

TWO-SIDED APPROXIMATIONS OF INVERSES, SQUARE ROOTS AND CHOLSKY FACTORS

J. W. SCHMIDT

*Department of Mathematics, Technical University of Dresden,
 Dresden, German Democratic Republic*

1. Introduction

For solving the problem $F(x) = 0$ in partially ordered spaces super-linearly convergent methods are known which give a monotonous enclosing of a zero x^* of F ,

$$(1) \quad x_0 \leq \dots \leq x_{n-1} \leq x_n \leq x^* \leq y_n \leq y_{n-1} \leq \dots \leq y_0,$$

see e.g. [1], [3], [5] and [10]. In general these methods are based on the well-known Newton method and on related methods.

This lecture is concerned with the following special matrix problems:

- (i) $F(x) = x^{-1} - a = 0 \quad (x^* = a^{-1}, \text{ inverse}),$
- (ii) $F(x) = x^2 - a = 0 \quad (x^* = a^{1/2}, \text{ square root}),$
- (iii) $F(x) = xx^T - a = 0 \quad (x^* = l, \text{ Cholesky factor}).$

Of course, it is much more appropriate to derive theorems about these concrete problems than to apply theorems referring to the general zero-problem. Furthermore, the original methods mostly have to be adapted to the problem in question.

The aim of the present lecture is to review some recent results in the mentioned matrix problems.

2. Enclosing of inverses

Let R be a complete normed ring with a unit element e and let $K \subset R$ be a cone by which R is partially ordered: $x \leq y$ for $x, y \in R$ means that

THEOREM 1. *Suppose that*

$$(4) \quad x_0 \leq y_0, \quad x_0 \leq x_1, \quad y_1 \leq y_0,$$

$$(5) \quad e - ax_0 \geq 0,$$

$$(6) \quad x_0^{-1} \text{ exists and } \|e - ax_0\| < 1.$$

Then the inverse of a is well-defined and the sequences produced by method (3) ensure the monotonous enclosing

$$(7) \quad x_0 \leq \dots \leq x_{n-1} \leq x_n \leq a^{-1} \leq y_n \leq y_{n-1} \leq \dots \leq y_0.$$

Furthermore they converge to a^{-1} at least with the R -order 2.

Proof. The enclosing (7) together with $e - ax_n \geq 0$ is verified by induction.

To prove the propositions for $n = 1$ set

$$Tx = x_0 + x(e - ax_0) \quad \text{for } x \in R.$$

Because of (6) the operator T is contractive. Furthermore, $x_1 \leq y_1$ holds, and the interval $[x_1, y_1]$ is closed because the cone K is so. As easily seen, $Tx \in [x_1, y_1]$ is valid for $x \in [x_1, y_1]$. Therefore, by Banach's fixed-point theorem, T has a unique fixed-point $x^* \in [x_1, y_1]$. Using (6), it follows that x^* is equal to a^{-1} . Thus $x_0 \leq x_1 \leq a^{-1} \leq y_1 \leq y_0$ is obtained.

For the step from n to $n+1$ it is seen that

$$x_{n+1} \leq x_n + a^{-1}(e - ax_n) = a^{-1},$$

$$y_{n+1} \geq x_n + a^{-1}(e - ax_n) = a^{-1}$$

and $e - ax_{n+1} = (e - ax_n)^2 \geq 0$. Because of

$$x_{n+1} - x_n = (x_n - x_{n-1})\{(e - ax_{n-1}) + (e - ax_{n-1})^2\},$$

$$y_{n+1} - y_n = (y_n - y_{n-1})\{(e - ax_{n-1}) + (e - ax_{n-1})^2\}$$

the inequalities $x_n \leq x_{n+1}$ and $y_{n+1} \leq y_n$ are fulfilled. Hence, the enclosing (7) is completely proved.

Finally, the convergence of (x_n) and (y_n) to a^{-1} with the R -order 2 is to be verified. Using $\|e - ax_{n+1}\| \leq \|e - ax_n\|^2$, the estimation

$$\|e - ax_n\| \leq \|e - ax_0\|^{2^n}$$

can be established. Since $y_{n+1} - x_{n+1} = (y_n - x_n)(e - ax_n)$ and $\|e - ax_n\| < 1$, the inequalities

$$\|y_n - x_n\| \leq \|y_0 - x_0\| \quad \text{and} \quad \|y_{n+1} - x_{n+1}\| \leq \|y_0 - x_0\| \|e - ax_0\|^{2^n}$$

hold, and hence the sequence $(y_n - x_n)$ converges to 0 with the R -order 2. Because of $0 \leq a^{-1} - x_n \leq y_n - x_n$, $0 \leq y_n - a^{-1} \leq y_n - x_n$ and the monotony of the norm the proof of the convergence properties is complete.

