# On the Quadratic Convergence of the Falk-Langemeyer Method

Ivan Slapničar[*]
Vjeran Hari[†]

## Abstract

The Falk–Langemeyer method for solving a real definite generalized eigenvalue problem, $Ax = \lambda Bx$, $x \neq 0$, is proved to be quadratically convergent under arbitrary cyclic pivot strategy if the eigenvalues of the problem are simple. The term "quadratically convergent" roughly means that the sum of squares of the off–diagonal elements of matrices from the sequence of matrix pairs generated by the method tends to zero quadratically per cycle.

**Key words:** generalized eigenvalue problem, Jacobi method, quadratic convergence, asymptotic convergence

**AMS(MOS) subject classification.** 65F15, 65F30

[*]Department of Mathematics, Faculty of Electrical Engineering, Mechanical Engineering and Naval Architecture, University of Split, R. Boškovića b.b., YU-58000 Split, Yugoslavia

[†]Department of Mathematics, University of Zagreb, pp.187, YU-41000 Zagreb, Yugoslavia

# 1  Introduction

In this paper we study the asymptotic convergence of the method established in 1960 by S. Falk and P. Langemeyer in [2]. Their method solves generalized eigenvalue problem

$$Ax = \lambda Bx , \qquad x \neq 0, \qquad (1)$$

where $A$ and $B$ are real symmetric matrices of order $n$ such that the pair $(A, B)$ is *definite*. By definition the pair $(A, B)$ is *definite* if the matrices $A$ and $B$ are hermitian or real symmetric and there exist real constants $a$ and $b$ such that the matrix $a\,A + b\,B$ is positive definite.

The Falk–Langemeyer method is the most commonly used Jacobi–type method for solving problem (1). Its advantages over other methods of solving problem (1) are that it applies to problem (1) for the widest class of starting pairs. Although it is not, in general, the fastest method for solving the given problem, in some cases it is the most appropriate. The QR method [11] is usually several times faster, at least on conventional computers, but it solves problem (1) only if matrix $B$ is positive definite (or positive definitizing shift for the pair is known in advance) and if matrix $B$ is well conditioned for Cholesky decomposition. The Falk–Langemeyer method is superior to the QR method in terms of numerical stability if matrix $B$ is badly conditioned for Cholesky decomposition. It is also superior to the QR method if approximate eigenvectors are known, i.e. if the matrices $A$ and $B$ are almost diagonal. This happens in the course of modeling the parameters of a system where a sequence of matrix pairs differing only slightly from each other has to be reduced. This also happens in various subspace iteration techniques (see [11]). Another reason why Jacobi–type methods have attracted attention recently is that they are adaptable for parallel processing (see [12], [10]).

The Jacobi–type method for solving problem (1) recently proposed by Veselić in [15] is somewhere in between previously mentioned methods in both, speed and requirements. Although Veselić's method works for definite matrix pairs, a linear combination $\rho A - \lambda B$ which is reasonably well conditioned for J–symmetric Cholesky decomposition must be known in advance. This method is one of the implicit methods, i.e. it works only on the eigenvectors matrix, and is therefore approximately two times faster than the Falk–Langemeyer method.

The Jacobi–type method considered by Zimmerman in [19] is closely re-

lated to the Falk–Langemeyer method (this is briefly described in Section 3) but requires positive definite matrix $B$. In [19] the convergence of this method is proved under the assumption that the starting matrices are almost diagonal. The same conclusion holds for the Falk–Langemeyer method as we shall show in this paper.

In [4] Hari studied the asymptotic convergence of complex extension of Zimmerman's method (also for positive definite $B$). He showed that his method converges quadratically under the cyclic pivot strategies if the eigenvalues of the problem are simple, while in the case of multiple eigenvalues the method can be modified so that the quadratic convergence persists. We are interested only in cyclic pivot strategies since some of them are amenable for parallel processing.

These results, the informal analysis of the convergence properties of the Falk–Langemeyer method performed by Hari in [7], and the numerical investigation suggested that the Falk–Langemeyer method behaves in the similar fashion. In this paper we prove that the Falk–Langemeyer method is quadratically convergent if the eigenvalues of the problem are simple and the pivot strategy is cyclic. The technique of the proof, originally established by the late J. H. Wilkinson in [16] (cf. [6]), is similar to that used in [4] .

Two main problems that had to be solved are that neither of the matrices $A$ and $B$ has to be positive definite and that the transformation matrices are not orthogonal and therefore difficult to estimate. Both problems were solved using the results about almost diagonal definite matrix pairs from [7].

The paper is organized as follows. In Section 2 we state the known results about almost diagonal definite matrix pairs from [7] to the extent necessary for understanding the rest of the paper. In Section 3 we describe the Falk–Langemeyer method, show that it always works for definite matrix pairs (without use of definitizing shifts), and give its algorithm. We also briefly describe Zimmerman method from [19] and [4] and relate it to the Falk–Langemeyer method. Section 4 is the central section of the paper. We first state the known result about the quadratic convergence of Zimmerman method from [4] and show to what extent can this result be applied to the Falk–Langemeyer method. We introduce measure $\widetilde{\varepsilon}_k$ which we use for defining and proving quadratic convergence. Then we prove the quadratic convergence of the Falk–Langemeyer method under the assumptions that the eigenvalues of the problem are simple, the pivot strategy is arbitrary cyclic and the matrices $A$ and $B$ are almost diagonal. At the end we show that the

quadratic convergence implies the convergence of Falk–Langemeyer method. In Section 5 we give the quadratic convergence results for parallel and serial strategies, briefly explain the possible modification of the Falk–Langemeyer method in case of multiple eigenvalues, and briefly discuss numerical experiments.

Most of the results presented in the paper are part of an M. S. thesis [13] done under the supervision of professor V. Hari.

We would like to thank professor K. Veselić from Fernuniversität Hagen for his helpful suggestions. We would also like to thank both reviewers for their comments which helped us clarify some important parts of the paper.

# 2   Almost Diagonal Definite Matrix Pairs

Here we consider the structure of almost diagonal definite matrix pair. We first state some properties of definite matrix pairs. Then we introduce chordal metric for measuring distance between eigenvalues of definite matrix pairs. We define the measures for the almost diagonality of the square matrix and of the pair of square matrices. At the end we state an important theorem from [7]. The theorem and its corollary reveal the structure of almost diagonal definite matrix pairs. All results are given for the general case of hermitian matrices even though in the rest of the paper we shall consider only the case of real symmetric matrices.

Definite matrix pair $(A, B)$ has some important properties:
a) There exists a nonsingular matrix $F$ such that

$$
\begin{aligned}
F^*AF &= \operatorname{diag}(a_1, \ldots, a_n) = D_A \\
F^*BF &= \operatorname{diag}(b_1, \ldots, b_n) = D_B.
\end{aligned}
\tag{2}
$$

The ratios $a_i/b_i$, $i = 1, \ldots, n$, of real numbers $a_i, b_i$ are the eigenvalues of the pair $(A, B)$ and are unique to the ordering. If $[f_1, \ldots, f_n]$ denotes the partition by columns of $F$, vectors $f_i$, $i = 1, \ldots, n$, are the corresponding eigenvectors. Matrices $D_A$ and $D_B$ are not uniquely determined by the pair $(A, B)$. In the real symmetric case $F^*$ can be changed to $F^T$ in the relation (2).
b) The Crawford constant $c(A, B)$,

$$
c(A, B) = \inf\{\mid x^*(A + iB)x \mid ; \ x \in \mathbf{C}^n, \parallel x \parallel = 1\}
\tag{3}
$$

4

is positive. Therefore, $A$ and $B$ share no common nul–subspace and $|a_i| + |b_i| > 0$, $i = 1, \ldots, n$, independently of the choice of $F$. Note that the choice $x = e_i$ (the i-th coordinate vector) in the relation (3) for $i = 1, \ldots, n$ implies

$$d_i \;=\; \sqrt[4]{(a_{ii})^2 + (b_{ii})^2} \;>\; 0 \;, \qquad i = 1, \ldots, n \;, \tag{4}$$

where $A = (a_{ij})$ and $B = (b_{ij})$. Hence the matrix

$$D = \operatorname{diag}\left(\frac{1}{d_1}, \ldots, \frac{1}{d_n}\right) \;, \tag{5}$$

is positive definite. In the real symmetric case for $n \neq 2$ only real vectors $x$ can be taken in the relation (3).

c) There exists a real number $\varphi$, such that the matrix $B_\varphi$ from the pair $(A_\varphi, B_\varphi)$,

$$
\begin{aligned}
A_\varphi &\;=\; A \cos \varphi \;-\; B \sin \varphi \\
B_\varphi &\;=\; A \sin \varphi \;+\; B \cos \varphi \;,
\end{aligned}
\tag{6}
$$

is positive definite. The matrices $A$ and $B$ can be simultaneously diagonalized if and only if the same holds for the matrices $A_\varphi$ and $B_\varphi$.

The proofs of the above properties are simple (see [14]). If some $f_i$ is a vector from the nul–subspace of $B$, the eigenvalue $\lambda_i$ is infinite. Such eigenvalues are not badly posed because they are zero eigenvalues of the pair $(B, A)$ counting their multiplicities. Hence, it is better to define eigenvalues as pairs of numbers $\lambda_i = [a_i, b_i]$, $i = 1, \ldots, n$. It is also necessary to choose a finite metric for measuring the distance between eigenvalues. Such is the *chordal metric*.

Let $\mathbf{R}^2 = \mathbf{R} \times \mathbf{R}$ and $\mathbf{R}_0^2 = \mathbf{R}^2 \setminus \{[0,0]\}$, where $\mathbf{R}$ is the set of real numbers. We say that the pairs $[a, b], [c, d] \in \mathbf{R}_0^2$ are *equivalent* if $ad - bc = 0$ and write $[a, b]\, \rho\, [c, d]$. It is easily seen that $\rho$ is an equivalence relation on $\mathbf{R}_0^2$. Let $\mathbf{R}_0^2\,|_\rho$ be the set of equivalence classes. Let $\lambda, \mu \in \mathbf{R}_0^2\,|_\rho$ and let $[a, b], [c, d]$ be their representatives, respectively. *Chordal distance* between $[a, b]$ and $[c, d]$ is defined with the formula

$$\chi([a, b], [c, d]) = \frac{|\, ad - bc \,|}{\sqrt{a^2 + b^2}\;\sqrt{c^2 + d^2}} \;.$$

It is easily seen that $\chi$ is constant when $[a, b]$ and $[c, d]$ vary over $\lambda$ and $\mu$, respectively. This defines metric $\widetilde{\chi} : \mathbf{R}_0^2 \mid_\rho \times \mathbf{R}_0^2 \mid_\rho \to \mathbf{R}$ by $\widetilde{\chi}(\lambda, \mu) = \chi([a, b], [c, d])$ where $[a, b]$ and $[c, d]$ are any representatives of $\lambda$ and $\mu$, respectively. However, for the sake of simplicity we shall use $\chi$ for the both functions $\chi$ and $\widetilde{\chi}$. We see that $\chi(\lambda, \mu) \leq 1$ for all $\lambda, \mu \in \mathbf{R}_0^2 \mid_\rho$. The proof of these and some other properties of the chordal metric can be found in [11], [14] and [13].

From now on, let $n$ denote the order of the matrices $A$ and $B$ and let $p$ denote the number of distinct eigenvalues of the pair $(A, B)$. We assume that

$$n \geq 3, \qquad p \geq 2,$$

and that the pair $(A, B)$ is definite. Note that if $p = 1$ then $A = \lambda B$, so $\lambda_1 = \lambda_2 = \cdots = \lambda_n = \lambda$ and *all* vectors are eigenvectors of the pair $(A, B)$.

The *off–norm of the square matrix* $A$ is the quantity

$$S(A) = \sqrt{\sum_{\substack{i,j=1 \\ i \neq j}}^{n} \mid a_{ij} \mid^2} = \| A - \operatorname{diag}(A) \|,$$

where $\| \cdot \|$ denotes the Euclidean matrix norm.

The *off–norm of the pair* $(A, B)$ is the quantity

$$\varepsilon(A, B) = \sqrt{S^2(A) + S^2(B)}. \tag{7}$$

Where no misunderstanding can arise, $\varepsilon$ shall be used instead of $\varepsilon(A, B)$. Let

$$\lambda_1 = \ldots = \lambda_{t_1}, \ \lambda_{t_1+1} = \ldots = \lambda_{t_2}, \ldots, \ \lambda_{t_{p-1}+1} = \ldots = \lambda_{t_p}, \tag{8}$$

where

$$\lambda_{t_i} = [s_i, c_i], \qquad s_i^2 + c_i^2 = 1, \qquad i = 1, \ldots, p, \tag{9}$$

be all eigenvalues of the pair $(A, B)$. Thus, we assume that the pair $(A, B)$ has $p$ distinct eigenvalues $\lambda_{t_1}, \ldots, \lambda_{t_p}$ with the appropriate multiplicities

$$n_i = t_i - t_{i-1}, \qquad i = 1, \ldots, p, \qquad t_0 = 0, \tag{10}$$

and the representatives which behave as sine and cosine are chosen. Since $p > 2$ we can define quantities

$$\delta_i = \frac{1}{3} \min_{\substack{1 \leq j \leq p \\ j \neq i}} \chi(\lambda_{t_i}, \lambda_{t_j}), \qquad \delta = \min_{1 \leq i \leq p} \delta_i. \tag{11}$$

6

Note that $\delta > 0$.

In the analysis we shall need matrices $\widetilde{A}$ and $\widetilde{B}$ defined as

$$\widetilde{A} = DAD \,, \qquad \widetilde{B} = DBD \,, \qquad (12)$$

where matrix $D$ is defined with relations (5) and (4). Since $D$ is positive definite, the matrices $\widetilde{A}$ and $\widetilde{B}$ are congruent to the matrices $A$ and $B$, respectively. Let us partition the matrices $\widetilde{A}$ and $\widetilde{B}$,

$$\widetilde{A} = \begin{bmatrix} \widetilde{A}_{11} & \cdots & \widetilde{A}_{1p} \\ \vdots & & \vdots \\ \widetilde{A}_{p1} & \cdots & \widetilde{A}_{pp} \end{bmatrix}, \qquad \widetilde{B} = \begin{bmatrix} \widetilde{B}_{11} & \cdots & \widetilde{B}_{1p} \\ \vdots & & \vdots \\ \widetilde{B}_{p1} & \cdots & \widetilde{B}_{pp} \end{bmatrix}, \qquad (13)$$

where $\widetilde{A}_{ii}$ and $\widetilde{B}_{ii}$ are diagonal blocks of order $n_i$ , $i = 1, \ldots, p$, and $n_i$'s are defined with relation (10). The relation (13) shall be written as $\widetilde{A} = (\widetilde{A}_{ij})$ and $\widetilde{B} = (\widetilde{B}_{ij})$.

Let the matrices $A$ and $B$ be partitioned according to the relation (13). *Departure from the block–diagonal form of the pair* $(A, B)$ *is the quantity*

$$\tau(A, B) \;=\; \sqrt{\tau^2(A) + \tau^2(B)} \,,$$

where

$$\tau^2(A) \;=\; \sum_{i=1}^{p} \sum_{\substack{j=1 \\ j \neq i}}^{p} \| A_{ij} \|^2 \,, \qquad \tau^2(B) \;=\; \sum_{i=1}^{p} \sum_{\substack{j=1 \\ j \neq i}}^{p} \| B_{ij} \|^2 \; .$$

THEOREM 1 *Let* $(A, B)$ *be a definite pair and let the matrices* $\widetilde{A}$ *and* $\widetilde{B}$ *be defined by the relations (12), (5) and (4). If*

$$\varepsilon(\widetilde{A}, \widetilde{B}) < \delta \,, \qquad (14)$$

*then there exists a permutation matrix* $P$ *such that for matrices* $\widetilde{A}' = P^T \widetilde{A} P$ *and* $\widetilde{B}' = P^T \widetilde{B} P$, *partitioned according to the relation (13), holds*

$$\| c_i \widetilde{A}'_{ii} - s_i \widetilde{B}'_{ii} \| \;\leq\; \frac{1}{\delta_i} \sum_{\substack{j=1 \\ j \neq i}}^{p} \| c_i \widetilde{A}'_{ij} - s_i \widetilde{B}'_{ij} \|^2 \,, \quad i = 1, \ldots, p \,. \qquad (15)$$

*On the both sides of the inequalities (15) the Euclidean matrix norm can be substituted with the spectral norm.*

PROOF: The proof of this theorem is found in [7]. Q.E.D.

