior with the harmonic Ritz values (of  $A - \mu I$ ) than with the Ritz values may be expected.

In [20] the harmonic Ritz values for symmetric matrices are discussed. The non-symmetric case has been considered in [30], [8]. However, in all these papers the discussion is restricted to harmonic Ritz values of  $A$  with respect to Krylov subspaces. In [14] the harmonic Ritz values are considered for more general subspaces associated with symmetric matrices. The approach is based on a generalized Rayleigh–Ritz procedure, and it is pointed out in [14] that the harmonic Ritz values are to be preferred for the Davidson method when aiming for interior eigenvalues.

In connection with the Jacobi–Davidson method for unsymmetric matrices, we propose a slightly more general approach based on projections. To this end, as well as for the introduction of notations that we will also need later, we discuss to some extent the harmonic Ritz values in §5.1.

### 5.1. Harmonic Ritz values on general subspaces.

*Ritz values.* If  $\mathcal{V}_k$  is a linear subspace of  $\mathbb{C}^n$  then  $\theta_k$  is a *Ritz value* of  $A$  with respect to  $\mathcal{V}_k$  with *Ritz vector*  $u_k$  if

$$(20) \quad u_k \in \mathcal{V}_k, u_k \neq 0, \quad Au_k - \theta_k u_k \perp \mathcal{V}_k.$$

How well the Ritz pair  $(\theta_k, u_k)$  approximates an eigenpair  $(\lambda, w)$  of  $A$  depends on the angle between  $w$  and  $\mathcal{V}_k$ .

In practical computations one usually computes Ritz values with respect to a growing sequence of subspaces  $\mathcal{V}_k$  (that is,  $\mathcal{V}_k \subset \mathcal{V}_{k+1}$  and  $\dim(\mathcal{V}_k) < \dim(\mathcal{V}_{k+1})$ ).

If  $A$  is normal, then any Ritz value is in the convex hull of the spectrum of  $A$ : any Ritz value is a mean (convex combination) of eigenvalues. For normal matrices, at least, this helps to explain the often regular convergence of extremal Ritz values with respect to extremal eigenvalues. For further discussions on the convergence behavior of Ritz values (for symmetric matrices), see [22], [29].

*Harmonic Ritz values.* A value  $\tilde{\theta}_k \in \mathbb{C}$  is a *harmonic Ritz value* of  $A$  with respect to some linear subspace  $\mathcal{W}_k$  if  $\tilde{\theta}_k^{-1}$  is a Ritz value of  $A^{-1}$  with respect to  $\mathcal{W}_k$  [20].

For normal matrices,  $\tilde{\theta}_k^{-1}$  is in the convex hull of the collection of  $\lambda^{-1}$ 's, where  $\lambda$  is an eigenvalue of  $A$ : any harmonic Ritz value is a harmonic mean of eigenvalues. This property explains their name and, at least for normal matrices, it explains why we may expect a more regular convergence behavior of harmonic Ritz values with respect to the eigenvalues that are closest to the origin.

Of course, we would like to avoid computing  $A^{-1}$  or solving linear systems involving  $A$ . The following theorem gives a clue about how that can be done.

**THEOREM 5.1.** *Let  $\mathcal{V}_k$  be some  $k$ -dimensional subspace with basis  $v_1, \dots, v_k$ . A value  $\tilde{\theta}_k \in \mathbb{C}$  is a harmonic Ritz value of  $A$  with respect to the subspace  $\mathcal{W}_k := A\mathcal{V}_k$  if and only if*

$$(21) \quad A\tilde{u}_k - \tilde{\theta}_k \tilde{u}_k \perp A\mathcal{V}_k \quad \text{for some} \quad \tilde{u}_k \in \mathcal{V}_k, \tilde{u}_k \neq 0.$$

If  $w_1, \dots, w_k$  span  $AV_k$  then,<sup>4</sup> with

$$V_k := [v_1 | \dots | v_k], \quad W_k := [w_1 | \dots | w_k], \quad \text{and} \quad \tilde{H}_k := (W_k^* V_k)^{-1} W_k^* A V_k,$$

property (21) is equivalent to

$$(22) \quad \tilde{H}_k s = \tilde{\theta}_k s \quad \text{for some} \quad s \in \mathbb{C}^k, s \neq 0 \quad (\text{and} \quad \tilde{u}_k = V_k s):$$

the eigenvalues of the  $k \times k$  matrix  $\tilde{H}_k$  are the harmonic Ritz values of  $A$ .

*Proof.* By (20),  $(\tilde{\theta}_k^{-1}, A\tilde{u}_k)$  is a Ritz pair of  $A^{-1}$  with respect to  $AV_k$  if and only if

$$(A^{-1} - \tilde{\theta}_k^{-1} I) A \tilde{u}_k \perp AV_k$$

for a  $\tilde{u}_k \in V_k$ ,  $\tilde{u}_k \neq 0$ .

Since  $(A^{-1} - \tilde{\theta}_k^{-1} I) A \tilde{u}_k = -\tilde{\theta}_k^{-1} (A \tilde{u}_k - \tilde{\theta}_k \tilde{u}_k)$  we have the first property of the theorem.

For the second part of the theorem, note that (21) is equivalent to

$$(23) \quad AV_k s - \tilde{\theta}_k V_k s \perp W_k \quad \text{for an} \quad s \neq 0,$$

which is equivalent to

$$W_k^* AV_k s - \tilde{\theta}_k (W_k^* V_k) s = 0$$

or  $\tilde{H}_k s - \tilde{\theta}_k s = 0$ .  $\square$

We will call the vector  $\tilde{u}_k$  in (21) the *harmonic Ritz vector* associated with the *harmonic Ritz value*  $\tilde{\theta}_k$  and  $(\tilde{\theta}_k, \tilde{u}_k)$  is a *harmonic Ritz pair*.

In the context of Krylov subspace methods (Arnoldi or Lanczos),  $V_k$  is the Krylov subspace  $\mathcal{K}_k(A; v_1)$ . The  $v_j$  are orthonormal and such that  $v_1, \dots, v_i$  span  $\mathcal{K}_i(A; v_1)$  for  $i = 1, 2, \dots$ . Then  $AV_k = V_{k+1} H_{k+1,k}$ , with  $H_{k+1,k}$  a  $(k+1) \times k$  upper Hessenberg matrix.