2.2. For constructing initial elements x_0 and y_0 according to Theorem 1 condition (5) is crucial. Therefore, the following method for solving (2) should be preferred, because no assumption of the type (5) is required,

$$(8) \quad \begin{aligned} m_{n+1} &= m_n + m_n(e - am_n), \\ r_{n+1} &= r_n |e - am_n|. \end{aligned}$$

Here, in addition, R is assumed to be a lattice where the absolute value is defined by $|x| = \sup(x, -x)$ for $x \in R$. With the abbreviations

$$(9) \quad x_n = m_n - r_n, \quad y_n = m_n + r_n$$

from (8) one obtains

$$\begin{aligned} x_{n+1} &= m_n + x_n(e - am_n)^+ - y_n(e - am_n)^-, \\ y_{n+1} &= m_n + y_n(e - am_n)^+ - x_n(e - am_n)^- \end{aligned}$$

where $x^+ = \sup(x, 0)$ and $x^- = \sup(-x, 0)$ for $x \in R$. Now, the following statement is valid.

THEOREM 2. *Assume that*

$$(10) \quad r_0 \geq 0, \quad |m_0 - m_1| \leq r_0 - r_1,$$

$$(11) \quad m_0^{-1} \text{ exists and } \|e - am_0\| < 1.$$

Then a can be inverted and the sequences (m_n) and (r_n) of method (8) provide the monotonous inclusion

$$(12) \quad |m_n - a^{-1}| \leq r_n \leq r_{n-1} \leq \dots \leq r_0$$

or, using (9),

$$(13) \quad x_0 \leq \dots \leq x_{n-1} \leq x_n \leq a^{-1} \leq y_n \leq y_{n-1} \leq \dots \leq y_0.$$

If, in addition,

$$(14) \quad \|e - am_0\| < 1,$$

the sequences (m_n) , (x_n) and (y_n) converge to a^{-1} at least with the R -order 2 and (r_n) converges to 0 at least with the same R -order.

A proof of Theorem 2 is given by J. W. Schmidt [9] and will not be repeated here. For related results with regard to the simplified method

$$(15) \quad \begin{aligned} m_{n+1} &= m_0 + m_n(e - am_0), \\ r_{n+1} &= r_n |e - am_0| \end{aligned}$$

see the earlier paper [8].

Since a^{-1} commutes with a , it is interesting that only the approximations m_n commute with a if m_0 does so.

2.3. In the case where R is the ring of the real (N, N) -matrices, initial matrices can be constructed if m_0 is a sufficiently good approximation to a^{-1} , see [8]. For this, let c be a matrix with positive elements,

$$c_{ik} = (c)_{ik} > 0 \quad \text{for } i, k = 1, \dots, N,$$

and assume

$$(c - c|e - am_0|)_{ik} > 0 \quad \text{for } i, k = 1, \dots, N.$$

Then, as is immediately seen, the conditions (10) are satisfied if

$$(16) \quad r_0 = \xi c \quad \text{with} \quad \xi = \max_{i,k=1,\dots,N} \frac{(|m_0(e - am_0)|)_{ik}}{(c - c|e - am_0|)_{ik}}.$$

Notice that (10) is equivalent to $x_0 \leq y_0$, $x_0 \leq x_1$ and $y_1 \leq y_0$.

2.4. The method (8) is closely related to the interval method of G. Alefeld and J. Herzberger [1]

$$(17) \quad X_{n+1} = X_n \cap \{m(X_n) + X_n(e - am(X_n))\}$$

where X_n is an interval matrix and $m(X_n)$ the corresponding midpoint matrix. Using the interval rules on midpoint and radius matrices, the method without intersection

$$X_{n+1} = m(X_n) + X_n(e - am(X_n))$$

is seen to be a special case of (8). If

$$m(X_\nu) + X_\nu(e - am(X_\nu)) \subset X_\nu \quad \text{for some } \nu,$$

with regard to Theorem 2, the intersection in method (17) can be avoided for all $n \geq \nu$.