COROLLARY 2 *Let the relation (14) hold for the definite pair $(A, B)$. Then there exists a permutation matrix $P$ such that for the matrices $\widetilde{A}' = P^T \widetilde{A} P = (\widetilde{a}'_{ij})$ , $\widetilde{B}' = P^T \widetilde{B} P = (\widetilde{b}'_{ij})$, $A' = P^T A P = (a'_{ij})$ and $B' = P^T B P = (b'_{ij})$, partitioned according to the relation (13), holds*

$$\sum_{i=1}^{p} \| c_i \widetilde{A}'_{ii} - s_i \widetilde{B}'_{ii} \|^2 \leq \frac{\tau^4(\widetilde{A}', \widetilde{B}')}{2\delta^2} , \tag{16}$$

$$\sum_{i=1}^{p} \sum_{j=t_{i-1}+1}^{t_i} \chi^2([s_i, c_i], [a'_{jj}, b'_{jj}]) = \sum_{i=1}^{p} \sum_{j=t_{i-1}+1}^{t_i} | c_i \widetilde{a}'_{jj} - s_i \widetilde{b}'_{jj} |^2$$

$$\leq \frac{\tau^4(\widetilde{A}', \widetilde{B}')}{2\delta^2} , \tag{17}$$

$$\chi([s_i, c_i], [a'_{jj}, b'_{jj}]) = | c_i \widetilde{a}'_{jj} - s_i \widetilde{b}'_{jj} | \leq \frac{\tau^2(\widetilde{A}', \widetilde{B}')}{2\delta} ,$$

$$j = t_{i-1} + 1, \ldots, t_i , \quad i = 1, \ldots, p . \tag{18}$$

PROOF: By Theorem 1 there exists a permutation matrix $P$ such that the relation (15) holds for the matrices $\widetilde{A}'$ and $\widetilde{B}'$. The Cauchy–Schwarz inequality implies

$$\| c_i \widetilde{A}'_{ij} - s_i \widetilde{B}'_{ij} \|^2 \leq (| c_i | \| \widetilde{A}'_{ij} \| + | s_i | \| \widetilde{B}'_{ij} \|)^2$$

$$\leq \| \widetilde{A}'_{ij} \|^2 + \| \widetilde{B}'_{ij} \|^2 , \quad i \neq j . \tag{19}$$

From the relations (15) and (19), the definition of $\tau(\widetilde{A}', \widetilde{B}')$ , and the symmetry of matrices $\widetilde{A}'$ and $\widetilde{B}'$ follows

$$\| c_i \widetilde{A}'_{ii} - s_i \widetilde{B}'_{ii} \| \leq \frac{1}{\delta_i} \sum_{\substack{j=1 \\ j \neq i}}^{p} (\| \widetilde{A}'_{ij} \|^2 + \| \widetilde{B}'_{ij} \|^2)$$

$$\leq \frac{1}{2\delta} \tau^2(\widetilde{A}', \widetilde{B}') , \quad i = 1, \ldots, p . \tag{20}$$

Finally, the relations (15), (19) and (20) and the definition of $\tau(\widetilde{A}', \widetilde{B}')$ imply

$$\sum_{i=1}^{p} \| c_i \widetilde{A}'_{ii} - s_i \widetilde{B}'_{ii} \|^2 \leq \frac{1}{2\delta} \tau^2(\widetilde{A}', \widetilde{B}') \sum_{i=1}^{p} \sum_{\substack{j=1 \\ j \neq i}}^{p} \frac{1}{\delta_i} (\| \widetilde{A}'_{ij} \|^2 + \| \widetilde{B}'_{ij} \|^2)$$

8

$$\leq \frac{1}{2\delta^2}\tau^4(\widetilde{A}', \widetilde{B}') \,,$$

which completes the proof of the relation (16).

The equalities in the relations (17) and (18) follow from the definition of the chordal metric and the fact that it does not depend upon the choice of the representatives. Inequality in the relation (17) now follows from the relation (16) and inequality in the relation (18) from the relation (20). Q.E.D.

Theorem 1 and Corollary 2 reveal the structure of almost diagonal definite matrix pairs in both the hermitian and the real symmetric case. The relation (18) implies that for $i = 1, \ldots, p$, pairs $[a'_{jj}, b'_{jj}]$, $j \in \{t_{i-1} + 1, \ldots, t_i\}$ approximate the eigenvalues $\lambda_{t_i}$ with an error of order of magnitude $\tau^2(\widetilde{A}', \widetilde{B}')$ in the chordal metric. The relation (16) implies that the blocks $\widetilde{A}'_{ii}$ and $\widetilde{B}'_{ii}$, $i = 1, \ldots, p$, are proportional with the proportionality constants being $\lambda_{t_i}$ also with the error of order $\tau^2(\widetilde{A}', \widetilde{B}')$. This proportionality becomes apparent when $\tau(\widetilde{A}', \widetilde{B}')$ is small enough compared to $\delta$. Note that the relations (15) and (18) do not imply that the off-diagonal elements of blocks $\widetilde{A}'_{ii}$ and $\widetilde{B}'_{ii}$ tend to zero together with $\tau(\widetilde{A}', \widetilde{B}')$. The relation (15) shows that for fixed $i$ the proportionality of the blocks $\widetilde{A}'_{ii}$ and $\widetilde{B}'_{ii}$ depends on the local separation $\delta_i$ of the eigenvalue $\lambda_{t_i}$ from other eigenvalues and on quantities $\| c_i \widetilde{A}'_{ij} - s_i \widetilde{B}'_{ij} \|^2$, $j = 1, \ldots, p$, $j \neq i$.

æ

# 3    The Falk–Langemeyer method

In this section we define the Falk–Langemeyer method, show that it always works for definite matrix pairs, and give its algorithm. At the end of the section we briefly describe the method of Zimmermann from [4] and [19], because it is closely related with the Falk–Langemeyer method. This relationship is also described.

The Falk–Langemeyer method solves problem (1) by constructing a sequence of "congruent" matrix pairs

$$(A^{(1)}, B^{(1)}), (A^{(2)}, B^{(2)}), \ldots \qquad (21)$$

where

$$A^{(1)} = A \,, \qquad\qquad B^{(1)} = B \,,$$

$$A^{(k+1)} = F_k^T A^{(k)} F_k , \qquad B^{(k+1)} = F_k^T B^{(k)} F_k , \qquad k \geq 1 . \quad (22)$$

Note that the transformation (22) with nonsingular matrix $F_k$ preserves the eigenvalues of the pair $(A^{(k)}, B^{(k)})$. This is a Jacobi–type method, hence the transformation matrices are chosen as nonsingular *elementary plane matrices*. An *elementary plane matrix* $F = (f_{ij})$ differs from the identity matrix only at the positions $(l, l)$, $(l, m)$, $(m, l)$ and $(m, m)$, where $1 \leq l < m \leq n$. The matrix

$$\widehat{F} = \begin{bmatrix} f_{ll} & f_{lm} \\ f_{ml} & f_{mm} \end{bmatrix}$$

is called $(l, m)$–*restriction* of the square matrix $F = (f_{ij})$.

For each $k \geq 1$, the $(l, m)$–restriction of the matrix $F_k$ has the form

$$\widehat{F}_k = \begin{bmatrix} 1 & \alpha_k \\ -\beta_k & 1 \end{bmatrix} , \qquad (23)$$

where real parameters $\alpha_k$ and $\beta_k$ are chosen to satisfy the condition

$$a_{lm}^{(k+1)} = 0 , \qquad b_{lm}^{(k+1)} = 0 . \qquad (24)$$

Here $A^{(k)} = (a_{ij}^{(k)})$ and $B^{(k)} = (b_{ij}^{(k)})$. Indices $l$ and $m$ are called *pivot indices* and the pair $(l, m)$ is called *pivot pair*. As $k$ varies the pivot pair also varies, hence $l = l(k)$ and $m = m(k)$. The transition from the pair $(A^{(k)}, B^{(k)})$ to the pair $(A^{(k+1)}, B^{(k+1)})$ is called the $k$–*th step* of the method. The manner in which we choose elements which are to be annihilated in the $k$–th step (or just the indices $(l, m)$ of these elements) is called *pivot strategy*. The pivot strategy is *cyclic* if every sequence of $N = n(n-1)/2$ successive pairs $(l, m)$ contains all pairs $(i, j), 1 \leq i < j \leq n$. A sequence of $N$ successive steps is referred to as a *cycle*. Two most common cyclic pivot strategies are the *column–cyclic strategy* and the *row–cyclic strategy*. The former is defined by the sequence of pairs

$$(1, 2), (1, 3), (2, 3), (1, 4), (2, 4), (3, 4), \ldots, (1, n), (2, n), \ldots, (n-1, n),$$

and the latt er by the sequence of pairs

$$(1, 2), (1, 3), \ldots, (1, n), (2, 3), (2, 4), \ldots, (2, n), (3, 4), \ldots, (n-1, n).$$

These two strategies are also called *serial strategies*. *Parallel cyclic strategies* are cyclic strategies which enable simultaneous execution of approximately

$n/2$ steps on parallel computers. These strategies have recently attracted considerable attention (see [12], [10]). We state the quadratic convergence results for serial and parallel strategies in Section 5.

Note that if the eigenvectors are needed, we must calculate the sequence of matrices $F^{(1)}$, $F^{(2)}$, ..., where

$$F^{(1)} = I, \qquad F^{(k+1)} = F^{(k)}F_k, \qquad k \geq 1. \qquad (25)$$

From the relations (22) and (25) we obtain for $k \geq 2$

$$\begin{aligned} F^{(k)} &= F_1 \cdots F_{k-1} \\ A^{(k)} &= (F^{(k)})^T A^{(1)} F^{(k)}, \qquad B^{(k)} = (F^{(k)})^T B^{(1)} F^{(k)}. \end{aligned}$$

We shall now derive one step of the algorithm. Note that only $(l, m)$–restrictions of the involved matrices are needed. Since (22) is the congruence transformation, the pairs $(A^{(k)}, B^{(k)})$ are definite for every $k \geq 1$ and the pairs of the corresponding $(l, m)$–restrictions are definite as well.

Let (index $k$ is omitted for the sake of simplicity)

$$\widehat{F}^T \widehat{A} \widehat{F} = \widehat{A}', \qquad \widehat{F}^T \widehat{B} \widehat{F} = \widehat{B}',$$

or respectively

$$\begin{bmatrix} 1 & -\beta \\ \alpha & 1 \end{bmatrix} \begin{bmatrix} a_{ll} & a_{lm} \\ a_{lm} & a_{mm} \end{bmatrix} \begin{bmatrix} 1 & \alpha \\ -\beta & 1 \end{bmatrix} = \begin{bmatrix} a'_{ll} & a'_{lm} \\ a'_{lm} & a'_{mm} \end{bmatrix},$$

$$\begin{bmatrix} 1 & -\beta \\ \alpha & 1 \end{bmatrix} \begin{bmatrix} b_{ll} & b_{lm} \\ b_{lm} & b_{mm} \end{bmatrix} \begin{bmatrix} 1 & \alpha \\ -\beta & 1 \end{bmatrix} = \begin{bmatrix} b'_{ll} & b'_{lm} \\ b'_{lm} & b'_{mm} \end{bmatrix}. \qquad (26)$$

Condition (24) now reads $a'_{lm} = b'_{lm} = 0$. From the relation (26) the system in unknowns $\alpha$ and $\beta$ is obtained

$$\begin{aligned} a'_{lm} &= \alpha a_{ll} + (1 - \alpha\beta)a_{lm} - \beta a_{mm} = 0, \\ b'_{lm} &= \alpha b_{ll} + (1 - \alpha\beta)b_{lm} - \beta b_{mm} = 0. \end{aligned} \qquad (27)$$

Eliminating nonlinear terms in both equations we obtain

$$\alpha = \frac{\Im_m}{\nu}, \qquad \beta = \frac{\Im_l}{\nu}, \qquad (28)$$

11

where $\nu$ is solution of the equation

$$\nu^2 - \Im_{lm}\nu - \Im_l \Im_m = 0 \tag{29}$$

and

$$
\begin{aligned}
\Im_l &= a_{ll}b_{lm} - b_{ll}a_{lm}\,, \\
\Im_m &= a_{mm}b_{lm} - b_{mm}a_{lm}\,, \\
\Im_{lm} &= a_{ll}b_{mm} - a_{mm}b_{ll}\,.
\end{aligned} \tag{30}
$$

Defining

$$\Im = (\Im_{lm})^2 + 4\Im_l \Im_m \tag{31}$$

we obtain

$$\nu_\pm = \frac{1}{2}\, \mathrm{sgn}\,(\Im_{lm})(\mid \Im_{lm} \mid \pm\sqrt{\Im})\,.$$

The algorithm is more stable if $\alpha$ and $\beta$ are smaller in modulus, so we take

$$\nu = \nu_+ = \frac{1}{2}\, \mathrm{sgn}\,(\Im_{lm})(\mid \Im_{lm} \mid +\sqrt{\Im})\,. \tag{32}$$

From the above fromula we see that the necessary condition for carrying out this step is $\Im \geq 0$. Let us show that this condition is fulfilled in each step due to the definitness of the pairs $((A^{(k)},{}^{(k)})$, $k \geq 1)$.

PROPOSITION 3 *Let the pair* $(A, B)$ *be definite. Then the following holds:*

(i) $\Im \geq 0\,,$

(ii) *The following statments are equivalent:*

    (a) $\Im = 0\,,$

    (b) $\Im_{lm} = \Im_l = \Im_m = 0\,,$

    (c) *There exist real constants* $s$ *and* $t$, $|s| + |t| \geq 0$, *such that*

$$s\widehat{A} + t\widehat{B} = O\,.$$

PROOF: The proof can be found in [4] and [13] but for the completeness of exposition we present it below.

Using the relation (6) we can define the pair $(A_\varphi, B_\varphi)$ such that the matrix $B_\varphi$ is positive definite. Let us calculate the quantities $(\Im_{lm})_\varphi, (\Im_l)_\varphi, (\Im_m)_\varphi$ and $(\Im)_\varphi$ from the pair $(A_\varphi, B_\varphi)$ using the relations (30) and (31). It is easy to verify that

$$\Im_l = (\Im_l)_\varphi, \qquad \Im_m = (\Im_m)_\varphi,$$
$$\Im_{lm} = (\Im_{lm})_\varphi, \qquad \Im = (\Im)_\varphi.$$

Therefore, without loss of generality, we can assume that the matrix $B$ from the pair $(A, B)$ is positive definite. The statement (c) is now equivalent to the statement $\widehat{A} = c\widehat{B}, c \in \mathbf{R}$ .

(i) With the notation

$$x = a_{ll}\sqrt{\frac{b_{mm}}{b_{ll}}}\,, \qquad y = a_{mm}\sqrt{\frac{b_{ll}}{b_{mm}}}\,, \qquad z = \frac{b_{lm}}{\sqrt{b_{ll}b_{mm}}}\,,$$

the following identity holds:

$$\Im = b_{ll}b_{mm}[(x-y)^2 + 4(xz - a_{lm})(yz - a_{lm})].$$

Since the right side of the above relation is the square polynomial in $a_{lm}$, we have

$$\begin{aligned}\Im &= b_{ll}b_{mm}P_2(a_{lm}) \geq b_{ll}b_{mm}P_2\left(\frac{x+y}{2}z\right)\\ &= b_{ll}b_{mm}(x-y)^2(1-z^2) = \Im_{lm}^2\left(1 - \frac{b_{lm}^2}{b_{ll}b_{mm}}\right) \geq 0 \ . \qquad (33)\end{aligned}$$

In the last inequality we have used the assumption that the matrix $B$, and therefore the matrix $\widehat{B}$, is positive definite.

(ii) Let (a) hold. The relation (33) implies that $\Im_{lm} = 0$. Matrix B is by the assumption positive definite. Therefore $b_{ll} > 0, b_{mm} > 0$, and the equality $\Im_{lm} = 0$ can be written as

$$\frac{b_{mm}}{b_{ll}}a_{ll} = a_{mm}.$$

Using this relation we can write

$$\Im_m = a_{mm}b_{lm} - b_{mm}a_{lm} = \frac{b_{mm}}{b_{ll}}(a_{ll}b_{lm} - b_{ll}a_{lm}) = \frac{b_{mm}}{b_{ll}}\Im_l \ ,$$

or $b_{ll}\Im_m = b_{mm}\Im_l$. From the definition of $\Im$, since $\Im_{lm} = 0$, we conclude that $\Im_l = \Im_m = 0$. This gives (b).

