# A norm-reducing Jacobi-like algorithm for the eigenvalues of non-normal matrices

I. Kiessling and A. Paulik \*

### **ABSTRACT**

A new norm decreasing Jacobi-like method for reducing a non-normal matrix to a normal one is described. The method is an improved version of the Huang-Gregory's procedure [6], in which certain norm reducing non-optimal steps are replaced by the optimal ones, which correspond no more to regular matrix transformations. The method renders itself particularly effective in dealing with defective matrices of special forms. Theory and experiments alike indicate that this algorithm is in all computational aspects — accuracy, convergence rate, computing time, complexity of the computer program — better or at least as good as the original one.

Both procedures, the original and the improved one, end in all of the authors computed examples with the diagonal matrix containing the eigenvalues of the non-normal initial matrix.

#### 1. INTRODUCTION

The aim of this paper is to present an improved version of the norm-decreasing Jacobi-like algorithm for reducing a non-normal matrix to a normal one, which was introduced by Huang and Gregory [6]. The above-mentioned algorithm, like all norm-decreasing Jacobi-like algorithms, cf. e.g. [1, 2, 3, 10, 11, 14, 15] originates in the following well known facts:

Let  $A = (a_{j,k} | j, k=1,...,n)$ ,  $a_{j,k} \in \mathcal{C}$  be some  $n \times n$ -matrix,  $\lambda_j$  (j=1,...,n) its - not necessary different - eigenvalues and

$$\|\mathbf{A}\| = \left(\sum_{j,k=1}^{n} |\mathbf{a}_{j,k}|^2\right)^{1/2} \tag{1.1}$$

its Euclidean norm. Then there holds not only the Schurinequality, cf. [13],

$$\sum_{k=1}^{n} |\lambda_k|^2 \le ||\mathbf{A}||^2, \tag{1.2}$$

with equality only iff A is normal, but – since the spectrum of A does not change by a similarity transformation  $A \rightarrow ZAZ^{-1}$  – also its generalisation

$$\sum_{k=1}^{n} |\lambda_k|^2 \leq \|ZAZ^{-1}\|^2, \text{ for all } Z \in R_n, \tag{1.3}$$

where  $R_n$  denotes the set of all regular complex-valued  $n \times n$ -matrices. Mirski [9] has shown that there holds

$$\sum_{k=1}^{n} |\lambda_k|^2 = \inf_{Z \in R_n} \|z A z^{-1}\|^2. \tag{1.4}$$

Let now  $(Z_p) \subseteq R_n$  be some sequence with

$$\sum_{k=1}^{n} |\lambda_{k}|^{2} = \lim_{p \to \infty} \|A_{p}\|^{2}, \quad A_{p} := Z_{p}AZ_{p}^{-1}, \quad (1.5)$$

then it can be shown that this is the case then and only then, if

$$\lim_{p \to \infty} \|A_p^* A_p - A_p A_p^*\| = 0.$$
 (1.6)

Hence, choosing p large enough, A<sub>p</sub> can be made arbitrarily close to a normal matrix with the same eigenvalues as A. But for normal matrices there exists an extension of the Jacobi algorithm [7,12] for symmetric matrices, cf. [4], so the following simple model algorithm for the computation of the eigenvalues of a general matrix may be proposed:

- i) Find the minimizing sequence of regular matrices  $(Z_p) \subseteq R_n$  so that (1.5) holds.
- ii) Choose  $\widetilde{p}$  large enough (or take some normal accumulation point  $A_{\infty} = \lim_{k \to \infty} A_{p_k}$ ) so that  $A_{\widetilde{p}}$  is nearly normal.
- iii) If A<sub>p</sub> (or A<sub>∞</sub>) is not diagonal, apply the sequence
   (U<sub>p</sub>) of the unitary matrices of the extended Jacobi algorithm [4] on A<sub>p</sub> in order to obtain the eigenvalues λ<sub>L</sub>:

$$\Lambda = \operatorname{diag}(\lambda_1, ..., \lambda_n) = \lim_{p \to \infty} \operatorname{U}_p^* A_{\widetilde{p}} \operatorname{U}_p (\lim_{p \to \infty} \operatorname{U}_p^* A_{\infty} \operatorname{U}_p)$$

As one can see, the crucial point in designing any algorithm for the eigenvalue computation of non-normal matrices is to find the sequence  $(Z_p)$ . The Jacobi-like algorithms [1,2,3,6,10,11,14,15] differ mainly in the way how one chooses the norm-minimizing sequence  $(Z_p)$ . In each algorithm ii) and iii) are done simultaneously. In the following we confine ourself to the description

<sup>\*</sup> I. Kiessling, A. Paulik, Institut für Numerische und Angewandte Mathematik der Georg August Universität, D-3400 Göttingen, BRD.

and analysis of the algorithm of Huang and Gregory [6].

