

SLOBODAN LAKIĆ

**TWO ITERATIVE METHODS FOR THE MATRIX
INVERSE SQUARE ROOT**

Abstract: This paper presents two fast iterative methods for finding the matrix inverse square root. The two iterative methods with high convergence rates are analyzed and their numerical stability properties are investigated. Numerical examples show that this methods can be superior then some earlier methods.

Subject Classification 65F30

Key words: matrix function, symmetric orthogonalization.

1. Introduction

Given a nonsingular matrix $A \times C^{n,n}$, a matrix X such that

$$A X^2 = I$$

is called an inverse square root of A and is denoted by $X = A^{-1/2}$. The inverse square root of a matrix has applications in the computation of an optimal symmetric orthogonalization of a set of vectors [6], theory of oscillations [5], etc. It is known that if the matrix A is real, then a real inverse square root which is a function of A exists if A has no negative real eigenvalues [4]. That X is a function of A means it is a polynomial in A , [5]. Numerical methods for finding the inverse square root have been proposed in many papers, for example [1], [2], [3], etc. In section 2 we give an iterative method to compute $A^{-1/2}$. This method we will show has $(2k+1)$ th - order convergence rate and improves the rate of convergence than the methods in [1], [2], [3]. Under some restriction on the spectrum of A this method is locally stable. In section 3, we proposed an alternative locally stable method to compute $A^{-1/2}$. In section 4, we illustrate the performance of the methods by numerical examples.

2. Comutation of $A^{-1/2}$

Lemma. Let w be a complex number such that $w \neq 0$ and $\arg(w) \neq \pi$. For some pre-chosen natural number k , we define the sequence $(z_{n,k})$ by

$$(2.1) \quad \begin{cases} z_{n,k} = 1, \\ z_{n+1,k} = z_{n,k} \frac{\sum_{m=0}^k \alpha_{m,k} (wz_{n,k}^2)^m}{\sum_{m=0}^k \beta_{m,k} (wz_{n,k}^2)^m}, \end{cases}$$

where $\alpha_{m,k} = \binom{2k+1}{2m+1}$, $\beta_{m,k} = \binom{2k+1}{2m}$.

Then $\lim_{n \rightarrow \infty} z_{n,k} = \frac{1}{\sqrt{w}}$, where \sqrt{w} is the principal square root of w .

(Without loss of generality, we will denote $z_{n,k}$, $\alpha_{m,k}$, $\beta_{m,k}$ by z_n , α_m , β_m respectively.)

Proof. First we prove that the sequence (z_n) is well defined. For some pre-chosen $k \in \mathbb{N}$ we define the function

$$h_k(z) = \frac{(1+z)^{2k+1} - (1-z)^{2k+1}}{(1+z)^{2k+1} + (1-z)^{2k+1}}, \quad z \in \mathbb{C}.$$

Let $S = C_+ \cup C_-$, where

$$C_+ = \{z \in \mathbb{C} : \operatorname{Re}(z) > 0\}, \quad C_- = \{z \in \mathbb{C} : \operatorname{Re}(z) < 0\}.$$

If we suppose that $(1+z)^{2k+1} + (1-z)^{2k+1} = 0$, then $z \neq 1$ and $(1+z)/(1-z) = e^{i\phi}$ where ϕ is not an odd multiple of π , else $z = \infty$. Now we have $z = \frac{i \sin \phi}{1 + \cos \phi}$ and the poles of h_k lie on the imaginary axis. Now

we prove that $h_k(S) \subseteq S$. Let $z \in S$ and we suppose that

$$(2.2) \quad h_k(z) = ai, \quad a \in \mathbb{R}$$

From (2.2) it follows

$$(2.3) \quad \frac{(1+ai)}{(1-ai)} = \left(\frac{(1+z)}{(1-z)} \right)^{2k+1}$$

Since $z \in S$ hence $|(1+z)/(1-z)| \neq 1$ and $|(1+z)/(1-z)|^{2k+1} \neq 1$ which is a contradiction with (2.3) because $|(1+ai)/(1-ai)| = 1$. So, $h_k(S) \subseteq S$. The function h_k is continuous on the set S and hence either $h_k(C_+) \subseteq C_+$ or $h_k(C_+) \subseteq C_-$ because the continuous image of a connected set is connected.