The elements of  $H_{k+1,k}$  follow from the orthogonalization procedure for the Krylov subspace basis vectors. In this situation, with  $H_{k,k}$  the upper  $k \times k$  block of  $H_{k+1,k}$ , we see that

$$\begin{aligned} (W_k^* V_k)^{-1} W_k^* A V_k &= (H_{k+1,k}^* V_{k+1}^* V_k)^{-1} H_{k+1,k}^* V_{k+1}^* V_{k+1} H_{k+1,k} \\ &= H_{k,k}^{*-1} H_{k+1,k}^* H_{k+1,k}. \end{aligned}$$

Since  $H_{k+1,k} = H_{k,k} + \beta e_{k+1} e_k^*$ , where  $\beta$  is equal to the element in position  $(k+1, k)$  of  $H_{k+1,k}$ , the harmonic Ritz values can be computed from a matrix which is a rank-one update of  $H_{k,k}$ :

$$\tilde{H}_k = H_{k,k}^{*-1} H_{k+1,k}^* H_{k+1,k} = H_{k,k} + |\beta|^2 H_{k,k}^{*-1} e_k e_k^*.$$

In [8], the author is interested in quasi-kernel polynomials (e.g., GMRES and quasi-minimal residual (QMR) polynomials). The zeros of these polynomials are harmonic Ritz values with respect to Krylov subspaces. This follows from Corollary 5.3 in [8], where these zeros are shown to be the eigenvalues of  $H_{k,k}^{*-1} H_{k+1,k}^* H_{k+1,k}$ . However,

<sup>4</sup> If  $AV_k$  has dimension less than  $k$ , then this subspace contains an eigenvector of  $A$ ; this situation is often referred to as a lucky breakdown. We do not consider this situation here.

in that paper these zeros are not interpreted as the Ritz values of  $A^{-1}$  with respect to some Krylov subspace.

In the context of Davidson's method we have more general subspaces  $\mathcal{V}_k$  and  $\mathcal{W}_k = A\mathcal{V}_k$ . According to Theorem 5.1 we have to construct the matrix  $\tilde{H}_k$  (which will not be Hessenberg in general), and this can be accomplished by either constructing an orthonormal basis for  $A\mathcal{V}_k$  (similar to Arnoldi's method) or by constructing bi-orthogonal bases for  $\mathcal{V}_k$  and  $A\mathcal{V}_k$  (similar to the bi-Lanczos method). We will consider this in more detail in §5.2.

## 5.2. The computation of the harmonic Ritz values.

**5.2.1. Bi-orthogonal basis construction.** In our algorithms, we expand the subspace  $\mathcal{V}_k$  by one vector in each sweep of the iteration. We proceed as follows.

Suppose that  $v_1, \dots, v_k$  span  $\mathcal{V}_k$  and that  $w_1, \dots, w_k$  span  $A\mathcal{V}_k$ , in such a way that, with  $V_k := [v_1 | \dots | v_k]$  and  $W_k := [w_1 | \dots | w_k]$ ,

$$AV_k = W_k$$

and

$$L_k := W_k^* V_k \quad \text{is lower triangular}$$

(in this case we say that  $W_k$  and  $V_k$  are bi-orthogonal).

According to Theorem 5.1 the harmonic Ritz values are the eigenvalues of

$$\tilde{H}_k := L_k^{-1} \hat{H}_k, \quad \text{where} \quad \hat{H}_k := W_k^* W_k.$$

Hence, if  $(\tilde{\theta}_k, s)$  is an eigenpair of  $\tilde{H}_k$  then  $(\tilde{\theta}_k, V_k s)$  is a harmonic Ritz pair.

Let  $t$  be the vector by which we want to expand the subspace  $\mathcal{V}_k$ .

First, we bi-orthogonalize  $t$  with respect to  $V_k$  and  $W_k$ :

$$(24) \quad \tilde{t} := t - V_k L_k^{-1} W_k^* t \quad \text{and} \quad v_{k+1} := \frac{\tilde{t}}{\|\tilde{t}\|_2}.$$

Then  $v_{k+1}$  is our next basis vector in the  $\mathcal{V}$ -space and we expand  $V_k$  by  $v_{k+1}$  to  $V_{k+1}$ :  $V_{k+1} := [V_k | v_{k+1}]$ .

Next, we compute  $w_{k+1} := Av_{k+1}$ , our next basis vector in the  $A\mathcal{V}$ -space, and we expand  $W_k$  to  $W_{k+1}$ .

Then the vector  $w_{k+1}^* V_{k+1}$  is computed and  $L_k$  is expanded to  $L_{k+1}$  by this  $k+1$  row vector. Finally, we compute  $W_{k+1}^* w_{k+1}$  as the new column of  $\hat{H}_{k+1}$ . By symmetry, we automatically have the new row of  $\hat{H}_{k+1}$ .

Since  $L_k^* = V_k^* AV_k$  and  $L_k^*$  is upper triangular, we see that  $L_k$  is diagonal if  $A$  is self-adjoint.

The formulation of the bi-orthogonalization step (24) does not allow for the use of modified Gram-Schmidt orthogonalization (due to the  $k \times k$  matrix  $L_k$ ). We can incorporate  $L_k$  into  $W_k$  by working with  $\tilde{W}_k := W_k L_k^{-*}$  instead of with  $W_k$ , and then modified Gram-Schmidt is possible:

$$v^{(1)} = t, \quad v^{(i+1)} = v^{(i)} - v_i \tilde{w}_i^* v^{(i)} \quad (i = 1, \dots, k-1), \quad \tilde{t} = v^{(k)}.$$

However, in this approach we have to update  $\tilde{W}_k$ , which would require  $k$  additional long vector updates.

**5.2.2. Orthogonal basis construction.** For stability reasons one might prefer to work with an orthogonal basis rather than with bi-orthogonal ones. In the context of harmonic Ritz values, an orthonormal basis for the image space  $AV_k$  is attractive:

$$(25) \quad \begin{aligned} w &= At, \quad \tilde{w} := w - W_k W_k^* w, \quad \text{and} \quad w_{k+1} := \frac{\tilde{w}}{\|\tilde{w}\|_2}, \\ \tilde{t} &:= t - V_k W_k^* w \quad \text{and} \quad v_{k+1} := \frac{\tilde{t}}{\|\tilde{t}\|_2}. \end{aligned}$$

Then  $W_k^* W_k = I$ ,  $W_k = AV_k$ , and (cf. Theorem 5.1)

$$(26) \quad \tilde{H}_k = (W_k^* V_k)^{-1} W_k^* A V_k = (W_k^* V_k)^{-1}.$$

It is not necessary to invert  $W_k^* V_k$  since the harmonic Ritz values are simply the inverses of the eigenvalues of  $W_k^* V_k$ . The construction of an orthogonal basis for  $AV_k$  can be done with modified Gram–Schmidt.