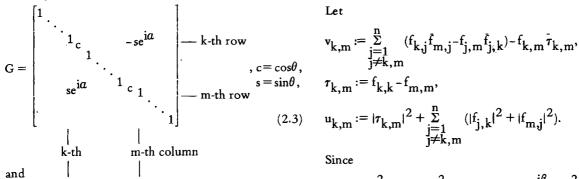
## 2. HUANG-GREGORY'S ALGORITHM

The algorithm [6] uses three kinds of similarity matrix transformation. Given a matrix  $A_p$  and the pivot pair (k, m), it generates the matrix  $A_{p+1}$  by the rule :

$$A_{p+1} = \begin{cases} D^{-1} A_p D, & \text{if } k=m \\ L^{-1} FL \text{ and } F=G^* A_p G, & \text{if } k \neq m, \end{cases}$$
 (2.1)

with

$$D = diag(1,...,1,d,1,...,1), d in the k-th row.$$
 (2.2)



The similarity transformations using matrices D and L are norm reducing, the similarity transformation with the unitary matrix G is diagonalising. The transformation parameters d,  $\theta$ ,  $\alpha$ ,  $\rho$ ,  $\beta$  are chosen in the following way:

Parameter d :

Let 
$$A_p = (b_{j,k} | j,k = 1,...,n)$$
 and

$$\mu_{k}^{2} = \sum_{\substack{j=1\\j\neq k}}^{n} |b_{k,j}|^{2}, \ \xi_{k}^{2} = \sum_{\substack{j=1\\j\neq k}}^{n} |b_{j,k}|^{2},$$

$$\mu_{k}^{2} = \mu_{k}^{2} + |b_{k,k}|^{2}, \ \hat{\xi}_{k}^{2} = \xi_{k}^{2} + |b_{k,k}|^{2}.$$
(2.5)

If  $\mu_{\mathbf{k}}$ ,  $\xi_{\mathbf{k}} \neq 0$ , set  $\mathbf{d} := (\mu_{\mathbf{k}}/\xi_{\mathbf{k}})^{1/2}$ , achieving the maximum possible norm reduction

$$\Delta = \|\mathbf{A}_{\mathbf{p}}\|^2 - \|\mathbf{A}_{\mathbf{p}+1}\|^2 = (\mu_{\mathbf{k}} - \xi_{\mathbf{k}})^2. \tag{2.6a}$$

If  $\mu_k \cdot \xi_k = 0$  and  $\hat{\mu}_k \cdot \hat{\xi}_k \neq 0$ , set  $\mathbf{d} := (\hat{\mu}_k / \hat{\xi}_k)^{1/2}$ , achieving the norm reduction

$$\Delta = \|\mathbf{A}_{\mathbf{p}}\|^2 - \|\mathbf{A}_{\mathbf{p}+1}\|^2 = (\hat{\mu}_{\mathbf{k}} - \hat{\xi}_{\mathbf{k}})^2 + |\mathbf{b}_{\mathbf{k},\mathbf{k}}|^2 (\mathbf{d} - \frac{1}{\mathbf{d}})^2. \tag{2.6b}$$

If 
$$\mu_k \cdot \hat{\xi}_k = 0$$
, i.e.  $b_{j,k} = b_{k,j}$  (j=1,...,n), deflate the matrix and continue (2.6c)

Parameters  $\theta$  and a:

Choose G so as to annihilate the element fm, k in  $F = (f_{i,k} | i,j = 1,...,n)$ , i.e. the parameters  $\theta$  and  $\alpha$  must satisfy the equation

$$tan \theta = \frac{2b_{m,k}e^{i\alpha}}{t_{k,m} + (t_{k,m}^2 + 4b_{k,m}b_{m,k})^{1/2}},$$
with  $t_{k,m} = b_{k,k} - b_{m,m}$ . (2.7)

Here a is chosen so that tan  $\theta$  is real and  $\theta$  is chosen to be the smaller (in absolute value) of the rotation angles.