2.5. The method (15) can be extended to the case where the element a is not given exactly. Now, it is only assumed that

$$(18) \quad \underline{a} \leq a \leq \bar{a}$$

with known $\underline{a}, \bar{a} \in R$. Let m be an approximation to all corresponding a^{-1} and suppose that elements $\underline{b}, \bar{b}, \underline{h}, \bar{h} \in R$ are available with

$$(19) \quad 0 \leq \underline{b} \leq (e - am)^+ \leq \bar{b}, \quad 0 \leq \underline{h} \leq (e - am)^- \leq \bar{h}$$

for all a with (18). In this situation the method

$$(20) \quad \begin{aligned} x_{n+1} &= m + x_n^+ \underline{b} - x_n^- \bar{b} - y_n^+ \bar{h} + y_n^- \underline{h}, \\ y_{n+1} &= m + y_n^+ \bar{b} - y_n^- \underline{b} - x_n^+ \underline{h} + x_n^- \bar{h}, \end{aligned}$$

which is due to J. W. Schmidt [8], has the following properties.

THEOREM 3. *Assume that*

$$(21) \quad x_0 \leq y_0, \quad x_0 \leq x_1, \quad y_1 \leq y_0,$$

$$(22) \quad m^{-1} \text{ exists and } \|\bar{b}\| + \|\bar{h}\| < 1.$$

Then all a with (18) are invertible and method (20) produces the monotonous enclosing

$$(23) \quad x_0 \leq \dots \leq x_{n-1} \leq x_n \leq a^{-1} \leq y_n \leq y_{n-1} \leq \dots \leq y_0,$$

valid for all elements a with (18).

For a proof of Theorem 3 see [8] where, in addition, the construction (16) of initial matrices is extended to method (20).

2.6. To demonstrate method (20) let

$$a = \begin{bmatrix} 99 & 99 & 98 \\ 99 & 98 & 98 \\ 98 & 98 + \varepsilon & 97 + \delta \end{bmatrix}, \quad \varepsilon, \delta \in [-0.001, +0.001].$$

Set $x^* = \lim x_n$ and $y^* = \lim y_n$. Then, by method (20) one gets the inclusion $x^* \leq a^{-1} \leq y^*$ for all matrices a in question, where now

$$\begin{aligned} x^* &= \begin{bmatrix} -108.77 \dots & 0.89 \dots & 87.23 \dots \\ 1 & -1 & 0 \\ 87.12 \dots & -0.11 \dots & -109.88 \dots \end{bmatrix}, \\ y^* &= \begin{bmatrix} -87.23 \dots & 1.11 \dots & 108.77 \dots \\ 1 & -1 & 0 \\ 108.88 \dots & 0.11 \dots & -88.12 \dots \end{bmatrix}. \end{aligned}$$

The optimal interval enclosing $\bar{x} \leq a^{-1} \leq \bar{y}$ by means of $\bar{x} = \inf \{a^{-1}\}$ and $\bar{y} = \sup \{a^{-1}\}$ is a little sharper because of

$$\bar{x} = \begin{bmatrix} -108.77 \dots & 0.89 \dots & 89.17 \dots \\ 1 & -1 & 0 \\ 89.08 \dots & -0.11 \dots & -109.88 \dots \end{bmatrix},$$

$$\bar{y} = \begin{bmatrix} -89.17 \dots & 1.11 \dots & 108.77 \dots \\ 1 & -1 & 0 \\ 108.88 \dots & 0.11 \dots & -90.08 \dots \end{bmatrix}.$$

However, here the diameter of the optimal interval $\|\bar{y} - \bar{x}\|$ is considerably greater than the loss of sharpness $\max\{\|\bar{x} - x^*\|, \|\bar{y} - y^*\|\}$.

3. Enclosing of square roots

Let a be a real symmetric, nonnegative definite (N, N) -matrix with the elements $a_{ik} = (a)_{ik}$. Then there exists a unique real symmetric, nonnegative definite matrix $a^{1/2}$ such that

$$a^{1/2} a^{1/2} = a,$$

and $a^{1/2}$ is called the square root of a . No finite method for computing $a^{1/2}$ is known.

3.1. In connection with the enclosing of square roots especially two kinds of partial orderings are of interest. The definite ordering denoted by " \leq " and generated by the cone of real symmetric, nonnegative definite matrices seems to be more convenient for theoretical manipulation than the natural ordering " \leq "; the first one is, however, of less numerical utility. Now, by an important theorem of W. Burmeister [2], starting from a given inclusion with respect to the definite ordering, it is always possible to construct an enclosing with respect to the natural ordering. The main step in this direction is the following result.