Finally, note that  $\tilde{H}_k^{-1} = W_k^* V_k = W_k^* A^{-1} W_k$ , and we explicitly have the matrix of the projection of  $A^{-1}$  with respect to an orthonormal basis of  $W_k$ . This again reveals how the harmonic Ritz values appear as inverses of Ritz values of  $A^{-1}$  with respect to  $AV_k$ .

**5.3. A restart strategy.** In the Jacobi–Davidson algorithm, Algorithm 1, it is suggested, just as for the original Davidson method, to restart simply by taking the Ritz vector  $u_m$  computed so far as a new initial guess. However, the process may construct a new search space that has considerable overlap with the previous one; this phenomenon is well known for the restarted power method and the restarted Arnoldi (without deflation) and it may lead to a reduced speed of convergence or even to stagnation. One may try to prevent this by retaining part of the search space, i.e., by returning to step 3 of Algorithm 1 with a well chosen  $\ell$ -dimensional subspace of the span of  $v_1, \dots, v_m$  for some  $\ell > 1$ . With our simple restart, we expect that the process will also construct vectors with strong components in directions of eigenvectors associated with eigenvalues close to the wanted eigenvalue. And this is just the kind of information that we have discarded at the point of restart.

This suggests a strategy of retaining  $\ell$  Ritz vectors associated with the Ritz values closest to this eigenvalue as well (including the Ritz vector  $u_m$  that is the approximation for the desired eigenvector). In Algorithm 1, these would be the  $\ell$  largest Ritz values. A similar restart strategy can be used for the harmonic Ritz values and, say, bi-orthogonalization: for the initial matrices  $V_\ell$  and  $W_\ell$  after restart we should take care that  $W_\ell = AV_\ell$  and the matrices should be bi-orthogonal (i.e.,  $W_\ell^* V_\ell$  should be lower triangular).

**5.4. The use of the harmonic Ritz values.** According to our approaches for the computation of harmonic Ritz values in §5.2, there are two variants for an algorithm that exploits the convergence properties of the harmonic Ritz values toward the eigenvalues closest to the origin. Of course, these algorithms can also be used to compute eigenvalues that are close to some  $\mu \in \mathbb{C}$ . In that case one should work with  $A - \mu I$  instead of  $A$ .

We start with the variant based on bi-orthogonalization.

**5.4.1. Jacobi–Davidson with bi-orthogonal basis.** When working with harmonic Ritz values, we have to be careful in applying Jacobi’s expansion technique. If  $(\theta_k, \tilde{u}_k)$  is a harmonic Ritz pair of  $A$  then  $r = A\tilde{u}_k - \tilde{\theta}_k \tilde{u}_k$  is orthogonal to  $A\tilde{u}_k$ , whereas in our discussion about the new Jacobi–Davidson method with regular Ritz



values (cf. §3) the vector  $r$  was orthogonal to  $u_k$ . However, we can follow Jacobi's approach here as well by using a skew basis or a skew projection. The update for  $\tilde{u}_k$  should be in the space orthogonal to  $\hat{u} := A\tilde{u}_k$ . If  $\lambda$  is the eigenvalue of  $A$  closest to the harmonic Ritz value  $\tilde{\theta}_k$ , then the optimal update is the solution  $v$  of

$$(27) \quad (B - \lambda I)v = -r \quad \text{where} \quad B := \left( I - \frac{\tilde{u}_k \hat{u}^*}{\hat{u}^* \tilde{u}_k} \right) A \left( I - \frac{\tilde{u}_k \hat{u}^*}{\hat{u}^* \tilde{u}_k} \right).$$

In our algorithm we will solve this equation (27) approximately. To be more precise, we solve approximately

$$(28) \quad \left( I - \frac{\tilde{u}_k \hat{u}^*}{\hat{u}^* \tilde{u}_k} \right) (A - \tilde{\theta}_k I) \left( I - \frac{\tilde{u}_k \hat{u}^*}{\hat{u}^* \tilde{u}_k} \right) t = -r.$$

Note that  $\hat{u}$  can be computed without an additional matrix vector product with  $A$  since  $\hat{u} = A\tilde{u}_k = AV_k s = W_k U_k s$  (if  $W_k U_k = AV_k$ , where  $U_k$  is a matrix of order  $k$ ).

The above considerations lead to Algorithm 2.

**5.4.2. Jacobi–Davidson with orthogonal basis.** If we want to work with an orthonormal basis for  $AV_k$  then we may proceed as follows.

Let  $v_1, v_2, \dots, v_k$  be the Jacobi–Davidson vectors obtained after  $k$  steps. Then we orthonormalize the set  $Av_1, Av_2, \dots, Av_k$  (as in (25)). The eigenvalues of  $\tilde{H}_k$  (cf. (26)) are the harmonic Ritz values of  $A$ , and let  $\tilde{\theta}_k$  be the one of interest, with corresponding harmonic Ritz vector  $\tilde{u}_k = V_k s$  ( $s$  is the eigenvector of  $\tilde{H}_k$ , normalized such that  $\|A\tilde{u}_k\|_2 = 1$ ).

Since  $(\tilde{\theta}_k^{-1}, A\tilde{u}_k)$  is a Ritz pair of  $A^{-1}$  with respect to  $AV_k$ , we have with  $z := A\tilde{u}_k$  that  $w := A^{-1}z - \tilde{\theta}_k^{-1}z$  is orthogonal to  $z$ , and although we do not have  $A^{-1}$  available, the vector  $w$  can be efficiently computed from

$$(29) \quad w = A^{-1}z - \tilde{\theta}_k^{-1}z = A^{-1}AV_k s - \tilde{\theta}_k^{-1}z = V_k s - \tilde{\theta}_k^{-1}z.$$

The orthonormal basis set for  $AV_k$  should be expanded by a suitable vector  $At$ , which, according to our Jacobi–Davidson approach, is computed approximately from

$$(30) \quad (I - zz^*)(A^{-1} - \tilde{\theta}_k^{-1}I)(I - zz^*)At = -w.$$

Also in this case we can avoid working with  $A^{-1}$ , since (30) reduces to

$$(31) \quad (I - zz^*)(A - \tilde{\theta}_k I)(I - V_k s z^* A)t = \tilde{\theta}_k w.$$

We may not expect any difference between the use of (28) and (31) when we use GMRES as the inner iteration, as we will see now.