Parameters  $\rho$  and  $\beta$ :

$$\mathbf{v}_{k,m} := \sum_{\substack{j=1\\j\neq k,m}}^{n} (f_{k,j}\tilde{f}_{m,j} - f_{j,m}\tilde{f}_{j,k}) - f_{k,m}\tilde{\tau}_{k,m},$$

$$\tau_{k,m} := f_{k,k} - f_{m,m},$$

$$\mathbf{u}_{k,m} := |\tau_{k,m}|^2 + \sum_{\substack{j=1\\ j \neq k, m}}^{n} (|f_{j,k}|^2 + |f_{m,j}|^2). \tag{2.8}$$

$$\Delta = \|\mathbf{A}_{\mathbf{p}}\|^2 - \|\mathbf{A}_{\mathbf{p}+1}\|^2 = 2\rho \operatorname{Re}(\mathbf{v}_{k,m} e^{-i\beta}) - \rho^2 \mathbf{u}_{k,m},$$

$$\mathbf{v}_{\mathbf{k},\mathbf{m}} e^{-\mathbf{i}\beta} = |\mathbf{v}_{\mathbf{k},\mathbf{m}}| \tag{2.9}$$

$$\rho = \begin{cases} \frac{|\mathbf{v}_{k,m}|}{\mathbf{u}_{k,m}}, & \text{for } 0 < |\mathbf{v}_{k,m}| \le \mathbf{u}_{k,m}, \\ 1, & \text{for } 0 < \mathbf{u}_{k,m} < |\mathbf{v}_{k,m}|, \\ 0, & \text{for } 0 = \mathbf{v}_{k,m}, \end{cases}$$
(2.10a)

$$\rho = \frac{1}{1}, \text{ for } 0 < u_{k,m} < |v_{k,m}|, \qquad (2.10b)$$

$$0, \text{ for } 0 = \mathbf{v}_{\mathbf{k}, \mathbf{m}},$$
 (2.10c)

achieving the norm reduction

$$\Delta = \begin{cases} \frac{|\mathbf{v}_{k,m}|^2}{\mathbf{u}_{k,m}}, & \text{for } 0 < |\mathbf{v}_{k,m}| \le \mathbf{u}_{k,m}, \\ 2|\mathbf{v}_{k,m}| \sim \mathbf{u}_{k,m}, & \text{for } 0 < \mathbf{u}_{k,m} < |\mathbf{v}_{k,m}|, \\ 0, & \text{for } 0 = \mathbf{v}_{k,m}. \end{cases}$$
(2.11a)

$$\Delta = \begin{cases} 2|\mathbf{v}_{k,m}| - \mathbf{u}_{k,m}, & \text{for } 0 < \mathbf{u}_{k,m} < |\mathbf{v}_{k,m}|, \\ 0 < \mathbf{u}_{k,m} < |\mathbf{v}_{k,m}|, \end{cases}$$
 (2.11b)

$$0, \text{ for } 0 = \mathbf{v_{k m}}.$$
 (2.11c)

#### 3. IMPROVED ALGORITHM

Let us now analyse the norm reducing transformations. Looking closely at the choice of the parameters in the non-optimal (not maximal) norm reduction transformations, one can see that letting drop the requirement that the norm reduction should happen by the similarity transformation there is a possibility to reach maximal norm reduction, cf. e.g. [8]. Not to get lost in details, we mention only the well known and in the following extensively used principle: in an upper or lower triangular block matrix one may set all the off-diagonal blocks

equal to zero, without changing eigenvalues and their algebraic multiplicities.

We improve now the choice of the transformation parameters d and  $\rho$  in (2.6b) and (2.10c) in the following way:

Proposition 1 (improves 2.6b).

If  $\mu_k \cdot \xi_k = 0$  then form  $A_{p+1}$  from  $A_p$  by the rule : set all off diagonal elements in the k-th row and the k-th column of the matrix A<sub>p</sub> equal to zero, achieving the maximum possible norm reduction

$$\Delta = \|\mathbf{A_p}\|^2 - \|\mathbf{A_{p+1}}\|^2 = (\mu_k - \xi_k)^2. \tag{3.1}$$

Proposition 2 (improves 2.10c)