Let (b) hold. Then

$$a_{mm} = b_{mm}\frac{a_{ll}}{b_{ll}} \; , \qquad\qquad a_{lm} = b_{lm}\frac{a_{ll}}{b_{ll}} \; .$$

Therefore, $\widehat{A} = c\widehat{B}$, where

$$c = \frac{a_{ll}}{b_{ll}} = \frac{a_{mm}}{b_{mm}} \; ,$$

and (c) holds.

Let (c) hold. Then obviously (b) holds, and if (b) holds, then (a) holds.

Q.E.D.

Now we see that the Falk–Langemeyer method can be applied to all definite matrix pairs. Note that definitizing shifts are not used and need not to be known.

We have two special cases in the algorithm. If $\Im = 0$, then the matrices $\widehat{A}$ and $\widehat{B}$ are proportional as shown in Proposition 3. Therefore, the two equations in the system (27) are linearly dependent and the system has a parametric solution in one of the following forms:

$$\begin{aligned}
(\alpha, \beta) &= \left(\frac{c\,b_{mm} - b_{lm}}{b_{ll} - c\,b_{lm}}, c\right) \; , \qquad (\alpha, \beta) = \left(\frac{c\,a_{mm} - a_{lm}}{a_{ll} - c\,a_{lm}}, c\right) \; , \\
(\alpha, \beta) &= \left(c, \frac{c\,b_{ll} + b_{lm}}{b_{mm} + c\,b_{lm}}\right) \; , \qquad (\alpha, \beta) = \left(c, \frac{c\,a_{ll} + a_{lm}}{a_{mm} + c\,a_{lm}}\right) \; ,
\end{aligned}$$

where $c$ is real. For every $c$ at least one of the quotients is well defined due to the definiteness of the pair $(\widehat{A}, \widehat{B})$. It is best to set $c = 0$ to ensure that $\alpha_k$ and $\beta_k$ tend to zero together with $\varepsilon(A^{(k)}, B^{(k)})$ as $k \to \infty$ (see step (5a) in Algorithm 4). Setting $c = 0$ also reduces the operation count. This choice yields four possibilities for $(\alpha, \beta)$:

$$\left(-\frac{b_{lm}}{b_{ll}} \; , \; 0\right) \; , \qquad\qquad \left(-\frac{a_{lm}}{a_{ll}} \; , \; 0\right) \; , \tag{34}$$

$$\left(0 \; , \; \frac{b_{lm}}{b_{mm}}\right) \; , \qquad\qquad \left(0 \; , \; \frac{a_{lm}}{a_{mm}}\right) \; . \tag{35}$$

14

Due to the definiteness of the pair $(A, B)$, we have

$$|a_{ii}| + |b_{ii}| > 0, \qquad i = 1, \ldots, n \; , \qquad\qquad (36)$$

so at least one quotient is defined in each of the relations (34) and (35). In order to obtain better condition of the transformation matrix, we choose the relation in which the defined quotient has smaller absolute value. If both quotients in the chosen relation are defined, then they are equal, and for numerical reasons it is better to choose one in which the sum of squares of the numerator and the denominator is greater.

The second special case is when $\Im > 0$ and $\Im_{lm} = 0$. This means that diagonals of the matrices $\widehat{A}$ and $\widehat{B}$ are proportional, while the matrices themselves are not. Then $\mathrm{sgn}(\Im_{lm})$ is not defined. Since $\Im_l \Im_m > 0$, we have $\mathrm{sgn}(\Im_l) = \mathrm{sgn}(\Im_m)$. Substituting $\mathrm{sgn}(\Im_{lm})$ with $\mathrm{sgn}(\Im_l)$ in the equation (32) gives

$$\nu = \mathrm{sgn}(\Im_l)\sqrt{\Im_l \Im_m} \; .$$

Inserting this in the equation (28) gives, after simple calculation,

$$\alpha = \sqrt{\frac{b_{mm}}{b_{ll}}} = \sqrt{\frac{a_{mm}}{a_{ll}}} \; , \qquad\qquad \beta = \frac{1}{\alpha} \; . \qquad\qquad (37)$$

The relation (36) implies that at least one of the quotients $b_{mm}/b_{ll}$ and $a_{mm}/a_{ll}$ is defined and different from zero. If both quotients are defined then they are equal and it is better to choose one in which the sum of squares of the numerator and the denominator is greater.

We can now define an algorithm of the method:

ALGORITHM 4 Definite matrix pair $(A, B)$ is given.

(1) Set $k = 1$, $A^{(1)} = A$, $B^{(1)} = B$, $F^{(1)} = I$ and choose the pivot strategy.

(2) Choose the pivot pair $(l, m) = (l(k), m(k))$ according to the strategy.

(3) If $a_{lm}^{(k)} = b_{lm}^{(k)} = 0$, then set $k = k + 1$, $A^{(k+1)} = A^{(k)}$, $B^{(k+1)} = B^{(k)}$, $F^{(k+1)} = F^{(k)}$ and go to step (2). Otherwise go to step (4).

15

(4) Calculate the quantities $\Im_l^{(k)}$, $\Im_m^{(k)}$, $\Im_{lm}^{(k)}$ and $\Im^{(k)}$ from formulas

$$\Im_l^{(k)} = a_{ll}^{(k)} b_{lm}^{(k)} - b_{ll}^{(k)} a_{lm}^{(k)} , \qquad \Im_m^{(k)} = a_{mm}^{(k)} b_{lm}^{(k)} - b_{mm}^{(k)} a_{lm}^{(k)} ,$$

$$\Im_{lm}^{(k)} = a_{ll}^{(k)} b_{mm}^{(k)} - a_{mm}^{(k)} b_{ll}^{(k)} , \qquad \Im^{(k)} = (\Im_{lm}^{(k)})^2 + 4\Im_l^{(k)} \Im_m^{(k)} .$$

(5)

(a) If $\Im^{(k)} = 0$ perform the following steps: If $\mid b_{ll}^{(k)} \mid \geq \mid a_{ll}^{(k)} \mid$ , then set $\alpha_k = -b_{lm}^{(k)}/b_{ll}^{(k)}$ ;

otherwise set $\alpha_k = -a_{lm}^{(k)}/a_{ll}^{(k)}$ .

If $\mid b_{mm}^{(k)} \mid \geq \mid a_{mm}^{(k)} \mid$ , then set $\beta_k = b_{lm}^{(k)}/b_{mm}^{(k)}$ ;

otherwise set $\beta_k = a_{lm}^{(k)}/a_{mm}^{(k)}$ .

Finally, if $\mid \alpha_k \mid \geq \mid \beta_k \mid$ , then set $\alpha_k = 0$ ; otherwise set $\beta_k = 0$ .

(b) If $\Im^{(k)} > 0$ perform the following steps:

(i) If $\Im_{lm}^{(k)} \neq 0$ , then calculate

$$\nu_k = \frac{1}{2}\,\mathrm{sgn}\,(\Im_{lm}^{(k)})(\mid \Im_{lm}^{(k)} \mid + \sqrt{\Im^{(k)}}),$$

$$\alpha_k = \frac{\Im_m^{(k)}}{\nu_k} , \qquad \beta_k = \frac{\Im_l^{(k)}}{\nu_k} .$$

(ii) If $\Im_{lm}^{(k)} = 0$ , then, according to the relation (37), calculate

$$\alpha_k = \sqrt{\frac{b_{mm}^{(k)}}{b_{ll}^{(k)}}} = \sqrt{\frac{a_{mm}^{(k)}}{a_{ll}^{(k)}}} , \qquad \beta_k = \frac{1}{\alpha_k}.$$

If both quotients for $\alpha_k$ are defined, then choose one in which the sum of squares of the numarator and the denominator is greater.

16

(6) Perform the transformation

$$A^{(k+1)} = F_k^T A^{(k)} F_k, \qquad\qquad B^{(k+1)} = F_k^T B^{(k)} F_k , \qquad\qquad (38)$$

$$F^{(k+1)} = F^{(k)} F_k . \qquad\qquad (39)$$

(7) Set $k = k + 1$ and move to step (2). $\blacksquare$

Since matrices $A^{(k)}$, $B^{(k)}$, $A^{(k+1)}$ and $B^{(k+1)}$ are symmetric, it is enough to store and to transform only upper triangles. In the transformation (38) only $l-$th and $m-$th row and column of the matrices $A^{(k)}$ and $B^{(k)}$ are changed and in the transformation (39) only $l-$th and $m-$th columns of the matrix $F^{(k)}$ are changed. Note that the eigenvalues can be found without calculating the matrices $F^{(k)}$, $k \geq 1$, and therefore the trasformation (39) can be omitted. This reduces the operational count about fifty percent.

Stopping criteria of the infinite iterative procedure defined with this algorithm are described in Section 5.

From now on, the term "Falk–Langemeyer method" denotes the method described by Algorithm 4.

**The Zimmerman method.** We shall now relate the Falk–Langemeyer method with another method for solving the generalized eigenvalue problem. This method is due to K. Zimmermann who roughly described it in her thesis [19]. Later on, in his thesis [4], Hari derived its algorithm and proved its quadratic convergence.

The Zimmermann method is defined for symmetric matrix pairs $(A, B)$ where matrix $B$ is positive definite. We shall denote this fact as $B > 0$. At the beginning of the iterative procedure the initial pair $(A, B)$ is normalized such that

$$A^{(1)} = DAD , \qquad\qquad B^{(1)} = DBD ,$$

where

$$D = \mathrm{diag} \left( \frac{1}{\sqrt{b_{11}}}, \ldots, \frac{1}{\sqrt{b_{nn}}} \right) .$$

Therefore, $b_{ii}^{(1)} = 1, i = 1, \ldots, n$. The Zimmerman method constructs a sequence of pairs $((A^{(k)}, B^{(k)}), k \geq 1)$ by the rule

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

17

The nonsingular matrices $Z_k$ are chosen to preserve the units on the diagonal of $B^{(k+1)}$ (automatic normalization at each step) and to annihilate the pivot elements. In [4] it is shown that for $k \geq 1$ holds

$$\widehat{Z}_k = \frac{1}{\sqrt{1 - \left(b_{lm}^{(k)}\right)^2}} \left[ \begin{array}{cc} \cos \varphi_k & \sin \varphi_k \\ -\sin \psi_k & \cos \psi_k \end{array} \right] ,$$

where

$$
\begin{array}{rcl}
\cos \varphi_k & = & \cos \theta_k + \xi_k (\sin \theta_k - \eta_k \cos \theta_k) , \\
\sin \varphi_k & = & \sin \theta_k - \xi_k (\cos \theta_k + \eta_k \sin \theta_k) , \\
\cos \psi_k & = & \cos \theta_k - \xi_k (\sin \theta_k + \eta_k \cos \theta_k) , \\
\sin \psi_k & = & \sin \theta_k + \xi_k (\cos \theta_k - \eta_k \sin \theta_k) ,
\end{array}
$$

$$
\begin{array}{rcl}
\xi_k & = & \dfrac{b_{lm}^{(k)}}{\sqrt{1 + b_{lm}^{(k)}} + \sqrt{1 - b_{lm}^{(k)}}} , \\[3ex]
\eta_k & = & \dfrac{b_{lm}^{(k)}}{(1 + \sqrt{1 + b_{lm}^{(k)}})(1 + \sqrt{1 - b_{lm}^{(k)}})} , \\[3ex]
\tan 2\theta_k & = & \dfrac{2a_{lm}^{(k)} - (a_{ll}^{(k)} + a_{mm}^{(k)})b_{lm}^{(k)}}{(a_{mm}^{(k)} - a_{ll}^{(k)})\sqrt{1 - \left(b_{lm}^{(k)}\right)^2}} , \\[3ex]
& & -\dfrac{\pi}{4} \leq \theta_k \leq \dfrac{\pi}{4} .
\end{array}
$$

If $a_{lm}^{(k)} = b_{lm}^{(k)} = 0$ we set $\theta_k = 0$. If the $(l, m)$–restrictions of $A^{(k)}$ and $B(k)$ are proportional and $a_{lm}^{(k)}$ and $b_{lm}^{(k)}$ are not both equal to zero, we set $\theta_k = \frac{\pi}{4}$.

If the matrix $B$ is not positive definite but the pair $(A, B)$ is, then there exists a definitizing shift $\mu$ such that the matrix $A - \mu B$ is positive definite. If this shift is known in advance, then the Zimmermann method can be applied to the pair $(A, B)$ in the sense that each $Z_k$ is computed from the pair $(B^{(k)}, A^{(k)} - \mu B^{(k)})$.

Although the Zimmermann method seems quite different from the Falk–Langemeyer method, the two methods are closely related. The following theorem gives precise formulation of this relationship. For this occasion only

we assume that in step (5a) of Algorithm 4 (that is when $\Im^{(k)} = 0$), parameters $\alpha_k$ and $\beta_k$ are computed according to the formulae (37). For this version of the Falk–Langemeyer method holds:

THEOREM 5 *Let $A$ and $B$ be symmetric matrices of order $n$ and let $B$ be positive definite. Let the sequences $((A^{(k)}, B^{(k)}), k \geq 1)$ and $((A^{(k)'}, B^{(k)'}), k \geq 1)$ be generated from the starting pair $(A, B)$ with the Falk–Langemeyer and the Zimmermann method, respectively. If the corresponding pivot strategies are the same, then*

$$A^{(k)'} = D^{(k)} A^{(k)} D^{(k)}, \qquad B^{(k)'} = D^{(k)} B^{(k)} D^{(k)}, \qquad k \geq 1 \,,$$

*where*

$$D^{(k)} = \mathrm{diag}\,\Big(\frac{1}{\sqrt{b_{11}^{(k)}}}, \ldots, \frac{1}{\sqrt{b_{nn}^{(k)}}}\Big) \,, \qquad k \geq 1 \ .$$

PROOF: The proof of this theorem is found in [4] Section 2.3. \qquad Q.E.D.

Let us suppose again that the matrix $B$ is not positive definite while the pair $(A, B)$ is, and that a positive definitizing shift $\mu$ is known in advance. Let us apply to the pair $(A, B)$ the Zimmermann method in the sence mentioned above and the version of the Falk–Langemeyer method which we used in Theorem 5. It is easy to see that the parameters $\alpha_k$ and $\beta_k$ from the Falk–Langemeyer method are invariant under the transformations $(A, B) \to (B, A - \mu B)$. Therefore, Theorem 5 holds in this case, as well, with

$$D^{(k)} \;=\; \mathrm{diag}\,\left(\frac{1}{\sqrt{a_{11}^{(k)} - \mu\, b_{11}^{(k)}}}, \ldots, \frac{1}{\sqrt{a_{nn}^{(k)} - \mu\, b_{nn}^{(k)}}}\right) \,, \qquad k \geq 1 \ .$$

We can conclude that *if the starting pair is positive definite or the definitizing shift is known in advance, then the Falk–Langemmeyer (Zimmermann) method is the fast scaled (normalized) version of the Zimmermann (Falk–Langemmeyer) method.*

19

# 4  Quadratic convergence

In this section we prove that the Falk–Langemeyer method is quadratically convergent if the starting definite pair has simple eigenvalues and the pivot strategy is cyclic. Definitizing shifts are not used and need not to be known. We first state the result about the quadratic convergence of the Zimmermann's method, and show to what extent can this result be applied to the Falk–Langemeyer method if the matrix $B$ is positive definite. Then we define the quadratic convergence for the Falk–Langemeyer method. In Subsection 4.1 we prove preliminary results which we use in the proof of the quadratic convergence of the Falk–Langemeyer method in Subsection 4.2.

The result about the quadratic convergence of Zimmermann method can be summarized as follows. Let the sequence $((A^{(k)}, B^{(k)}), k \geq 1)$ be generated by the Zimmerman method from the pair $(A, B)$, $B > 0$, and let $\varepsilon_k = \varepsilon(A^{(k)}, B^{(k)})$, where $\varepsilon$ is defined with the relation (7). Note that $\varepsilon_k$ is natural measure for convergence of the Zimmerman method since each matrix $B^{(k)}$ has units along the diagonal.