Since  $\|z\|_2 = 1$  and  $w \perp z = A\tilde{u}_k$ , we see that  $1 = z^* z = \tilde{\theta}_k z^* \tilde{u}_k$ . Furthermore, we have that  $\tilde{\theta}_k w = \tilde{\theta}_k \tilde{u}_k - A\tilde{u}_k = -r$ .

It can be shown that the operator in the left-hand side of (28) and the one in the left-hand side of (31) differ only by a rank-one matrix of the form  $r(2(A - \tilde{\theta}_k I)^* \tilde{u}_k)^*$ . Therefore, the operators generate identical Krylov subspaces if, in both cases,  $r$  is the first Krylov subspace vector: Krylov subspace methods like GMRES( $m$ ) with initial approximation  $x_1 = 0$  lead, in exact arithmetic, to identical approximate solutions when used to solve equations (28) and (31).

It is not yet clear which of the two approaches will be more efficient in practical situations. Much depends upon how the projected systems are approximately solved.

ALGORITHM 2. *The Jacobi–Davidson algorithm with harmonic Ritz values and bi-orthogonalization.*

1. **Start:** Choose an initial nontrivial vector  $v$ .

- Compute  $v_1 = v/\|v\|_2$ ,  $w_1 = Av_1$ ,  $l_{11} = w_1^* v_1$ ,  $h_{11} = w_1^* w_1$ , set  $\ell = 1$ ,  $V_1 = [v_1]$ ,  $W_1 = [w_1]$ ,  $L_1 = [l_{11}]$ ,  $\hat{H}_1 = [h_{11}]$ ,  $\tilde{u} = v_1$ ,  $\hat{u} = w_1$ ,  $\tilde{\theta} = h_{11}/l_{11}$ , compute  $r = \hat{u} - \tilde{\theta}\tilde{u}$ .

2. **Iterate:** Until convergence do:

3. **Inner loop:** For  $k = \ell, \dots, m-1$  do:

- Solve approximately  $t \perp \hat{u}$ ,

$$\left(I - \frac{\tilde{u}\tilde{u}^*}{\tilde{u}^*\tilde{u}}\right)(A - \tilde{\theta}I)\left(I - \frac{\tilde{u}\tilde{u}^*}{\tilde{u}^*\tilde{u}}\right)t = -r.$$

- Bi-orthogonalize  $t$  against  $V_k$  and  $W_k$   
( $\tilde{t} = t - V_k L_k^{-1} W_k^* t$ ,  $v_{k+1} = \tilde{t}/\|\tilde{t}\|_2$ )  
and expand  $V_k$  with this vector to  $V_{k+1}$ .
- Compute  $w_{k+1} := Av_{k+1}$   
and expand  $W_k$  with this vector to  $W_{k+1}$ .
- Compute  $w_{k+1}^* V_{k+1}$ , the last row vector of  $L_{k+1} := W_{k+1}^* V_{k+1}$ ,  
compute  $w_{k+1}^* W_{k+1}$ , the last row vector of  $\hat{H}_{k+1} := W_{k+1}^* W_{k+1}$ ,  
its adjoint is the last column of  $\hat{H}_{k+1}$ .  $\tilde{H}_{k+1} := L_{k+1}^{-1} \hat{H}_{k+1}$ .
- Compute the smallest eigenpair  $(\tilde{\theta}, s)$  of  $\tilde{H}_{k+1}$ .
- Compute the harmonic Ritz vector  $\tilde{u} := V_{k+1}s/\|V_{k+1}s\|_2$ ,  
compute  $\hat{u} := A\tilde{u}$  ( $= W_{k+1}s/\|W_{k+1}s\|_2$ ), and  
the associated residual vector  $r := \hat{u} - \tilde{\theta}\tilde{u}$ .
- Test for convergence. Stop if satisfied.

4. **Restart:** Choose an appropriate value for  $\ell < m$ .

- Compute the smallest  $\ell$  eigenvalues of  $\tilde{H}_m$  and  
the matrix  $Y$  with columns the associated eigenvectors.  
Orthogonalize  $Y$  with respect to  $L_m$  (cf. §5.2.1):  
 $Y = ZR$  with  $R$  upper triangular and  $Z^* L_m Z$  lower triangular.
- Set  $V_\ell := V_m Z$ ,  $W_\ell := W_m Z$ ,  $L_\ell := Z^* L_m Z$ ,  $\hat{H}_\ell := Z^* \hat{H}_m Z$ , and goto 3.

**6. Numerical examples.** The new algorithm has already been successfully applied in applications from chemistry (in which Davidson's method was the preferred one before) and magnetohydrodynamics (MHD) models. For reports on these experiences, see [2], [28].

The simple examples that we will present here should be seen as illustrations only of the new approach. The examples have been coded in MATLAB and have been executed on a SUN SPARC workstation in about 15 decimals working precision. Most of our examples are for symmetric matrices, since it was easier for us to check the behavior of the Ritz values in this case, but our codes did not take advantage of this fact for the generation of the matrices  $H_k$ .

*Example 1.* We start with a simple tridiagonal diagonally dominant matrix  $A$ ,

with diagonal elements 2.4 and off-diagonal elements 1, of order  $n = 100$ . The starting vector is taken to be the vector with all 1's, scaled to unit length.

In Davidson's method, the suggested approach is to approximate the inverse of  $A - \theta_k I$ , and in this example we take just  $(A - \theta_k I)^{-1}$ , which is the best one can do if one wants to approximate well. Note that the standard choice  $(D_A - \theta_k I)$  leads to (almost) Arnoldi's method, which converges only very slowly (for  $i \ll n$ ) in this case.

In Figure 1 we have plotted the log of  $|\lambda - \theta_k|$  as the dashed curve, and we see that this indeed almost leads to stagnation (some progress is made, since we have computed the inverse in floating point arithmetic). If, however, we use the Jacobi-Davidson method (as in §3), again with exact inversion of the projected operator  $B - \theta_k I$ , then we observe (the lower curve) very fast (cubical?) convergence, just as expected.

Of course, solving the involved linear systems exactly is usually too expensive, and it might be more realistic to investigate what happens if we take more crude approximations for the operators involved. For the Davidson method we take 5 steps of GMRES for the approximation of the solution of  $(A - \theta_k I)v = r$ , and for Jacobi-Davidson we also take 5 steps of GMRES for the approximate solution of the projected system (13). The results are given in Figure 2 (the dashed curve represents the results for the Davidson method).

