

# Globally convergent Jacobi methods for positive definite matrix pairs

Vjeran Hari

Received: date / Accepted: date

**Abstract** The paper derives and investigates the Jacobi methods for the generalized eigenvalue problem  $Ax = \lambda Bx$ , where  $A$  is a symmetric and  $B$  is a symmetric positive definite matrix. The methods first “normalize”  $B$  to have the unit diagonal and then maintain that property during the iterative process. The global convergence is proved for all such methods. That result is obtained for the large class of generalized serial strategies from [10]. Preliminary numerical tests confirm a high relative accuracy of some of those methods, provided that the both matrices are positive definite, and the spectral condition numbers of  $\Delta_A A \Delta_A$  and  $\Delta_B B \Delta_B$  are small, for some nonsingular diagonal matrices  $\Delta_A$  and  $\Delta_B$ .

**Keywords** Generalized eigenvalue problem · Jacobi method · Global convergence

**Mathematics Subject Classification (2000)** 65F15

## 1 Introduction

In this paper we consider the global convergence of the Jacobi-type algorithms for the positive definite generalized eigenvalue problem (PGEP)

$$Ax = \lambda Bx, \quad x \neq 0,$$

where  $A$  and  $B$  are the real symmetric matrices of order  $n$  and  $B$  is positive definite. We first derive four algorithms, of which three are new, and one was derived in [4] but has not been published since. These algorithms have a common property: they require that  $B$  has ones along the diagonal. That property is maintained during the iteration. Since the methods simultaneously diagonalize the pivot submatrices, we call them briefly the Jacobi methods for PGEP.

---

This work has been fully supported by Croatian Science Foundation under the project: IP.09.2014.3670.

Department of Mathematics, Faculty of Science, University of Zagreb, Bijenička cesta 30, 10000 Zagreb, Croatia. E-mail: hari@math.hr

The most well-known Jacobi method for PGEP is the Falk–Langemeyer (*FL*) method [2], which is well defined for a definite initial pair  $(A, B)$  (see [17]), i.e., for the case when  $\alpha A + \beta B$  is positive definite for some real  $\alpha$  and  $\beta$ . The asymptotic quadratic convergence and accuracy of that method were considered in [17] and [12], respectively. Since the spectral norm of the transformation matrix used by the *FL* method is larger than one, the norms of the iteration matrices  $A^{(k)}$ ,  $B^{(k)}$ , and of the matrix of accumulated transformations  $F^{(k)}$ , gradually increase. Thus, occasionally, these matrices have to be “normalized”. Furthermore, the stopping criterion uses the normalized matrices. In [13] it was shown that on parallel computers these tasks can be demanding and it is better to work with the “normalized version” of the *FL* method, which they call the Hari–Zimmermann variant of the Falk–Langemeyer method. The method considered in [13] is a one-sided block version of the method derived in [4], implemented as the one-sided block Jacobi method for the generalized singular value problem. Then, to quote from [13], “it is almost perfectly parallelizable, so parallel shared memory versions of the algorithm are highly scalable, and their speedup almost solely depends on the number of cores used”. It compares favorably to the LAPACK’s DTGSJA algorithm. Hence, our first goal is to derive the core algorithm for that block method.

The idea of the original element-wise, two-sided method had been briefly outlined by Zimmermann [20] and the algorithm was later derived and analyzed by Hari [4]. We call it simply the *HZ* method. Here, we consider only the element-wise, two-sided methods.

The *FL* and *HZ* methods are extensions of the Jacobi method for symmetric matrices, because they diagonalize the pivot submatrices in each step. The *FL* method uses simplified transformation matrices, with the ones along the diagonal. The *HZ* method uses simplified iteration matrices  $B^{(k)}$ , with the ones along the diagonal. So, for the *HZ* method, a preliminary step is employed to make the diagonal elements of  $B$  equal to one. In that step, the congruence transformation with the diagonal matrix is used:

$$A \rightarrow DAD, \quad B \rightarrow DBD, \quad D = \text{diag} \left( b_{11}^{-1/2}, b_{22}^{-1/2}, \dots, b_{nn}^{-1/2} \right). \quad (1.1)$$

Then  $(DAD, DBD)$  is the starting pair for the method. Actually, the transformation (1.1) has a double effect. First, it normalizes and balances  $B$  in such a way that it obtains the unit diagonal. Second, it nicely preconditions  $B$ . Namely,  $B$  now has almost the optimal condition that can be obtained by a symmetric diagonal scaling [18].

Each one of these two methods, *FL* and *HZ*, has both some advantages and some shortcomings. The advantage of the *FL* method is that it is well defined for a more general initial matrix pair, and the transformations are somewhat cheaper to apply, because the transformation matrices have the unit diagonal. The shortcoming lies in the fact that the elements of  $A^{(k)}$ ,  $B^{(k)}$ , and  $F^{(k)}$  increase. On the contrary, the *HZ* method has no need for renormalizations. It is a proper generalization of the standard Jacobi method for the matrix  $A$  since the *HZ* method reduces to it when  $B = I$ . Its weaker side lies in the fact that at least one of the matrices,  $A$  or  $B$ , has to be positive definite. If  $A$  is such, then the method can be applied to the pair  $(B, A)$ , which has the same eigenvectors as  $(A, B)$ , and reciprocal eigenvalues.

The *HZ* method has a significance on its own, since it can be employed for solving the PGEP with smaller matrices (say, for  $n \leq 2000$ ) on standard PCs. Even so, its best role nowadays is to serve as a kernel algorithm for a block method, like the one from [13]. That is a natural application, because the *HZ* method is very fast (just a few sweeps are needed) and highly accurate (cf., [11]) on the pairs of almost diagonal matrices.

In this paper we derive three similar Jacobi methods for the PGEP. The first two are called *LL<sup>T</sup>J* and *RR<sup>T</sup>J* methods because, in step  $k$ , their algorithms are based on the *LL<sup>T</sup>* and *RR<sup>T</sup>* factorizations of  $\hat{B}^{(k)}$ , the pivot submatrix of  $B^{(k)}$ , followed by the Jacobi transformation that diagonalizes the updated  $\hat{A}^{(k)}$ . The third method is a combination of them. At each step it chooses the algorithm that is more accurate. In this way a special hybrid method is defined that we call the *CJ* method. All these methods are identical except in their method of transforming each  $\hat{B}^{(k)}$  to  $I_2$ .

Furthermore, one can combine all three methods, *HZ*, *LL<sup>T</sup>J* and *RR<sup>T</sup>J*. Hence, we introduce a *hybrid method* that employs in each step any of the three algorithms. This leads us to a further generalization that we call a *general Jacobi method* for the PGEP. Roughly speaking, it is any method that simultaneously diagonalizes the both pivot submatrices, while maintaining the unit diagonal of  $B$ . All those methods can serve as the kernel algorithms for the block Jacobi methods.

A kernel algorithm should at least be provably globally convergent. Hence, this paper is devoted to proving the global convergence of all those methods. Recall that the global convergence problem for the one-sided Jacobi methods reduces to the one for the corresponding two-sided methods, so the global convergence is always linked to the two-sided methods. The global convergence is proved for the large class of the generalized serial strategies from [10]. In addition, the numerical tests have been performed in MATLAB to inspect the high relative accuracy of the new methods.

The paper is divided into 6 sections. In Section 2 we derive a detailed algorithm for the *HZ* method. In Section 3, we derive algorithms for the *LL<sup>T</sup>J*, *RR<sup>T</sup>J* and *CJ* methods. We also define a hybrid and a general Jacobi method. In Section 4 we prove the global convergence of all those methods. In Section 5 we provide preliminary numerical tests in MATLAB which indicate the high relative accuracy of the *HZ* and *CJ* methods. Finally, Section 6 gives a brief summary of this research and outlines the future work.

## 2 The *HZ* Method

Here we derive the algorithm of the *HZ* method. In [4] the algorithm had been derived for the complex Hermitian matrices, and then it was simplified by assuming that the matrices were real symmetric. Here we derive the algorithm directly for the real symmetric matrices  $A$  and  $B$  such that  $B$  is positive definite. The method has the form

$$A^{(k+1)} = Z_k^T A^{(k)} Z_k, \quad B^{(k+1)} = Z_k^T B^{(k)} Z_k, \quad k \geq 0,$$

where  $A^{(0)} = DAD$ ,  $B^{(0)} = DBD$ , and  $D$  is the diagonal matrix from the relation (1.1). Each transformation matrix  $Z_k$  is a nonsingular elementary plane matrix that defers from the identity  $I_n$  in one principal submatrix of order 2. That submatrix is denoted

by  $\hat{Z}_k$  and is called a *pivot submatrix* of  $Z_k$ . We assume that it lies on the intersection of the rows and columns  $i$  and  $j$ , so we can write

$$\hat{Z}_k = \begin{bmatrix} z_{ii}^{(k)} & z_{ij}^{(k)} \\ z_{ji}^{(k)} & z_{jj}^{(k)} \end{bmatrix}, \quad k \geq 0.$$

The indices  $i, j, i < j$ , both depending on  $k$ , are the *pivot indices*,  $(i, j)$  is a pivot pair, and a way how  $(i, j)$  is selected in each step is a *pivot strategy*. The role of  $\hat{Z}_k$  is to diagonalize the corresponding (pivot) submatrices of  $A^{(k)}$  and  $B^{(k)}$  and to maintain the unit diagonal of  $B^{(k)}$ .

The method is *globally convergent* if, for every initial pair of the symmetric matrices  $(A, B)$  such that  $B$  is positive definite, the sequence of matrices  $(B^{(k)}, k \geq 0)$  tends to the identity matrix, while  $(A^{(k)}, k \geq 0)$  tends to the diagonal matrix of the eigenvalues of  $(A, B)$ . In such a case, the columns of  $Z^{(k)} = DZ_0Z_1 \cdots Z_k$  approach the set of the eigenvectors of the pair  $(A, B)$ .

## 2.1 The derivation of the HZ algorithm

Here we derive the formulas related to one step of the method. To simplify the notation, we omit the indices  $k$  and  $k+1$ , and assume  $b_{11} = \cdots = b_{mm} = 1$ . Then one step of the method has the form  $A' = Z^T A Z$ ,  $B' = Z^T B Z$ . For the pivot submatrices we have

$$\hat{A}' = \hat{Z}^T \hat{A} \hat{Z}, \quad \hat{B}' = \hat{Z}^T \hat{B} \hat{Z}.$$

We first derive the algorithm for computing  $\hat{Z}$ . It is sought in the form of a product of two Jacobi rotations and one diagonal matrix. We have two possibilities:

$$\begin{aligned} \text{(a)} \quad \hat{Z} &= \begin{bmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{1+b_{ij}}} & 0 \\ 0 & \frac{1}{\sqrt{1-b_{ij}}} \end{bmatrix} \begin{bmatrix} \cos(\theta - \frac{\pi}{4}) & -\sin(\theta - \frac{\pi}{4}) \\ \sin(\theta - \frac{\pi}{4}) & \cos(\theta - \frac{\pi}{4}) \end{bmatrix}, \\ \text{(b)} \quad \hat{Z} &= \begin{bmatrix} \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{1-b_{ij}}} & 0 \\ 0 & \frac{1}{\sqrt{1+b_{ij}}} \end{bmatrix} \begin{bmatrix} \cos(\theta + \frac{\pi}{4}) & -\sin(\theta + \frac{\pi}{4}) \\ \sin(\theta + \frac{\pi}{4}) & \cos(\theta + \frac{\pi}{4}) \end{bmatrix}. \end{aligned}$$

The leftmost rotation is just the Jacobi rotation for  $\hat{B}$ . The second transformation reestablishes the unit diagonal of  $\hat{B}$ . The rightmost rotation is the Jacobi rotation for the transformed  $\hat{A}$ .

The both approaches yield the same matrix  $\hat{Z}$ , so we consider only the case (a). The derivation of  $\hat{Z}$  in the case (b) is quite similar to the derivation presented below.

(a) Let us first consider how the matrices  $\hat{B}$  and  $\hat{A}$  are transformed. Recall that  $\hat{Z} = \hat{R}(\frac{\pi}{4}) D_+ \hat{R}(\theta - \frac{\pi}{4})$  and the diagonal elements of  $\hat{B}$  are ones. If we write

$$\hat{B}_1 = \hat{R}(\frac{\pi}{4})^T \hat{B} \hat{R}(\frac{\pi}{4}), \quad \hat{B}_2 = D_+^T \hat{B}_1 D_+, \quad \hat{B}' = \hat{R}(\theta - \frac{\pi}{4})^T \hat{B}_2 \hat{R}(\theta - \frac{\pi}{4}),$$

then a simple calculation yields

$$\hat{B}_1 = \begin{bmatrix} 1+b_{ij} & 0 \\ 0 & 1-b_{ij} \end{bmatrix}, \quad \hat{B}_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_2, \quad \hat{B}' = I_2.$$

Now, let us consider the transformations on  $\hat{A}$ . If we write

$$\hat{A}_1 = \hat{R}\left(\frac{\pi}{4}\right)^T \hat{A} \hat{R}\left(\frac{\pi}{4}\right), \quad \hat{A}_2 = D_+^T \hat{A}_1 D_+, \quad \hat{A}' = \hat{R}\left(\theta - \frac{\pi}{4}\right)^T \hat{A}_2 \hat{R}\left(\theta - \frac{\pi}{4}\right),$$

we obtain

$$\begin{aligned} \hat{A}_1 &= \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} a_{ii} & a_{ij} \\ a_{ij} & a_{jj} \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} a_{ii} + 2a_{ij} + a_{jj} & a_{jj} - a_{ii} \\ a_{jj} - a_{ii} & a_{jj} - 2a_{ij} + a_{ii} \end{bmatrix}, \\ \hat{A}_2 &= \begin{bmatrix} \frac{1}{\sqrt{1+b_{ij}}} & 0 \\ 0 & \frac{1}{\sqrt{1-b_{ij}}} \end{bmatrix} \hat{A}_1 \begin{bmatrix} \frac{1}{\sqrt{1+b_{ij}}} & 0 \\ 0 & \frac{1}{\sqrt{1-b_{ij}}} \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \frac{a_{ii}+2a_{ij}+a_{jj}}{1+b_{ij}} & \frac{a_{jj}-a_{ii}}{\sqrt{1-(b_{ij})^2}} \\ \frac{a_{jj}-a_{ii}}{\sqrt{1-(b_{ij})^2}} & \frac{a_{jj}-2a_{ij}+a_{ii}}{1-b_{ij}} \end{bmatrix}, \\ \hat{A}' &= \frac{1}{2} \begin{bmatrix} \frac{a_{ii}+2a_{ij}+a_{jj}}{1+b_{ij}} + \tan\left(\theta - \frac{\pi}{4}\right) \frac{a_{jj}-a_{ii}}{\sqrt{1-(b_{ij})^2}} & 0 \\ 0 & \frac{a_{jj}-2a_{ij}+a_{ii}}{1-b_{ij}} - \tan\left(\theta - \frac{\pi}{4}\right) \frac{a_{jj}-a_{ii}}{\sqrt{1-(b_{ij})^2}} \end{bmatrix}, \end{aligned}$$

where

$$\tan\left(2\left(\theta - \frac{\pi}{4}\right)\right) = \frac{2 \frac{a_{jj}-a_{ii}}{\sqrt{1-(b_{ij})^2}}}{\frac{a_{ii}+2a_{ij}+a_{jj}}{1+b_{ij}} - \frac{a_{jj}-2a_{ij}+a_{ii}}{1-b_{ij}}} = \frac{\sqrt{1-(b_{ij})^2} (a_{jj} - a_{ii})}{2a_{ij} - (a_{ii} + a_{jj})b_{ij}}.$$