Set p according to

$$\rho = \begin{cases} \frac{|\mathbf{v}_{k,m}|}{\mathbf{u}_{k,m}}, & \text{for } 0 < |\mathbf{v}_{k,m}| \le \mathbf{u}_{k,m}, \\ 1, & \text{for } 0 < \mathbf{u}_{k,m} < |\mathbf{v}_{k,m}|, \end{cases}$$
(3.2a)

$$\rho = \begin{cases} 1, & \text{for } 0 < u_{k,m} < |v_{k,m}|, \\ 0, & \text{for } 0 = v_{k,m} \text{ and } u_{k,m} \neq |\tau_{k,m}|^2, \end{cases}$$
(3.2b)

achieving, except for (3.2b), the maximum possible

$$\Delta = \begin{cases} \frac{|\mathbf{v}_{k,m}|^2}{\mathbf{u}_{k,m}}, & \text{for } 0 < |\mathbf{v}_{k,m}| \le \mathbf{u}_{k,m}, \\ 2|\mathbf{v}_{k,m}| - \mathbf{u}_{k,m}, & \text{for } 0 < \mathbf{u}_{k,m} < |\mathbf{v}_{k,m}|, \\ 0, & \text{for } 0 = \mathbf{v}_{k,m} \text{ and } \mathbf{u}_{k,m} \ne |\tau_{k,m}|^2. \end{cases}$$
(3.3a)

$$\Delta = \begin{cases} 2|\mathbf{v}_{k,m}| - \mathbf{u}_{k,m}, & \text{for } 0 < \mathbf{u}_{k,m} < |\mathbf{v}_{k,m}|, \\ 0, & \text{for } 0 = \mathbf{v}_{k,m} \text{ and } \mathbf{u}_{k,m} \neq |\tau_{k,m}|^2. \end{cases}$$
(3.3b)

If  $v_{k,m} = 0$  and  $u_{k,m} = |\tau_{k,m}|^2$ , then form  $A_{p+1}$  from F by setting all off-diagonal elements of the k-th and m-th rows and columns equal to zero, achieving the maximum possible norm reduction

$$\Delta = \|A_{p}\|^{2} - \|A_{p+1}\|^{2}$$

$$= |f_{k,m}|^{2} + \sum_{\substack{j=1\\ j \neq k \ m}}^{n} |f_{k,j}|^{2} + |f_{j,m}|^{2}.$$
(3.4)

Proof of the proposition 1

Let  $\mu_k$ ,  $\xi_k = 0$ . Since this fact implies that  $A_p$  has one of

k-th column
$$\begin{bmatrix}
x \\
\vdots \\
x \\
0...0a_{k,k}0...0
\end{bmatrix} \text{ or } \begin{bmatrix}
0 \\
\vdots \\
0 \\
x...xa_{k,k}x...x
\end{bmatrix}, (3.5)$$

we may set all the remaining off-diagonal elements of the k-th column or k-th row, respectively, equal to zero, without changing the eigenvalues of A<sub>D</sub>. Hence (3.1) follows.

Proof of the proposition 2

Since we have chosen  $\beta$  according to (2.9) we obtain for ∆ the expression

$$\Delta = \|A_{\mathbf{p}}\|^2 - \|A_{\mathbf{p}+1}\|^2 = 2\rho |\mathbf{v}_{\mathbf{k},\mathbf{m}}| - \rho^2 \mathbf{u}_{\mathbf{k},\mathbf{m}}, \qquad (3.6)$$

which is for  $u_{k,m} \neq 0$  a quadratic function of  $\rho$ . If  $u_{k,m} \neq 0$ , it is positive for

$$0 < \rho < \frac{2|\mathbf{v}_{\mathbf{k},\mathbf{m}}|}{\mathbf{u}_{\mathbf{k},\mathbf{m}}} =: \hat{\rho} \tag{3.7}$$

and achieves a maximum (3.3a) for  $\rho$  chosen according to (3.2a). If  $\hat{\rho} > 1$ , then we choose for computational purposes  $\rho = 1$  according to (3.2b) and obtain the (not maximum possible) norm reduction (3.3b). If  $v_{k,m} = 0$ , then (3.6) reduces to

$$\Delta = -\rho^2 u_{\mathbf{k},\mathbf{m}},\tag{3.8}$$

implying that the only possible regular transformation, which does not increase the norm of A<sub>p+1</sub> is to choose  $\rho = 0$ . However, there is an important exception for  $u_{k,m} = |\tau_{k,m}|^2$ . In this case we obtain from (2.8)  $f_{i,k} = f_{m,i} = 0 \ (j=1,...,n; j \neq k)$ 

$$\begin{bmatrix}
0 & x \\
\vdots & \vdots \\
0 & x
\end{bmatrix}$$

$$x...x f_{k,k} x...x f_{k,m} x...x$$

$$0 & x \\
\vdots & \vdots \\
0 & x
\end{bmatrix}$$

$$0...0 0 0...0 f_{m,m} 0...0$$

$$0 & x \\
\vdots & \vdots \\
0 & x
\end{bmatrix}$$

$$0...0 x$$

$$0 & x \\
\vdots & \vdots \\
0 & x
\end{bmatrix}$$

Hence, we may set all the remaining off-diagonal elements of the k-th and m-th rows and columns equal to zero, without changing the eigenvalues of An, obtaining the maximum possible norm reduction (3.4).