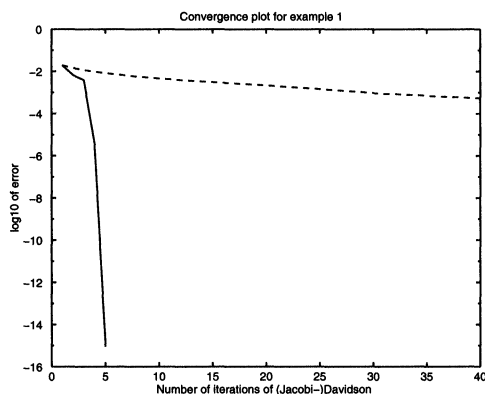


FIG. 1. Convergence of Ritz values with exact inverses.

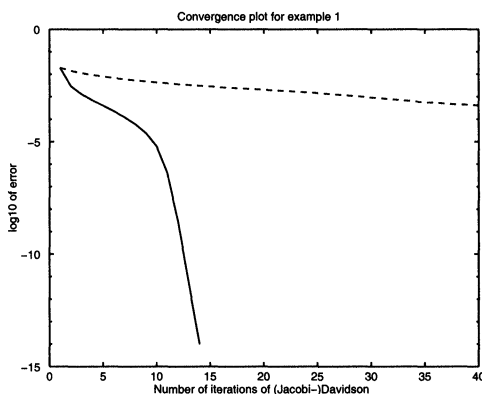


FIG. 2. Convergence of Ritz values with approximate inverses.

Again we see that for this moderately diagonally dominant matrix it is attractive to work with the Jacobi-Davidson method. Note that the 5 GMRES steps for the approximate inversion step are not sufficient to have quadratical convergence, but the linear convergence takes place with a very small convergence factor.

We also learn from this example that Krylov subspace methods may not be expected to make good preconditioners for Davidson's method: with a few steps one suffers from the fact that  $A - \theta_k I$  has a very small eigenvalue, and if carried out to high precision (almost) stagnation is to be expected. The Jacobi-Davidson method does not have these problems, since the projected operator  $B - \theta_k I$  does not have a small eigenvalue (unless the eigenvalue  $\lambda$  is close to some other eigenvalue of  $A$ ).

*Example 2.* Our second example is still highly artificial, but here we try to mimic more or less what happens when a matrix is not diagonally dominant. The matrix  $A$

is constructed as  $B = Q_1 A Q_1$ , with

$$A = \text{tridiag}(-1, 2, -1),$$

and  $Q_1$  is a Householder transformation. The order of the matrices is 100. Note that the distribution of the eigenvalues at the upper and lower ends of the spectrum is not particularly favorable for Krylov subspace methods since they are not well separated in a relative sense.

For those who wish to repeat our experiments we add that the Householder vector  $h$  was chosen with elements  $h_j = \sqrt{j + .45}$ ,  $j = 1, 2, \dots, 100$ . The starting vector for the iteration algorithms was chosen as a vector with all elements equal to 1. Furthermore, we restarted the outer iterations after each 20 steps, which represents a usual strategy in practical situations.

The Davidson algorithm (with  $D_A - \theta_k I$ ) needed 565 iteration steps to find the largest eigenvalue

$$\lambda = 3.9990325 \dots$$

to almost working precision.

In the Jacobi–Davidson algorithm we did the inner iterations, necessary for solving (13) approximately, with 5 steps of GMRES. This time we needed 65 outer iterations (i.e., 320 inner iteration steps). The inner iteration method (GMRES), as well as the number of steps (5 steps), has been chosen arbitrarily. In actual computations one may choose any appropriate means to approximate the solution of the projected system (13), such as, e.g., the incomplete LU (ILU).

*Example 3.* This example illustrates that our new algorithms (in §§3 and 5) may also be used for the computation of interior eigenvalues. In this example we compute an approximation for the eigenvalue of smallest absolute value. For this purpose the Jacobi–Davidson algorithm that uses Ritz values (§3) is modified; instead of computing the largest eigenpair  $(\theta_k, u_k)$  of  $H_k$  we compute the one with smallest absolute value for the Ritz value.

Again we take a simple matrix:  $A$  is the  $100 \times 100$  diagonal matrix with spectrum

$$\left\{ t^2 - 0.8 \mid t = \frac{j}{100}, j = 1, \dots, 100 \right\}.$$

All coordinates of the starting vector  $v_1$  are equal to 1. We solve the projected equations approximately using 8 steps of GMRES. We do not use restarts for the (outer) iterations.

Note that our algorithms, like any Krylov subspace method, do not take advantage of the fact that  $A$  is diagonal as long as we do not use diagonal preconditioning. Indeed, with  $D = Q^T A Q$  and  $y = Qv$ , we see that  $\mathcal{K}_i(D; v) = Q^T \mathcal{K}_i(A; y)$ , so that except for an orthogonal transformation the same subspaces are generated. In particular, the Rayleigh quotients, of which the Ritz values are the local extrema for symmetric matrices, are the same for both subspaces. This means that if the starting vectors are properly related as  $y = Qv$ , then the Krylov method with  $D$  and  $v$  leads to the same convergence history as the method with  $A$  and  $y$ .

We have also used the harmonic Ritz value variant of Jacobi–Davidson as in §5. Figures 3 and 5 show the convergence history of  $\log_{10} \|r_k\|_2$ , where the residual  $r_k$ , at step  $k$ , is either  $Au_k - \theta_k u_k$  for the Ritz pair with Ritz value of smallest absolute value, or  $A\tilde{u}_k - \tilde{\theta}_k \tilde{u}_k$  for the harmonic Ritz pair with harmonic Ritz value of smallest

absolute value. Along the vertical axis we have  $\log_{10} \|r_k\|_2$  and we have the iteration number  $k$  along the horizontal axis. Figures 4 and 6 show the convergence history of all (harmonic) Ritz values (indicated by  $\cdot$ ) in the interval indicated along the vertical axis. Again, we have the number of iteration steps along the horizontal axis. We have marked the positions of the eigenvalues (with  $+$ ) at the right of Figures 4 and 6.