We say that the Zimmerman method is *quadratically convergent* on the pair $(A, B)$ if $\varepsilon_k \to 0$ as $k \to \infty$ and there exist a constant $c_0 > 0$ and an integer $r_0$ such that for $r \geq r_0$ holds

$$\varepsilon_{(r+1)N+1} \; \leq \; c_0 \, \varepsilon_{rN+1}^2 \; .$$

Hence of special importance are conditions under which the above relation holds for $r = 1$. We call them *asymptotic assumptions*. Let

$$\sigma \; = \; \mathrm{spr}\,(A, B) \; = \; \max_{1 \leq i \leq n} \, |\lambda_i| \; , \qquad \delta \; = \; \frac{1}{3} \min_{i \neq j} \mid \lambda_i - \lambda_j \mid \; .$$

THEOREM 6 *Let the sequence $((A^{(k)}, B^{(k)}, k \geq 1)$ be generated by the Zimmerman method from the starting pair $(A, B)$, $B > 0$, and let the asymptotic assumptions*

$$S(B^{(1)}) \leq \frac{1}{2N} \; , \qquad 2\sqrt{1 + \sigma^2}\, \varepsilon_1 < \delta \; , \tag{40}$$

*hold. If the eigenvalues of the pair $(A, B)$ are simple and the pivot strategy is cyclic, then*

$$\varepsilon_{N+1} \; \leq \; \sqrt{N(1 + \sigma^2)}\, \frac{\varepsilon_1^2}{\delta} \; . \tag{41}$$

PROOF: The proof of this theorem is found in [4] Section 3.3.     Q.E.D.

In Theorem 6 the term $\sigma$ appears in the assumption (40) and in the assertion (41) because matrix $B$ is not diagonal and matrix $A$ is not normalized. From Theorem 5 we see that Theorem 6 holds for the Falk–Langemeyer method provided that the step (5a) of Algorithm 4 is appropriately changed, the matrix $B$ is positive definite, and the pairs $(A^{(k)}, B^{(k)})$ generated by the Falk–Langemeyer method are normalized so that $b_{ii}^{(k)} = 1, i = 1 \ldots n, k \geq 1$.

In the rest of this section we prove that the Falk–Langemeyer method defined with Algorithm 4 is quadratically convergent on definite matrix pairs with simple eigenvalues if the pivot strategy is cyclic. We first have to define the measure for the quadratic convergence.

Let $(A, B)$ be a definite pair. We shall use the measure $\widetilde{\varepsilon} = \widetilde{\varepsilon}(A, B)$ defined by

$$\widetilde{\varepsilon}(A, B) = \varepsilon(\widetilde{A}, \widetilde{B}),$$

where $\widetilde{A}$ and $\widetilde{B}$ are given by the relations (12),(5) and (4). The measure $\widetilde{\varepsilon}$ enables us to use results of Corollary 2 and it takes into account the diagonal elements of matrices $A$ and $B$. Note that the measure $\varepsilon(A, B)$ is generally not the proper measure for almost diagonality of the pair $(A, B)$ since it takes no account of the diagonals of matrices $A$ and $B$.

Let the sequence of pairs

$$(A^{(1)}, B^{(1)}), (A^{(2)}, B^{(2)}), \ldots \tag{42}$$

be generated by the Falk–Langemeyer method from the starting definite pair $(A, B)$. For $k \geq 1$ we set

$$\widetilde{\varepsilon}_k = \widetilde{\varepsilon}(A^{(k)}, B^{(k)}) = \varepsilon(\widetilde{A}^{(k)}, \widetilde{B}^{(k)}) , \tag{43}$$

$$\widetilde{A}^{(k)} = D_k A^{(k)} D_k , \qquad \widetilde{B}^{(k)} = D_k B^{(k)} D_k , \tag{44}$$

$$D_k = \operatorname{diag}\left(\frac{1}{d_1^{(k)}} , \ldots , \frac{1}{d_n^{(k)}}\right) , \tag{45}$$

$$d_i^{(k)} = \sqrt[4]{(a_{ii}^{(k)})^2 + (b_{ii}^{(k)})^2} , \ i = 1, \ldots, n . \tag{46}$$

From the relations (44), (45) and (46) we see that the pairs $(\widetilde{A}^{(k)}, \widetilde{B}^{(k)})$ are normalized in the sence that

$$(\widetilde{a}_{ii}^{(k)})^2 + (\widetilde{b}_{ii}^{(k)})^2 = 1 , \qquad i = 1, \ldots, n . \tag{47}$$

21

DEFINITION **7** *The Falk-Langemeyer method is* quadratically convergent *on the pair* $(A, B)$ *if* $\widetilde{\varepsilon}_k \to 0$ *as* $k \to \infty$ *and there exist a constant* $c_0 > 0$ *and an integer* $r_0$ *such that for* $r \geq r_0$ *holds*

$$\widetilde{\varepsilon}_{(r+1)N+1} \leq c_0 \widetilde{\varepsilon}_{rN+1}^2.$$ (48)

From Definition 7 we see that ultimately $\widetilde{\varepsilon}_k$ decreases quadratically per cycle. At the end of Subsection 4.2 we shall show that the quadratic convergence implies the convergence of the sequence (42) towards the pair of diagonal matrices $(D_A, D_B)$, where

$$D_A = \text{diag}\,(a_1, \ldots, a_n), \qquad D_B = \text{diag}\,(b_1, \ldots, b_n)\,.$$ (49)

Here $\lambda_i = [a_i, b_i], i = 1, \ldots, n$, are the eigenvalues of the pair $(A, B)$. Finally, we shall show that ultimately the quadratic reduction of $\widetilde{\varepsilon}_{rN+1}$ implies the quadratic reduction of $\varepsilon_{rN+1}$ and vice versa[1].

In order to be able to observe the measure $\widetilde{\varepsilon}$ we must solve one more problem. The transformation matrices $F_k$ are calculated from unnormalized pairs $(A^{(k)}, B^{(k)})$ and are therefore difficult to estimate. To solve this problem we shall observe the sequence obtained from the pair $(A, B)$ with following process:

> *normalization, step of the method, normalization, step of the method,...*

This sequence reads

$$(\overline{\widetilde{A}}^{(1)}, \overline{\widetilde{B}}^{(1)}), (\overline{A}^{(2)}, \overline{B}^{(2)}), (\overline{\widetilde{A}}^{(2)}, \overline{\widetilde{B}}^{(2)}), (\overline{A}^{(3)}, \overline{B}^{(3)}), (\overline{\widetilde{A}}^{(3)}, \overline{\widetilde{B}}^{(3)}), \ldots,$$ (50)

where

$$(\overline{\widetilde{A}}^{(1)}, \overline{\widetilde{B}}^{(1)}) = (\widetilde{A}^{(1)}, \widetilde{B}^{(1)})\,,$$ (51)

and for $k \geq 1$ holds

$$\overline{A}^{(k+1)} = \widetilde{\overline{F}}_k^T \widetilde{\overline{A}}^{(k)} \widetilde{\overline{F}}_k\,, \qquad \overline{B}^{(k+1)} = \widetilde{\overline{F}}_k^T \widetilde{\overline{B}}^{(k)} \widetilde{\overline{F}}_k\,,$$ (52)

---

[1] Here $\varepsilon_k$ measures off–diagonal elements of the pairs from sequence (42) and should not be confused with the quantity used in connection with Zimmerman method.

$$\widetilde{\overline{A}}^{(k+1)} \;=\; \overline{D}_{k+1}\overline{A}^{(k+1)}\overline{D}_{k+1}\,, \qquad\qquad \widetilde{\overline{B}}^{(k+1)} = \overline{D}_{k+1}\overline{B}^{(k+1)}\overline{D}_{k+1}\,, \quad (53)$$

$$\overline{D}_{k+1} \;=\; \operatorname{diag}\big(\tfrac{1}{\overline{d}_1^{(k+1)}}\,,\ldots,\ \tfrac{1}{\overline{d}_n^{(k+1)}}\big)\,, \qquad\qquad\qquad\qquad (54)$$

$$\overline{d}_i^{(k+1)} \;=\; \sqrt[4]{(\overline{a}_{ii}^{(k+1)})^2 + (\overline{b}_{ii}^{(k+1)})^2}\,, \qquad\qquad i = 1,\ldots,n\,. \quad (55)$$

Of course, the sequences (42) and (50) are generated using the same pivot strategy. The matrices $\widetilde{F}_k$ are calculated according to Algorithm 4, but from the pairs $(\widetilde{\overline{A}}^{(k)}, \widetilde{\overline{B}}^{(k)})$. Since in the transition from $(\widetilde{\overline{A}}^{(k)}, \widetilde{\overline{B}}^{(k)})$ to $(\overline{A}^{(k+1)}, \overline{B}^{(k+1)})$ of all diagonal elements only those at positions $(l,l)$ and $(m,m)$ are being changed, we conclude that

$$\overline{d}_i^{(k+1)} = \sqrt[4]{(\overline{a}_{ii}^{(k+1)})^2 + (\overline{b}_{ii}^{(k+1)})^2} = \sqrt[4]{(\widetilde{a}_{ii}^{(k)})^2 + (\widetilde{b}_{ii}^{(k)})^2} = 1,\ i = 1,\ldots,n,\ i \neq l,m\,.$$
$$(56)$$

We will now show that the operations of normalization and of carrying out one step of the algorithm commute. This is equivalent to showing that $\widetilde{\overline{A}}^{(k)} = \widetilde{A}^{(k)}$ and $\widetilde{\overline{B}}^{(k)} = \widetilde{B}^{(k)}$ for $k \geq 1$.

Let $\widetilde{F}_k$ be the transformation matrices obtained according to Algorithm 4 from the pairs $(\widetilde{A}^{(k)}, \widetilde{B}^{(k)})$, $k \geq 1$. The following proposition shows that the matrices $F_k$ and $\widetilde{F}_k$ are simply related.

PROPOSITION 8 *For $k \geq 1$ holds* $\widetilde{F}_k \;=\; D_k^{-1} F_k D_k$.

PROOF: Because of the relations (44), (45) and (46) we have

$$\widehat{\widetilde{A}}^{(k)} = \begin{bmatrix} \dfrac{a_{ll}^{(k)}}{(d_l^{(k)})^2} & \dfrac{a_{lm}^{(k)}}{d_l^{(k)} d_m^{(k)}} \\[2ex] \dfrac{a_{lm}^{(k)}}{d_l^{(k)} d_m^{(k)}} & \dfrac{a_{mm}^{(k)}}{(d_m^{(k)})^2} \end{bmatrix}\,, \qquad \widehat{\widetilde{B}}^{(k)} = \begin{bmatrix} \dfrac{b_{ll}^{(k)}}{(d_l^{(k)})^2} & \dfrac{b_{lm}^{(k)}}{d_l^{(k)} d_m^{(k)}} \\[2ex] \dfrac{b_{lm}^{(k)}}{d_l^{(k)} d_m^{(k)}} & \dfrac{b_{mm}^{(k)}}{(d_m^{(k)})^2} \end{bmatrix}\,. \quad (57)$$

The assertion is now obtained by simply using the relation (57) in Algorithm 4 and calculating the matrix $\widetilde{F}_k$. \hfill Q.E.D.

PROPOSITION 9 *For $k \geq 1$ the following holds:*

$$(i) \qquad\qquad (\widetilde{\overline{A}}^{(k)}, \widetilde{\overline{B}}^{(k)}) \;=\; (\widetilde{A}^{(k)}, \widetilde{B}^{(k)})\,,$$

*(ii)* $\qquad D_k \;=\; D_1 \overline{D}_2 \overline{D}_3 \cdots \overline{D}_k \;.$

PROOF: The proof is by induction in respect to $k$.

$(i)$ For $k = 1$ the assertion holds due to the relation (51). Suppose that the assertion holds for some $k \geq 1$. This means that

$$\widetilde{\overline{A}}^{(k)} = \widetilde{A}^{(k)} \;, \qquad \widetilde{\overline{B}}^{(k)} = \widetilde{B}^{(k)} \;, \qquad \widetilde{\overline{F}}_k = \widetilde{F}_k \;. \qquad (58)$$

From the relation (52) it follows that $\overline{A}^{(k+1)} = \widetilde{\overline{F}}_k^{\,T} \widetilde{\overline{A}}^{(k)} \widetilde{\overline{F}}_k$ , which, because of the relation (58), implies that $\overline{A}^{(k+1)} = \widetilde{F}_k^{\,T} \widetilde{A}^{(k)} \widetilde{F}_k$. Since the relation (44) and Proposition 8 imply

$$
\begin{aligned}
\overline{A}^{(k+1)} &= D_k F_k^T D_k^{-1} D_k A^{(k)} D_k D_k^{-1} F_k D_k = D_k F_k^T A^{(k)} F_k D_k \\
&= D_k A^{(k+1)} D_k \;,
\end{aligned} \qquad (59)
$$

we conclude that normalizations of the matrices $A^{(k+1)}$ and $\overline{A}^{(k+1)}$ give the same matrix. Now we use the same argument to show that $\widetilde{\overline{B}}^{(k+1)} = \widetilde{B}^{(k+1)}$ for $k \geq 1$ and to prove $(i)$.

$(ii)$ For $k = 1$ the assertion is trivially fulfilled. Let the assertion hold for some $k \geq 1$. From the relations (59), (53) and the assertion $(i)$ we obtain

$$
\begin{aligned}
\overline{D}_{k+1} D_k A^{(k+1)} D_k \overline{D}_{k+1} &= \overline{D}_{k+1} \overline{A}^{(k+1)} \overline{D}_{k+1} = \\
= \widetilde{\overline{A}}^{(k+1)} &= \widetilde{A}^{(k+1)} = D_{k+1} A^{(k+1)} D_{k+1} \;.
\end{aligned}
$$

It is obvious that $D_{k+1} = D_k \overline{D}_{k+1}$ and inserting the induction assumption we conclude that $(ii)$ holds. $\qquad$ Q.E.D.

From Proposition 9 we see that the relations (50), (52) and (53) can be written as

$$(\widetilde{A}^{(1)}, \widetilde{B}^{(1)}), \; (\overline{A}^{(2)}, \overline{B}^{(2)}), \; (\widetilde{A}^{(2)}, \widetilde{B}^{(2)}), \; (\overline{A}^{(3)}, \overline{B}^{(3)}), \; (\widetilde{A}^{(3)}, \widetilde{B}^{(3)}), \; \ldots \quad (60)$$

$$\overline{A}^{(k+1)} = \widetilde{F}_k^T \widetilde{A}^{(k)} \widetilde{F}_k \;, \qquad \overline{B}^{(k+1)} = \widetilde{F}_k^T \widetilde{B}^{(k)} \widetilde{F}_k \;, \qquad (61)$$

$$\widetilde{A}^{(k+1)} = \overline{D}_{k+1} \overline{A}^{(k+1)} \overline{D}_{k+1} \;, \qquad \widetilde{B}^{(k+1)} = \overline{D}_{k+1} \overline{B}^{(k+1)} \overline{D}_{k+1} \;. \quad (62)$$

The relations (60), (61), (62), (54) and (55) define the normalized Falk–Langemeyer method. We use the normalized method only as an aid to obtain information about the quantity $\widetilde{\varepsilon}_k$. æ

24

## 4.1 Preliminaries

Here we define asymptotic assumptions and prove several lemmas which are used later in the proof of the quadratic convergence of the Falk–Langemeyer method. The quadratic convergence proof is based on the idea of Wilkinson (see [18]) which consists in estimating the growth of already annihilated elements in the current cycle. To this end we must estimate the transformation parametars $\widetilde{\alpha}_k$ and $\widetilde{\beta}_k$ and also the growth of all off-diagonal elements in the current cycle. These two tasks are solved in Lemma 11, Lemma 13, Lemma 14 and Lemma 15. Lemma 10 gives us two numeric relations which are used in the proof. Lemma 11 and Lemma 12 estimate the transformation parametars $\widetilde{\alpha}_k$ and $\widetilde{\beta}_k$, and the measure $\widetilde{\varepsilon}_k$ in one step. Lemma 13, Lemma 14 and Lemma 15 estimate the growth of $\widetilde{\alpha}_k$, $\widetilde{\beta}_k$ and $\widetilde{\varepsilon}_k$ during $N$ consecutive steps. Lemma 15 is the most important for the proof of the quadratic convergence. In this subsection we do not assume that the pivot strategy is cyclic. Therefore the results of this subsection hold for any pivot strategy. However, if the pivot strategy is cyclic, then Lemma 13, Lemma 14 and Lemma 15 explain the behaviour of $\widetilde{\alpha}_k$, $\widetilde{\beta}_k$ and $\widetilde{\varepsilon}_k$ during one cycle.