THEOREM 4. *Let the matrices $r \geq 0$ and m be given. Then, the interval of matrices relative to the definite ordering*

$$(24) \quad I = \{x: m - r \leq x \leq m + r\}$$

is optimally enclosed by an interval with respect to the natural ordering whose bounds are given by

$$(25) \quad (\leq)\text{-inf } I = m - s, \quad (\leq)\text{-sup } I = m + s$$

where

$$(26) \quad s_{ik} = \sqrt{r_{ii}r_{kk}} \quad \text{for } i, k = 1, \dots, N.$$

For a proof of this theorem see [2]. Now, Theorem 4 immediately leads to the announced result, see [2].

THEOREM 5. *The monotone enclosing with respect to the definite ordering*

$$(27) \quad x_{n-1} \leq x_n \leq x^* \leq y_n \leq y_{n-1}$$

implies the monotone inclusion relative to the natural ordering

$$(28) \quad u_{n-1} \leq u_n \leq x^* \leq v_n \leq v_{n-1}$$

where

$$(29) \quad m_n = \frac{1}{2}(x_n + y_n),$$

$$(s_n)_{ik} = \frac{1}{2} \sqrt{(y_n - x_n)_{ii}(y_n - x_n)_{kk}} \quad \text{for } i, k = 1, \dots, N$$

and

$$(30) \quad u_n = m_n - s_n, \quad v_n = m_n + s_n.$$

Proof. Let

$$I_n = \{x: x_n \leq x \leq y_n\} = \{x: m_n - r_n \leq x \leq m_n + r_n\}$$

with $r_n = (y_n - x_n)/2$. Then Theorem 4 gives $(\leq)\text{-inf } I_n = u_n$ and $(\leq)\text{-sup } I_n = v_n$. Since $x^* \in I_n$, it follows that $u_n \leq x^* \leq v_n$, and because of $I_n \subset I_{n-1}$, using the definition of infimum and supremum, $u_{n-1} \leq u_n$ and $v_n \leq v_{n-1}$ are obtained.

3.2. It is well known that for a real number $a > 0$ the method

$$(31) \quad y_{n+1} = \frac{1}{2}(y_n + a/y_n)$$

provides a quadratically convergent sequence (y_n) with the limit \sqrt{a} for every initial value $y_0 > 0$. Furthermore, using

$$(32) \quad x_n = a/y_n,$$

the monotonous enclosing

$$(33) \quad x_1 \leq \dots \leq x_{n-1} \leq x_n \leq \sqrt{a} \leq y_n \leq y_{n-1} \leq \dots \leq y_1$$

holds.

This result can be generalized to the computation of square roots of matrices if the definite ordering is employed, see [2].

THEOREM 6. Let a and y_0 be real symmetric, positive definite matrices with $ay_0 = y_0a$. Then the method

$$(34) \quad x_n = ay_n^{-1}, \quad y_{n+1} = \frac{1}{2}(x_n + y_n)$$

is well-defined, and by the sequences (x_n) and (y_n) the square root $a^{1/2}$ is monotonously enclosed with respect to the definite ordering,

$$(35) \quad x_{n-1} \leq x_n \leq a^{1/2} \leq y_n \leq y_{n-1} \quad \text{for } n \geq 2.$$

In addition, the sequences converge to $a^{1/2}$ at least with the R-order 2.

Theorem 6, inclusive enclosing (35), is also valid for the method without inversion

$$(36) \quad \begin{aligned} x_{n+1} &= x_n + z_n(a - x_n^2), & y_{n+1} &= y_n + z_n(a - y_n^2), \\ z_{n+1} &= z_n + z_n(e - 2z_n y_{n+1}) \end{aligned}$$

if $0 \leq a \leq e$. But now the initial matrices have to be restricted, e.g. to $x_0 = y_0 = e$ and $z_0 = e/2$.

With regard to Theorem 5 the definite enclosing (35) can be transformed into a natural one. The new sequences (u_n) and (v_n) converge to $a^{1/2}$ likewise at least with R -order 2.