From the trigonometric identities

$$\tan\left(2\left(\theta - \frac{\pi}{4}\right)\right) = \tan\left(2\theta - \frac{\pi}{2}\right) = -\cot(2\theta) = \frac{1}{-\tan(2\theta)},$$

we obtain

$$\tan(2\theta) = \frac{2a_{ij} - (a_{ii} + a_{jj})b_{ij}}{\sqrt{1-(b_{ij})^2} (a_{ii} - a_{jj})}, \quad -\frac{\pi}{4} \leq \theta \leq \frac{\pi}{4}. \quad (2.1)$$

Since  $B$  is positive definite with the unit diagonal, such is also  $\hat{B}$ , and we have  $|b_{ij}| < 1$ . This implies  $\sqrt{1-(b_{ij})^2} > 0$ . Next, we derive a compact form for  $\hat{Z}$ . To this end, let  $c = \cos \theta$ ,  $s = \sin \theta$ . We have

$$\cos\left(\theta - \frac{\pi}{4}\right) = \frac{\sqrt{2}}{2} (c+s), \quad \sin\left(\theta - \frac{\pi}{4}\right) = \frac{\sqrt{2}}{2} (s-c),$$

$$\begin{aligned} \hat{Z} &= \frac{\sqrt{2}}{2} \begin{bmatrix} \frac{1}{\sqrt{1+b_{ij}}} & -\frac{1}{\sqrt{1-b_{ij}}} \\ \frac{1}{\sqrt{1+b_{ij}}} & \frac{1}{\sqrt{1-b_{ij}}} \end{bmatrix} \frac{\sqrt{2}}{2} \begin{bmatrix} c+s & c-s \\ s-c & c+s \end{bmatrix} = \frac{1}{\sqrt{1-(b_{ij})^2}} \\ &\cdot \frac{1}{2} \begin{bmatrix} \sqrt{1-b_{ij}}(c+s) - \sqrt{1+b_{ij}}(s-c) & \sqrt{1-b_{ij}}(c-s) - \sqrt{1+b_{ij}}(c+s) \\ \sqrt{1-b_{ij}}(c+s) + \sqrt{1+b_{ij}}(s-c) & \sqrt{1-b_{ij}}(c-s) + \sqrt{1+b_{ij}}(c+s) \end{bmatrix}. \end{aligned}$$

We would like to obtain  $\hat{Z}$  in the form

$$\hat{Z} = \frac{1}{\sqrt{1-(b_{ij})^2}} \begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \psi & \cos \psi \end{bmatrix}.$$

To this end we use the following identities, which hold for  $|x| \leq 1$ :

$$\begin{aligned}\sqrt{1+x} - \sqrt{1-x} &= \frac{2x}{\sqrt{1+x} + \sqrt{1-x}}, \\ \sqrt{1+x} + \sqrt{1-x} &= 2 - \frac{2x^2}{(1 + \sqrt{1+x})(1 + \sqrt{1-x})(\sqrt{1+x} + \sqrt{1-x})}.\end{aligned}$$

Let

$$\xi = \frac{b_{ij}}{\sqrt{1+b_{ij}} + \sqrt{1-b_{ij}}}, \quad \eta = \frac{b_{ij}}{(1 + \sqrt{1+b_{ij}})(1 + \sqrt{1-b_{ij}})}, \quad (2.2)$$

$$\rho = 1 - \xi\eta = \xi + \sqrt{1-b_{ij}} = \frac{1}{2}(\sqrt{1+b_{ij}} + \sqrt{1-b_{ij}}). \quad (2.3)$$

A straightforward computation shows that we have  $\xi + \xi\eta^2 - 2\eta = 0$ . This is equivalent to  $\rho^2 + \xi^2 = 1$ . Now we have

$$\begin{aligned}\cos \phi &= \frac{1}{2} \left[ c(\sqrt{1-b_{ij}} + \sqrt{1+b_{ij}}) - s(\sqrt{1+b_{ij}} - \sqrt{1-b_{ij}}) \right] = \\ &= c - \frac{cb_{ij}^2}{(1 + \sqrt{1+b_{ij}})(1 + \sqrt{1-b_{ij}})(\sqrt{1+b_{ij}} + \sqrt{1-b_{ij}})} - \frac{sb_{ij}}{\sqrt{1+b_{ij}} + \sqrt{1-b_{ij}}} \\ &= c - \xi(s + \eta c) = \rho \cos \theta - \xi \sin \theta.\end{aligned}$$

In the similar way, we obtain

$$\left. \begin{aligned}\cos \phi &= \cos \theta - \xi(\sin \theta + \eta \cos \theta) = \rho \cos \theta - \xi \sin \theta, \\ \sin \phi &= \sin \theta + \xi(\cos \theta - \eta \sin \theta) = \rho \sin \theta + \xi \cos \theta, \\ \cos \psi &= \cos \theta + \xi(\sin \theta - \eta \cos \theta) = \rho \cos \theta + \xi \sin \theta, \\ \sin \psi &= \sin \theta - \xi(\cos \theta + \eta \sin \theta) = \rho \sin \theta - \xi \cos \theta.\end{aligned}\right\} \quad (2.4)$$

From the relation (2.4) one obtains  $\cos^2 \phi + \sin^2 \phi = 1$  and  $\cos^2 \psi + \sin^2 \psi = 1$ . Using  $2\rho\xi = b_{ij}$ , one easily obtains the following identities

$$\cos \phi \sin \psi = \cos \theta \sin \theta - \rho\xi = 0.5(\sin 2\theta - b_{ij}), \quad (2.5)$$

$$\cos \psi \sin \phi = \cos \theta \sin \theta + \rho\xi = 0.5(\sin 2\theta + b_{ij}), \quad (2.6)$$

$$\cos \phi \cos \psi = \rho^2 \cos^2 \theta - \xi^2 \sin^2 \theta, \quad (2.7)$$

$$\sin \phi \sin \psi = \rho^2 \sin^2 \theta - \xi^2 \cos^2 \theta. \quad (2.8)$$

Using the bounds for  $|\xi|$ ,  $|\eta|$ , and  $\rho$ ,

$$|\xi| \leq \sqrt{2}/2, \quad |\eta| \leq \sqrt{2} - 1, \quad \sqrt{2}/2 \leq \rho \leq 1, \quad (2.9)$$

we obtain

$$\min\{\cos \phi, \cos \psi\} \geq \rho \cos \theta - \frac{|b_{ij}|}{2\rho} |\sin \theta| \geq \left(\rho - \frac{|b_{ij}|}{2\rho}\right) \cos \theta > 0, \quad (2.10)$$

$$\max\{\cos \phi, \cos \psi\} = \rho \cos \theta + |\xi \sin \theta| = \cos(|\theta| - \zeta) \geq \cos(|\theta|) \geq \frac{\sqrt{2}}{2}. \quad (2.11)$$

Here,  $\zeta$  is defined by  $\cos \zeta = \rho$ ,  $\sin \zeta = |\xi|$ . Thus  $\tan \zeta = |\xi|/\rho \leq 1$ , and therefore  $0 \leq \zeta \leq \pi/4$ .

## 2.2 The HZ algorithm

Here we collect the obtained formulas and write down a detailed algorithm for one step of the method. In step  $k$ , input to the algorithm consists of the matrices  $A$ ,  $B$ , and the pivot pair  $(i, j)$ . The pivot submatrices are given by the relation (2.12) below

$$\hat{A} = \begin{bmatrix} a_{ii} & a_{ij} \\ a_{ij} & a_{jj} \end{bmatrix}, \quad \hat{B} = \begin{bmatrix} 1 & b_{ij} \\ b_{ij} & 1 \end{bmatrix}, \quad |b_{ij}| < 1. \quad (2.12)$$

If  $a_{ij} = 0$  and  $b_{ij} = 0$ , we set  $Z = I$ , and continue with the next step. If this is not the case, we compute the pivot submatrix of  $Z$ , i.e.,

$$\hat{Z} = \frac{1}{\sqrt{1-(b_{ij})^2}} \begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \psi & \cos \psi \end{bmatrix} = \begin{bmatrix} c1 & -s1 \\ s2 & c2 \end{bmatrix}. \quad (2.13)$$

Here  $\cos \phi$ ,  $\sin \phi$ ,  $\cos \psi$ ,  $\sin \psi$  are given by the relations (2.4), (2.2), (2.3).

If  $a_{ii} = a_{jj}$  and  $2a_{ij} = (a_{ii} + a_{jj})b_{ij}$ , then the expression for the angle  $\theta$  in the relation (2.1) reduces to the form  $0/0$ . In that case  $\hat{A}$  and  $\hat{B}$  are proportional, and we choose  $\theta = 0$ . Then we have

$$\hat{Z} = \frac{1}{\sqrt{1-(b_{ij})^2}} \begin{bmatrix} \rho & -\xi \\ -\xi & \rho \end{bmatrix} = \frac{1}{\sqrt{1-(b_{ij})^2}} \begin{bmatrix} \rho & -\frac{b_{ij}}{2\rho} \\ -\frac{b_{ij}}{2\rho} & \rho \end{bmatrix}.$$

It is easy to show that in this case  $a'_{ii} = a_{ii}$  and  $a'_{jj} = a_{jj}$ .

In the general case, the diagonal elements of  $\hat{B}'$  are ones, while the diagonal elements of  $\hat{A}'$  can be computed as follows:

$$\begin{aligned} a'_{ii} &= a_{ii} + \frac{1}{1-b_{ij}^2} [(b_{ij}^2 - \sin^2 \phi) a_{ii} + 2 \cos \phi \sin \psi a_{ij} + \sin^2 \psi a_{jj}] \\ &= a_{ii} + [(\frac{b_{ij}^2}{1-b_{ij}^2} - s1^2) a_{ii} + 2 c1 s2 a_{ij} + s2^2 a_{jj}] \end{aligned} \quad (2.14)$$

$$\begin{aligned} a'_{jj} &= a_{jj} - \frac{1}{1-b_{ij}^2} [(\sin^2 \psi - b_{ij}^2) a_{jj} + 2 \cos \psi \sin \phi a_{ij} - \sin^2 \phi a_{ii}] \\ &= a_{jj} - [(s2^2 - \frac{b_{ij}^2}{1-b_{ij}^2}) a_{jj} + 2 c2 s1 a_{ij} - s1^2 a_{ii}] \end{aligned} \quad (2.15)$$

Here, the relations (2.5) and (2.6) can be used. The pivot elements are set to zero:  $b_{ij} = 0$ ,  $b'_{ji} = 0$ ,  $a_{ij} = 0$ ,  $a'_{ji} = 0$ , while the off-diagonal elements are transformed using the formulas

$$\begin{aligned} a'_{ki} &= c1 \cdot a_{ki} + s2 \cdot a_{kj}, & b'_{ki} &= c1 \cdot b_{ki} + s2 \cdot b_{kj}, & a'_{ik} &= a'_{ki}, & b'_{ik} &= b'_{ki}, & k &\neq i, j, \\ a'_{kj} &= c2 \cdot a_{kj} - s1 \cdot a_{ki}, & b'_{kj} &= c2 \cdot b_{kj} - s1 \cdot b_{ki}, & a'_{jk} &= a'_{kj}, & b'_{jk} &= b'_{kj}, & k &\neq i, j. \end{aligned}$$

It can immediately be seen that in the case  $b_{ij} = 0$ , the quantities  $\xi$  and  $\eta$  become zero. Then the transformation  $Z$  becomes the Jacobi rotation for the matrix  $A$ . Therefore this method is a proper generalization of the simple Jacobi method for  $A$  to the positive definite pair  $(A, B)$ , under the constraint that  $B$  has ones on the diagonal.

When considering the two-sided methods it is advantageous to express the updated diagonal elements in the form  $a'_{ii} = a_{ii} + \delta a_{ii}$ ,  $a'_{jj} = a_j + \delta a_{jj}$ . For the *HZ* method,  $\delta a_{ii}$  and  $\delta a_{jj}$  are the expressions within brackets in the relations (2.14) and (2.15), respectively. Then contributions to the diagonal elements coming from all steps within the current sweep can be accumulated separately, and at the end of the sweep, added to the diagonal elements that are saved prior to the sweep [15]. Another set of formulas for the diagonal elements is derived in Section 4 (see (4.9), (4.10)). Next, we have to decide whether it is better to set the pivot elements zero or to compute them. Numerical tests with badly conditioned matrices have shown that it is better to compute  $a'_{ij}$  while  $b'_{ij}$  can be set zero. This leads to the algorithm which is displayed below.

There are many open problems connected with a way to implement the algorithm and they will be addressed elsewhere. Say, if  $|b_{ij}| = 1 - \varepsilon$  for a tiny  $\varepsilon > 0$ , then the spectral condition of  $F$  is as large as  $1/\sqrt{\varepsilon}$  and it may have bad impact on accuracy of the computed eigenvalues and eigenvectors. Such a case indicates a huge condition number of  $B^{(0)}$ . Recall that  $\kappa_2(X) = \|X\|_2 \|X^{-1}\|_2$  is the spectral condition number of  $X$ , where  $\|X\|_2 = \sqrt{\text{spr}(X^*X)}$  is the spectral norm of  $X$ .