As we said in Section 1, the quadratic convergence can always be expected if the eigenvalues of problem (1) are simple. We will therefore use two quadratic convergence assumptions:

(A1) The eigenvalues of the pair $(A, B)$ are simple, i. e.

$$p = n \geq 3 \,.$$

(A2) The pair $(A, B)$ is almost diagonal, i. e.

$$\frac{\widetilde{\varepsilon}_1}{\delta} < \frac{1}{2N} \,.$$

Asymptotic assumption (A2) is similar to the assumptions used in Theorem 6 and in convergence results of some other Jacobi–type methods (see [4], [1]). Assumption (A1) implies

$$N \geq 3 \tag{63}$$

and

$$\varepsilon_k = \tau_k \,, \qquad \widetilde{\varepsilon}_k = \widetilde{\tau}_k \,, \qquad k \geq 1 \,, \tag{64}$$

25

where $\tau_k = \tau(A^{(k)}, B^{(k)})$ and $\widetilde{\tau}_k = \tau(\widetilde{A}^{(k)}, \widetilde{B}^{(k)})$. We shall use the notation

$$\widetilde{a}_k \;=\; |\widetilde{a}_{lm}^{(k)}| \,, \qquad\qquad \widetilde{b}_k \;=\; |\widetilde{b}_{lm}^{(k)}| \,, \qquad\qquad k \geq 1 \,. \qquad (65)$$

LEMMA **10** *Let $r$ be an integer such that $r \geq 3$ and let $x$ be a nonnegative real number satisfying $2xr < 1$. Then the following inequalities hold:*

$$(1 - x)^{-r} \leq 1 + \frac{12}{7} \cdot r \cdot x \,, \qquad\qquad (1 + x)^r \leq 1 + \frac{4}{3} \cdot r \cdot x \,.$$

PROOF: The proof of this lemma is elementary and can be found in [4]. Q.E.D.

The following lemma shows how are the transformation parameters $\widetilde{\alpha}_k$ and $\widetilde{\beta}_k$ from matrices $\widetilde{F}_k$ bounded with $\widetilde{\epsilon}_k$.

LEMMA **11** *Let the assumption (A1) hold. If for some $k \geq 1$ holds*

$$\widetilde{\varepsilon}_k \;<\; \frac{2}{3N}\delta \,, \qquad\qquad (66)$$

*then*

$$\max\{|\widetilde{\alpha}_k|, |\widetilde{\beta}_k|\} \;\leq\; 0.34 \cdot \frac{\sqrt{(\widetilde{a}_k)^2 + (\widetilde{b}_k)^2}}{\delta} \,. \qquad (67)$$

PROOF: Suppose that for some $k \geq 1$ the relation (66) holds. Then Theorem 1 and Corollary 2 hold for the pair $(\widetilde{A}^{(k)}, \widetilde{B}^{(k)})$ as well. The assumption (A1) and the relations (63), (64) and (18) imply that there exists an ordering of the eigenvalues of the pair $(A, B)^2$ such that

$$\chi(\lambda_i, [\widetilde{a}_{ii}^{(k)}, \widetilde{b}_{ii}^{(k)}]) \leq \frac{\widetilde{\varepsilon}_k^2}{2\delta} \;<\; \frac{4\delta^2}{9N^2} \cdot \frac{1}{2\delta} \;<\; \frac{2}{81} \cdot \delta \;<\; 0.025 \cdot \delta \,,$$
$$i \;=\; 1, \ldots, n \,. \qquad (68)$$

Applying twice the triangle inequality and using the definition (11) and the relation (68), we obtain

$$\begin{aligned}
| \, \widetilde{\mathfrak{S}}_{lm}^{(k)} \, | &\;=\; | \, \widetilde{a}_{ll}^{(k)}\widetilde{b}_{mm}^{k)} - \widetilde{a}_{mm}^{(k)}\widetilde{b}_{ll}^{(k)} \, | \;=\; \chi([\widetilde{a}_{ll}^{(k)}, \widetilde{b}_{ll}^{(k)}], [\widetilde{a}_{mm}^{(k)}, \widetilde{b}_{mm}^{(k)}]) \\
&\;\geq\; \chi(\lambda_l, \lambda_m) - \chi(\lambda_l, [\widetilde{a}_{ll}^{(k)}, \widetilde{b}_{ll}^{(k)}]) \,-\, \chi(\lambda_m, [\widetilde{a}_{mm}^{(k)}, \widetilde{b}_{mm}^{(k)}]) \\
&\;>\; 3 \cdot \delta - 2 \cdot 0.025 \cdot \delta = 2.95 \cdot \delta \,. \qquad (69)
\end{aligned}$$

_____

[2]Since $p = n$, the eigenvalues can be ordered so that the matrix $P$ from Theorem 1 and Corollary 2 is identity matrix.

It is obvious that $|\widetilde{\mathfrak{S}}_{lm}^{(k)}| \neq 0$. This excludes cases (5a) and (5bii) of Algorithm 3. Therefore, we have

$$\max\{|\widetilde{\alpha}_k|,\ |\widetilde{\beta}_k|\} \ \leq\ \frac{2}{|\widetilde{\mathfrak{S}}_{lm}^{(k)}| + \sqrt{\widetilde{\mathfrak{S}}^{(k)}}} \cdot \max\{|\widetilde{\mathfrak{S}}_l^{(k)}|\,,\ |\widetilde{\mathfrak{S}}_m^{(k)}|\}\ . \qquad (70)$$

From the Cauchy–Schwarz inequality and the relations (47) and (65) we have

$$\begin{aligned}
|\widetilde{\mathfrak{S}}_l^{(k)}| &=\ |\widetilde{a}_{ll}^{(k)}\widetilde{b}_{lm}^{(k)} - \widetilde{b}_{ll}^{(k)}\widetilde{a}_{lm}^{(k)}| \ \leq\ \sqrt{(\widetilde{a}_{ll}^{(k)})^2 + (\widetilde{b}_{ll}^{(k)})^2}\sqrt{(\widetilde{a}_{lm}^{(k)})^2 + (\widetilde{b}_{lm}^{(k)})^2}\\
&=\ \sqrt{(\widetilde{a}_k)^2 + (\widetilde{b}_k)^2}\ .
\end{aligned}$$

The same estimate holds for $\widetilde{\mathfrak{S}}_m^{(k)}$ and therefore

$$\max\{|\widetilde{\mathfrak{S}}_l^{(k)}|\,,\ |\widetilde{\mathfrak{S}}_m^{(k)}|\} \ \leq\ \sqrt{(\widetilde{a}_k)^2 + (\widetilde{b}_k)^2}\ . \qquad (71)$$

Since

$$(\widetilde{a}_k)^2 + (\widetilde{b}_k)^2 \ \leq\ \frac{1}{2} \cdot \widetilde{\varepsilon}_k^2\ , \qquad (72)$$

the relations (69), (71),and (72) imply

$$\begin{aligned}
\sqrt{\widetilde{\mathfrak{S}}^{(k)}} &=\ \sqrt{(\widetilde{\mathfrak{S}}_{lm}^{(k)})^2 + 4\,\widetilde{\mathfrak{S}}_l^{(k)}\widetilde{\mathfrak{S}}_m^{(k)}} \ \geq\ \sqrt{(2.95\,\delta)^2 - 4 \cdot \frac{\widetilde{\varepsilon}_k^2}{2}}\\
&\geq\ \sqrt{(2.95\,\delta)^2 - 2 \cdot \frac{4}{9\,N^2}\delta^2} \ \geq\ 2.933\,\delta\ .
\end{aligned}$$

The assertion (67) now follows from the relations (70), (69), (71) and the above relation. Q.E.D.

The following Lemma gives the relation between $\widetilde{\varepsilon}_k$ and $\widetilde{\varepsilon}_{k+1}$. It is used later in the proof of of Lemma 15.

LEMMA **12** *Let the assumption (A1) hold. If for some* $k \geq 1$ *the relation (66) holds, then*

$$\widetilde{\varepsilon}_{k+1}^2 \ \leq\ \frac{1 + 0.494 \cdot \frac{\widetilde{\varepsilon}_k}{\delta}}{1 - 0.077 \cdot \frac{\widetilde{\varepsilon}_k}{\delta}}\,[\widetilde{\varepsilon}_k^2 - 2\,(\widetilde{a}_k^2 + \widetilde{b}_k^2)]\,, \qquad (73)$$

PROOF: Suppose that the relation (66) holds for some $k \geq 1$. The relation (62), together with the definition of $\tilde{\varepsilon}_k$ , implies

$$\tilde{\varepsilon}_{k+1}^2 = S^2(\overline{D}_{k+1}\overline{A}^{(k+1)}\overline{D}_{k+1}) + S^2(\overline{D}_{k+1}\overline{B}^{(k+1)}\overline{D}_{k+1}) . \tag{74}$$

If

$$m_{k+1} = \min\{\overline{d}_1^{(k+1)}, \ldots, \overline{d}_n^{(k+1)}\} ,$$

then the relation (74) implies

$$\tilde{\varepsilon}_{k+1}^2 \leq S^2(\frac{1}{(m_{k+1})^2}\overline{A}^{(k+1)}) + S^2(\frac{1}{(m_{k+1})^2}\overline{B}^{(k+1)}) = \frac{1}{(m_{k+1})^4} \cdot \overline{\varepsilon}_{k+1}^2 . \tag{75}$$

Let us define vectors

$$
\begin{aligned}
\tilde{a}^l &= (\tilde{a}_{l,1}^{(k)}, \tilde{a}_{l,2}^{(k)}, \ldots, \tilde{a}_{l,l-1}^{(k)}, \tilde{a}_{l,l+1}^{(k)}, \ldots, \tilde{a}_{l,m-1}^{(k)}, \tilde{a}_{l,m+1}^{(k)}, \ldots, \tilde{a}_{l,n}^{(k)}) , \\
\tilde{a}^m &= (\tilde{a}_{m,1}^{(k)}, \tilde{a}_{m,2}^{(k)}, \ldots, \tilde{a}_{m,l-1}^{(k)}, \tilde{a}_{m,l+1}^{(k)}, \ldots, \tilde{a}_{m,m-1}^{(k)}, \tilde{a}_{m,m+1}^{(k)}, \ldots, \tilde{a}_{m,n}^{(k)}) , \\
\tilde{a}_l^T &= (\tilde{a}_{1,l}^{(k)}, \tilde{a}_{2,l}^{(k)}, \ldots, \tilde{a}_{l-1,l}^{(k)}, \tilde{a}_{l+1,l}^{(k)}, \ldots, \tilde{a}_{m-1,l}^{(k)}, \tilde{a}_{m+1,l}^{(k)}, \ldots, \tilde{a}_{n,l}^{(k)}) , \\
\tilde{a}_m^T &= (\tilde{a}_{1,m}^{(k)}, \tilde{a}_{2,m}^{(k)}, \ldots, \tilde{a}_{l-1,m}^{(k)}, \tilde{a}_{l+1,m}^{(k)}, \ldots, \tilde{a}_{m-1,m}^{(k)}, \tilde{a}_{m+1,m}^{(k)}, \ldots, \tilde{a}_{n,m}^{(k)}) ,
\end{aligned}
$$

where generally $a^T$ denotes the transposed vector $a$. Let $\overline{a}^l, \overline{a}^m, \overline{a}_l$ and $\overline{a}_m$ be row and column vectors defined in the same way, but from elements of the matrix $\overline{A}^{(k+1)}$. Relation (61) implies that

$$\begin{bmatrix} \overline{a}^l \\ \overline{a}^m \end{bmatrix} = \hat{\tilde{F}}_k^T \begin{bmatrix} \tilde{a}^l \\ \tilde{a}^m \end{bmatrix} , \qquad [\overline{a}_l, \overline{a}_m] = [\tilde{a}_l, \tilde{a}_m]\hat{\tilde{F}}_k .$$

Therefore,

$$\left\| \begin{bmatrix} \overline{a}^l \\ \overline{a}^m \end{bmatrix} \right\|^2 \leq \| \hat{\tilde{F}}_k^T \|_2^2 \left\| \begin{bmatrix} \tilde{a}^l \\ \tilde{a}^m \end{bmatrix} \right\|^2 , \quad \| [\overline{a}_l, \overline{a}_m] \|^2 \leq \| \hat{\tilde{F}}_k \|_2^2 \| [\tilde{a}_l, \tilde{a}_m] \|^2 .$$

The off–diagonal elements of the matrix $\tilde{A}^{(k)}$ which are changed in the transformation (61), are exactly the elements of vectors $\tilde{a}^l, \tilde{a}^m, \tilde{a}_l$ and $\tilde{a}_m$ with the exception of $\tilde{a}_{lm}^{(k)}$ and $\tilde{a}_{ml}^{(k)}$ which are annihilated. Since $\| \hat{\tilde{F}}_k^T \|_2 = \| \hat{\tilde{F}}_k \|_2$, we conclude that

$$S^2(\overline{A}^{(k+1)}) \leq S^2(\tilde{A}^{(k)}) + (\| \hat{\tilde{F}}_k \|_2^2 - 1)(\| \tilde{a}^l \|^2 + \| \tilde{a}^m \|^2 + \| \tilde{a}_l \|^2 + \| \tilde{a}_m \|^2) - 2\tilde{a}_k^2 .$$

28

Since $\| \widehat{\widetilde{F}}_k \|_2 \geq 1$ (see further in the proof), we conclude that

$$S^2(\overline{A}^{(k+1)}) \leq \| \widehat{\widetilde{F}}_k \|_2^2 \, (S^2(\widetilde{A}^{(k)}) - 2\widetilde{a}_k^2)\,.$$

By applying the similar analysis to matrix $\widetilde{B}^{(k)}$, we obtain

$$S^2(\overline{B}^{(k+1)}) \leq \| \widehat{\widetilde{F}}_k \|_2^2 \, (S^2(\widetilde{B}^{(k)}) - 2\widetilde{b}_k^2)\,.$$

Adding two previous inequalities and using the definitions of $\overline{\varepsilon}_{k+1}$ and $\widetilde{\varepsilon}_k$, gives

$$\overline{\varepsilon}_{k+1}^2 \leq \| \widehat{\widetilde{F}}_k \|_2^2 \, [\widetilde{\varepsilon}_k^2 - 2(\widetilde{a}_k^2 + \widetilde{b}_k^2)]\,.$$

Inserting this inequality into relation (75), we obtain

$$\widetilde{\varepsilon}_{k+1}^2 \leq \frac{\| \widehat{\widetilde{F}}_k \|_2^2}{(m_{k+1})^4} \cdot [\widetilde{\varepsilon}_k^2 - 2(\widetilde{a}_k^2 + \widetilde{b}_k^2)]\,. \tag{76}$$

To complete the proof we must find the upper bound for $\| \widehat{\widetilde{F}}_k \|_2^2$ and the lower bound for $m_{k+1}$.