3.3. The following numerical example is taken from [2]. Let a be the matrix with the known square root

$$(a^{1/2})_{ik} = (1 + 2|i - k|)^{-1} \quad \text{for } i, k = 1, \dots, N,$$

and let the initial matrix be $y_0 = 5e$. In general, due to the rounding errors the theoretically assured commutativity of x_n and y_n with a and consequently the enclosing (35) are not preserved during a larger number of steps. Therefore the iteration (34) is stopped if the condition

$$(x_n)_{ii} \leq (y_n)_{ii} \quad \text{for } i = 1, \dots, N$$

necessary for $x_n \leq y_n$ is violated. For $N = 30$ the elements in the fifth row and first column are

n	$(x_n)_{51}$	$(y_n)_{51}$	$(u_n)_{51}$	$(v_n)_{51}$
0	0.0893 ...	0.0	—	—
1	0.1279 ...	0.0446 ...	-1.0003 ...	1.1703 ...
2	0.1312 ...	0.0863 ...	-0.3148 ...	0.5324 ...
3	0.1136 ...	0.1087 ...	0.0016 ...	0.2207 ...
4	0.1110 1863 ...	0.1111 9731 ...	0.1000 1492 ...	0.1222 0103 ...
5	0.1111 1425 8912	0.1111 0797 3684	0.1109 4305 9024	0.1112 7917 3572
6	0.1111 1110 5937	0.1111 1111 6299	0.1111 1105 0143	0.1111 1117 2092

Indeed, the elements of x_n and y_n give only an approximate information about the corresponding element of $a^{1/2}$. An inclusion, however, is provided by the elements of u_n and v_n .

4. Enclosing of Cholesky factors

A lower-triangular matrix l with positive diagonal elements is called the *Cholesky factor* of a real symmetric, positive definite matrix a if

$$ll^T = a.$$

The Cholesky factor uniquely exists and can be obtained by the Cholesky method. Here the aim is to give an error estimation for the computed factor by means of a monotonous enclosing with respect to the natural ordering, at first for Stieltjes matrices and then for arbitrary matrices.

4.1. A matrix $a = (a_{ik})$ is defined to be an M -matrix if $a_{ik} \leq 0$ for $i \neq k$ and $a^{-1} \geq 0$ relative to the natural ordering. A Stieltjes matrix is a real symmetric M -matrix. Stieltjes matrices a are positive definite and the corresponding Cholesky factors l are lower-triangular M -matrices uniquely characterized by

$$l_{ii} > 0, \quad l_{ik} \leq 0 \quad \text{for } k < i \quad (i, k = 1, \dots, N).$$

Applying the Newton type method

$$(37) \quad \begin{aligned} F(y_n) + F'(y_n)(y_{n+1} - y_n) &= 0, \\ F(x_n) + F'(y_n)(x_{n+1} - x_n) &= 0 \end{aligned}$$

to $F(x) = xx^T - a = 0$ the following method is obtained:

$$(38) \quad \begin{aligned} y_n(y_{n+1} - y_n)^T + (y_{n+1} - y_n)y_n^T &= a - y_n y_n^T, \\ y_n(x_{n+1} - x_n)^T + (x_{n+1} - x_n)y_n^T &= a - x_n x_n^T. \end{aligned}$$

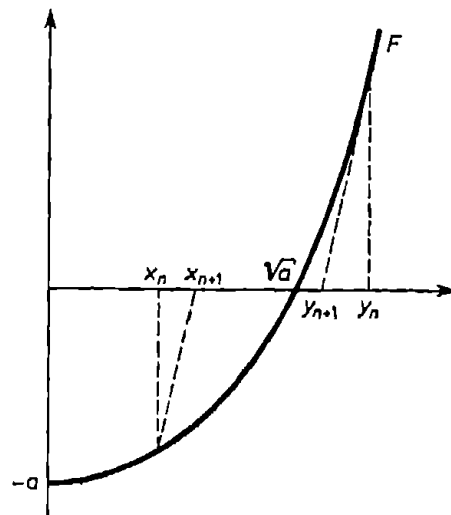


Fig. 2. Method (38) for single-valued functions (iii)

From these linear equations the elements of the lower-triangular matrices y_{n+1} and x_{n+1} can easily be determined, e.g. row by row if $(y_n)_{ii} \neq 0$ for $i = 1, \dots, N$.

The following result about method (38) is due to J. W. Schmidt and U. Patzke [13].

THEOREM 7. *Let a be a Stieltjes matrix and x_0, y_0 be lower-triangular M -matrices such that*

$$(39) \quad x_0 x_0^T \leq a \leq y_0 y_0^T.$$

Then, the sequences (x_n) and (y_n) according to (38) are well-defined. They converge to the Cholesky factor l of a at least with the R -order 2 and

$$(40) \quad x_0 \leq \dots \leq x_{n-1} \leq x_n \leq l \leq y_n \leq y_{n-1} \leq \dots \leq y_0$$

holds.