#### Remark 1

The norm reduction in proposition 2 occurs also in the important case of the defective eigenvalues: the difference between the algebraic and the geometric multiplicity of the eigenvalues  $f_{k,k}$ ,  $f_{m,m}$  in the case  $f_{k,k} = f_{m,m}$ is at least 1, if  $\Delta > 0$ .

# 4. NUMERICAL RESULTS

All calculations were performed with a UNIVAC Series 1100 computer using 60 bits for mantissa. The authors of [6] performed their calculations with an IBM 360/65 using 56 bits for mantissa. Our implementation of the Huang-Gregory's algorithm reproduced therefore the results for the matrices given in [6] with little higher

Let us now consider the conditions which must be satis-

TABLE 1

ex. 4 of [6]	stopping criterion	A	В		С	
			$\epsilon = 0$	$\epsilon = 10^{-12}$	$\epsilon = 0$	$\epsilon = 10^{-12}$
$  A_p - \operatorname{diag}(A_p)  $	i)	.90.10 <sup>-21</sup>	.20.10 <sup>-35</sup>		.20.10 <sup>-35</sup>	0
	ii)		.66.10 <sup>-144</sup>	.17.10 <sup>-3</sup>	0	0
A <sub>p</sub> A <sub>p</sub> *-A <sub>p</sub> *A <sub>p</sub>	i)		.23.10 <sup>-38</sup>		.23.10 <sup>-88</sup>	0
	ii)		.4.10 <sup>-144</sup>	.45.10 <sup>-7</sup>	0	0
A-A <sub>p</sub>      A	i)	.40.10 <sup>-7</sup>	.41.10 <sup>-8</sup>		.41.10 <sup>-8</sup>	.41.10 <sup>-8</sup>
	ii)		.41.10 <sup>-8</sup>	.23.10 <sup>-4</sup>	.41.10 <sup>-8</sup>	.41.10 <sup>-8</sup>
Number of	i)	9	9		9	8
sweeps p	ii)		30	30	11	8
CPU [s]	i)		.650		.697	.542
	ii)		1.899	2.139	.799	.542

fied applying the improved algorithm. Consider  $\xi_k \cdot \mu_k = 0$  and  $\hat{\xi}_k \cdot \hat{\mu}_k \neq 0$  (cf. (2.6b)). This case occurs for instance, if  $\xi_k = 0$ ;  $\mu_k$ ,  $\hat{\xi}_k$ ,  $\hat{\mu}_k \neq 0$ . Due to the rounding errors the evaluation of  $\xi_k$  may give a small value > 0 instead of zero. Then proposition 1 giving maximal norm reduction would not be applied. To suppress this influence of rounding errors we define a parameter  $\epsilon > 0$  and evaluate the relational expression  $\xi_k > \epsilon$  instead of  $\xi = 0$ . If a matrix occurs having the form (3.9) we proceed in an analogous way.

The numerous examples computed by us can be divided into three classes.

- 1. General class, containing examples 2 to 7 from [6] for instance. The improved algorithm is at least as good as the original one.
- 2. Special matrices having essentially triangular or Jordan's normal form. These structures are recognized by the improved algorithm giving the eigenvalues after the first sweep (1). Curiously, the Huang-Gregory's algorithm needs a multiple number of sweeps or it fails.
- 3. In exceptional cases, see ex. 1 [6], both algorithms give a normal, but not a diagonal matrix. Performing the algorithms in such a way that small perturbations of the matrices A<sub>p</sub> due to the rounding errors are admitted does not change the eigenvalues. However, the computation algorithms are instable in that sense that after slightly disturbing the matrices A<sub>p</sub> converging against some normal matrix, convergence occurs against a diagonal matrix. So we obtain in ex. 1 [6] the eigenvalues with high accuracy after 6 sweeps.