The relation (56) implies that

$$m_{k+1} = \min\{1, \overline{d}_l^{(k+1)}, \overline{d}_m^{(k+1)}\}\,. \tag{77}$$

Relation (61) implies that

$$\overline{a}_{ll}^{(k+1)} = \widetilde{a}_{ll}^{(k)} - 2\widetilde{\beta}_k \widetilde{a}_{lm}^{(k)} + \widetilde{\beta}_k^2 \, \widetilde{a}_{mm}^{(k)}\,, \qquad \overline{b}_{ll}^{(k+1)} = \widetilde{b}_{ll}^{(k)} - 2\widetilde{\beta}_k \widetilde{b}_{lm}^{(k)} + \widetilde{\beta}_k^2 \, \widetilde{b}_{mm}^{(k)}\,,$$

$$\overline{a}_{mm}^{(k+1)} = \widetilde{\alpha}_k^2 \, \widetilde{a}_{ll}^{(k)} + 2\widetilde{\alpha}_k \widetilde{a}_{lm}^{(k)} + \widetilde{a}_{mm}^{(k)}\,, \qquad \overline{b}_{mm}^{(k+1)} = \widetilde{\alpha}_k^2 \, \widetilde{b}_{ll}^{(k)} + 2\widetilde{\alpha}_k \widetilde{b}_{lm}^{(k)} + \widetilde{b}_{mm}^{(k)}\,.$$

Therefore

$$(\overline{d}_l^{(k+1)})^4 = (\overline{a}_{ll}^{(k+1)})^2 + (\overline{b}_{ll}^{(k+1)})^2 = 1 + 4\beta_k^2 \, [(\widetilde{a}_{lm}^{(k)})^2 + (\widetilde{b}_{lm}^{(k)})^2] \, +$$

$$+ \, \widetilde{\beta}_k^4 \, - \, 4\widetilde{\beta}_k(\widetilde{a}_{ll}^{(k)}\widetilde{a}_{lm}^{(k)} + \widetilde{b}_{ll}^{(k)}\widetilde{b}_{lm}^{(k)}) - 4\widetilde{\beta}_k^3(\widetilde{a}_{mm}^{(k)}\widetilde{a}_{lm}^{(k)} + \widetilde{b}_{mm}^{(k)}\widetilde{b}_{lm}^{(k)}) \, +$$

$$+ \, 2\widetilde{\beta}_k^2(\widetilde{a}_{ll}^{(k)}\widetilde{a}_{mm}^{(k)} + \widetilde{b}_{ll}^{(k)}\widetilde{b}_{mm}^{(k)})$$

$$\geq \, 1 - 4\,|\,\widetilde{\beta}_k\,|\,|\,\widetilde{a}_{ll}^{(k)}\widetilde{a}_{lm}^{(k)} + \widetilde{b}_{ll}^{(k)}\widetilde{b}_{lm}^{(k)}\,| - 4\,|\,\widetilde{\beta}_k\,|^3\,|\,\widetilde{a}_{mm}^{(k)}\widetilde{a}_{lm}^{(k)} + \widetilde{b}_{mm}^{(k)}\widetilde{b}_{lm}^{(k)}\,| \, -$$

$$- \, 2\,\widetilde{\beta}_k^2\,|\,\widetilde{a}_{ll}^{(k)}\widetilde{a}_{mm}^{(k)} + \widetilde{b}_{ll}^{(k)}\widetilde{b}_{mm}^{(k)}\,|\,. \tag{78}$$

29

Using the relation (47) and the Cauchy–Schwarz inequality in the relation (78), we obtain

$$(\overline{d}_l^{(k+1)})^4 \geq 1 - 4 \mid \tilde{\beta}_k \mid (1 + \mid \tilde{\beta}_k \mid^2)\sqrt{(\tilde{a}_{lm}^{(k)})^2 + (\tilde{b}_{lm}^{(k)})^2} - 2\tilde{\beta}_k^2 . \qquad (79)$$

Using similar argument we obtain

$$(\overline{d}_m^{(k+1)})^4 \geq 1 - 4 \mid \tilde{\alpha}_k \mid (1 + \mid \tilde{\alpha}_k \mid^2)\sqrt{(\tilde{a}_{lm}^{(k)})^2 + (\tilde{b}_{lm}^{(k)})^2} - 2\tilde{\alpha}_k^2 . \qquad (80)$$

Relations (77), (79), (80), (72) and Lemma 11 now imply

$$m_{k+1}^4 \geq 1 - 4 \cdot \frac{0.34}{\sqrt{2}} \cdot \frac{\tilde{\varepsilon}_k}{\delta} \cdot \left[ 1 + \left( \frac{0.34}{\sqrt{2}} \cdot \frac{\tilde{\varepsilon}_k}{\delta} \right)^2 \right] \cdot \frac{\tilde{\varepsilon}_k}{\sqrt{2}} - 2 \cdot \left( \frac{0.34}{\sqrt{2}} \cdot \frac{\tilde{\varepsilon}_k}{\delta} \right)^2 . \qquad (81)$$

Since $\chi$ is chordal metric, from the relation (11) we see that $\delta \leq 1/3$. The relations (66) and (63) therefore imply

$$\tilde{\varepsilon}_k < \frac{2}{9N} < \frac{2}{27} . \qquad (82)$$

Inserting the relation (82) and the assumption (66) into the relation (81) we obtain

$$m_{k+1}^4 > 1 - \frac{0.34}{\sqrt{2}} \cdot \frac{\tilde{\varepsilon}_k}{\delta} \left[ 4 \cdot \left( 1 + \left( \frac{0.34}{\sqrt{2}} \cdot \frac{2}{3N} \right)^2 \right) \cdot \frac{2}{27\sqrt{2}} + 0.34 \cdot \sqrt{2} \cdot \frac{2}{3N} \right] .$$

Finally, taking into account that $N \geq 3$ we obtain

$$m_{k+1}^4 \geq 1 - 0.077 \cdot \frac{\tilde{\varepsilon}_k}{\delta} \qquad (83)$$

We shall now estimate $\| \widehat{\tilde{F}}_k \|_2^2$. Since

$$\| \widehat{\tilde{F}}_k \|_2^2 \leq \| \widehat{\tilde{F}}_k \|_1 \cdot \| \widehat{\tilde{F}}_k \|_\infty ,$$

where $\|A\|_1 = \max_j \sum_i |a_{ij}|$, $\|A\|_\infty = \max_i \sum_j |a_{ij}|$ for $A = (a_{ij})$, we obtain

$$\| \widehat{\tilde{F}}_k \|_2^2 \leq (1 + \max\{|\tilde{\alpha}_k|, |\tilde{\beta}_k|\})^2 .$$

30

Lemma 11 and the relation (66) now imply

$$\| \widehat{\widetilde{F}}_k \|_2^2 \;\leq\; \left(1 + \frac{0.34}{\sqrt{2}}\,\frac{\widetilde{\varepsilon}_k}{\delta}\right)^2 \;\leq\; 1 + \frac{0.34}{\sqrt{2}}\,\frac{\widetilde{\varepsilon}_k}{\delta}\left(2 + \frac{0.34}{\sqrt{2}}\,\frac{2}{3N}\right)$$

$$\leq\; 1 + 0.494\,\frac{\widetilde{\varepsilon}_k}{\delta}\; . \tag{84}$$

The relation (73) now follows from the relations (76), (83) and (84). Q.E.D.æ

We shall now prove that if the assumptions (A1) and (A2) hold, then Lemma 11 and Lemma 12 hold during $N$ consecutive steps.

LEMMA **13** *Let the asymptotic assumptions (A1) and (A2) hold. Then for each* $k \in \{1, \ldots, N\}$ *holds*

$$\widetilde{\varepsilon}_k \;\leq\; \frac{1}{1 - 0.3 \cdot (k-1)\frac{\widetilde{\varepsilon}_1}{\delta}} \cdot \widetilde{\varepsilon}_1 \;, \qquad\qquad \frac{\widetilde{\varepsilon}_k}{\delta} \;<\; \frac{2}{3N} \; .$$

PROOF: The proof is by induction. For $k = 1$ lemma is trivially fulfilled. Suppose that lemma holds for some $k \in \{1, \ldots, N-1\}$. From the second inequality in the induction assumption we conclude that, for the chosen $k$, Lemma 11 and Lemma 12 hold. From Lemma 12 it follows that

$$\widetilde{\varepsilon}_{k+1}^2 \;\leq\; \frac{1 + 0.494 \cdot \frac{\widetilde{\varepsilon}_k}{\delta}}{1 - 0.077 \cdot \frac{\widetilde{\varepsilon}_k}{\delta}} \cdot \widetilde{\varepsilon}_k^2 \;\leq\; \frac{1}{(1 - 0.494\frac{\widetilde{\varepsilon}_k}{\delta})(1 - 0.077\frac{\widetilde{\varepsilon}_k}{\delta})}\widetilde{\varepsilon}_k^2$$

$$\leq\; \frac{1}{\left(1 - 0.3\frac{\widetilde{\varepsilon}_k}{\delta}\right)^2}\widetilde{\varepsilon}_k^2 \; . \tag{85}$$

Hence

$$\widetilde{\varepsilon}_{k+1} \;\leq\; \frac{1}{1 - 0.3 \cdot \frac{\widetilde{\varepsilon}_k}{\delta}} \cdot \widetilde{\varepsilon}_k \; .$$

Inserting the induction assumption in this inequality we obtain

$$\widetilde{\varepsilon}_{k+1} \;\leq\; \frac{1}{1 - 0.3\frac{1}{1 - 0.3(k-1)\frac{\widetilde{\varepsilon}_1}{\delta}}\frac{\widetilde{\varepsilon}_1}{\delta}} \cdot \frac{1}{1 - 0.3(k-1)\frac{\widetilde{\varepsilon}_1}{\delta}} \cdot \widetilde{\varepsilon}_1$$

$$\leq\; \frac{1}{1 - 0.3(k-1)\frac{\widetilde{\varepsilon}_1}{\delta} - 0.3\frac{\widetilde{\varepsilon}_1}{\delta}} \cdot \widetilde{\varepsilon}_1 \;=\; \frac{1}{1 - 0.3 \cdot k \cdot \frac{\widetilde{\varepsilon}_1}{\delta}} \cdot \widetilde{\varepsilon}_1 \;,$$

and the first assertion of the lemma is proved. From this assertion for $k + 1$, because of the asymptotic assumption (A2), we now have

$$\frac{\widetilde{\varepsilon}_{k+1}}{\delta} \leq \frac{1}{1 - 0.3 \cdot k \cdot \frac{\widetilde{\varepsilon}_1}{\delta}} \cdot \frac{\widetilde{\varepsilon}_1}{\delta} < \frac{1}{1 - 0.3 \, (N-1)\frac{1}{2N}} \cdot \frac{1}{2N} = \frac{1}{2 - 0.3} \cdot \frac{1}{N} < \frac{2}{3N}$$

which completes the proof. Q.E.D.

LEMMA **14** *If the asymptotic assumptions (A1) and (A2) hold, then the assertions (67) and (73) of Lemma 11 and Lemma 12 hold for every $k \in \{1, \ldots, N\}$.*

PROOF: The assertion follows imidiately from second assertion of Lemma 13.

Q.E.D.

The next lemma explains behaviour of $S(\widetilde{A}^{(k)})$, $S(\widetilde{B}^{(k)})$ and $\widetilde{\varepsilon}_k$, and of the transformation parameters $\widetilde{\alpha}_k$ and $\widetilde{\beta}_k$ during $N$ consecutive steps. Let us define the quantity

$$c_N = \frac{1 + 0.494 \cdot \frac{2}{3N}}{1 - 0.077 \cdot \frac{2}{3N}} \, . \tag{86}$$

LEMMA **15** *Let the asymptotic assumptions (A1) and (A2) hold. Then:*

*(i) For $k = 1, \ldots, N$ holds*

$$\begin{bmatrix} S^2(\widetilde{A}^{(k+1)}) \\ S^2(\widetilde{B}^{(k+1)}) \\ \widetilde{\varepsilon}_{k+1}^2 \end{bmatrix} \leq (c_N)^k \begin{bmatrix} S^2(\widetilde{A}^{(1)}) \\ S^2(\widetilde{B}^{(1)}) \\ \widetilde{\varepsilon}_1^2 \end{bmatrix} \leq 1.566 \begin{bmatrix} S^2(\widetilde{A}^{(1)}) \\ S^2(\widetilde{B}^{(1)}) \\ \widetilde{\varepsilon}_1^2 \end{bmatrix} \, .$$

*(ii) For any choice $\widetilde{\omega}_k \in \{\widetilde{\alpha}_k, \widetilde{\beta}_k\}$, $1 \leq k \leq N$, holds*

$$\sum_{k=1}^{N} \widetilde{\omega}_k^2 \leq 0.426 \cdot \frac{\widetilde{\varepsilon}_1^2}{\delta^2} \, .$$

PROOF: $(i)$ Because of Lemma 12 and Lemma 14, for $k = 1, \ldots, N$ holds

$$
\begin{aligned}
\widetilde{\varepsilon}_{k+1}^2 &\leq c_N(\widetilde{\varepsilon}_k^2 - 2\,(\widetilde{a}_k^2 + \widetilde{b}_k^2)) \\
&\leq c_N\{c_N[\widetilde{\varepsilon}_{k-1}^2 - 2\,(\widetilde{a}_{k-1}^2 + \widetilde{b}_{k-1}^2)] - 2\,(\widetilde{a}_k^2 + \widetilde{b}_k^2)\} \\
&\leq \ldots \leq (c_N)^k \widetilde{\varepsilon}_1^2 - 2\sum_{j=1}^{k}(c_N)^{k-j+1}(\widetilde{a}_j^2 + \widetilde{b}_j^2)\,.
\end{aligned}
\tag{87}
$$

From the relation (87) immediately follows

$$
\widetilde{\varepsilon}_{k+1}^2 \ \leq \ (c_N)^k \widetilde{\varepsilon}_1^2 \ \leq \ (c_N)^N \widetilde{\varepsilon}_1^2\,, \qquad k = 1, \ldots, N\,.
\tag{88}
$$

Using Lemma 10 we obtain

$$
(c_N)^N \ < \ (1 + \frac{4}{3} \cdot 0.494 \cdot \frac{2}{3N} \cdot N)(1 + \frac{12}{7} \cdot 0.077 \cdot \frac{2}{3N} \cdot N) \ < \ 1.566\,.
\tag{89}
$$

Inserting this inequality into relation (88), we obtain

$$
\widetilde{\varepsilon}_{k+1}^2 \ \leq \ 1.566 \cdot \widetilde{\varepsilon}_1^2\,, \qquad k = 1, \ldots, N\,.
$$

From the proof of Lemma 12 we se that the above estimates hold for the quantities $S^2(\widetilde{A}^{(k+1)})$ and $S^2(\widetilde{B}^{(k+1)})$, as well. Therefore $(i)$ is proved.

$(ii)$ Since $c_N > 1$, from the relation (87) for $k = N$ we have

$$
\widetilde{\varepsilon}_{N+1}^2 \ \leq \ (c_N)^N \widetilde{\varepsilon}_1^2 - 2\sum_{k=1}^{N}(\widetilde{a}_k^2 + \widetilde{b}_k^2)\,.
$$

Since $\widetilde{\varepsilon}_{N+1}^2 \geq 0$, this inequality implies

$$
\sum_{k=1}^{N}(\widetilde{a}_k^2 + \widetilde{b}_k^2) \ \leq \ \frac{1}{2} \cdot (c_N)^N \widetilde{\varepsilon}_1^2 \ \leq \ 0.783 \cdot \widetilde{\varepsilon}_1^2\,.
$$

The above inequality together with Lemma 11 and Lemma 14 imply

$$
\begin{aligned}
\sum_{k=1}^{N} \widetilde{\omega}_k^2 &\leq \sum_{k=1}^{N} \max\{\widetilde{\alpha}_k^2, \widetilde{\beta}_k^2\} \leq \sum_{k=1}^{N} 0.34^2 \cdot (\widetilde{a}_k^2 + \widetilde{b}_k^2) \cdot \frac{1}{\delta^2} \\
&\leq 0.1156 \cdot 0.783 \cdot \frac{\widetilde{\varepsilon}_1^2}{\delta^2} \leq 0.091 \cdot \frac{\widetilde{\varepsilon}_1^2}{\delta^2}
\end{aligned}
$$

and the lemma is proved. $\hspace{3cm}$ Q.E.D.

## 4.2 The proof

Here we prove that the Falk–Langemeyer method is quadratically convergent if the assumptions (A1) and (A2) are fulfilled and the pivot strategy is cyclic. Then we prove that the quadratic convergence implies the convergence of the sequence of pairs (42) towards the pair of diagonal matrices. At the end we prove that the measures $\widetilde{\varepsilon}_k$ and $\varepsilon_k$ are equivalent in the sense that ultimately the quadratic reduction of $\widetilde{\varepsilon}_{kN+1}$ implies the quadratic reduction of $\varepsilon_{kN+1}$ and vice versa.

We can now prove our paper's central theorem.