If  $A$  is positive definite, then a possible remedy is to work with the pair  $(B, A)$  instead of  $(A, B)$ . If  $A$  is indefinite, one can try with the *FL* method or with the *J*-Jacobi method [19, 7, 8] applied to  $(B, A)$ .

#### **% % % Algorithm HZ**

```

select the pivot pair (i, j)
if  $a_{ij} \neq 0$  or  $b_{ij} \neq 0$  then
     $\rho = 0.5 * (\text{sqrt}(1 + b_{ij}) + \text{sqrt}(1 - b_{ij}))$ ;  $\xi = b_{ij} / (2 * \rho)$ ;  $\tau = \text{sqrt}((1 + b_{ij}) * (1 - b_{ij}))$ ;
     $t2 = 2 * a_{ij} - (a_{ii} + a_{jj}) * b_{ij}$ ;
    if  $t2 = 0$  then  $t = 0$ ;
        else  $ct2 = \tau * (a_{ii} - a_{jj}) / t2$ ;  $t = \text{sign}(ct2) / (\text{abs}(ct2) + (1 + \text{sqrt}(1 + ct2^2)))$ ;
    end;
     $cs = 1 / \text{sqrt}(1 + t^2)$ ;  $sn = t / \text{sqrt}(1 + t^2)$ ;  $c1 = (\rho * cs - \xi * sn) / \tau$ ;
     $c2 = (\rho * cs + \xi * sn) / \tau$ ;  $s1 = (\rho * sn + \xi * cs) / \tau$ ;  $s2 = (\rho * sn - \xi * cs) / \tau$ ;
     $a'_{ij} = (c1 * c2 - s1 * s2) * a_{ij} + (c2 * s2 * a_{jj} - c1 * s1 * a_{ii})$ ;  $a'_{ji} = a'_{ij}$ ;  $b'_{ij} = 0$ ;  $b'_{ji} = b'_{ij}$ ;
     $\delta_i = (b_{ij} / \tau - s1) * (b_{ij} / \tau + s1) * a_{ii} + (2 * c1 * a_{ij} + s2 * a_{jj}) * s2$ ;  $a'_{ii} = a_{ii} + \delta_i$ ;
     $\delta_j = (s2 - b_{ij} / \tau) * (s2 + b_{ij} / \tau) * a_{jj} + (2 * c2 * a_{ij} - s1 * a_{ii}) * s1$ ;  $a'_{jj} = a_{jj} - \delta_j$ ;
    for  $k = 1, \dots, n, k \neq i, j$  do
         $a'_{ki} = c1 * a_{ki} + s2 * a_{kj}$ ;  $b'_{ki} = c1 * b_{ki} + s2 * b_{kj}$ ;  $a'_{ik} = a'_{ki}$ ;  $b'_{ik} = b'_{ki}$ ;
         $a'_{kj} = c2 * a_{kj} - s1 * a_{ki}$ ;  $b'_{kj} = c2 * b_{kj} - s1 * b_{ki}$ ;  $a'_{jk} = a'_{kj}$ ;  $b'_{jk} = b'_{kj}$ ;
    end
end
end

```

### **3 The $LL^T J$ , $RR^T J$ and $CJ$ Algorithms**

Looking at the derivation of the *HZ* algorithm, one can see that instead of applying the Jacobi step followed by a diagonal transformation to  $\hat{B}$ , one can apply a congruence

transformation with the inverse of the Cholesky factor of  $\hat{B}$ . Actually, we have two possibilities: (a) using the  $LL^T$  and (b) using the  $RR^T$  factorization of  $\hat{B}$ . After each of these transformations the algorithm proceeds with the Jacobi step, as in the  $HZ$  algorithm. This approach yields two algorithms which we call the  $LL^T J$  and the  $RR^T J$  algorithm. We denote the transformation matrix in these algorithms by  $F$ .

Numerical investigation has shown that neither of these two algorithms has favorable accuracy properties. Fortunately, there is a combination of these algorithms with excellent accuracy properties. We call it the *Cholesky-Jacobi*, or shorter, the *CJ* algorithm.

### 3.1 The $LL^T J$ algorithm

Let us write the Cholesky factorization of  $\hat{B}$  by elements,

$$\begin{bmatrix} 1 & b_{ij} \\ b_{ij} & 1 \end{bmatrix} = \hat{B} = \hat{L}\hat{L}^T = \begin{bmatrix} 1 & 0 \\ a & c \end{bmatrix} \begin{bmatrix} 1 & a \\ 0 & c \end{bmatrix} = \begin{bmatrix} 1 & a \\ a & a^2 + c^2 \end{bmatrix}.$$

Assuming positive  $c$ , one immediately obtains  $a = b_{ij}$ ,  $c = \sqrt{1 - b_{ij}^2}$ , hence

$$\hat{L} = \begin{bmatrix} 1 & 0 \\ b_{ij} & \sqrt{1 - b_{ij}^2} \end{bmatrix} \quad \text{and} \quad \hat{L}^{-1} = \frac{1}{\sqrt{1 - b_{ij}^2}} \begin{bmatrix} \sqrt{1 - b_{ij}^2} & 0 \\ -b_{ij} & 1 \end{bmatrix}.$$

If we write  $\hat{F}_1 = \hat{L}^{-T}$ , then  $\hat{F}_1^T \hat{B} \hat{F}_1 = I_2$ , and we have

$$\begin{aligned} \hat{F}_1^T \hat{A} \hat{F}_1 &= \begin{bmatrix} 1 & 0 \\ f_{ij} & f_{jj} \end{bmatrix} \begin{bmatrix} a_{ii} & a_{ij} \\ a_{ij} & a_{jj} \end{bmatrix} \begin{bmatrix} 1 & f_{ij} \\ 0 & f_{jj} \end{bmatrix} = \begin{bmatrix} a_{ii} & f_{ij}a_{ii} + f_{jj}a_{ij} \\ f_{ij}a_{ii} + f_{jj}a_{ij} & f_{ij}^2a_{ii} + 2f_{ij}f_{jj}a_{ij} + f_{jj}^2a_{jj} \end{bmatrix} \\ &= \begin{bmatrix} a_{ii} & (a_{ij} - b_{ij}a_{ii})/\sqrt{1 - b_{ij}^2} \\ (a_{ij} - b_{ij}a_{ii})/\sqrt{1 - b_{ij}^2} & a_{jj} - \frac{2a_{ij} - (a_{ii} + a_{jj})b_{ij}}{1 - b_{ij}^2} b_{ij} \end{bmatrix}, \end{aligned}$$

where  $f_{ij} = -b_{ij}/\sqrt{1 - b_{ij}^2}$ ,  $f_{jj} = 1/\sqrt{1 - b_{ij}^2}$ . The final transformation  $\hat{F}$  has the form  $\hat{F} = \hat{F}_1 \hat{R}$ , where  $\hat{R}$  is a Jacobi rotation that annihilates the  $(1,2)$ -element of  $\hat{F}_1^T \hat{A} \hat{F}_1$ . Its angle  $\vartheta_1$  is determined by the formula

$$\tan(2\vartheta_1) = \frac{2(a_{ij} - b_{ij}a_{ii})\sqrt{1 - b_{ij}^2}}{a_{ii} - a_{jj} + 2(a_{ij} - b_{ij}a_{ii})b_{ij}}, \quad -\frac{\pi}{4} \leq \vartheta_1 \leq \frac{\pi}{4}. \quad (3.1)$$

The transformation formulas for the diagonal elements of  $A$  read

$$a'_{ii} = a_{ii} + \tan \vartheta_1 \cdot \frac{a_{ij} - a_{ii}b_{ij}}{\sqrt{1 - b_{ij}^2}}, \quad (3.2)$$

$$a'_{jj} = a_{jj} - \frac{2a_{ij} - b_{ij}(a_{ii} + a_{jj})}{1 - b_{ij}^2} b_{ij} - \tan \vartheta_1 \cdot \frac{a_{ij} - a_{ii}b_{ij}}{\sqrt{1 - b_{ij}^2}}. \quad (3.3)$$

In the case  $a_{ii} = a_{jj}$ ,  $a_{ij} = a_{ii}b_{ij}$  the expression for  $\tan(2\vartheta_1)$  has the form  $0/0$ , and then we choose  $\vartheta_1 = 0$ . In that case  $a'_{ii}$  and  $a'_{jj}$  reduce to  $a_{ii}$  and  $a_{jj}$ , respectively.

The transformation matrix has a simpler form than in the *HZ* method. We have

$$\begin{aligned}\hat{F} &= \frac{1}{\sqrt{1-b_{ij}^2}} \begin{bmatrix} \sqrt{1-b_{ij}^2} & -b_{ij} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} c_{\vartheta_1} & -s_{\vartheta_1} \\ s_{\vartheta_1} & c_{\vartheta_1} \end{bmatrix} = \frac{1}{\sqrt{1-b_{ij}^2}} \begin{bmatrix} c_{\tilde{\vartheta}_1} & -s_{\tilde{\vartheta}_1} \\ s_{\tilde{\vartheta}_1} & c_{\tilde{\vartheta}_1} \end{bmatrix} \\ &= \begin{bmatrix} c1 & -s1 \\ s2 & c2 \end{bmatrix}, \quad \begin{aligned} c_{\tilde{\vartheta}_1} &= c_{\vartheta_1} \sqrt{1-b_{ij}^2} - s_{\vartheta_1} b_{ij}, \\ s_{\tilde{\vartheta}_1} &= c_{\vartheta_1} b_{ij} + s_{\vartheta_1} \sqrt{1-b_{ij}^2}, \end{aligned} \quad (3.4) \\ c1 &= c_{\vartheta_1} - s_{\vartheta_1} b_{ij} / \sqrt{1-b_{ij}^2}, \quad c2 = c_{\vartheta_1} / \sqrt{1-b_{ij}^2}, \\ s1 &= c_{\vartheta_1} b_{ij} / \sqrt{1-b_{ij}^2} + s_{\vartheta_1}, \quad s2 = s_{\vartheta_1} / \sqrt{1-b_{ij}^2}.\end{aligned}$$

It is easy to verify that  $c_{\tilde{\vartheta}_1}^2 + s_{\tilde{\vartheta}_1}^2 = 1$ .

### 3.2 The $RR^TJ$ algorithm

Instead of the  $LL^T$  one can use the  $RR^T$  factorization of  $\hat{B}$ . Then we have

$$\begin{bmatrix} 1 & b_{ij} \\ b_{ij} & 1 \end{bmatrix} = \hat{B} = \hat{R}\hat{R}^T = \begin{bmatrix} c & a \\ 0 & 1 \end{bmatrix} \begin{bmatrix} c & 0 \\ a & 1 \end{bmatrix} = \begin{bmatrix} a^2 + c^2 & a \\ a & 1 \end{bmatrix}.$$

Assuming positive  $c$ , one obtains  $a = b_{ij}$ ,  $c = \sqrt{1-b_{ij}^2}$ . If we write  $\hat{F}_2 = \hat{R}^{-T}$ , then  $\hat{F}_2^T \hat{B} \hat{F}_2 = \hat{R}^{-1} \hat{B} \hat{R}^{-T} = I_2$ , and we have

$$\hat{F}_2^T \hat{A} \hat{F}_2 = \begin{bmatrix} a_{ii} - \frac{2a_{ij} - (a_{ii} + a_{jj})b_{ij}}{1-b_{ij}^2} b_{ij} & (a_{ij} - b_{ij}a_{jj}) / \sqrt{1-b_{ij}^2} \\ (a_{ij} - b_{ij}a_{jj}) / \sqrt{1-b_{ij}^2} & a_{jj} \end{bmatrix}. \quad (3.5)$$

The final transformation  $\hat{F}$  is given by  $\hat{F} = \hat{F}_2 \hat{R}$ , where  $\hat{R}$  is a Jacobi transformation that annihilates  $(1, 2)$ -element of  $\hat{F}_2^T \hat{A} \hat{F}_2$ . Its angle  $\vartheta_2$  is determined from

$$\tan(2\vartheta_2) = \frac{2(a_{ij} - b_{ij}a_{jj})\sqrt{1-b_{ij}^2}}{a_{ii} - a_{jj} - 2(a_{ij} - b_{ij}a_{jj})b_{ij}}, \quad -\frac{\pi}{4} \leq \vartheta_2 \leq \frac{\pi}{4}. \quad (3.6)$$

The transformation formulas for the diagonal elements of  $A$  read

$$a'_{ii} = a_{ii} - \frac{2a_{ij} - (a_{ii} + a_{jj})b_{ij}}{1-b_{ij}^2} b_{ij} + \tan \vartheta_2 \cdot \frac{a_{ij} - a_{jj}b_{ij}}{\sqrt{1-b_{ij}^2}}, \quad (3.7)$$

$$a'_{jj} = a_{jj} - \tan \vartheta_2 \cdot \frac{a_{ij} - a_{jj}b_{ij}}{\sqrt{1-b_{ij}^2}}. \quad (3.8)$$

In the case  $a_{ii} = a_{jj}$ ,  $a_{ij} = a_{jj}b_{ij}$  the angle  $\vartheta_2$  is chosen to be 0. In that case  $a'_{ii}$  and  $a'_{jj}$  are read from the relation (3.5), and they reduce to  $a_{ii}$  and  $a_{jj}$ , respectively.

This leads to the pivot submatrix  $\hat{F}$  of the  $RR^T J$  algorithm:

$$\begin{aligned}\hat{F} &= \frac{1}{\sqrt{1-b_{ij}^2}} \begin{bmatrix} 1 & 0 \\ -b_{ij} & \sqrt{1-b_{ij}^2} \end{bmatrix} \begin{bmatrix} c_{\vartheta_2} & -s_{\vartheta_2} \\ s_{\vartheta_2} & c_{\vartheta_2} \end{bmatrix} = \frac{1}{\sqrt{1-b_{ij}^2}} \begin{bmatrix} c_{\vartheta_2} & -s_{\vartheta_2} \\ s_{\tilde{\vartheta}_2} & c_{\tilde{\vartheta}_2} \end{bmatrix} \\ &= \begin{bmatrix} c1 & -s1 \\ s2 & c2 \end{bmatrix}, \quad \begin{aligned} c_{\tilde{\vartheta}_2} &= c_{\vartheta_2} \sqrt{1-b_{ij}^2} + s_{\vartheta_2} b_{ij}, \\ s_{\tilde{\vartheta}_2} &= s_{\vartheta_2} \sqrt{1-b_{ij}^2} - c_{\vartheta_2} b_{ij}, \end{aligned} \quad (3.9) \\ c1 &= c_{\vartheta_2} / \sqrt{1-b_{ij}^2}, \quad c2 = c_{\vartheta_2} + s_{\vartheta_2} b_{ij} / \sqrt{1-b_{ij}^2}, \\ s1 &= s_{\vartheta_2} / \sqrt{1-b_{ij}^2}, \quad s2 = s_{\vartheta_2} - c_{\vartheta_2} b_{ij} / \sqrt{1-b_{ij}^2}.\end{aligned}$$

Again, it is easy to prove that  $c_{\tilde{\vartheta}_2}^2 + s_{\tilde{\vartheta}_2}^2 = 1$ .