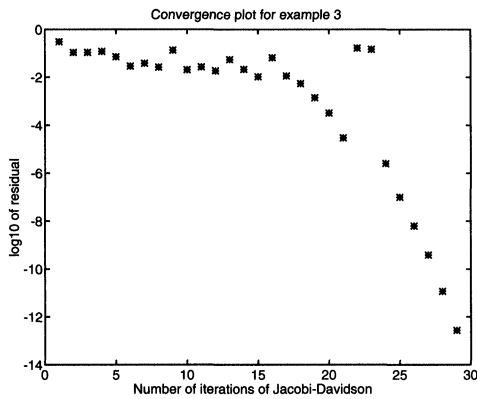


FIG. 3. *Convergence residuals using Ritz values.*

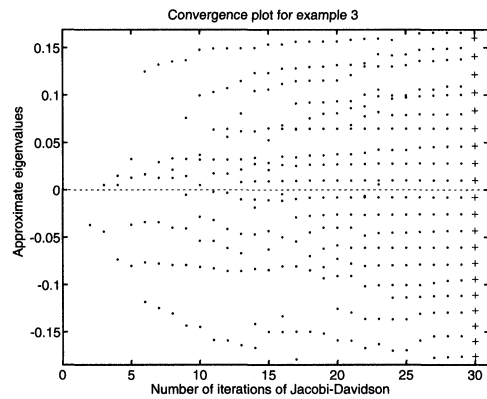


FIG. 4. *The Ritz values.*

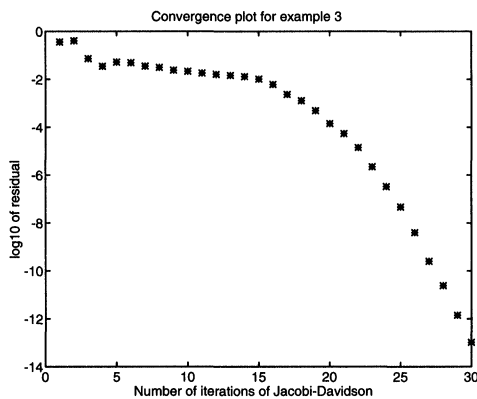


FIG. 5. *Convergence residuals using harmonic Ritz values.*

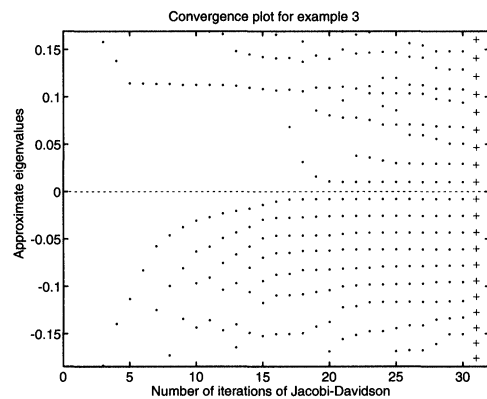


FIG. 6. *The harmonic Ritz values.*

For this example, the algorithm based on Ritz values (see Figures 3 and 4) converges about as fast as the algorithm based on harmonic Ritz values (see Figures 5 and 6), but the convergence history with harmonic Ritz values is much smoother. The difference in smoothness seems to also be typical for other examples. This fact can be exploited for the construction of restart strategies and stopping criteria. From experiments, we have learned that restarting for the “Ritz value” algorithm can be quite problematic; see also [14] for similar observations.

*Example 4.* In our previous examples the matrices  $A$  are symmetric. However, our algorithms are not restricted to the symmetric case and may also be used for the approximation of nonreal eigenvalues. In this example, we use complex arithmetic.

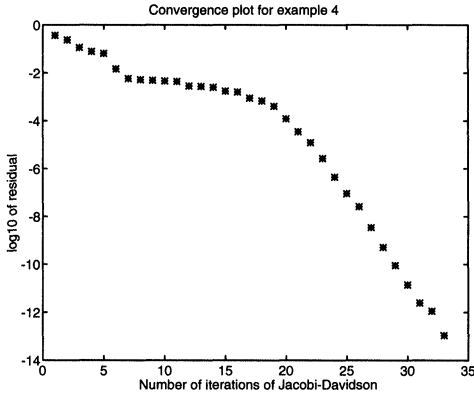


FIG. 7. *Convergence residuals using harmonic Ritz values.*

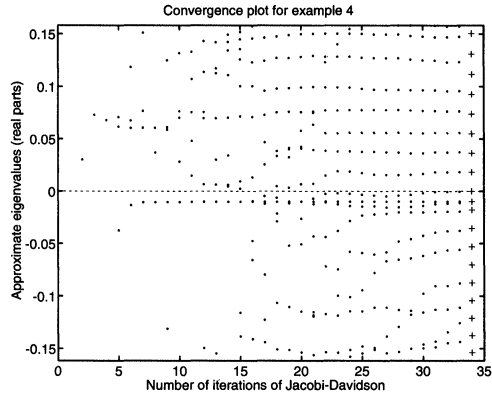


FIG. 8. *Real parts of harmonic Ritz values.*

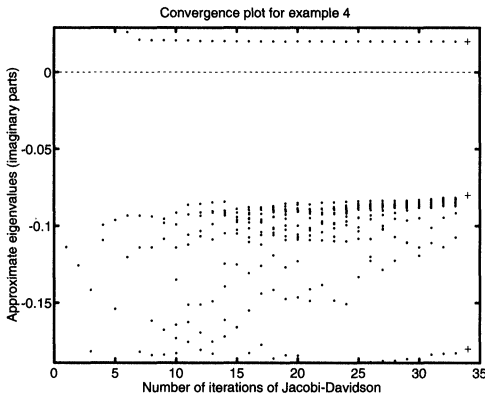


FIG. 9. *Imaginary parts of harmonic Ritz values.*

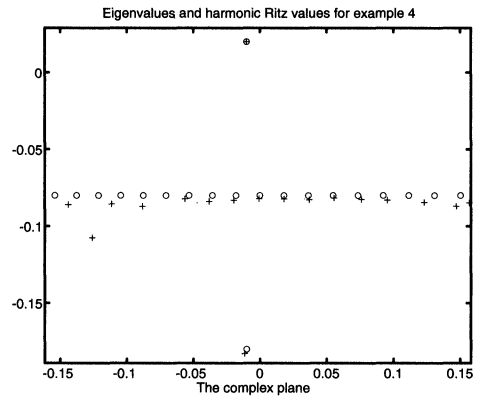


FIG. 10. *Harmonic Ritz values in  $\mathbb{C}$  at step 33.*

For this example we have simply augmented the diagonal matrix of Example 3 by the two complex diagonal elements  $0.8 + 0.1i$  and  $0.8 - 0.1i$  to a matrix of order 102, and we have applied the harmonic Ritz value variant of our algorithm (§5) to the matrix  $A - \mu I$  with shift  $\mu = 0.81 + 0.08i$ . Again, all coordinates of the starting vector are 1. Now we solve the projected equations approximately using 10 steps of GMRES.