THEOREM **16** *Let the asymptotic assumptions (A1) and (A2) hold and let the sequence $((A^{(k)}, B^{(k)}),\ k \geq 1)$ be generated with the Falk–Langemeyer method from the pair $(A, B)$. Then for any cyclic strategy holds*

$$\widetilde{\varepsilon}_{N+1} \ \leq \ \sqrt{N} \cdot \frac{\widetilde{\varepsilon}_1^2}{\delta}\,.$$

PROOF: Let us fix some $k \in \{1, \ldots, N\}$. Then the pivot pair $(l, m)$ is also fixed. We want to know what happens with the element on this position till the end of cycle. Therefore, we will observe the elements $\widetilde{a}_{lm}^{(r)}$, $r = k + 1, \ldots, N$. We know that $\widetilde{a}_{lm}^{(k+1)} = 0$ and that the elements $\widetilde{a}_{lm}^{(r)}$ actually change at most $2(n - 2)$ times. Let $r_1, \ldots, r_s$, $s \leq 2n - 4$, denote those values of $r$ for which $\widetilde{a}_{lm}^{(r)}$ changes in the $r$–th step. Let us introduce the notation:

$$
\begin{aligned}
\overline{z}_i &= \overline{a}_{lm}^{(r_i+1)}, & \widetilde{z}_i &= \widetilde{a}_{lm}^{(r_i+1)}\,, \\
\overline{d}_j^{(i)} &= \sqrt[4]{(\overline{a}_{jj}^{(r_i+1)})^2 + \overline{b}_{jj}^{(r_i+1)})^2}\,, & j &\in \{l, m\}\,, \\
\overline{m}_{lm}^{(i)} &= \min\{\overline{d}_l^{(i)}, \overline{d}_m^{(i)}\}\,, & d_N &= \sqrt{1 - 0.077 \cdot \frac{2}{3N}}\,. \qquad (90)
\end{aligned}
$$

Performing the $r_1$–th step according to Algorithm 4, gives

$$\overline{z}_1 = \left(0 \cdot 1 \pm \widetilde{a}^{(r_1)} \widetilde{\omega}_{r_1}\right),$$

34

where $\widetilde{\omega}_{r_1} \in \{\widetilde{\alpha}_{r_1}, \widetilde{\beta}_{r_1}\}$ and $\widetilde{a}^{(r_1)}$ is some off–diagonal element of the matrix $\widetilde{A}^{(r_1)}$. Since

$$\widetilde{z}_i = \frac{\overline{z}_i}{\overline{d}_l^{(i)}\,\overline{d}_m^{(i)}} , \qquad i = 1,\ldots,s , \tag{91}$$

from the relations (90), (83) and Lemma 12 follows that

$$\mid \widetilde{z}_1 \mid \ \leq \ \frac{1}{(\overline{m}_{lm}^{(i)})^2} \mid \widetilde{a}^{(r_1)} \mid\mid \widetilde{\omega}_{r_1} \mid \ \leq \ \frac{1}{d_N} \mid \widetilde{a}^{(r_1)} \mid\mid \widetilde{\omega}_{r_1} \mid \ . \tag{92}$$

Further, in the $r_2$–th step, we have

$$\overline{z}_2 = (1 \cdot \widetilde{z}_1 \pm \widetilde{a}^{(r_2)}\widetilde{\omega}_{r_2}) , \tag{93}$$

where $\widetilde{\omega}_{r_2} \in \{\widetilde{\alpha}_{r_2}, \widetilde{\beta}_{r_2}\}$, and $\widetilde{a}^{(r_2)}$ is some off–diagonal element of the matrix $\widetilde{A}^{(r_2)}$. The relations (93), (92) and (91) imply

$$\mid \widetilde{z}_2 \mid \ \leq \ \frac{1}{d_N} \cdot (\frac{1}{d_N} \mid \widetilde{a}^{(r_1)} \mid\mid \widetilde{\omega}_{r_1} \mid + \mid \widetilde{a}^{(r_2)} \mid\mid \widetilde{\omega}_{r_2} \mid) .$$

By induction we obtain

$$\mid \widetilde{z}_j \mid \ \leq \ \sum_{i=1}^{j} \frac{1}{(d_N)^{j-i+1}} \mid \widetilde{a}^{(r_i)} \mid\mid \widetilde{\omega}_{r_i} \mid , \qquad j = 1,\ldots,s . \tag{94}$$

For $k = 1,\ldots,N+1$ following notation is introduced:

$$\widetilde{A}^{(k)} = \widetilde{D}_A^{(k)} + \widetilde{E}^{(k)} , \qquad \widetilde{D}_A^{(k)} = \mathrm{diag}\,(\widetilde{a}_{ii}^{(k)}) . \tag{95}$$

Matrix $\widetilde{E}^{(N+1)}$ obviously consists of elements which have undergone the maximal number of changes. If $s(i,j)$ denotes the number of changes of the element on position $(i,j)$, then

$$s(i,j) \ \leq \ 2n - 4 , \qquad i,j \in \{1,\ldots,n\}, i \neq j .$$

The quantity $s(i,j)$ depends upon $(i,j)$ and the pivot strategy. Elements of the matrix $\widetilde{E}^{(N+1)}$ can therefore be denoted as $\widetilde{z}_{s(i,j)}$.

Having in mind relation (94), we can now write

$$\mid \widetilde{E}^{(N+1)} \mid \ \leq \ \frac{1}{(d_N)^{2n-4}}(\mid \widetilde{P}^{(2)} \mid\mid \widetilde{\omega}_2 \mid + \mid \widetilde{P}^{(3)} \mid\mid \widetilde{\omega}_3 \mid + \ldots + \mid \widetilde{P}^{(N)} \mid\mid \widetilde{\omega}_N \mid) . \tag{96}$$

Here the notation $\mid C \mid = (\mid c_{ij} \mid)$ for $C = (c_{ij})$ is used. Matrix $\widetilde{P}^{(k)}$ consists precisely of those elements of $l(k)$–th and $m(k)$–th row and column of the matrix $\widetilde{E}^{(k)}$ which already were pivot elements[3], i. e. of elements which contribute to the final estimate. All other elements of the matrix $\widetilde{P}^{(k)}$ are zeros.

Assertion $(i)$ of Lemma 15 gives us

$$\| \mid \widetilde{P}^{(k)} \mid \| = \| \widetilde{P}^{(k)} \| \leq S(\widetilde{A}^{(k)}) \leq \sqrt{1.566} \cdot S(\widetilde{A}^{(1)}), \qquad k = 2, \ldots, N \, . \tag{97}$$

From the relations (96) and (97), Lemma 15 and the Cauchy–Schwarz inequality we obtain

$$S(\widetilde{A}^{(N+1)}) = \| \widetilde{E}^{(N+1)} \| = \| \mid \widetilde{E}^{(N+1)} \mid \| \leq \frac{1}{(d_N)^{2n-4}} \sqrt{1.566} \cdot S(\widetilde{A}^{(1)}) \sum_{k=2}^{N} \mid \widetilde{\omega}_k \mid$$

$$\leq \frac{1.252}{(d_N)^{2n-4}} \cdot S(\widetilde{A}^{(1)}) \cdot [(N-1) \sum_{k=2}^{N} \widetilde{\omega}_k^2]^{\frac{1}{2}} \leq \frac{1.252}{(d_N)^{2n-4}} \cdot S(\widetilde{A}^{(1)}) \cdot [N \sum_{k=1}^{N} \widetilde{\omega}_k^2]^{\frac{1}{2}} \, . \tag{98}$$

Since $N \geq 3$, from Lemma 10 follows that

$$\frac{1}{(d_N)^{2n-4}} = \frac{1}{(1 - 0.077\frac{2}{3N})^{n-2}} < 1 + \frac{12}{7} \cdot 0.077 \cdot \frac{2}{3N}(n-2)$$

$$\leq 1 + \frac{12}{7} \cdot 0.077 \cdot \frac{4}{3} \cdot \frac{1}{n} \leq 1.059 \, .$$

Finally, inserting this inequality and assertion $(ii)$ of Lemma 15 into relation (98), we obtain

$$S(\widetilde{A}^{(N+1)}) \leq 0.4 \cdot S(\widetilde{A}^{(1)}) \sqrt{N} \cdot \frac{\widetilde{\varepsilon}_1}{\delta} \, .$$

Applying a similar analysis to matrices $\widetilde{B}^{(k)}$ yields

$$S(\widetilde{B}^{(N+1)}) \leq 0.4 \cdot S(\widetilde{B}^{(1)}) \sqrt{N} \cdot \frac{\widetilde{\varepsilon}_1}{\delta} \, .$$

From the last two inequalities and the definitions of $\widetilde{\varepsilon}_{N+1}$ and $\widetilde{\varepsilon}_1$ follows

$$\widetilde{\varepsilon}_{N+1} \leq 0.4 \cdot \sqrt{N} \cdot \frac{\widetilde{\varepsilon}_1^2}{\delta} \, ,$$

---

[3]Here $(l(k), m(k))$ denotes pivot pair in the $k$–th step so this $k$ should not be confused with the $k$ that was fixed at the beginning of the proof.

and the theorem is proved. Q.E.D.

Note that in the proof of Theorem 16 it is not necessary to assume that the affiliation is preserved, i.e. that the pairs $[a_{ii}^{(k)}, b_{ii}^{(k)}]$ approximate the eigenvalues $\lambda_i$ for $i = 1, \ldots, n$, $k = 1, \ldots, N$. However, for large enough $k$ this fact follows from Theorem 17.

æ From Theorem 16 and the assumptions (A1) and (A2) follows that

$$\widetilde{\varepsilon}_{N+1} \; < \; \sqrt{N} \cdot \frac{1}{2N} \cdot \widetilde{\varepsilon}_1 \; = \; \frac{1}{2\sqrt{N}} \cdot \widetilde{\varepsilon}_1 \; < \; 0.3 \cdot \widetilde{\varepsilon}_1 \,. \tag{99}$$

Applying inductively the relation (99) we obtain

$$\widetilde{\varepsilon}_{rN+1} \; \leq \; (0.3)^r \cdot \widetilde{\varepsilon}_1 \,, \qquad r \geq 1 \,. \tag{100}$$

Therefore,

$$\lim_{r \to \infty} \widetilde{\varepsilon}_{rN+1} \; = \; 0 \,. \tag{101}$$

From the relation (101) and the assertion (i) of Lemma 15 we conclude that

$$\lim_{k \to \infty} \widetilde{\varepsilon}_k \; = \; 0 \,. \tag{102}$$

*The relation (102) and Theorem 16 imply the quadratic convergence of the Falk–Langemeyer method according to Definition 7 if the eigenvalues are simple and the pivot strategy is cyclic.*

Next we prove that under assumptions of Theorem 16 the sequences of matrices $(A^{(k)}, k \geq 1)$ and $(B^{(k)}, k \geq 1)$, generated by the Falk–Langemeyer method, converge towards diagonal matrices.

THEOREM **17** *Let the assumptions of Theorem 16 hold. Then*

$$\lim_{k \to \infty} A^{(k)} \; = \; D_A \,, \qquad \lim_{k \to \infty} B^{(k)} \; = \; D_B \,,$$

*where $D_A$ and $D_B$ are diagonal matrices.*

**Proof.** The relation (44) implies that

$$A^{(k)} \; = \; (D_k)^{-1} \widetilde{A}^{(k)} (D_k)^{-1} \,, \qquad B^{(k)} \; = \; (D_k)^{-1} \widetilde{B}^{(k)} (D_k)^{-1} \,, \tag{103}$$

37

where diagonal matrices $D_k$ are defined with the relations (45) and (46). It is therefore sufficient to proove that the sequences $(\widetilde{A}^{(k)}, \ k \geq 1)$, $(\widetilde{B}^{(k)}, \ k \geq 1)$ and $((D_k)^{-1}, \ k \geq 1)$ converge towards diagonal matrices. The relation (102) implies that the off-diagonal elements of matrices $\widetilde{A}^{(k)}$ and $\widetilde{B}^{(k)}$ tend to zero as $k \to \infty$. Therefore, it remains to proove that for $i = 1, \ldots, n$ the sequences $(\widetilde{a}_{ii}^{(k)}, \ k \geq 1)$ and $(\widetilde{b}_{ii}^{(k)}, \ k \geq 1)$ converge. The relation (18) and the assumption (A1) imply that for each $k \geq 1$ there exists an ordering of the eigenvalues $\lambda_i = [s_i, c_i]$, $i = 1, \ldots, n$, such that

$$| c_i \widetilde{a}_{ii}^{(k)} - s_i \widetilde{b}_{ii}^{(k)} | \leq \frac{\widetilde{\varepsilon}_k^2}{2\delta} , \qquad i = 1, \ldots, n . \qquad (104)$$

Let us consider unit vectors $[s_i, c_i]^T$ and $[\widetilde{a}_{ii}^{(k)}, \widetilde{b}_{ii}^{(k)}]^T$ in $\mathbf{R}^2$. The left-hand side of the inequality (104) is $| \sin \varphi_i^{(k)} |$ where $\varphi_i^{(k)}$ is the angle between these two vectors. The relations (102) and (104) imply

$$\lim_{k \to \infty} \sin \varphi_i^{(k)} = 0 , \qquad i =, 1 \ldots, n .$$

Hence, for each $i$ the sequence of vectors $([\widetilde{a}_{ii}^{(k)}, \widetilde{b}_{ii}^{(k)}]^T, k \geq 1)$ has only finite number of accumulation points in $\mathbf{R}^2$. Therefore, it suffices to show that for large enough $k$ the changes in $\widetilde{a}_{ii}^{(k)}$ and $\widetilde{b}_{ii}^{(k)}$ are arbitrary small. From the relation (102) and Lemma 11 we see that $\widetilde{\alpha}_k \to 0$ and $\widetilde{\beta}_k \to 0$ as $k \to \infty$. Therefore, the changes in $\widetilde{a}_{ii}^{(k)}$ and $\widetilde{b}_{ii}^{(k)}$ tend to zero as $k \to \infty$. This prooves that for each $i \in \{1, \ldots, n\}$ limits $\lim_{k \to \infty} \widetilde{a}_{ii}^{(k)}$ and $\lim_{k \to \infty} \widetilde{b}_{ii}^{(k)}$ exist.

We shall now prove that $((D_k)^{-1}, k \geq 1)$ is a convergent sequence. Looking at the definition of $D_k$ (relation (45)) we see that it suffices to prove that for each $i \in \{1, \ldots, n\}$ the sequence $(d_i^{(k)}, k \geq 1)$, converges to a nonzero number. From Proposition 9 we have

$$d_i^{(k)} = d_i^{(1)} \overline{d}_i^{(2)} \cdots \overline{d}_i^{(k)} , \qquad i = 1, \ldots, n , \qquad k \geq 2 .$$

From the definiteness of pairs $(A^{(1)}, B^{(1)})$ and $(\overline{A}^{(k)}, \overline{B}^{(k)})$ we conclude that $d_i^{(1)}$ and $\overline{d}_i^{(k)}$ are different from zero for all $i$ and $k$. Hence it suffices to prove that the infinite product $\prod_{k=2}^{\infty} \overline{d}_i^{(k)}$ converges[4]. This product converges if and only if the product $\prod_{k=2}^{\infty} \left( \overline{d}_i^{(k)} \right)^4$ converges. Therefore, it suffices to show that

---

[4]Since all factors in the product are nonzero the limit, if exists, is also nonzero.

the later product is absolutely convergent. From the relation (78) we see that for $i \in \{1, \ldots, n\}$ and $k \geq 2$ we can write $\left(\overline{d}_i^{(k)}\right)^4 = 1 + u_i^{(k)}$, so it suffices to show that the series $\sum_{k=2}^{\infty} u_i^{(k)}$ are absolutely convergent for all $i \in \{1, \ldots, n\}$.

The relation (83) of of Lemma 12 implies that

$$(\overline{d}_i^{(k+1)})^4 \geq 1 - 0.077 \cdot \frac{\widetilde{\varepsilon}_k}{\delta} , \qquad 1 = 1, \ldots, n , \qquad k \geq 1 .$$

Looking for upper bound instead of lower bound in the relation (78) and making similar estimates as in the relation (83), we obtain

$$(\overline{d}_i^{(k+1)})^4 \leq 1 + 0.077 \cdot \frac{\widetilde{\varepsilon}_k}{\delta} , \qquad 1 = 1, \ldots, n , \qquad k \geq 1 .$$

Therefore,

$$\mid u_i^{(k+1)} \mid = \mid (\overline{d}_i^{(k+1)})^4 - 1 \mid \leq 0.077 \cdot \frac{\widetilde{\varepsilon}_k}{\delta} , \qquad 1 = 1, \ldots, n , \qquad k \geq 1 .$$

Hence it suffices to show that the series $\sum_{k=1}^{\infty} \widetilde{\varepsilon}_k$ converges. From the assertion (i) of Lemma 15 we have

$$\widetilde{\varepsilon}_{rN+i} \leq 1.3 \cdot \widetilde{\varepsilon}_{rN+1} , \qquad 1 \leq i \leq N , \qquad r \geq 1 ,$$

hence it suffices to prove the convergence of the sequence $\sum_{r=1}^{\infty} \widetilde{\varepsilon}_{rN+1}$. From the relation (100) we see that the later series is majorized by the convergent series $\sum_{r=1}^{\infty} (0.3)^r \cdot \widetilde{\varepsilon}_1$. This proves the absolute convergence of the series $\sum_{k=2}^{\infty} u_i^{(k)}$ for $i \in \{1, \ldots, n\}$ and therefore the convergence of the sequence $((D_k)^{-1}, k \geq 1)$.