### 3.3 The $CJ$ algorithm

The third algorithm is a combination of the  $LL^T J$  and the  $RR^T J$  algorithms, and we describe it briefly as follows:

select the pivot pair  $(i, j)$ ;  
if  $a_{ii} \leq a_{jj}$  then choose the  $LL^T J$  algorithm, else choose the  $RR^T J$  algorithm.

Collecting the formulas (3.1)–(3.9) we can write it in more detail. One simple version is given below. In section 5 we shall justify the importance of the  $CJ$  algorithm.

#### **% % % Algorithm CJ**

```
select the pivot pair (i, j)
if  $a_{ij} \neq 0$  or  $b_{ij} \neq 0$  then
     $\beta = b_{ij}$ ;  $\tau = \text{sqrt}((1+\beta)*(1-\beta))$ ;
    if  $a_{ii} \leq a_{jj}$  then  $\sigma = 1$ ;  $\alpha_1 = a_{ii}$ ;  $\alpha_2 = a_{jj}$ ;
    else  $\sigma = -1$ ;  $\alpha_1 = a_{jj}$ ;  $\alpha_2 = a_{ii}$ ;
end;  $e = a_{ij} - \beta * \alpha_1$ ;  $ct2 = (0.5 * (\alpha_1 - \alpha_2) + e * \beta) / (\sigma * e * \tau)$ ;
 $t = \text{sign}(ct2) / (\text{abs}(ct2) + \text{sqrt}(1 + ct2^2))$ ;  $cs = 1 / \text{sqrt}(1 + t^2)$ ;  $sn = t / \text{sqrt}(1 + t^2)$ ;
 $\delta_1 = \sigma * t * e / \tau$ ;  $\delta_2 = \delta_1 + (\beta / \tau) * (2 * a_{ij} - (a_{ii} + a_{jj}) * \beta) / \tau$ ;  $\alpha_1 = \alpha_1 + \delta_1$ ;  $\alpha_2 = \alpha_2 - \delta_2$ ;
if  $\sigma > 0$  then  $c2 = cs / \tau$ ;  $s2 = sn / \tau$ ;  $c1 = cs - s2 * \beta$ ;  $s1 = sn + c2 * \beta$ ;  $a'_{ii} = \alpha_1$ ;  $a'_{jj} = \alpha_2$ ;
    else  $c1 = cs / \tau$ ;  $s1 = sn / \tau$ ;  $c2 = cs + s1 * \beta$ ;  $s2 = sn - c1 * \beta$ ;  $a'_{ii} = \alpha_2$ ;  $a'_{jj} = \alpha_1$ ;
end;
 $a'_{ij} = (c1 * c2 - s1 * s2) * a_{ij} + (c2 * s2 * a_{jj} - c1 * s1 * a_{ii})$ ;  $a'_{ji} = a'_{ij}$ ;
 $b'_{ij} = (c1 * c2 - s1 * s2) * b_{ij} + (c2 * s2 - c1 * s1)$ ;  $b'_{ji} = b'_{ij}$ ;
for  $k = 1, \dots, n, k \neq i, j$  do
     $a'_{ki} = c1 * a_{ki} + s2 * a_{kj}$ ;  $b'_{ki} = c1 * b_{ki} + s2 * b_{kj}$ ;  $a'_{ik} = a'_{ki}$ ;  $b'_{ik} = b'_{ki}$ ;
     $a'_{kj} = c2 * a_{kj} - s1 * a_{ki}$ ;  $b'_{kj} = c2 * b_{kj} - s1 * b_{ki}$ ;  $a'_{jk} = a'_{kj}$ ;  $b'_{jk} = b'_{kj}$ ;
end
end
```

### 3.4 The hybrid and the general PGEP Jacobi method

Recall that the *CJ* algorithm has been obtained by combining the  $LL^T J$  and the  $RR^T J$  algorithms. Actually, at each step we can combine more algorithms, say by including the *HZ* algorithm in the selection list of algorithms. Yet we cannot make decision what is the best combination of algorithms in each step, because we need to inspect the high relative accuracy and the asymptotic behavior of so obtained methods. Such research can be complex if it includes the case of the close and the multiple eigenvalues of the initial pair [5,6]. In that case, in the later stage of the process, the algorithms operate on the pair of almost diagonal symmetric matrices with a special intrinsic structure. That structure can cause the formulas for the angles to become unstable, and some modifications of the algorithms could be needed [4].

For that reason, in this paper we focus on just one important property of the methods, which is the global convergence. This property is not influenced by the stability of the algorithms. Even more, the global convergence can be proved for the pretty general Jacobi methods, which combine all known algorithms in the most general way.

**Definition 3.1** Let  $A, B$  be the symmetric matrices of the same dimension such that  $B$  is positive definite with the ones on the diagonal. Let  $\mathcal{H}$  denote a collection of the Jacobi methods for the generalized eigenvalue problem  $Ax = \lambda Bx$ , which satisfy in each step  $k$  the following two rules:

1. the pivot submatrix  $\hat{A}^{(k)}$  is diagonalized and  $\hat{B}^{(k)}$  is transformed to  $I_2$ ,
2. at least one of the two diagonal elements of  $\hat{F}_k$  is not smaller than  $\sqrt{2}/2$ .

An element of  $\mathcal{H}$  is called a *general PGEP Jacobi method*. A *hybrid Jacobi method* is any method from  $\mathcal{H}$  that uses in each step either the *HZ*, the  $LL^T J$ , or the  $RR^T J$  algorithm.

In Definition 3.1 the pivot strategy is not specified, hence any can be used. If the Jacobi method uses only the *HZ* ( $LL^T J$ ,  $RR^T J$ , *CJ*) algorithm, it will be called the *HZ* ( $LL^T J$ ,  $RR^T J$ , *CJ*) method.

In Section 4.1 we prove that the *HZ*, the  $LL^T J$ , the  $RR^T J$ , and the *CJ* methods are in  $\mathcal{H}$ . They are the special cases of the hybrid (Jacobi) method. The general PGEP Jacobi method can choose, in any step, any conceivable algorithm satisfying the above two rules. We shall briefly call it a *general Jacobi method*.

## 4 The Global Convergence

In this section we prove the global convergence of the four methods which have been derived so far, plus of the hybrid and of the general Jacobi method. The global convergence is proved under the class of generalized serial strategies from [10, Definition 3.7]. That class of cyclic strategies includes the weakly wavefront strategies from [16], and many more.

For each method we use the same notation for the iteration matrices:  $A^{(k)} = (a_{rs}^{(k)})$ ,  $B^{(k)} = (b_{rs}^{(k)})$ ,  $k \geq 0$ . The first rule in Definition 3.1 implies that all diagonal elements of each  $B^{(k)}$ ,  $k \geq 0$ , are equal to 1.

#### 4.1 The hybrid method belongs to $\mathcal{H}$

Here we also examine the general form of the transformation matrix  $\hat{F}$  which simultaneously diagonalizes the pivot submatrices  $\hat{A}$  and  $\hat{B}$ . This form is used in the global convergence proof from Subsection 4.4.

We denote the elements of  $\hat{A}$  and  $\hat{B}$  as in the relation (2.12). The elements of  $\hat{F}$  are denoted as the elements of  $\hat{Z}$  in the relation (2.13). Furthermore, let

$$\gamma = \phi - \psi.$$

We first examine the matrix  $\hat{F}$  from the  $HZ$  algorithm, i.e., the matrix  $\hat{Z}$  from the relation (2.13). The relations (2.5) – (2.8) and (2.9) imply

$$\sin \gamma = \sin(\phi - \psi) = \sin \phi \cos \psi - \sin \psi \cos \phi = 2\xi\rho = b_{ij}, \quad (4.1)$$

$$\cos \gamma = \cos(\phi - \psi) = \cos \phi \cos \psi + \sin \psi \sin \phi = \rho^2 - \xi^2 \geq 0. \quad (4.2)$$

Hence  $\cos \gamma = \sqrt{1 - b_{ij}^2}$ ,  $\gamma$  has the same sign as  $b_{ij}$ , and since  $|b_{ij}| < 1$ , we have  $-\pi/2 < \gamma < \pi/2$ . The relations (2.5)–(2.8) imply

$$\phi + \psi = 2\theta, \quad \text{hence} \quad \phi = \theta + \gamma/2, \quad \psi = \theta - \gamma/2. \quad (4.3)$$

The matrix  $\hat{F}$  from the  $HZ$  algorithm has the form

$$\hat{F} = \frac{1}{\cos \gamma} \begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \psi & \cos \psi \end{bmatrix}, \quad -\frac{\pi}{2} < \gamma = \phi - \psi < \frac{\pi}{2}, \quad \sin \gamma = b_{ij}. \quad (4.4)$$

From the relations (2.10), (2.11) we see that  $\cos \phi > 0$ ,  $\cos \psi > 0$  and the smaller angle is from the segment  $[-\pi/4, \pi/4]$ . Thus, the  $HZ$  method is an element of  $\mathcal{H}$ .

Let us examine  $\hat{F}$  from the  $LL^T J$  algorithm. From the relation (3.4), we conclude that the angle  $\psi$  equals to the angle  $\vartheta_1$  from the relation (3.1), while  $\phi = \psi + \gamma$ , where  $\cos \gamma = \sqrt{1 - b_{ij}^2}$ ,  $\sin \gamma = b_{ij}$ . Hence, the form of  $\hat{F}$  is as in the relation (4.4), with the same  $\gamma$ . We have  $\cos \psi \geq \sqrt{2}/2$ , while  $\cos \phi \geq 0$  if and only if  $\tan \psi b / \sqrt{1 - b_{ij}^2} \leq 1$ .

In particular, if  $|b_{ij}| \leq \sqrt{2}/2$  then  $\cos \phi \geq 0$ . So, the  $LL^T J$  method is also from  $\mathcal{H}$ .

The matrix  $\hat{F}$  from the  $RR^T J$  algorithm also has the form as in the relation (4.4). From the relation (3.9), we see that  $\phi$  is just  $\vartheta_2$  from the relation (3.6). Hence  $\cos \phi \geq \sqrt{2}/2$ , while  $\cos \psi \geq 0$  if and only if  $\tan \phi b / \sqrt{1 - b_{ij}^2} \geq -1$ . In particular, if  $|b_{ij}| \leq \sqrt{2}/2$  then  $\cos \psi \geq 0$ . It is also seen that  $\psi = \phi - \gamma$ , with  $\cos \gamma = \sqrt{1 - b_{ij}^2}$ ,  $\sin \gamma = b_{ij}$ , so  $\gamma$  is as in the relation (4.4). Hence the  $RR^T J$  method is also from  $\mathcal{H}$ .

We immediately conclude that *every hybrid method belongs to the class  $\mathcal{H}$* . In particular, the  $CJ$  method is an element of  $\mathcal{H}$ . Furthermore, every pivot submatrix  $\hat{F}$  appearing in the hybrid method satisfies the relation (4.4).

All these results are in accordance with the Gose's result [3], which in the case when the positive definite matrix  $\hat{B}$  has ones on the diagonal, takes the following form:

every matrix  $\hat{F}$  satisfying  $\hat{F}^T \hat{B} \hat{F} = I_2$  has the form as in the relation (4.4).

Gose's result has important implications, which we state as a remark.

*Remark 4.1* Note that  $b_{ij} = 0$  (i.e.,  $\hat{B} = I_2$ ) if and only if  $\gamma = 0$  (i.e.,  $\hat{F}$  from the hybrid method is just the Jacobi rotation for  $\hat{A}$ ). Hence,  $b_{ij} \neq 0$  if and only if  $\gamma \neq 0$ , i.e.,  $0 < \cos \gamma < 1$ . In that case  $|\cos \phi| + |\cos \psi| = |\cos(\psi + \gamma)| + |\cos \psi| > 0$ , hence at least one diagonal element of  $\hat{F}$  must be nonzero. This fact is notable but it does not warrant that the second rule in Definition 3.1 is satisfied.

Suppose that  $\hat{F}$  is such that only the first rule in Definition 3.1 is satisfied. Then  $\hat{F}$  has the form as in the relation (4.4), but  $\cos \phi$  and  $\cos \psi$  can be negative, zero or smaller than  $\sqrt{2}/2$ . If so,  $\hat{F}$  can be updated to additionally satisfy the second rule. Indeed, either from

$$\hat{\tilde{F}} = \hat{F} \begin{bmatrix} 0 & \zeta_1 \\ -\zeta_1 & 0 \end{bmatrix} \quad \text{or from} \quad \hat{\tilde{F}} = \hat{F} \begin{bmatrix} \zeta_1 & 0 \\ 0 & \zeta_1 \end{bmatrix}, \quad \zeta_1 \in \{-1, 1\},$$

we see that the value of  $\zeta_1$  can be chosen in such a way that at least one diagonal element of  $\hat{\tilde{F}}$  is not smaller than  $\sqrt{2}/2$ . So updated  $\hat{F}$  still diagonalizes  $\hat{A}$  and  $\hat{B}$ , maintains the ones on the diagonal of  $\hat{B}$  and maintains the difference between the new  $\phi$  and  $\psi$  (one can easily check that the angle shifts can be  $\pm\pi/2$ ,  $\pm\pi$  and 0). This justifies the appearance of  $\sqrt{2}/2$  in the second rule in Definition 3.1.

For example, if  $\hat{F}$  satisfies  $\hat{F}^T \hat{B} \hat{F} = I_2$  and  $\hat{F}^T \hat{A} \hat{F} = \text{diag}(a'_{ii}, a'_{jj})$ , then  $-\hat{F}$  satisfies the same relations. This change corresponds to the shift of the angles  $\phi$  and  $\psi$  by  $\pi$ , which does not change the difference  $\phi - \psi = \gamma$ , but it can make the diagonal element(s) of  $\hat{F}$  positive.