The proof [13] is based on three lemmas.

LEMMA 1. *Let x and y be lower-triangular M -matrices. Then*

$$xx^T \leq yy^T \text{ implies } x \leq y.$$

LEMMA 2. *Let x be a lower-triangular M -matrix and y be a lower-triangular matrix. Then*

$$xy^T + yx^T \geq 0 \text{ implies } y \geq 0.$$

LEMMA 3. *Let x be a lower-triangular M -matrix, y be a lower-triangular matrix and z be a symmetric matrix. Then*

$$xy^T + yx^T \leq z, y \geq 0 \text{ implies } \|y\| \leq \alpha \|z\|,$$

where α is independent of y and z . The norm is defined by

$$\|z\| = \max \{\|z\|_\infty, \|z^T\|_\infty\}$$

where $\|\cdot\|_\infty$ denotes the row sum norm.

Now, a sketch of a proof of Theorem 7 shall be given. At first the enclosing (40) together with $x_n x_n^T \leq a \leq y_n y_n^T$ is verified by induction.

The assumption (39) in view of Lemma 1 implies the inclusion $x_0 \leq l \leq y_0$.

In the step from n to $n+1$, because of $l \leq y_n \leq y_0$ the lower-triangular matrix l_n is an M -matrix. Therefore, the two relations (38) and

$$y_n(y_{n+1}-l)^T + (y_{n+1}-l)y_n^T = (y_n-l)(y_n-l)^T \geq 0,$$

$$y_n(x_{n+1}-l)^T + (x_{n+1}-l)y_n^T = (y_n-x_n)(x_n-l)^T + (x_n-l)(y_n-l)^T \leq 0$$

by means of Lemma 2 yield the inequalities $y_{n+1} \leq y_n$, $x_{n+1} \geq x_n$, $y_{n+1} \geq l$, $x_{n+1} \leq l$. Finally, the remaining inequalities $y_{n+1}y_{n+1}^T - a \geq 0$, $x_{n+1}x_{n+1}^T - a \leq 0$ are proved by simple estimations. Therefore (40) holds, and because the Cholesky factor is unique, $\lim x_n = \lim y_n = l$ is true. Using

$$\begin{aligned} & l(y_{n+1}-x_{n+1})^T + (y_{n+1}-x_{n+1})l^T \\ & \leq y_n(y_{n+1}-x_{n+1})^T + (y_{n+1}-x_{n+1})y_n^T = (y_n-x_n)(y_n-x_n)^T, \end{aligned}$$

Lemma 3 ensures that

$$\|y_{n+1} - x_{n+1}\| \leq \alpha \|y_n - x_n\|^2.$$

Hence the sequences $(y_n - x_n)$ and consequently (x_n) and (y_n) converge at least with the R -order 2.

4.2. The matrix

$$a = \begin{bmatrix} 2 & & -1 & & & \\ -1 & & 2 & & -1 & \\ & \ddots & & \ddots & & \\ & & & & & \ddots \end{bmatrix} \} N \text{ rows}$$

is a Stieltjes matrix, and $\text{cond}(a) \rightarrow \infty$ for $N \rightarrow \infty$.

Let

$$y_0 = \begin{bmatrix} \alpha & & & \\ \beta & & a & \\ & \ddots & & \ddots \end{bmatrix}, \quad x_0 = \begin{bmatrix} \gamma & & & \\ \delta & & \gamma & \\ & \ddots & & \ddots \end{bmatrix}.$$

Then $x_0 x_0^T \leq a \leq y_0 y_0^T$ holds if $\alpha^2 \geq 2$, $\alpha\beta \geq -1$, $\alpha > 0$, $\beta \leq 0$ and $\gamma^2 + \delta^2 \leq 2$, $\gamma\delta \leq -1$, $\gamma > 0$, $\delta \leq 0$. For $N = 8$ and for $\alpha = 1.42$, $\beta = -0.7$, $\gamma = 1$ and $\delta = -1$ the element of l , e.g. in the eighth row and column, is enclosed by method (38) as follows, see [13]:

n	$(x_n)_{88}$	$(y_n)_{88}$
0	1	1.42
1	1.001 ...	1.181 ...
2	1.020 ...	1.091 ...
3	1.052 ...	1.062 ...
4	1.060 5019 ...	1.060 6654 ...
5	1.060 6601 3921	1.060 6601 7180
6	1.060 6601 7178	1.060 6601 7178

A more general method for the construction of starting matrices according to (39) is described in [13]. There it is assumed that a sufficiently good approximation to l is known.