<div align="right">Q.E.D.</div>

Note that the global convergence (i.e. the convergence for all definite pairs $(A, B)$ ) of the Falk–Langemeyer method in the case of cyclic pivot strategies is not yet proved.

We end this section by showing that our asymptotic assumptions also imply ultimate quadratic reduction of $\varepsilon_{rN+1}$. Indeed, for $r \geq 1$ the relation (103) implies

$$\varepsilon_{rN+1} \leq (d_{max}^{(rN+1)})^2 \cdot \widetilde{\varepsilon}_{rN+1} , \qquad \widetilde{\varepsilon}_{rN+1} \leq \frac{1}{(d_{min}^{(rN+1)})^2} \varepsilon_{rN+1} ,$$

where

$$d_{max}^{(rN+1)} = \max\{d_1^{(rN+1)}, \dots, d_n^{(rN+1)}\} \ ,$$
$$d_{min}^{(rN+1)} = \min\{d_1^{(rN+1)}, \dots, d_n^{(rN+1)}\} \ .$$

Theorem 16 implies

$$\varepsilon_{(r+1)N+1} \leq (d_{max}^{((r+1)N+1)})^2 \widetilde{\varepsilon}_{(r+1)N+1} \leq (d_{max}^{((r+1)N+1)})^2 \frac{\sqrt{N}}{\delta} \widetilde{\varepsilon}_{rN+1}^2$$

$$\leq \left[ \frac{d_{max}^{((r+1)N+1)}}{(d_{min}^{(rN+1)})^2} \right]^2 \frac{\sqrt{N}}{\delta} \varepsilon_{rN+1}^2 \leq c \cdot \frac{\sqrt{N}}{\delta} \varepsilon_{rN+1}^2 \ , \qquad r \geq 1 \ ,$$

where $c$ is an upper bound of the convergent sequence ($[d_{max}^{((r+1)N+1)}/(d_{min}^{(rN+1)})^2]^2$, $r \geq 1$). In a similar way one can prove that quadratic reduction of $\varepsilon_{rN+1}$ ultimately implies quadratic reduction of $\widetilde{\varepsilon}_{rN+1}$.

The techniques described in this section can be used for studying asymptotic convergence properties of various different Jacobi–type algorithms.

# 5 Concluding remarks

In Algorithm 4 only $(l,m)-$restrictions of the pair $(A^{(k)}, B^{(k)})$ are used in each step. Therefore, parallel strategies are in fact cyclic (see [10]) and Theorem 16 and Theorem 17 hold for them as well.

In [13] it is proved that if the assumptions of Theorem 16 hold and the pivot strategy is serial, then

$$\widetilde{\varepsilon}_{N+1} \leq \frac{\widetilde{\varepsilon}_1^2}{\delta} \ .$$

**Modified method.** If the problem (1) has multiple eigenvalues, the method can fail to be quadratically convergent. This failure occurs because when pairs $[a_{ll}^{(k)}, b_{ll}^{(k)}]$ and $[a_{mm}^{(k)}, b_{mm}^{(k)}]$ (here $(l,m)$ is the pivot pair in the $k$–th step) approximate the same eigenvalue, then parameters $\widetilde{\alpha}_k$ and $\widetilde{\beta}_k$ can be of order $O(1)$ and, therefore, some previously annihilated elements can become of order $O(\widetilde{\varepsilon}_k)$ again. This situation is described in detail in [7] and [13]. Simple

omitting of these critical steps does not yield to the quadratic convergence, even though the measure $\widetilde{\tau}_k = \tau(\widetilde{A}^{(k)}, \widetilde{B}^{(k)})$, $k \geq 1$, from Corollary 2 tends to zero. The relation (16) does not imply that the off-diagonal elements of diagonal blocks tend to zero together with $\widetilde{\tau}_k$, but merely that the diagonal blocks become more and more proportional. Therefore, $\widetilde{\varepsilon}_k$ does not have to tend to zero at all and the convergence of $\widetilde{\tau}_k$ can considerably slow down. If we modify the method so that in such cases we use triangular transformation matrices similar to the matrix from step (5a) of Algorithm 4, the quadratic convergence persists.

Modification of the Falk–Langemeyer method and the proof of quadratic convergence of the modified method will be topics of our subsequent paper.

**Numerical results.** Our test program is written in FORTRAN in double precision. Test pairs were generated in the manner that $A = G^T D_A G$ and $B = G^T D_B G$, where diagonal matrices $D_A$ and $D_B$ are being read and $G$ is random. For elements of matrix $G$ only numbers which are sums of the powers of 2 were used, so the test pairs were stored as accurately as possible.

The iterative process is terminated when, after some cycle $r$, inequality

$$\varepsilon_{rN+1} \;<\; eps \cdot \sqrt{\|A\|^2 + \|B\|^2} \cdot 2N$$

is fulfilled, where $eps$ is machine precision. After the end of the process, the maximal error of the residual

$$\max_{1 \leq i \leq n} \left\{ \frac{\| b_i' A f_i' - a_i' B f_i' \|_{max}}{\sqrt{(a_i')^2 + (b_i')^2}\sqrt{\| A f_i' \|^2 + \| B f_i' \|^2}} \right\} ,$$

is calculated. Here $[a_i', b_i']$ are the calculated eigenvalues of the pair $(A, B)$ and $f_i'$ are the corresponding eigenvectors. Also the maximal absolute values of the off-diagonal elements of matrices $(F')^T A F'$ and $(F')^T B F'$ are calculated. Those three quantities were usually of order. Infinite eigenvalues were represented with numbers of order of magnitude $O(1/\text{machine precision})$.

We observed the convergence of both measures $\varepsilon_k$ and $\widetilde{\varepsilon}_k$. Observations confirmed all theoretical results. For starting pairs that were not almost diagonal, convergence was in the beginning linear and several cycles were needed before quadratic convergence started. The asymptotic assumption (A2) appears to be very adequate because in almost all cases quadratic convergence

41

started after it was fulfilled. Algorithm behaved very regularly in the sense that the condition $\Im^{(k)} \geq 0$, $k \geq 1$, (see assertion (i) of Proposition 3) was always fulfilled for definite starting pairs. This condition was fulfilled even in some cases when the starting pair was semidefinite, or slightly indefinite.

Average number of cycles for smaller matrices ($n \leq 15$) was around ten and for larger matrices ($n \leq 100$) around fifteen. Last cycles were usually empty, i.e. not all $N$ steps were executed. For orientation, the approximate duration of the process is five minutes for $n = 40$ and one and a half hour for $n = 100$ on IBM PC/AT with a coprocessor, and about thirty times shorter on IBM 4371.

In the presence of very close eigenvalues several additional cycles were usually needed because the quadratic convergence was delayed. The existance of additional cycles does not disagree with theoretical results since the quantity $\delta$ from the asymptotic assumption (A2) is in this case very small.

We observed that the results are generally better if increasing or decreasing order of numbers defined with diagonal pairs $[a_{ii}^{(k)}, b_{ii}^{(k)}]$ is preserved by interchanging pivot rows and columns if necessary. However, interchanging must be stopped after the asymptotic assumption (A2) is fulfilled. Otherwise some off–diagonal element which was not yet annihilated can "run away" from annihilation and therefore terminate quadratic convergence.

æ **Example.** We give an example of the pair of order 10 generated in the previously described manner. Elements of the matrices $D_A$ and $D_B$ are

$$-2, \ 1, \ 10, \ 0, \ -0.001, \ 10, \ 1, \ 5, \ 5, \ 4$$

and

$$-1, \ 0.1, \ -1, \ -100, \ -100, \ 0, \ -1, \ 0.1, \ 1, \ 1,$$

respectively, so the exact eigenvalues of the problem are

$$2, \ 10, \ -10, \ 0, \ 0.00001, \ \infty, \ -1, \ 50, \ 5, \ 4 \,.$$

Elements of the matrix $G$ are uniformly distributed integers from the interval $[-10, 10]$. Note that both matrices $A$ and $B$ are indefinite, while the pair $(A, B)$ itself is definite (for example $A - 3B > 0$). In order to increase the stability of the computation, the process started from the normalized pair $(\widetilde{A}, \widetilde{B})$.

Only upper triangles of the matrices $A$ and $B$ are displayed. Each row begins with the diagonal element. Asymptotic convergence is described as

follows: in column CYC is the number of cycle; in column ROT is the number of rotations performed in the cycle; columns SUMA, SUMB, SUM and SUMT display values of $S(A^{(k)})$, $S(B^{(k)})$, $\varepsilon_k$ and $\widetilde{\varepsilon}_k$ after the cycle, respectively.

```
ORDER OF MATRICES N =  10
COLUMN CYCLIC PIVOT STRATEGY
STOPPING CRITERION: SUM(K) <  .49D-13

MATRIX A
ROW
 1     .21350D+04     .41900D+03     .11600D+03    -.11430D+04    -.10490D+04
       .44002D+03    -.13750D+04     .20027D+02     .51903D+03    -.60802D+03
 2     .14310D+04    -.32700D+03    -.34200D+03    -.26100D+03    -.10390D+04
      -.29600D+03    -.43200D+03     .50000D+03    -.13100D+03
 3     .18320D+04     .28000D+03     .93300D+03     .64100D+03    -.91000D+03
      -.38500D+03    -.74500D+03    -.58200D+03
 4     .11860D+04     .85799D+03     .34099D+03     .56100D+03    -.40001D+03
      -.49101D+03     .57401D+03
 5     .99295D+03     .68695D+03     .22402D+03    -.53006D+03    -.69606D+03
      -.90965D+02
 6     .13360D+04    -.57298D+03    -.43106D+03    -.97606D+03     .43035D+02
 7     .13470D+04     .29703D+03     .10027D+02     .67399D+03
 8     .87292D+03     .32992D+03    -.49550D+01
 9     .88792D+03     .40105D+03
10     .82798D+03
```

```
MATRIX B
ROW
 1   -.72420D+04   -.81550D+04    .40130D+04   -.20630D+04   -.39800D+03
      .31460D+04   -.73170D+04    .83650D+04    .27080D+04   -.78930D+04
 2   -.99425D+04    .50853D+04   -.28814D+04   -.28007D+04    .93320D+03
     -.79655D+04    .68774D+04   -.11530D+03   -.79841D+04
 3   -.26020D+04    .14878D+04    .14297D+04   -.63790D+03    .39712D+04
     -.34409D+04   -.16500D+02    .39654D+04
 4   -.10848D+04   -.15558D+04   -.55260D+03   -.20896D+04    .11306D+04
     -.83980D+03   -.18888D+04
 5   -.59091D+04   -.46836D+04   -.23130D+03   -.41948D+04   -.61837D+04
      .10755D+04
 6   -.50607D+04    .28013D+04   -.69463D+04   -.62826D+04    .43479D+04
 7   -.72920D+04    .82531D+04    .26803D+04   -.78484D+04
 8   -.12938D+05   -.81476D+04    .10050D+05
 9   -.82005D+04    .46091D+04
10   -.88056D+04
```

```
A S Y M P T O T I C   C O N V E R G E N C E
   CYC   ROT    SUMA      SUMB      SUM       SUMT
    1    45    .60D+00   .35D+01   .35D+01   .18D+01
    2    45    .72D+00   .32D+00   .79D+00   .12D+01
    3    45    .57D+00   .28D+00   .64D+00   .72D+00
    4    45    .31D+00   .21D+00   .38D+00   .24D+00
    5    45    .18D-01   .23D-01   .30D-01   .25D-01
    6    45    .39D-02   .13D-01   .13D-01   .93D-02
    7    45    .56D-04   .15D-03   .16D-03   .92D-04
    8    44    .18D-09   .61D-09   .63D-09   .19D-09
    9    29    .26D-20   .19D-20   .32D-20   .93D-20
TOTAL NO. OF ROTATIONS   388          TIME(sec)   5.68
```

```
C A L C U L A T E D    E I G E N V A L U E S
 I       A(I,I)            B(I)                D(I)
 1     .36774301D+01    .91935754D+00     .40000000000001D+01
 2     .29605919D+01    .59211837D-01     .50000000000001D+02
 3     .24951699D+00    .24951699D-01     .10000000000015D+02
 4     .62005903D+00    .12401181D+00     .50000000000011D+01
 5    -.16044487D+00   -.80222433D-01     .19999999999999D+01
 6    -.94464581D-04   -.94464581D+01     .10000000000004D-04
 7    -.43510239D-16   -.23229057D+01     .18730953468577D-16
 8     .62365226D-01   -.62365226D-01    -.99999999999990D+00
 9     .42815492D+01   -.25722206D-13    -.16645341957511D+15
10     .32339485D+01   -.32339485D+00    -.99999999999999D+01

   MAXIMAL(relative) ERROR =   .15D-13   FOR  I =    3
   MAXIMAL OFF-DIAGONAL ELEMENTS:
   Ft A F =   .35D-14         Ft B F =   .30D-13
```

# References

[1] Charlier, J.-P., Van Dooren, P., *A Jacobi–like algorithm for computing the generalized Schur form of a regular pencil*, Philips Research Laboratory, technical report,1988.

[2] Falk, S., Langemeyer, P., *Das Jacobische Rotations–Verfahren für reelsymmetrische Matrizen–Paare I, II*, Elektronische Datenverarbeitung (1960) 30–43.

[3] Gose, G., *Das Jacobi Verfahren für  $Ax = \lambda Bx$* , ZAMM 59(1979) 93–101.

[4] Hari, V., *On Cyclic Jacobi Methods for the Positive Definite Generalized Eigenvalue Problem,* Disertation, Fernuniversität Hagen, 1984.

[5] Hari, V., *On the Convergence of cyclic Jacobi–Like Processes,* Lin. Alg. and Its Appl. 81(1986) 105–127.

[6] Hari, V., *On the Quadratic Convergence of Jacobi Algorithms,* Radovi matematički 2(1986) 127–146.

[7] Hari, V., *On Pairs of Almost Diagonal Matrices,* proposed for publication in Lin. Alg. and Its Appl.

[8] Henrici, P., Zimmermann, K., *An Estimate for the Norms of Certain Cyclic Jacobi Operators,* Lin. Alg. and Its Appl. 1(1968) 489–501.

[9] Kempen, H.P.M. van, *On the Quadratic Convergence of the Serial Cyclic Jacobi Method,* Numer. Math. 9(1966) 19–22.

[10] Luk, F., Park, H., *On the equivalence and convergence of parallel Jacobi algorithms,* Proceedings SPIE conference on Advanced Algorithms and Architecture II, 1987.

[11] Parlett, B.N., *Symmetric Eigenvalue Problem,* Prentice Hall Inc., Englewood Cliffs, N.J., 1980.

[12] Sameh, A. H., *On Jacobi and Jacobi–like Algorithms for a Parallel Computer,* Mathematics of Computation, Vol. 25 No. 9(1971) 579–590.

[13] Slapničar, I., *Quadratic Convergence of the Falk–Langemeyer Method,* M. S. Thesis, University of Zagreb, 1988. (in Croatian language)

[14] Stewart, G.W., *Perturbation Bounds for the Definite Generalized Eigenvalue Problem,* Lin. Alg. and Its Appl. 23(1960) 69–85.

[15] Veselić, K., *An eigenreduction algorithm for definite matrix pairs and its applications to overdamped linear systems,* to appear in SIAM J. Sci. Stat. Comp.

[16] Wilkinson, J.H., *Note on the Quadratic Convergence of the Cyclic Jacobi Processes,* Numer. Math. 4(1962) 296–300.

[17] Wilkinson, J.H., *Almost Diagonal Matrices with Multiple or Close Eigenvalues,* Lin. Alg. and Its Appl. 1(1968) 1–12.

[18] Wilkinson, J.H., *The Algebraic Eigenvalue Problem,* Oxford University Press, Oxford,1965.

[19] Zimmermann, K., *On the Convergence of a Jacobi Process for Ordinary and Generalized Eigenvalue Problems,* Disertation No 4305(1965), Eidgenossische Technische Hochschule Zürich.

æ