#### 4.2 The results for the general Jacobi method

Let us now consider the general Jacobi method, i.e., an arbitrary element from  $\mathcal{H}$ . In each step we have

$$\hat{F}^T \hat{A} \hat{F} = \text{diag}(a'_{ii}, a'_{jj}), \quad \hat{F}^T \hat{B} \hat{F} = I_2, \quad b_{ii} = b_{jj} = 1. \quad (4.5)$$

From (4.5) one obtains  $\hat{A} \hat{F} = \hat{F}^{-T} \text{diag}(a'_{ii}, a'_{jj})$  and  $\hat{B} \hat{F} = \hat{F}^{-T}$ , where  $\hat{F}$  is (by the Gose's result) from the relation (4.4). Writing these relations by elements, we have

$$\frac{1}{\cos \gamma} \begin{bmatrix} a_{ii} & a_{ij} \\ a_{ij} & a_{jj} \end{bmatrix} \begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \psi & \cos \psi \end{bmatrix} = \begin{bmatrix} \cos \psi & -\sin \psi \\ \sin \phi & \cos \phi \end{bmatrix} \begin{bmatrix} a'_{ii} \\ a'_{jj} \end{bmatrix}, \quad (4.6)$$

$$\frac{1}{\cos \gamma} \begin{bmatrix} 1 & b_{ij} \\ b_{ij} & 1 \end{bmatrix} \begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \psi & \cos \psi \end{bmatrix} = \begin{bmatrix} \cos \psi & -\sin \psi \\ \sin \phi & \cos \phi \end{bmatrix}. \quad (4.7)$$

The relations (4.7), (4.6) and (4.5) imply

$$\cos \gamma = \frac{\cos \phi}{\cos \psi} + b_{ij} \tan \psi = \frac{\cos \psi}{\cos \phi} - b_{ij} \tan \phi, \quad (4.8)$$

$$a'_{ii} = \frac{1}{\cos \gamma} \left( a_{ii} \frac{\cos \phi}{\cos \psi} + a_{ij} \tan \psi \right) = \frac{a_{ii} + a_{ij} \frac{\sin \psi}{\cos \phi}}{1 + b_{ij} \frac{\sin \psi}{\cos \phi}}, \quad (4.9)$$

$$a'_{jj} = \frac{1}{\cos \gamma} \left( a_{jj} \frac{\cos \psi}{\cos \phi} - a_{ij} \tan \phi \right) = \frac{a_{jj} - a_{ij} \frac{\sin \phi}{\cos \psi}}{1 - b_{ij} \frac{\sin \phi}{\cos \psi}}, \quad (4.10)$$

$$2 \cos(\phi + \psi) a_{ij} = a_{ii} \sin(2\phi) - a_{jj} \sin(2\psi). \quad (4.11)$$

If  $b_{ij} = 0$ , we have  $\gamma = 0$ , hence  $\phi = \psi$ . Then the relations (4.9)–(4.11) take the familiar form

$$a'_{ii} = a_{ii} + a_{ij} \tan \phi, \quad a'_{jj} = a_{jj} - a_{ij} \tan \phi, \quad \tan(2\phi) = 2a_{ij}/(a_{ii} - a_{jj}),$$

which is associated with the standard Jacobi method. If  $b_{ij} \neq 0$ , then at least one of the two expressions for each diagonal element is well defined (the denominator is not zero), and it can be used in computation. Actually, the formulas (4.9)–(4.10) are very attractive alternatives to the formulas (2.14)–(2.15), (3.2)–(3.3) and (3.7)–(3.8).

Next, we derive two auxiliary relations which are used in the global convergence proof. Using  $\phi = \psi + \gamma$  and some trigonometric identities, one straightforwardly obtains

$$\begin{aligned} \cos(\phi + \psi) &= \cos(2\psi + \gamma) = \cos(2\psi) - 2 \sin(\phi + \psi - \frac{\gamma}{2}) \sin(\frac{\gamma}{2}) \\ &= \cos(2\phi - \gamma) = \cos(2\phi) + 2 \sin(\phi + \psi + \frac{\gamma}{2}) \sin(\frac{\gamma}{2}), \\ \sin(2\phi) &= \sin(2\psi + 2\gamma) = \sin(2\psi) + 2 \cos(\phi + \psi) \sin \gamma. \end{aligned}$$

These identities, together with the relation (4.11), yield

$$2a_{ij} \cos(2\psi) - (a_{ii} - a_{jj}) \sin(2\psi) = 4a_{ij} \sin(2\psi + \frac{\gamma}{2}) \sin(\frac{\gamma}{2}) + 2a_{ii} \cos(2\psi + \gamma) \sin \gamma, \quad (4.12)$$

$$2a_{ij} \cos(2\phi) - (a_{ii} - a_{jj}) \sin(2\phi) = -4a_{ij} \sin(2\phi - \frac{\gamma}{2}) \sin(\frac{\gamma}{2}) + 2a_{jj} \cos(2\phi - \gamma) \sin \gamma. \quad (4.13)$$

We end this subsection by showing that the matrices generated by the general Jacobi method are bounded. To this end we need the notion of *spectral radius*.

The spectral radius of the positive definite pair  $(A, B)$  is defined by the formula

$$\mu = \max_{\lambda \in \sigma(A, B)} |\lambda|, \quad \sigma(A, B) \text{ is the spectrum of } (A, B).$$

We obviously have  $\sigma(A, B) = \sigma(A^{(k)}, B^{(k)})$ ,  $k \geq 0$ . Using the Cholesky factorization of  $B^{(k)}$  and the min-max theorem for the eigenvalues of symmetric matrices, one obtains

$$\mu = \max_{x \neq 0} \frac{|x^T A^{(k)} x|}{x^T B^{(k)} x} = \max_{\|x\|_2=1} \frac{|x^T A^{(k)} x|}{x^T B^{(k)} x}, \quad k \geq 0. \quad (4.14)$$

**Lemma 4.1** *The sequences of matrices  $(A^{(k)}, k \geq 0)$  and  $(B^{(k)}, k \geq 0)$  generated by the general Jacobi method from  $\mathcal{H}$  are bounded. In particular, we have*

$$\|B^{(k)}\|_2 < n, \quad \|A^{(k)}\|_2 \leq \mu \|B^{(k)}\|_2 < n\mu, \quad k \geq 0, \quad (4.15)$$

where  $\mu$  is the spectral radius of  $(A^{(0)}, B^{(0)})$ .

*Proof* For the sequence  $(B^{(k)}, k \geq 0)$  the assertion (4.15) is obvious since

$$\|B^{(k)}\|_2 \leq \|B^{(k)}\|_\infty < n, \quad k \geq 0.$$

The first inequality follows from the fact that each  $B^{(k)}$  is symmetric. The second one follows from the fact that all off-diagonal elements of  $B^{(k)}$  lie in the open interval  $(-1, 1)$ . From the relation (4.14) we have, for each  $k$ ,

$$|x^T A^{(k)} x| \leq \mu x^T B^{(k)} x \leq \mu \|B^{(k)}\|_2 \quad \text{for any unit vector } x.$$

Since each  $A^{(k)}$  is symmetric, there is some  $z_k$ ,  $\|z_k\|_2 = 1$ , such that  $|z_k^T A^{(k)} z_k| = \|A^{(k)}\|_2$ . This proves the second assertion of Lemma 4.1.  $\square$

### 4.3 Some auxiliary and preparatory results

The global convergence proof of the methods from  $\mathcal{H}$  is based on the following proposition whose proof is just a repetition of the corresponding proof from [3]. However, for the sake of completeness of the exposition, we present it here.

**Proposition 4.1** *Let  $B$  be a symmetric positive definite matrix with the unit diagonal. Consider the Jacobi-type iterative process of the form*

$$B^{(k+1)} = F_k^T B^{(k)} F_k, \quad k \geq 0; \quad B^{(0)} = B, \quad (4.16)$$

where  $F_k$ ,  $k \geq 0$ , are the nonsingular elementary plane matrices. Suppose the pivot submatrices  $\hat{F}_k$  of  $F_k$  are chosen to satisfy the requirement  $\hat{F}_k^T \hat{B}^{(k)} \hat{F}_k = I_2$ ,  $k \geq 0$ . Then under any pivot strategy the following assertions hold:

- (i)  $\lim_{k \rightarrow \infty} b^{(k)} = 0$ , where  $b^{(k)} = b_{i(k)j(k)}^{(k)}$  is the pivot element of  $B^{(k)}$
- (ii) There is a sequence of plane rotations  $(R_k, k \geq 0)$  such that  $F_k - R_k \rightarrow 0$  as  $k \rightarrow \infty$ .

*Proof* The assumptions imply that each  $B^{(k)}$  is a symmetric positive definite with the ones on the diagonal. Hence the absolute value of each off-diagonal element of  $B^{(k)}$  is smaller than 1. The Gose's result [3] implies that each  $\hat{F}_k$  has the form

$$\hat{F}_k = \frac{1}{\cos \gamma_k} \begin{bmatrix} \cos(\psi_k + \gamma_k) & -\sin(\psi_k + \gamma_k) \\ \sin \psi_k & \cos \psi_k \end{bmatrix}, \quad \sin \gamma_k = b^{(k)}, \quad -\frac{\pi}{2} < \gamma_k < \frac{\pi}{2}. \quad (4.17)$$

Let  $B^{(k)} = (b_{rs}^{(k)}), k \geq 0$  and

$$H(B^{(k)}) = \frac{\det(B^{(k)})}{b_{11}^{(k)} b_{22}^{(k)} \cdots b_{nn}^{(k)}} = \det(B^{(k)}), \quad k \geq 0.$$

By the Hadamard's inequality we have

$$0 < H(B^{(k)}) \leq 1, \quad k \geq 0. \quad (4.18)$$

Since  $\det(F_k) = \det(\hat{F}_k) = 1/\cos \gamma_k$ , the relation (4.16) implies

$$H(B^{(k+1)}) = \det(B^{(k+1)}) = \det^2(F_k) \det(B^{(k)}) = \frac{1}{\cos^2 \gamma_k} H(B^{(k)}), \quad k \geq 0. \quad (4.19)$$

From the relations (4.19) and (4.18) we see that  $H(B^{(k)})$  is a nondecreasing sequence of positive real numbers, bounded above by 1. Hence it is convergent with limit  $\zeta$ ,  $0 < \zeta \leq 1$ . By taking the limit on the both sides of the equation (4.19), after cancelation with  $\zeta$ , one obtains

$$1 = \lim_{k \rightarrow \infty} \cos^2 \gamma_k = \lim_{k \rightarrow \infty} (1 - (b^{(k)})^2) = 1 - \lim_{k \rightarrow \infty} (b^{(k)})^2.$$

This proves the assertion (i).

(ii) Since  $\gamma_k \rightarrow 0$  as  $k \rightarrow \infty$ , the relation (4.17) implies

$$\hat{F}_k - \hat{R}_k = \frac{1}{\cos \gamma_k} \begin{bmatrix} \cos(\psi_k + \gamma_k) & -\sin(\psi_k + \gamma_k) \\ \sin \psi_k & \cos \psi_k \end{bmatrix} - \begin{bmatrix} \cos \psi_k & -\sin \psi_k \\ \sin \psi_k & \cos \psi_k \end{bmatrix} \rightarrow 0$$

as  $k \rightarrow \infty$ . This implies the second assertion and completes the proof.  $\square$

In the following corollaries we show how Proposition 4.1 applies to the  $HZ$ , the  $LL^T J$ , the  $RR^T J$ , the hybrid and the general Jacobi method. To simplify notation we use  $a_{ii}^{(k)}$ ,  $a_{ii}^{(k+1)}$  instead of  $a_{i(k)i(k)}^{(k)}$ ,  $a_{i(k+1)i(k+1)}^{(k+1)}$ , respectively, and similar for  $a_{jj}^{(k)}$ ,  $a_{jj}^{(k+1)}$ ,  $a_{ij}^{(k)}$ .

**Corollary 4.1** *Let  $\xi^{(k)}$ ,  $\eta^{(k)}$ ,  $\rho^{(k)}$ ,  $\phi^{(k)}$ ,  $\psi^{(k)}$ ,  $\gamma^{(k)}$  and  $\theta^{(k)}$  denote the quantities  $\xi$ ,  $\eta$ ,  $\rho$ ,  $\phi$ ,  $\psi$ ,  $\gamma$  and  $\theta$ , respectively, related to the  $HZ$  method in step  $k$ ,  $k \geq 0$ . Then under any pivot strategy, the following relations hold*

$$\lim_{k \rightarrow \infty} \xi^{(k)} = 0, \quad \lim_{k \rightarrow \infty} \eta^{(k)} = 0, \quad \lim_{k \rightarrow \infty} \rho^{(k)} = 1 \quad \lim_{k \rightarrow \infty} \gamma^{(k)} = 0$$

$$\lim_{k \rightarrow \infty} |\phi^{(k)} - \theta^{(k)}| = 0 \quad \lim_{k \rightarrow \infty} |\psi^{(k)} - \theta^{(k)}| = 0 \quad (4.20)$$

$$\lim_{k \rightarrow \infty} |a_{ii}^{(k+1)} - a_{ii}^{(k)} - a_{ij}^{(k)} \tan \theta^{(k)}| = 0 \quad (4.21)$$

$$\lim_{k \rightarrow \infty} |a_{jj}^{(k+1)} - a_{jj}^{(k)} + a_{ij}^{(k)} \tan \theta^{(k)}| = 0 \quad (4.22)$$

$$\lim_{k \rightarrow \infty} |(a_{ii}^{(k)} - a_{jj}^{(k)}) \sin(2\theta^{(k)}) - 2a_{ij}^{(k)} \cos(2\theta^{(k)})| = 0. \quad (4.23)$$

We also have  $-\pi/4 \leq \theta^{(k)} \leq \pi/4$  for all  $k \geq 0$ .

*Proof* The assertions hold because Proposition 4.1(i) applies to the  $HZ$  method. For the pivot element of  $B^{(k)}$  we have  $b^{(k)} \rightarrow 0$  as  $k \rightarrow \infty$ . In particular, the first two assertions are implied by the relations (2.2) – (2.3) and (4.1) – (4.3). The assertions (4.21) and (4.22) are implied by the relations (4.9) – (4.10), (4.20), and Lemma 4.1. The assertion (4.23) is implied by the relation (2.1) and Lemma 4.1. The last assertion is implied by the relation (2.1).  $\square$

The angle  $\theta^{(k)}$  in the assertions (4.20)–(4.23) can be replaced by  $\phi^{(k)}$  or  $\psi^{(k)}$ .