4.3. Let now a be an arbitrary real symmetric, positive definite (N, N) -matrix. The first equation of the iteration (38)

$$y_n y_{n+1}^T + y_{n+1} y_n^T = a + y_n y_n^T$$

leads to the following method of successive displacements approximating the Cholesky factor l of a :

$$(41) \quad \begin{aligned} (y_{n+1})_{kk} &= \frac{1}{2} \left\{ (y_n)_{kk} + \frac{a_{kk} - \sum_{\nu=1}^{k-1} (y_{n+1})_{k\nu}^2}{(y_n)_{kk}} \right\}, \\ (y_{n+1})_{ik} &= \frac{1}{(y_{n+1})_{kk}} \left\{ a_{ik} - \sum_{\nu=1}^{k-1} (y_{n+1})_{i\nu} (y_{n+1})_{k\nu} \right\} \quad \text{for } i > k. \end{aligned}$$

Here the elements $(y_{n+1})_{ik}$ of the lower-triangular matrix y_{n+1} are computed, e.g. column by column, $k = 1, \dots, N$.

U. Patzke [6] modified method (41) under the assumption that a starting matrix y_0 is known with

$$(42) \quad \text{sign}((y_0)_{ik}) = \text{sign}((l)_{ik}) \quad \text{for } i, k = 1, \dots, N, \quad i > k$$

in order to guarantee monotonous enclosing:

$$(43) \quad \text{For } k = 1, \dots, N \text{ do}$$

$$\begin{aligned} (y_{n+1})_{kk} &= \frac{1}{2(y_n)_{kk}} \left\{ a_{kk} - \sum_{\nu=1}^{k-1} (\mu_{k\nu}^{n+1})^2 + (y_n)_{kk}^2 \right\}, \\ (x_{n+1})_{kk} &= \frac{1}{(y_{n+1})_{kk}} \left\{ a_{kk} - \sum_{\nu=1}^{k-1} (\lambda_{k\nu}^{n+1})^2 \right\}, \\ (y_{n+1})_{ik} &= \frac{1}{a_{ik}^{n+1}} \left\{ a_{ik} - \sum_{\nu=1}^{k-1} \varrho_{i\nu k}^{n+1} \varrho_{k\nu i}^{n+1} \right\} \quad \text{for } i > k, \\ (x_{n+1})_{ik} &= \frac{1}{\beta_{ik}^{n+1}} \left\{ a_{ik} - \sum_{\nu=1}^{k-1} \sigma_{i\nu k}^{n+1} \sigma_{k\nu i}^{n+1} \right\} \quad \text{for } i > k, \end{aligned}$$

where

$$\begin{aligned} \mu_{k\nu}^{n+1} &= (x_{n+1})_{k\nu}, & \lambda_{k\nu}^{n+1} &= (y_{n+1})_{k\nu} & \text{for } \text{sign}((l)_{k\nu}) \geq 0, \\ \mu_{k\nu}^{n+1} &= (y_{n+1})_{k\nu}, & \lambda_{k\nu}^{n+1} &= (x_{n+1})_{k\nu} & \text{for } \text{sign}((l)_{k\nu}) < 0 \end{aligned}$$

and

$$\begin{aligned} \alpha_{ik}^{n+1} &= (x_{n+1})_{kk}, & \beta_{ik}^{n+1} &= (y_{n+1})_{kk} & \text{for } \text{sign}((l)_{ik}) \geq 0, \\ \alpha_{ik}^{n+1} &= (y_{n+1})_{kk}, & \beta_{ik}^{n+1} &= (x_{n+1})_{kk} & \text{for } \text{sign}((l)_{ik}) < 0 \end{aligned}$$

and

$$\begin{aligned} \varrho_{ivk}^{n+1} &= (x_{n+1})_{iv}, & \sigma_{ivk}^{n+1} &= (y_{n+1})_{iv} & \text{for } \text{sign}((l)_{kv}) \geq 0, \\ \varrho_{ivk}^{n+1} &= (y_{n+1})_{iv}, & \sigma_{ivk}^{n+1} &= (x_{n+1})_{iv} & \text{for } \text{sign}((l)_{kv}) < 0. \end{aligned}$$

The following theorem concerning method (43) is proved in [6].