But $1 \in C_+$ and $h_k(1) = 1 \in C_+$, so $h_k(C_+) \subseteq C_+$. Thus if $s_0 \in C_+$ then $s_{n+1} = h_k(s_n) \in C_+$. Using (2.1) one can prove that

$$(2.4) \quad z_{n+1} \sqrt{w} = \frac{(1+z_n \sqrt{w})^{2k+1} - (1-z_n \sqrt{w})^{2k+1}}{(1+z_n \sqrt{w})^{2k+1} + (1-z_n \sqrt{w})^{2k+1}}$$

Let $s_n = z_n \sqrt{w}$. Now $s_0 = \sqrt{w} \in C_+$ because $\arg(w) \neq \pi$. Consequently $s_n \in C_+$. If we suppose that for some n it holds $(1+s_n)^{2k+1} + (1-s_n)^{2k+1} = 0$, then s_n lie on the imaginary axis which is a contradiction. So for each n it holds

$$(1+z_n \sqrt{w})^{2k+1} + (1-z_n \sqrt{w})^{2k+1} \neq 0$$

and the sequence (z_n) is well defined.

It is easy to show that

$$(2.5)$$

$$\frac{z_{n+1} \pm 1}{\sqrt{w}} = \frac{2}{\sqrt{w} \left(\left(z_n + \frac{1}{\sqrt{w}} \right)^{2k+1} + \left(z_n - \frac{1}{\sqrt{w}} \right)^{2k+1} \right)} \left(z_n \pm \frac{1}{\sqrt{w}} \right)^{2k+1}$$

$$(2.6) \quad \frac{z_n - \frac{1}{\sqrt{w}}}{z_n + \frac{1}{\sqrt{w}}} = \left(\frac{1 - \sqrt{w}}{1 + \sqrt{w}} \right)^{(2k+1)^n}$$

Consequently it follows that

$$(2.7) \quad z_n - \frac{1}{\sqrt{w}} = \frac{2d_n}{\sqrt{w}(1-d_n)}$$

where

$$d_n = \left(\frac{1 - \sqrt{w}}{1 + \sqrt{w}} \right)^{(2k+1)^n}$$

Since \sqrt{w} is the principal square root of w , hence $\operatorname{Re} \sqrt{w} > 0$ and $|(1 - \sqrt{w}) / (1 + \sqrt{w})| < 1$.

So we have

$$\lim_{n \rightarrow \infty} d_n = 0 \quad \text{and consequently} \quad \lim_{n \rightarrow \infty} z_n = \frac{1}{\sqrt{w}}.$$

Now we define a matrix sequence (X_n) to compute $A^{-1/2}$.

$$(1) \quad \begin{cases} X_0 = 1 \\ X_{n+1} = X_n \left(\sum_{m=0}^k \alpha_m (AX_n^2)^m \right) \left(\sum_{m=0}^k \beta_m (AX_n^2)^m \right)^{-1} \end{cases}$$

where α_m , β_m and k are as in Lemma.

For our analysis we assume that A is diagonalizable, that is there exists a nonsingular matrix V such that

$$V^{-1}AV = D$$

where $D = \operatorname{diag}(a_1, \dots, a_n)$ and a_1, \dots, a_n are the eigenvalues of A .

Theorem 2.1. Let $A \in C^{n,n}$ be nonsingular and diagonalizable and assume A has no negative real eigenvalues. Then the matrix sequence (1) converges to X , where X is the principal inverse square root of A and

$$(2.9) \quad \|X_{n+1} - X\| = 0 \quad (\|X_n - X\|^{2k+1}).$$

Proof. We define $L_n = V^{-1} X_n V$ where V is as in (2.8). From iteration (1) we have

$$(2.10) \quad \begin{cases} L_0 = 1 \\ L_{n+1} = L_n \left(\sum_{m=0}^k \alpha_m (DL_n^2)^m \right) \left(\sum_{m=0}^k \beta_m (DL_n^2)^m \right)^{-1} \end{cases}$$

The sequence L_i is a sequence of diagonal matrices $L_i = \text{diag}(l_1^{(i)}, \dots, l_n^{(i)})$. The equation (2.10) is equivalent to n scalar equations

$$l_j^{(0)} = 1, \quad l_j^{(n+1)} = l_j^{(n)} \frac{\sum_{m=0}^k \alpha_m (a_j (l_j^{(n)})^2)^m}{\sum_{m=0}^k \beta_m (a_j (l_j^{(n)})^2)^m}, \quad j = 1, \dots, n.$$

Application of lemma yields

$$\lim_{n \rightarrow \infty} L_n = D^{-1/2} = \text{diag} \left(\frac{1}{\sqrt{a_1}}, \dots, \frac{1}{\sqrt{a_n}} \right).$$

Consequently $\lim_{n \rightarrow \infty} X_n = A^{-1/2}$.