**Corollary 4.2** *Let  $\phi^{(k)} = \tilde{\vartheta}_1^{(k)}$ ,  $\psi^{(k)} = \vartheta_1^{(k)}$  ( $\phi^{(k)} = \vartheta_2^{(k)}$ ,  $\psi^{(k)} = \tilde{\vartheta}_2^{(k)}$ ) denote the angles  $\tilde{\vartheta}_1$ ,  $\vartheta_1$  ( $\vartheta_2$ ,  $\tilde{\vartheta}_2$ ), respectively, related to the  $LL^T J$  ( $RR^T J$ ) method in step  $k$ ,  $k \geq 0$ . Then under any pivot strategy, the following relations hold*

$$\lim_{k \rightarrow \infty} \gamma^{(k)} = 0, \quad \text{where} \quad \gamma^{(k)} = \phi^{(k)} - \psi^{(k)} \quad (4.24)$$

$$\lim_{k \rightarrow \infty} |a_{ii}^{(k+1)} - a_{ii}^{(k)} - a_{ij}^{(k)} \tan \psi^{(k)}| = 0 \quad (4.25)$$

$$\lim_{k \rightarrow \infty} |a_{jj}^{(k+1)} - a_{jj}^{(k)} + a_{ij}^{(k)} \tan \psi^{(k)}| = 0 \quad (4.26)$$

$$\lim_{k \rightarrow \infty} |(a_{ii}^{(k)} - a_{jj}^{(k)}) \sin(2\psi^{(k)}) - 2a_{ij}^{(k)} \cos(2\psi^{(k)})| = 0. \quad (4.27)$$

In the relations (4.25)–(4.27), the angle  $\psi^{(k)}$  can be replaced by  $\phi^{(k)}$ . In addition, for the  $LL^T J$  ( $RR^T J$ ) method we have  $-\pi/4 \leq \psi^{(k)} \leq \pi/4$ ,  $k \geq 0$  ( $-\pi/4 \leq \phi^{(k)} \leq \pi/4$ ,  $k \geq 0$ ).

*Proof* Again, the assertions hold because Proposition 4.1(i) applies to these two methods. Hence the pivot element  $b^{(k)}$  of  $B^{(k)}$  tends to 0 as  $k \rightarrow \infty$ .

To prove the first assertion (4.24), one uses the relation (3.4) (resp. (3.9)).

For the  $LL^T J$  ( $RR^T J$ ) method, the assertions (4.25) and (4.26) are implied by the relations (3.2)–(3.3) (resp. (3.7)–(3.8)), the assertion (4.24), and Lemma 4.1.

The assertion (4.27) is implied by the relation (3.1) (resp. (3.6)), (4.24), and Lemma 4.1.

Because of the assertion (4.24) and Lemma 4.1, the angle  $\psi^{(k)}$  in the assertions (4.25)–(4.27) can be replaced by  $\phi^{(k)}$ .  $\square$

The next corollary summarizes the common assertions of the preceding results.

**Corollary 4.3** *Suppose that for the hybrid method from  $\mathcal{H}$ , the angles  $\phi^{(k)}$  and  $\psi^{(k)}$  have the same meaning as those in Corollary 4.1 and Corollary 4.2, depending on the method that is chosen in step  $k$ . Then the relations (4.24)–(4.27) hold and every appearance of  $\psi^{(k)}$  ( $\phi^{(k)}$ ) in the relations (4.25)–(4.27) can be replaced by  $\phi^{(k)}$  ( $\psi^{(k)}$ ). In addition, for every  $\varepsilon > 0$  there is  $k_0 \geq 0$  such that*

$$\max\{|\phi^{(k)}|, |\psi^{(k)}|\} \leq \pi/4 + \varepsilon, \quad k \geq k_0. \quad (4.28)$$

*Proof* The assertion (4.28) is implied by the second rule in the definition of  $\mathcal{H}$  and  $|\phi^{(k)} - \psi^{(k)}| \rightarrow 0$  as  $k \rightarrow \infty$ .  $\square$

In particular, Corollary 4.3 implies that the relations (4.24)–(4.27) hold for the  $CJ$  method. It remains to prove the same result for the general Jacobi method. To make a difference between the hybrid and the general method, the angles appearing in the general method shall be subscripted.

**Corollary 4.4** *Let  $\phi_k$ ,  $\psi_k$ , and  $\gamma_k$  denote the angles  $\phi$ ,  $\psi$ , and  $\gamma$  from the relation (4.4) that are obtained in step  $k$ , when the general Jacobi method from  $\mathcal{H}$  is applied to  $(A, B)$ . Then the relations (4.24)–(4.28) hold, provided that the angle superscripts are replaced by the subscripts. In the relations (4.25)–(4.27) the angles  $\psi_k$  and  $\phi_k$  can be interchanged.*

*Proof* The first assertion (4.24) is implied by the Gose's result [3] and Proposition 4.1(i), because  $\sin \gamma_k = b^{(k)}$ ,  $k \geq 0$ , where  $b^{(k)}$  is the pivot element of  $B^{(k)}$ .

To prove the assertion (4.28) note that the second rule in the definition of  $\mathcal{H}$  ensures that at least one of the angles  $\phi_k$  or  $\psi_k$  is from  $[-\pi/4, \pi/4]$ . By the first assertion we have  $\phi_k - \psi_k \rightarrow 0$  as  $k \rightarrow \infty$ . Hence, the both angles must satisfy the relation (4.28).

To prove (4.25) and (4.26), we use the following trigonometric identities

$$\begin{aligned} \frac{1}{\cos \gamma_k} \frac{\cos \phi_k}{\cos \psi_k} &= \frac{\cos(\psi_k + \gamma_k)}{\cos \gamma_k \cos \psi_k} = 1 - \tan \psi_k \tan \gamma_k, & \frac{1}{\cos \gamma_k} &= 1 + \frac{2 \sin^2(\gamma_k/2)}{\cos \gamma_k}, \\ \frac{1}{\cos \gamma_k} \frac{\cos \psi_k}{\cos \phi_k} &= \frac{\cos(\phi_k - \gamma_k)}{\cos \gamma_k \cos \phi_k} = 1 + \tan \phi_k \tan \gamma_k, & \tan \phi_k &= \tan \psi_k + \frac{(1 + \tan^2 \psi_k) \tan \gamma_k}{1 - \tan \psi_k \tan \gamma_k}. \end{aligned}$$

Now, from the relations (4.9), (4.10), and the assertion (4.15) of Lemma 4.1, we have

$$\left. \begin{aligned} |a_{ii}^{(k+1)} - a_{ii}^{(k)} - a_{ij}^{(k)} \tan \psi_k| &= |a_{ij}^{(k)} \frac{2 \sin^2(\gamma_k/2)}{\cos \gamma_k} - a_{ii}^{(k)} \tan \gamma_k| |\tan \psi_k| \leq v_k |\tan \psi_k| \\ |a_{jj}^{(k+1)} - a_{jj}^{(k)} + a_{ij}^{(k)} \tan \phi_k| &= |a_{ij}^{(k)} \frac{2 \sin^2(\gamma_k/2)}{\cos \gamma_k} - a_{jj}^{(k)} \tan \gamma_k| |\tan \phi_k| \leq v_k |\tan \phi_k| \end{aligned} \right\} \quad (4.29)$$

where

$$v_k = 2n\mu |\sin(\gamma_k/2)| / \cos(\gamma_k), \quad k \geq 0. \quad (4.30)$$

Here we have used the Cauchy-Schwarz inequality and the estimates

$$\max\{\sqrt{(a_{ij}^{(k)})^2 + (a_{ii}^{(k)})^2}, \sqrt{(a_{ij}^{(k)})^2 + (a_{jj}^{(k)})^2}\} \leq \|\hat{A}^{(k)}\|_2 \leq \|A^{(k)}\|_2 \leq n\mu, \quad k \geq 0.$$

By the Gose's result and Proposition 4.1(i) we have  $\phi_k - \psi_k = \gamma_k \rightarrow 0$  as  $k \rightarrow \infty$ . Hence  $v_k \rightarrow 0$  as  $k \rightarrow \infty$  and  $\lim_{k \rightarrow \infty} (\tan \phi_k - \tan \psi_k) = 0$ . Since  $\min\{|\tan \phi_k|, |\tan \psi_k|\} \leq 1$ , this proves the assertions (4.25) and (4.26) in which  $\phi^{(k)}$  and  $\psi^{(k)}$  are replaced by  $\phi_k$  and  $\psi_k$ , respectively.

To prove the remaining assertion, we use the relations (4.12) and (4.13). In step  $k$  we obtain

$$|2a_{ij}^{(k)} \cos(2\psi_k) - (a_{ii}^{(k)} - a_{jj}^{(k)}) \sin(2\psi_k)| \leq 2n\mu |2 \sin(\frac{\gamma_k}{2}) + \sin \gamma_k|, \quad (4.31)$$

$$|2a_{ij}^{(k)} \cos(2\phi_k) - (a_{ii}^{(k)} - a_{jj}^{(k)}) \sin(2\phi_k)| \leq 2n\mu |2 \sin(\frac{\gamma_k}{2}) + \sin \gamma_k|. \quad (4.32)$$

Here, we have used  $\max_{r,s} |a_{rs}^{(k)}| \leq \|A^{(k)}\|_2$ , Lemma 4.1 and the fact that  $\sin(\frac{\gamma_k}{2})$  and  $\sin \gamma_k$  have the same sign. Since  $\lim_{k \rightarrow \infty} \gamma_k = 0$ , the relation (4.27) is proved.  $\square$

The assertions (4.23), (4.27) and the relations (4.31), (4.32) imply that both angles  $\phi_k$  and  $\psi_k$  "asymptotically" satisfy the known equation for the angle of the Jacobi rotation. The same holds for the diagonal elements updates since  $v_k \rightarrow 0$  as  $k \rightarrow \infty$ .

From the relation (4.30) and the two relations preceding it, we have for  $k \geq 0$

$$\max\{|a_{ii}^{(k+1)} - a_{ii}^{(k)}|, |a_{jj}^{(k+1)} - a_{jj}^{(k)}|\} \leq \max\{|\tan \phi_k|, |\tan \psi_k|\} (|a_{ij}^{(k)}| + v_k). \quad (4.33)$$

#### 4.4 The convergence theorem

In the global convergence considerations we shall use the quantity  $S(A^{(k)}, B^{(k)})$ , which measures how close is the pair  $(A^{(k)}, B^{(k)})$  from the set of pairs of diagonal matrices. Here

$$S(X, Y) = [S^2(X) + S^2(Y)]^{\frac{1}{2}}, \quad S(X) = \frac{\sqrt{2}}{2} \|X - \text{diag}(X)\|_F; \quad X = X^T, Y = Y^T,$$

where  $\|X\|_F = \sqrt{\text{trace}(X^T X)}$  stands for the Frobenius norm of  $X$ . The condition  $\lim_{k \rightarrow \infty} S(A^{(k)}, B^{(k)}) = 0$  means that  $A^{(k)}$  approaches the set of diagonal matrices and  $B^{(k)}$  tends to  $I_n$ .

To prove the global convergence of the *HZ*, the *LL<sup>T</sup>*, the *RR<sup>T</sup>*, the *CJ*, the hybrid and the general Jacobi method, one can use either [9, Corollary 5.8] or [10, Corollary 5.3]. The proof is the same except at the place where one of these two references is invoked. The first reference presumes that the pivot strategy is weakly equivalent to the row-cyclic one (the so-called weakly wavefront strategy from [16]). The second reference presumes that any generalized serial strategy is used. We shall choose the latter choice because the set of generalized serial strategies is much larger and it includes the set of weakly wavefront strategies. We note that Corollary 5.3 from [10] holds for the block Jacobi methods, but if the partition which defines the matrix block-partition is specified to be  $(1, 1, \dots, 1)$ , then it holds for the element-wise methods and has the following form.

**Lemma 4.2 ([10], Corollary 5.3)** *Let  $H \neq 0$  be a symmetric matrix of order  $n$  and let the sequence  $H^{(0)} = H, H^{(1)}, \dots$  be generated by the Jacobi-type process*

$$H^{(k+1)} = F_k^T H^{(k)} F_k + E^{(k)}, \quad k \geq 0,$$

where  $F_k$  are nonsingular elementary plane matrices. Let the sequence  $(H^{(k)}, k \geq 0)$  be bounded,  $\lim_{k \rightarrow \infty} S(E^{(k)}) = 0$  and let the following three assumptions hold:

- A1** the pivot strategy of the process is generalized serial  
**A2** there is a sequence of orthogonal elementary plane matrices  $(U_k, k \geq 0)$ , such that
- $$\lim_{k \rightarrow \infty} (F_k - U_k) = 0$$
- A3** the diagonal element  $f_{ii}^{(k)}$  of  $F_k$  satisfies the relation  $\liminf_{k \rightarrow \infty} |f_{ii}^{(k)}| > 0$ .

Then the following two conditions are equivalent

- (i)  $\lim_{k \rightarrow \infty} h_{ij}^{(k+1)} = 0$ ,  
(ii)  $\lim_{k \rightarrow \infty} S(H^{(k)}) = 0$ .

Here  $H^{(k)} = (h_{rs}^{(k)})$ ,  $k \geq 0$ , and  $(i, j)$  is the pivot pair in step  $k$ .

Note that in **A3** the diagonal element  $f_{ii}^{(k)}$  can be replaced by  $f_{jj}^{(k)}$ . This is true because the assumption **A2** implies  $\lim_{k \rightarrow \infty} (|f_{ii}^{(k)}| - |f_{jj}^{(k)}|) = 0$ . We recall that the pivot indices  $i$  and  $j$  are functions of  $k$ , i.e.  $i = i(k)$ ,  $j = j(k)$ ,  $k \geq 0$ .

While Lemma 4.2 is used to prove that  $S(A^{(k)}, B^{(k)})$  tends to zero, the following lemma is used to prove the convergence of the diagonal elements. To formulate it, we need more notation.