To give a numerical example of the 1-st class we select ex. 4 of [6]. This is a defective 5\*5-matrix. In the following table 1 column A contains values taken from [6],

column B values of the Huang-Gregory's algorithm implemented by us, column C values of the improved algorithm, both for  $\epsilon = 0$  and  $\epsilon = 10^{-12}$ , respectively. All rows marked by i) contain the data obtained by that sweep after which condition i)  $\|A_p - \operatorname{diag}(A_p)\| \le 10^{-8}$ was satisfied for the first time. An empty field in a row marked by i) means that condition i) was not satisfied after the 30-th sweep. All rows marked by ii) contain data obtained by that sweep after which  $||A_p - diag(A_p)||$ or  $\|A_p A_p^* - A_p^* A_p\|$  became zero for the first time or those computed by the 30-th sweep. An empty field in a row marked by ii) means that condition i) was satisfied for the first time by  $||A_p - diag(A_p)|| = 0$ . One can see that for  $\epsilon=0$  the two algorithms do not differ remarkably. This behaviour is not surprising since the nonregular transformations do not occur in the first 9 sweeps due to the fact that under the influence of the rounding errors generally  $\mu_k$ ,  $\xi_k$ ,  $u_{k,m} \neq 0$ .

For  $\epsilon=10^{-12}$  the number of sweeps p of the improved algorithm needed to satisfy stopping criterion i) decreases to p = 8 whereas the Huang-Gregory's algorithm cannot satisfy this condition.

We set in the table  $\Lambda := \operatorname{diag}(\lambda_1,...,\lambda_n)$  with exact eigenvalues  $\lambda_k$  of the matrix A.

# REFERENCES

- EBERLEIN P. J.: 'A Jacobi-like method for the automatic computation of eigenvalues and eigenvectors of an arbitrary matrix', J. SIAM, 10 (1962) 74-88.
- EBERLEIN P. J., BOOTHROYD J.: 'Solution to the eigenproblem by a norm-reducing Jacobi-type method', Numer. Math., 11 (1968) 1-12.
- EBERLEIN P. J.: 'Solution to the complex eigenproblem by a norm-reducing Jacobi-type method', Numer. Math., 14 (1970) 232-245.

<sup>(1)</sup> A sweep is by definition a sequence of  $\frac{n}{2}(n+1)$  transformations (2.1), in which no pivot pair appears more than once.

- GOLDSTINE H. H., HORWITZ L. P.: 'A procedure for the diagonalisation of normal matrices', J. ACM, 6 (1959), 176-195.
- GREGORY R. T., KARNEY D. L.: A collection of matrices for testing computational algorithms, Wiley-Interscience, 1969.
- HUANG C. P., GREGORY R. T.: 'A norm-reducing Jacobilike algorithm for the eigenvalues of non-normal matrices, Colloquia Mathematica Societatis Janos Bolyai, Keszthely (Hungary), (1977) 365-393.
- JACOBI C. G. J.: 'Über ein leichtes Verfahren, die in der Theorie der Säkularstörungen vorkommenden Gleichungen numerisch aufzulösen', J. Reine Angew. Math., 30 (1864) 51-94.
- 8. KRESS R., DE VRIES H. L., WEGMANN R.: 'On nonnormal matrices', Linear Algebra and Appl., 8 (1974) 108-120.
- MIRSKY L.: 'On the minimization of matrix norms', Amer. Math. Monthly, 65 (1958) 106-107.
- RUHE A.: 'On the quadratic convergence of a generalization of the Jacobi method to arbitrary matrices', BIT 8 (1968) 210-231.
- 11. RUTISHAUSER H. 'Une méthode pour le calcul des valeurs propres des matrices non symétriques', Comptes Rendus, 259 (1964) 2758.
- RUTISHAUSER H.: 'The Jacobi method for real symmetric matrices', in J. H. Wilkinson; C. Reinsch, Handbook for automatic computation, Vol. II, Linear algebra, Springer Verlag, 1971
- SCHUR I.: 'Über die charakteristischen Wurzeln einer linearen Substitution mit einer Anwendung auf die Theorie der Integralgleichungen', Math. Ann., 66 (1909) 488-510.
- 14. VESELIČ K. 'On a class of Jacobi-like procedures for diagonalizing arbitrary real matrices', Numer. Math. 33, (1979) 157-172.
- 15. VESELIČ K., WENZEL H. J.: 'A quadratically convergent Jacobi-like method for real matrices with complex eigenvalues', Numer. Math. 33 (1979) 425-435.