Figures 7–9 show the convergence history of the residual vector (Figure 7), the real parts of the harmonic Ritz values (Figure 8), and the imaginary parts of the harmonic Ritz values (Figure 9). Figure 10 shows the harmonic Ritz values of order 33. In Figures 7–9, we have used the same symbols as in the previous example. In Figure 10, the harmonic Ritz values of step 33 of the Jacobi–Davidson iteration are represented by a +, while the eigenvalues of  $A$  are represented by  $\circ$ .

Clearly, the algorithm finds the eigenvalue  $\lambda = 0.8 + 0.1i$  of  $A$  close to the shift  $\mu$ , but also other ones, such as the conjugate of  $\lambda$  (which is quite far from the shift). This is typical for other experiments as well; usually a large number of (harmonic) Ritz values are converging in the Jacobi–Davidson method.

*Example 5.* In this example we experimentally compare the performances of the Davidson method, the Jacobi–Davidson method, and the accelerated shift and inexact invert (ASII) variant of observation 4 in §4.1, i.e., we expand our search space by (the orthogonal complement of) the approximate solution  $\tilde{t}$  of

$$\begin{aligned} (A - \theta_k I)t &= -r & (D), \\ (I - u_k u_k^*)(A - \theta_k I)(I - u_k u_k^*)t &= -r & (JD), \\ (A - \theta_k I)t &= u_k & (SI), \end{aligned}$$

respectively (cf. §4). We solve these equations approximately by  $m$  steps of full GMRES (with 0 as an initial guess). Since we are interested in the absolute smallest eigenvalue we take for  $\theta_k$  the absolute smallest eigenvalue of  $H_k = V_k^* A V_k$ . The preconditioner  $M$  for GMRES is kept fixed throughout the iteration process. The systems (D) and (SI) are preconditioned by  $M^{-1}$ , while the projection  $M_d := (I - u_k u_k^*)M(I - u_k u_k^*)$  is used as preconditioner for (JD). This means that for (JD) we have to solve equations of the form  $M_d z = y$ , where  $y$  is a given vector orthogonal to  $u_k$ . We follow the approach as indicated in §4.1 and we solve this equation by  $z = \alpha M^{-1} u_k - M^{-1} y$  with  $\alpha = u_k^* M^{-1} y / u_k^* M^{-1} u_k$ .

We have applied the three methods—Davidson, Jacobi–Davidson, and SI—for a matrix from the Harwell–Boeing set of test matrices:  $A$  is the SHERMAN4 matrix shifted by 0.5 (we wish to compute the eigenvalue of the SHERMAN4 matrix that is closest to 0.5).  $A$  is of order  $n = 1104$ . All eigenvalues of  $A$  appear to be real and are in the interval  $[0.030726, 66.497]$ . The smallest eigenvalues are (in 5 decimal places): 0.030726, 0.084702, 0.27757, 0.39884, 0.43154, 0.58979, 0.63480, ..., so that with the given shift we are aiming at the fifth eigenvalue.

For  $M$  we have selected the ILU(2) decomposition of  $A$ . We have plotted the  $\log_{10}$  of the norm of the residual versus the number of outer iteration steps (which is the dimension of the search space  $V_k$ ): Figures 11, 12, and 13 show the results for, respectively, 5, 10, and 25 steps of GMRES.

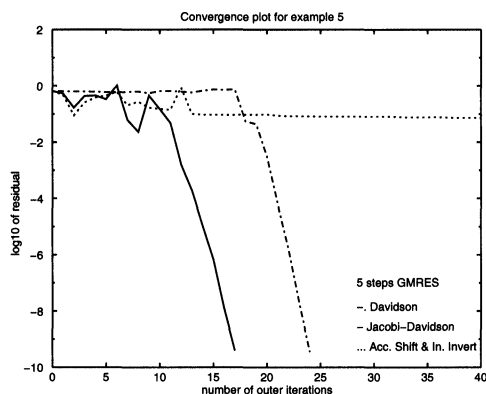


FIG. 11. Using 5 steps of preconditioned GMRES.

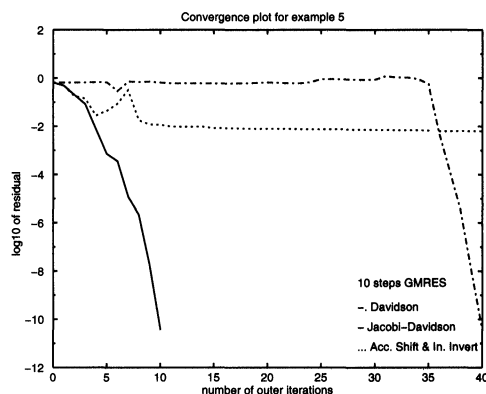


FIG. 12. Using 10 steps of preconditioned GMRES.

Larger values of  $m$  imply more accurate approximate solutions of the “expansion equations” (D), (JD), and (SI). In line with our discussions in §3 and our results in Example 1, we see that improving the approximation in Davidson’s method slows the speed of convergence and it may even lead to stagnation (cf. the dash-dotted curves

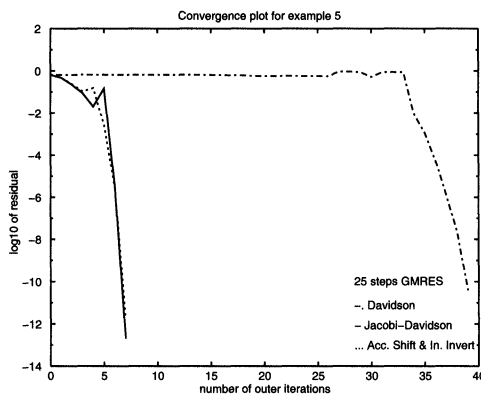


FIG. 13. Using 25 steps of preconditioned GMRES.

—). As might be anticipated, for ASII we observe the opposite effect (cf. the dotted curves ...); the more precisely we solve (SI), the faster the method converges, while stagnation may occur if (SI) is not solved accurately enough. The speed of convergence of our Jacobi–Davidson method does not depend so much upon the accuracy of the approximate solutions of (JD) (cf. the solid curves —): the method converges faster than Davidson and ASII.

As argued in §4.1, ASII may be rather sensitive to rounding errors, especially if the expanding vector  $\tilde{t}$  has a large component in the direction of  $u$ . For ASII, but also for Davidson, we had to apply modified Gram–Schmidt (mod-GS) twice to maintain sufficient orthogonality of  $V_k$ , while in Jacobi–Davidson this was not necessary. By doing mod-GS only once, the angle between the expansion vector  $\tilde{t}$  and the already available search space may become too small to allow an accurate computation of the orthogonal component. In such a situation, it may help to apply mod-GS more often [23]. For the present example, twice was enough (but other examples, not reported here, required more mod-GS sweeps).