Let the eigenvalues of the pair  $(A, B)$  be nonincreasingly ordered,

$$\lambda_1 = \cdots = \lambda_{s_1} > \lambda_{s_1+1} = \cdots = \lambda_{s_2} > \cdots > \lambda_{s_{p-1}+1} = \cdots = \lambda_{s_p}. \quad (4.34)$$

The case  $p = 1$  is not interesting, because it implies  $A = \lambda_1 B$ . Then every nonzero vector is an eigenvector belonging to the only eigenvalue  $\lambda_1$ . So, let  $p > 1$ .

If we set  $s_0 = 0$  in (4.34), we conclude that  $n_r = s_r - s_{r-1}$  is the multiplicity of  $\lambda_{s_r}$ . Furthermore, if we set  $\lambda_{s_0} = \lambda_0 = \infty$ ,  $\lambda_{s_{p+1}} = -\infty$ , then  $3\delta_r$ , where

$$3\delta_t = \min\{\lambda_{s_{t-1}} - \lambda_{s_t}, \lambda_{s_t} - \lambda_{s_{t+1}}\}, \quad 1 \leq t \leq p,$$

is the absolute gap in the spectrum of  $(A, B)$  associated with  $\lambda_{s_r}$ . Let

$$\delta = \min_{1 \leq r \leq p} \delta_r, \quad \delta_0 = \frac{\delta}{1 + \mu^2}. \quad (4.35)$$

We see that  $3\delta$  is the minimum absolute gap in the spectrum and  $\delta_0 < \delta$ .

**Lemma 4.3** *Let  $A, B$  be the symmetric matrices of order  $n$  such that  $B$  is positive definite with the unit diagonal. Let the eigenvalues of  $(A, B)$  be ordered as in the relation (4.34) and  $\delta_0$  be as in the relation (4.35). If*

$$S(A, B) < \delta_0,$$

*then there is a permutation matrix  $P$  such that for the matrix  $A' = P^T A P = (a'_{rl})$  we have*

$$2 \sum_{l=1}^n |a'_{ll} - \lambda_l|^2 \leq \frac{S^4(A, B)}{\delta_0^2}. \quad (4.36)$$

*Proof* The proof is just a reformulation of [5, Corollary 3.3], using  $\delta_0$ .  $\square$

Now, we can formulate and prove the main convergence theorem.

**Theorem 4.1** *The HZ, the  $LL^T J$ , the  $RR^T J$ , and the CJ, are globally convergent under the class of generalized serial strategies. The hybrid and the general Jacobi method are globally convergent under the same class of cyclic strategies.*

*Proof* Let us prove that  $S(A^{(k)}, B^{(k)}) \rightarrow 0$  as  $k \rightarrow \infty$  for each method. To this end, we apply Lemma 4.2 to the sequences of matrices  $(A^{(k)}, k \geq 0)$  and  $(B^{(k)}, k \geq 0)$ .

For all considered methods the pivot elements are annihilated, i.e.,  $a_{ij}^{(k+1)} = 0$ ,  $b_{ij}^{(k+1)} = 0$ . Hence the condition (i) of Lemma 4.2 is fulfilled. We also have  $E^{(k)} = 0$ ,  $k \geq 0$ . Next, we know that Lemma 4.1 holds for all considered methods. Hence the both sequences,  $(A^{(k)}, k \geq 0)$  and  $(B^{(k)}, k \geq 0)$ , are bounded.

It remains to check the assumptions **A1–A3**. The first assumption is just a selection of the pivot strategy. The second assumption is implied by Proposition 4.1(ii), and it holds for all considered methods. Finally, the assumption **A3** holds because the relation  $\phi_k - \psi_k = \gamma_k \rightarrow 0$  as  $k \rightarrow \infty$  and the assumption **A2** hold for all considered

methods. In addition, the relation (4.28) also holds for all considered methods, and it implies  $\liminf_{k \rightarrow \infty} |f_{ii}^{(k)}| \geq \sqrt{2}/2$ . Hence  $S(A^{(k)}, B^{(k)}) \rightarrow 0$  as  $k \rightarrow \infty$ . Since each  $B^{(k)}$  has the unit diagonal, this implies  $\lim_{k \rightarrow \infty} B^{(k)} = I_n$ .

It remains to prove that the diagonal elements of  $A^{(k)}$  converge. If  $p = 1$ , the proof is completed because then  $A^{(k)} = \lambda_1 B^{(k)}$  for each  $k$  and  $B^{(k)} \rightarrow I_n$  as  $k \rightarrow \infty$ .

So, let  $p > 1$ . We have to show that for large enough  $k$  the diagonal elements of  $A^{(k)}$  cannot change their affiliation to the eigenvalues of the initial pair  $(A, B)$ . It is sufficient to prove it for the general Jacobi method.

Since  $\lim_{k \rightarrow \infty} \gamma_k = 0$  and  $\lim_{k \rightarrow \infty} S(A^{(k)}, B^{(k)}) = 0$ , the relations (4.29), (4.30), and (4.28) with  $\phi_k, \psi_k$ , imply that there is an integer  $k_1$  such that

$$v_k < \frac{\delta}{2}, \quad \max\{|\tan \phi_k|, |\tan \psi_k|\} < 1.05, \quad S(A^{(k)}, B^{(k)}) < \min\{\delta_0, \frac{\delta}{2}\}, \quad k \geq k_1. \quad (4.37)$$

Here  $\delta$  and  $\delta_0$  are from the relation (4.35). By Lemma 4.3, for  $k \geq k_1$ , all diagonal elements of  $A^{(k)}$  are contained in the union of open intervals  $\cup_{r=1}^p \mathcal{D}_r$ , where

$$\mathcal{D}_r = \{x; |x - \lambda_{s_r}| < (\sqrt{2}/4)\delta\}, \quad 1 \leq r \leq p.$$

This is implied by the relations (4.36) and (4.37). Indeed, we have

$$\frac{S^2(A^{(k)}, B^{(k)})}{\sqrt{2}\delta_0} < \frac{1}{\sqrt{2}}S(A^{(k)}, B^{(k)}) < \frac{1}{\sqrt{2}}\frac{1}{2}\delta = \frac{\sqrt{2}}{4}\delta, \quad k \geq k_1.$$

Let  $a_{ii}^{(k)} \in \mathcal{D}_r, a_{jj}^{(k)} \in \mathcal{D}_t, r \neq t$ . Using the relations (4.33) and (4.37), we have

$$\begin{aligned} |a_{ii}^{(k+1)} - \lambda_{s_r}| &\leq |a_{ii}^{(k+1)} - a_{ii}^{(k)}| + |a_{ii}^{(k)} - \lambda_{s_r}| < 1.05(|a_{ij}^{(k)}| + v_k) + \frac{\sqrt{2}}{4}\delta \\ &< 1.05\left(\frac{\delta}{2} + \frac{\delta}{2}\right) + \frac{\sqrt{2}}{4}\delta = (1.05 + \frac{\sqrt{2}}{4})\delta < 1.404\delta, \end{aligned}$$

and in the similar way we obtain  $|a_{jj}^{(k+1)} - \lambda_{s_t}| < 1.404\delta$ . Hence for  $k \geq k_1$  the diagonal elements of  $A^{(k)}$  cannot change their affiliation to the eigenvalues.  $\square$

## 5 The High Relative Accuracy Experiments

Here we present several experiments in MATLAB which deal with high relative accuracy of the methods derived in earlier sections. The tests have been made on a PC with Intel(R) Core(TM) i7-2620M CPU and with 8 GiB of installed memory, under the 64-bit operating system Windows 8.1 Enterprise, using MATLAB R2016a.

Our goal is to check numerically whether some of the derived methods compute the eigenvalues of the positive definite pairs with high relative accuracy. First, we have to find a class of ‘‘well-behaved’’ matrix pairs. Roughly speaking, *a well behaved (positive definite) pair of matrices is the pair that allows only small relative perturbations of the eigenvalues and eigenvectors if the perturbation matrices are sufficiently small in some norm*. Such a pair of matrices has additional properties and

its perturbations are also somewhat special. Once we find a well-behaved pair, we can apply to it the methods and see how accurately the eigenvalues are computed. A method has high relative accuracy on that pair if it generates in finite arithmetic (in each step and cumulatively), errors that belong to that special kind of perturbations.

Our choice of well-behaved pairs is based on the following result of Drmač.

**Theorem 5.1** [1, Theorem 3.2] *Let  $A$  and  $B$  be symmetric positive definite matrices of order  $n$  and let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  be the eigenvalues of the pair  $(A, B)$ . Let  $A_S = D_A^{-1/2} A D_A^{-1/2}$ ,  $B_S = D_B^{-1/2} B D_B^{-1/2}$ , where  $D_A = \text{diag}(A)$ ,  $D_B = \text{diag}(B)$ . Let  $\delta A$  and  $\delta B$  be symmetric perturbations and  $\tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_n$  be the eigenvalues of the pair  $(A + \delta A, B + \delta B)$ . Let  $(\delta A)_S = D_A^{-1/2} \delta A D_A^{-1/2}$ ,  $\varepsilon_{A_S} = \|(\delta A)_S\|_2 / \|A_S\|_2$  and  $(\delta B)_S = D_B^{-1/2} \delta B D_B^{-1/2}$ ,  $\varepsilon_{B_S} = \|(\delta B)_S\|_2 / \|B_S\|_2$ . If*

$$\varepsilon_{A_S} \kappa_2(A_S) = \|(\delta A)_S\|_2 \|A_S^{-1}\|_2 < 1 \quad \text{and} \quad \varepsilon_{B_S} \kappa_2(B_S) = \|(\delta B)_S\|_2 \|B_S^{-1}\|_2 < 1,$$

then

$$\max_{1 \leq i \leq n} \frac{|\tilde{\lambda}_i - \lambda_i|}{\lambda_i} \leq \frac{\varepsilon_{A_S} \kappa_2(A_S) + \varepsilon_{B_S} \kappa_2(B_S)}{1 - \varepsilon_{B_S} \kappa_2(B_S)}. \quad (5.1)$$

From the theorem it follows that our class of well-behaved pairs is comprised of the pairs of symmetric positive definite matrices that can be well scaled symmetrically, i.e., for which  $\kappa_2(A_S)$  and  $\kappa_2(B_S)$  are small numbers. In addition, if the perturbation matrices can be well scaled symmetrically, i.e., if  $\varepsilon_{A_S}$  and  $\varepsilon_{B_S}$  are small, then the relative perturbations in all eigenvalues will be small.

Next, we have to find what methods generate small  $\varepsilon_{A_S^{(k)}} \kappa_2(A_S^{(k)})$  and  $\varepsilon_{B_S^{(k)}} \kappa_2(B_S^{(k)})$  in each step. Such a proof requires a detailed rounding error analysis which is a demanding task. The rounding error analysis is also used to show that all errors appearing in the process can be moved in some way to  $A^{(0)}$  and  $B^{(0)}$ . Then Theorem 5.1 can be applied just once to so perturbed starting pair  $(A^{(0)}, B^{(0)})$ .

Let us apply the Cauchy-Schwarz inequality to the numerator on the right-hand side of (5.1). We obtain  $\varepsilon_{A_S} \kappa_2(A_S) + \varepsilon_{B_S} \kappa_2(B_S) \leq \sqrt{\kappa_2^2(A_S) + \kappa_2^2(B_S)} \sqrt{\varepsilon_{B_S}^2 + \varepsilon_{B_S}^2}$ .

Recall that for all considered methods the starting matrix  $B^{(0)}$  is just  $B_S^{(0)}$ .

If some Jacobi method has high relative accuracy, then the relation

$$\rho_{(A,B)} = \max_{1 \leq i \leq n} \frac{|\tilde{\lambda}_i - \lambda_i|}{\lambda_i} / \sqrt{\kappa_2^2(A_S^{(0)}) + \kappa_2^2(B^{(0)})} \leq f(n) \mathbf{u}, \quad (5.2)$$

should hold for the pairs from our class of well-behaved pairs. In the relation (5.2),  $\tilde{\lambda}_i$ ,  $1 \leq i \leq n$  are the computed eigenvalues of the starting pair  $(A^{(0)}, B^{(0)})$ ,  $f(n)$  is a slowly growing function of  $n$  and  $\mathbf{u}$  is the machine round-off. A *strong indication that the method has high relative accuracy* can be obtained from a larger sample of pairs from our class of well-behaved pairs. We shall call it  $\Upsilon$ .

The relation (5.2) should hold regardless of the condition number  $\kappa_2(A^{(0)})$ . Therefore, it makes sense to investigate how  $\rho_{(A,B)}$  behaves with respect to  $\chi_{(A,B)}$ , where

$$\chi_{(A,B)} = \sqrt{\kappa_2^2(A^{(0)}) + \kappa_2^2(B^{(0)})}.$$

For the given sample of pairs  $Y$ , we shall make for each method its “graph of relative errors”  $\mathcal{E}$ , which is defined by

$$\mathcal{E} = \{(\chi_{(A,B)}, \rho_{(A,B)}) : (A, B) \in Y\}.$$

In MATLAB we can compute “nearly exact” eigenvalues  $\lambda_i$  using the variable precision arithmetic (vpa). The eigenvalues  $\hat{\lambda}_i$  are computed by the method that is tested with the standard double precision. Hence it is easy to compute the quantities  $\rho_{(A,B)}$  and  $\chi_{(A,B)}$ . The graph  $\mathcal{E}$  will be displayed using the MATLAB `scatter(x,y,3)` function. *The method will be indicated as high relative accurate if the y-values of the points on the graph are scattered around the machine epsilon,  $\mathbf{u} \approx 2.2 \cdot 10^{-16}$ .*

Next, we describe how to generate the pairs of symmetric positive definite matrices for numerical tests. They are determined by 4 diagonal matrices with positive diagonal elements:  $\Delta_A$ ,  $\Delta_B$ ,  $\Sigma$ ,  $\Delta$  and two orthogonal matrices  $U$ ,  $V$  of order  $n$ . The starting pair  $(A^{(0)}, B^{(0)})$  is computed in two steps:

- (1)  $F = U\Sigma V^T$ ,  $A = F^T \Delta_A F$ ,  $B = F^T \Delta_B F$ ,
- (2)  $B^{(0)} = D_B^{-1/2} B D_B^{-1/2}$  ( $= B_S$ ),  $A^{(0)} = \Delta D_A^{-1/2} A D_A^{-1/2} \Delta$  ( $= \Delta A_S \Delta$ ),

where  $D_A$  and  $D_B$  are the diagonal parts of  $A$  and  $B$ , as is defined in Theorem 5.1. The magnitudes of  $\kappa_2(A_S^{(0)})$  and  $\kappa_2(B^{(0)})$  can be controlled by the magnitudes of the diagonal entries of  $\Delta_A$ ,  $\Delta_B$ ,  $\Sigma$ . Indeed, by [18] we have  $\kappa_2(A_S^{(0)}) \leq n\kappa_2^2(\Sigma)\kappa_2(\Delta_A)$ ,  $\kappa_2(B^{(0)}) \leq n\kappa_2^2(\Sigma)\kappa_2(\Delta_B)$ , and almost always  $\kappa_2(A_S^{(0)})$  and  $\kappa_2(B^{(0)})$  are much smaller than these bounds. To simplify the construction, we have set  $\Delta_B = I_n$ .