THEOREM 8. *Let a be real symmetric, positive definite and let y_0 be lower-triangular with the property (42). If*

$$(44) \quad \begin{aligned} (x_1)_{kk} &> 0 & \text{for } k = 1, \dots, N, \\ (x_1)_{ik} &\geq 0 & \text{for } i > k \text{ and } \text{sign}((l)_{ik}) \geq 0, \\ (y_1)_{ik} &\leq 0 & \text{for } i > k \text{ and } \text{sign}((l)_{ik}) \leq 0, \end{aligned}$$

then the sequences (x_n) and (y_n) are well-defined by (43). They converge to the Cholesky factor l of a at least with the R -order 2 and l is enclosed monotonously,

$$(45) \quad x_1 \leq \dots \leq x_{n-1} \leq x_n \leq l \leq y_n \leq y_{n-1} \leq \dots \leq y_1.$$

To show the R -order 2, at first,

$$\|y_{n+1} - x_{n+1}\|_F \leq \alpha \|y_n - l\|_F^2$$

with the Frobenius norm $\|\cdot\|_F$ is verified. Then, using $x_{n+1} \leq l \leq y_{n+1}$, it follows immediately that

$$\|x_{n+1} - l\|_F \leq \alpha \|y_n - l\|_F^2,$$

$$\|y_{n+1} - l\|_F \leq \alpha \|y_n - l\|_F^2.$$

These estimations imply the R -order 2 of both sequences (x_n) and (y_n) , see e.g. [11].

4.4. Let

$$a = \begin{bmatrix} 25 & 5 & -5 & 5 \\ 5 & 17 & 3 & -3 \\ -5 & 3 & 18 & -6 \\ 5 & -3 & -6 & 28 \end{bmatrix}, \quad y_0 = \begin{bmatrix} 3 & 0 & 0 & 0 \\ + & 3 & 0 & 0 \\ - & + & 3 & 0 \\ + & - & - & 3 \end{bmatrix}.$$

Then by the method (43) e.g. for the element in the fourth row and column the following approximations are obtained [6]:

n	$(x_n)_{44}$	$(y_n)_{44}$
1	4.1294 ...	5.7947 ...
2	4.9321 ...	5.0584 ...
3	4.9996 ...	5.0003 ...
4	4.9999 9998	5.0000 0001 ...
5	5.0000 0000 000	5.0000 0000 000

References

- [1] G. Alefeld and J. Herzberger, *Einführung in die Intervallrechnung*, Bibliographisches Institut, Mannheim-Wien-Zürich 1974.
- [2] W. Burmeister, *Optimal interval enclosing of certain sets of matrices with application to monotone enclosing of square roots*, Computing 25 (1980), 283-295.
- [3] N. S. Kurpel and T. S. Kurchenko, *Two-sided methods for the solution of systems of equations* (in Russian), Naukova dumka, Kiev 1975.
- [4] W. Mönch, *Monotone Einschließung von positiven Inversen*, Z. Angew. Math. Mech. 53 (1973), 207-208.
- [5] J. M. Ortega and W. C. Rheinboldt, *Iterative solution of nonlinear equations in several variables*, Academic Press, New York and London 1970.
- [6] U. Patzke, *Ein monoton einschließendes Einzelschrittverfahren für Cholesky-Faktoren*, Beitr. Numer. Math. 11 (1983), 139-146.
- [7] J. W. Schmidt, *Monotone Einschließung von Inversen positiver Elemente durch Verfahren vom Schulz-Typ*, Computing 16 (1976), 211-216.
- [8] —, *Einschließung inverser Elemente durch Fixpunktverfahren*, Numer. Math. 31 (1978), 313-320.
- [9] —, *Monotone Eingrenzung von inversen Elementen durch ein quadratisch konvergentes Verfahren ohne Durchschnittsbildung*, Z. Angew. Math. Mech. 60 (1980), 202-204.
- [10] —, *Monotone Einschließungsverfahren bei nichtlinearen Gleichungen*, Mitteilungen Math. Gesellschaft, Heft 1/2 (1982), 62-87.
- [11] —, *On the R-order of coupled sequences*, Computing 26 (1981), 333-342.
- [12] J. W. Schmidt and H. Leonhardt, *Eingrenzung von Lösungen mit Hilfe der Regula falsi*, Computing 6 (1970), 318-329.
- [13] J. W. Schmidt and U. Patzke, *Iterative Nachverbesserung mit Fehlerengrenzung der Cholesky-Faktoren von Stieltjes-Matrizen*, J. Reine Angew. Math. 327 (1981), 81-92.

Presented to the Semester
Computational Mathematics
February 20 — May 30, 1980