**Acknowledgments.** We thank Albert Booten for his careful reading of an early version of the manuscript. He was also helpful in deriving formula (31). We gratefully acknowledge helpful discussions with Diederik Fokkema on the subject of §4. He also provided the numerical data for Example 5. Axel Ruhe helped us to improve our presentation.

We are also grateful to the referees who made many valuable suggestions. One of the referees drew our attention to the references [16], [18], [27], which contain information that should have been more widely known.

## REFERENCES

- [1] W. E. ARNOLDI, *The principle of minimized iterations in the solution of the matrix eigenvalue problem*, Quart. Appl. Math., 9 (1951), pp. 17–29.
- [2] J. G. L. BOOTEN, H. A. VAN DER VORST, P. M. MEIJER, AND H. J. J. TE RIELE, *A preconditioned Jacobi–Davidson method for solving large generalized eigenvalue problems*, Report NM-R9414, Department of Numerical Mathematics, CWI, Amsterdam, The Netherlands, 1994.
- [3] M. CROUZEIX, B. PHILIPPE, AND M. SADKANE, *The Davidson method*, SIAM J. Sci. Comput., 15 (1994), pp. 62–76.



- [4] E. R. DAVIDSON, *The iterative calculation of a few of the lowest eigenvalues and corresponding eigenvectors of large real symmetric matrices*, J. Comput. Phys., 17 (1975), pp. 87–94.
- [5] ———, *Matrix eigenvector methods in quantum mechanics*, in Methods in Computational Molecular Physics, G. H. F. Diercksen and S. Wilson, eds., Reidel, Dordrecht, The Netherlands, 1983, pp. 95–113.
- [6] ———, *Monster matrices: Their eigenvalues and eigenvectors*, Computers in Physics, 7 (1993), pp. 519–522.
- [7] A. DEN BOER, *De Jacobi methode van 1845 tot 1990*, Master's thesis, Department of Mathematics, University of Utrecht, Utrecht, The Netherlands, 1991. (In Dutch.)
- [8] R. W. FREUND, *Quasi-kernel polynomials and their use in non-Hermitian matrix iterations*, J. Comput. Appl. Math., 43 (1992), pp. 135–158.
- [9] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 2nd ed., The John Hopkins University Press, Baltimore and London, 1989.
- [10] Y. HUANG AND H. A. VAN DER VORST, *Some observations on the convergence behaviour of GMRES*, Tech. report 89-09, Delft University of Technology, Faculty of Tech. Math., Delft, The Netherlands, 1989.
- [11] C. G. J. JACOBI, *Ueber eine neue Auflösungsart der bei der Methode der kleinsten Quadrate vorkommende linearen Gleichungen*, Astronom. Nachr., (1845), pp. 297–306.
- [12] ———, *Ueber ein leichtes Verfahren, die in der Theorie der Säcularstörungen vorkommenden Gleichungen numerisch aufzulösen*, J. Reine und Angew. Math., (1846), pp. 51–94.
- [13] C. LANCZÔS, *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators*, J. Res. Nat. Bur. Standards, 45 (1950), pp. 255–282.
- [14] R. B. MORGAN, *Computing interior eigenvalues of large matrices*, Linear Algebra Appl., 154/156 (1991), pp. 289–309.
- [15] ———, *Generalizations of Davidson's method for computing eigenvalues of large nonsymmetric matrices*, J. Comput. Phys., 101 (1992), pp. 287–291.
- [16] R. B. MORGAN AND D. S. SCOTT, *Generalizations of Davidson's method for computing eigenvalues of sparse symmetric matrices*, SIAM J. Sci. Statist. Comput., 7 (1986), pp. 817–825.
- [17] ———, *Preconditioning the Lanczos algorithm for sparse symmetric eigenvalue problems*, SIAM J. Sci. Comput., 14 (1993), pp. 585–593.
- [18] J. OLSEN, P. JØRGENSEN, AND J. SIMONS, *Passing the one-billion limit in full configuration-interaction (FCI) calculations*, Chemical Physics Letters, 169 (1990), pp. 463–472.
- [19] A. M. OSTROWSKI, *On the convergence of the Rayleigh quotient iteration for the computation of characteristic roots and vectors. V*, Arch. Rational Mech. Anal., 3 (1959), pp. 472–481.
- [20] C. C. PAIGE, B. N. PARLETT, AND H. A. VAN DER VORST, *Approximate solutions and eigenvalue bounds from Krylov subspaces*, Num. Lin. Alg. with Appl., 2 (1995), pp. 115–134.
- [21] B. N. PARLETT, *The Rayleigh quotient iteration and some generalizations*, Math. Comp., 28 (1974), pp. 679–693.
- [22] ———, *The Symmetric Eigenvalue Problem*, Prentice-Hall, Englewood Cliffs, NJ, 1980.
- [23] A. RUHE, *Numerical aspects of Gram-Schmidt orthogonalization of vectors*, Linear Algebra Appl., 52/53 (1983), pp. 591–602.
- [24] ———, *Rational Krylov algorithms for nonsymmetric eigenvalue problems. II. Matrix pairs*, Linear Algebra Appl., 197/198 (1994), pp. 283–295.
- [25] Y. SAAD, *Numerical Methods for Large Eigenvalue Problems*, Manchester University Press, Manchester, UK, 1992.
- [26] Y. SAAD AND M. H. SCHULTZ, *GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM J. Sci. Statist. Comput., 7 (1986), pp. 856–869.
- [27] A. STATHOPOULOS, Y. SAAD, AND C. F. FISCHER, *Robust preconditioning of large sparse symmetric eigenvalue problems*, J. Comput. Appl. Math., 64 (1995), pp. 197–215.
- [28] H. J. J. VAN DAM, J. H. VAN LENTHE, G. L. G. SLEIJPEN, AND H. A. VAN DER VORST, *An improvement of Davidson's iteration method; Applications to MRCI and MRCEPA calculations*, J. Comput. Chem., to appear.
- [29] A. VAN DER SLUIS AND H. A. VAN DER VORST, *The convergence behavior of Ritz values in the presence of close eigenvalues*, Linear Algebra Appl., 88/89 (1987), pp. 651–694.
- [30] M. B. VAN GIJZEN, *Iterative solution methods for linear equations in finite element computations*, Ph.D. thesis, Delft University of Technology, Delft, The Netherlands, 1994.
- [31] J. H. VAN LENTHE AND P. PULAY, *A space-saving modification of Davidson's eigenvector algorithm*, J. Comput. Chem., 11 (1990), pp. 1164–1168.