Note that  $\kappa_2(A^{(0)}) \leq \kappa_2(A_S^{(0)})\kappa_2(\Delta)$ . If a method has high relative accuracy,  $\rho_{(A,B)}$  from the relation (5.2) should not depend on  $\kappa_2(A^{(0)})$ , which is controlled by  $\kappa_2(\Delta)$ .

If we set  $\Delta = I_n$  and  $(A^{(0)}, B^{(0)}) = (D_B^{-1/2} A D_B^{-1/2}, B_S)$ , then we know the eigenvalues of  $(A^{(0)}, B^{(0)})$  in advance. They are the quotients  $(\Delta_A)_{ll}/(\Delta_B)_{ll}$ ,  $1 \leq l \leq n$ . This can be used when considering the matrix pairs with multiple eigenvalues.

The diagonal matrices are constructed via the MATLAB function `diag(d)`, where  $d$  is a vector. The vectors are constructed by the MATLAB function `logspace`. We use it to make the diagonal matrices  $\Sigma$  and  $\Delta_A$ . For the construction of  $\Delta$  we use our m-function `scalvec(k1,k2,k3,n,k)` that generates vector  $d$  of length  $n$ ,  $d = [10^{k_1}, \dots, 10^{k_2}, \dots, 10^{k_3}]$ . Here  $k$  determines position of  $10^{k_2}$  among the components of  $d$ . To compute  $\Delta$ , `scalvec` is used within a 3-level loop, controlled by  $k_1$ ,  $k_2$ , and  $k_3$ . Altogether our main m-file uses a 7-level loop, 3 for computing  $\Delta$ , 2 for  $\Sigma$ , and 2 for  $\Delta_A$ . The orthogonal matrices  $U$  and  $V$  are computed using the QR factorization of the random matrices of order  $n$ . For example, for computing  $U$  the command `[Q,~]=qr(rand(n))` is used.

Once we have obtained  $A^{(0)}, B^{(0)}$ , we convert their copies to symbolic type. Then we apply the *variable precision arithmetic* (vpa) to those copies. We use vpa with 80 decimal digits to compute the reference eigenvalues and eigenvectors.

We have made tests for the following methods: the MATLAB `eig` function, the *HZ* method (m-function `dsyhz`), the *LL<sup>T</sup>J* method (`dsyllt`), the *RR<sup>T</sup>J* method (`dsyrrt`), and the hybrid *CJ* method (`dsylrt`). As a control method, we have used

MATLAB `eig` function employing `vpa` with 80 decimal digits. We have considered only the accuracy of the computed eigenvalues.

On input all those m-functions accept a pair  $(A, B)$  of the symmetric matrices such that  $B$  is positive definite. The m-functions use only the upper-triangles of the matrices  $A$  and  $B$ . On output each m-function yields the eigenvector matrix  $F$ , the diagonal matrix of eigenvalues and number of sweeps needed to terminate the process. We consider output of the control method accurate, and use it to compute the relative errors of the eigenvalues obtained by other methods.

Altogether, we have generated 18900 pairs of positive definite matrices of order 10. These pairs make the sample  $\mathcal{Y}$  for testing high relative accuracy of the  $HZ$ , the  $LL^T J$ , the  $RR^T J$ , and the  $CJ$  method.

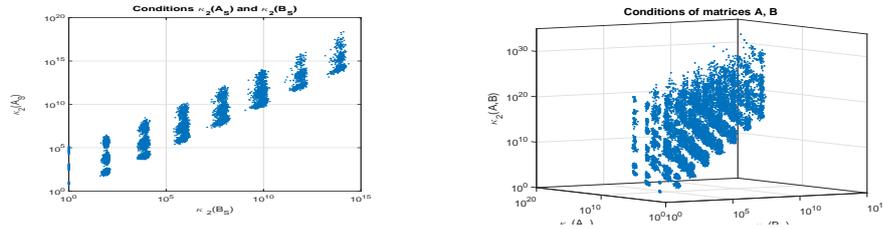
### 5.1 The results

Instead of displaying the m-files, which are used in the tests, we display figures. First, we use the MATLAB `scatter(x, y, 3)` and `scatter3(x, y, z, 3)` functions to display the sets

$$\{(\kappa_2(B_S), \kappa_2(A_S)) : (A, B) \in \mathcal{Y}\} \quad \text{and} \quad \{(\kappa_2(B_S), \kappa_2(A_S), \kappa_2(A, B)) : (A, B) \in \mathcal{Y}\},$$

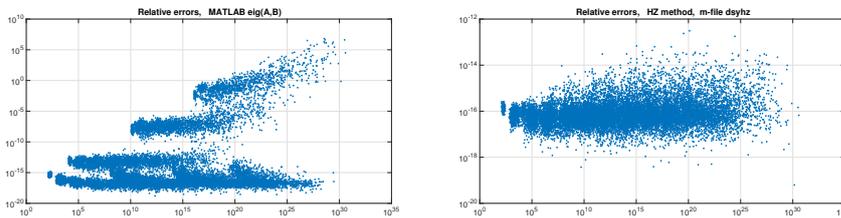
where  $\kappa_2(A, B) = \sqrt{\kappa_2^2(A^{(0)}) + \kappa_2^2(B^{(0)})}$  is the ‘‘spectral condition’’ of the pair  $(A, B)$ .

To this end, the vectors  $x, y, z$  of length 18900 hold the values of the scaled condition numbers of the starting matrices and of  $\kappa_2(A, B)$ . The  $i$ th entry of these vectors is as follows:  $x(i) = \kappa_2(B_S) = \kappa_2(B^{(0)})$ ,  $y(i) = \kappa_2(A_S)$ ,  $z(i) = \kappa_2(A, B)$ , where  $(A, B)$  is the  $i$ th matrix pair in the experiment. The following two figures are obtained:



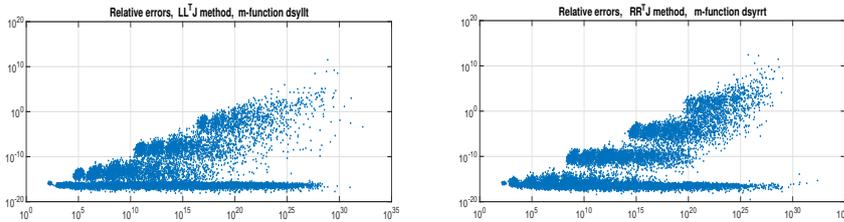
The 7 ‘‘regions’’ in the figures correspond to the 7-level loop that is used for generating the matrix pairs.

To display  $\mathcal{E}$  for the given method we use `scatter(x, y, 3)` with  $x(i) = \kappa_2(A, B)$ ,  $y(i) = \rho_{(A, B)}$ , where  $\rho_{(A, B)}$  is from the relation (5.2). The first two figures display the graphs of the results of MATLAB `eig(A, B)` function and of the  $HZ$  method:



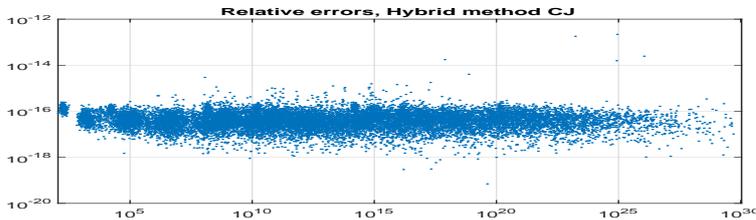
The figures indicate high relative accuracy of the  $HZ$  method, an important property that is not shared with the  $\text{eig}(A,B)$  function. Since  $A$  and  $B$  are symmetric positive definite,  $\text{eig}(A,B)$  computes by default the generalized eigenvalues using the Cholesky factorization of  $B$ .

For testing the  $LL^T J$  and the  $RR^T J$  methods we have used almost the same main and auxiliary m-scripts. The difference (in respect to the  $HZ$  method) comes from invoking different m-functions (`dsyllt` and `dsyrrt` instead of `dsyhz`). Here are the figures of the graph  $\mathcal{E}$  for these methods:



We see that the relative accuracy of the  $LL^T J$  and the  $RR^T J$  methods is similar to that of the MATLAB  $\text{eig}(A,B)$  function. We have made many modifications to `dsyllt` and `dsyrrt` m-functions in hope to enhance the relative accuracy of the computed eigenvalues. The only modification which made a huge difference was when we applied deRijk [14] pivot strategy. That strategy tries to ensure that at each step the diagonal element  $a_{ii}^{(k)}$  is larger than  $a_{jj}^{(k)}$ . We have discovered that the  $LL^T J$  and the  $RR^T J$  algorithms have to be combined in the special way, which yielded the hybrid  $CJ$  method.

We end this section with a figure of the graph  $\mathcal{E}$  of the hybrid  $CJ$  method:



It indicates that the  $CJ$  method has high relative accuracy property.

## 6 Conclusion and Future Work

In this paper we have derived several new element-wise, two-sided Jacobi methods for the PGEP. For all considered methods, the global convergence has been proved under a large class of the generalized serial strategies from [10]. The numerical tests indicate that two of them, the  $HZ$  and the  $CJ$  method, have high relative accuracy property on a sample of the well-behaved pairs of positive definite matrices. This makes them competitors to the  $FL$  method for the role of the best kernel algorithm for the PGEP block Jacobi method.

As a continuation of this research, future work will try to rigorously prove high relative accuracy of these methods. Then one-sided versions of these methods can be compared with the one-sided method from [1] which, by intention, firstly transforms the initial GSVD problem to the standard SVD problem. A very interesting problem is the asymptotic rate of convergence of the Jacobi methods in the case of multiple eigenvalues. However, the most important immediate problems are proving the global convergence and high relative accuracy of the block Jacobi methods for PGEP.

Finally, all those methods, element-wise and block, can be extended to the case of complex matrices. The above mentioned problems for real methods are also opened for the complex methods.

**Acknowledgements** The author is indebted to an anonymous reviewer for providing insightful comments and remarks. He is also thankful to V. Novaković for reading and improving the text of the manuscript.

## References

1. Drmač, Z.: A Tangent Algorithm for Computing the Generalized Singular Value Decomposition. *SIAM J. Numer. Anal.* 35 (5), 1804–1832 (1998)
2. Falk, S., Langemeyer, P.: Das Jacobische Rotations-Verfahren für reel symmetrische Matrizen-Paare I, II. *Elektronische Datenverarbeitung* 30-43 (1960)
3. Gose, G.: Das Jacobi Verfahren für  $Ax = \lambda Bx$ . *Zamm* 59, 93–101 (1979)
4. Hari, V.: On Cyclic Jacobi Methods for the Positive Definite Generalized Eigenvalue Problem. Ph.D. thesis, University of Hagen (1984)
5. Hari, V.: On Pairs of Almost Diagonal Matrices. *Linear Algebra and Its Appl.* 148, 193–223 (1991)
6. Hari, V., Drmač, Z.: On Scaled Almost Diagonal Hermitian Matrix Pairs. *SIAM J. Matrix Anal. Appl.* 18 (4), 1000-1012 (1997)
7. Hari, V., Singer, S., Singer, S.: Block-oriented  $J$ -Jacobi Methods for Hermitian Matrices. *Lin. Alg. and Its Appl.* 433 (8-10), 1491-1512 (2010)
8. Hari, V., Singer, S., Singer, S.: Full Block  $J$ -Jacobi Method for Hermitian Matrices. *Lin. Alg. and Its Appl.* 444, 1-27 (2014)
9. Hari, V.: Convergence to Diagonal Form of Block Jacobi-type Methods. *Numer. Math.* 129 (3), 449–481 (2015)
10. Hari, V., Begović Kovač, E.: Convergence of the Cyclic and Quasi-cyclic Block Jacobi Methods. To appear in *Electron. T. Numer. Ana. (ETNA)*
11. Matejaš, J.: Accuracy of the Jacobi Method on Scaled Diagonally Dominant Symmetric Matrices. *SIAM J. Matrix Anal. Appl.* 31(1), 133-153 (2009)
12. Matejaš, J.: Accuracy of one step of the Falk-Langemeyer method. *Numerical Algorithms* 68(4), 645-670 (2015)
13. Novaković, V., Singer, S., Singer, S.: Blocking and Parallelization of the Hari–Zimmermann Variant of the Falk–Langemeyer Algorithm for the Generalized SVD. *Parallel Comput.* 49, 136-152 (2015)
14. de Rijk, P. P. M.: A one-sided Jacobi algorithm for computing the singular value decomposition on a vector computer. *SIAM J. Sci. Stat. Comp.* 10, 359–371 (1989)
15. Rutishauser, H.: The Jacobi method for real symmetric matrices. *Handbook for Automatic Computation Series, Volum 2, Linear Algebra*, 202–211 (1969)  
*Numer. Math.* 9(1), 1–10 (1966) doi:10.1007/BF02165223. MR 1553948.
16. Shroff, G., Schreiber, R.: On the Convergence of the Cyclic Jacobi Method for Parallel Block Orderings. *SIAM J. Matrix Anal. Appl.* 10 (3), 326–346 (1989)
17. Slapničar, I., Hari, V.: On the Quadratic Convergence of the Falk-Langemeyer Method for Definite Matrix Pairs. *SIAM J. Matrix Anal. Appl.* 12 (1), 84-114 (1991)
18. van der Sluis, A.: Condition numbers and equilibration of matrices. *Numer. Math.* 14 (1), 14–23 (1969)
19. Veselić, K.: A Jacobi eigenreduction algorithm for definite matrix pairs, *Numer. Math.* 64 (1) 241–269 (1993).
20. Zimmermann, K.: On the Convergence of the Jacobi Process for Ordinary and Generalized Eigenvalue Problems. Ph.D. Thesis, Dissertation No. 4305 ETH, Zürich (1965)