Now, $X_{n+1} - X = 2X((X + X_n)^{2k+1} + (X - X_n)^{2k+1})^{-1}(X_n - X)^{2k+1}$

follows upon using equation (2.5).

From the above equation by taking norms we have (2.9).

2. Stability Analysis

We now investigate the numerical local stability of the method (I). By local stability we mean that in a neighbourhood of the solution $A^{-1/2}$, a perturbation due to rounding error will not be magnified in the succeeding steps.

Let a rounding error $E_n = O(\varepsilon)$ be introduced in the n -th step. Denote by \tilde{X}_n the computed solution at this step, that is $\tilde{X}_n = X_n + E_n$. Using the perturbation result in [7]

$$(2.11) \quad (A + B)^{-1} = A^{-1} - A^{-1}BA^{-1} + O(\|B\|^2)$$

we have

$$\tilde{X}_{n+1} = \tilde{X}_n \left(\sum_{m=0}^k \alpha_m (A\tilde{X}_n^2)^m \right) \left(\sum_{m=0}^k \beta_m (A\tilde{X}_n^2)^m \right)^{-1},$$

$$\begin{aligned} \tilde{X}_{n+1} &= (X_n + E_n) \times \\ &\times \left(\sum_{m=0}^k \alpha_m (AX_n^2)^m + \sum_{m=0}^k \alpha_m \sum_{i=0}^{m-1} (AX_n^2)^i A(X_n E_n + E_n X_n) (AX_n^2)^{m-i-1} \right) \times \\ &\times \left(\sum_{m=0}^k \beta_m (AX_n^2)^m + \sum_{m=0}^k \beta_m \sum_{i=0}^{m-1} (AX_n^2)^i A(X_n E_n + E_n X_n) (AX_n^2)^{m-i-1} \right)^{-1} + \\ &+ O(\varepsilon^2) \end{aligned}$$

$$\begin{aligned} \tilde{X}_{n+1} &= (X_n + E_n) \times \\ &\times \left(\sum_{m=0}^k \alpha_m (AX_n^2)^m + \sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} (AX_n^2)^i A(X_n E_n + E_n X_n) (AX_n^2)^{m-i-1} \right) \times \\ &\times \left(\left[\sum_{m=0}^k \beta_m (AX_n^2)^m \right]^{-1} - \left[\sum_{m=0}^k \beta_m (AX_n^2)^m \right]^{-1} \left[\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} (AX_n^2)^i \times \right. \right. \\ &\times \left. \left. A(X_n E_n + E_n X_n) (AX_n^2)^{m-i-1} \right] \left[\sum_{m=0}^k \beta_m (AX_n^2)^m \right]^{-1} \right) + O(\varepsilon^2). \end{aligned}$$

We define $\tilde{E}_n = V^{-1}E_n V$.

Direct calculations give

$$\tilde{E}_{n+1} = L_n \left(\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} (DL_n^2)^i D(L_n \tilde{E}_n + \tilde{E}_n L_n) (DL_n^2)^{m-i-1} \right) \times$$

$$\begin{aligned} & \times \left(\left[\sum_{m=0}^k \beta_m (DL_n^2)^m \right]^{-1} - \left[\sum_{m=0}^k \alpha_m (DL_n^2)^m \right] \left[\sum_{m=0}^k \beta_m (DL_n^2)^m \right]^{-1} \right. \\ & \times \left. \left[\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} (DL_n^2)^i A(L_n \tilde{E}_n + \tilde{E}_n L_n) (DL_n^2)^{m-i-1} \left[\sum_{m=0}^k \beta_m (DL_n^2)^m \right]^{-1} \right] \right) \\ & \times \tilde{E}_n \left[\sum_{m=0}^k \alpha_m (DL_n^2)^m \right] \left[\sum_{m=0}^k \beta_m (DL_n^2)^m \right]^{-1} + O(\varepsilon^2). \end{aligned}$$

Since $I_j^{(n)} = \varepsilon + \frac{1}{\sqrt{\alpha_j}}$, $\alpha_m = \beta_{k-m}$ and $\sum_{m=0}^k \alpha_m = \sum_{m=0}^k \beta_m$ we obtain

$$\begin{aligned} \tilde{E}_{n+1} &= \left(1 - \frac{\sum_{m=1}^k m(\beta_m - \beta_{k-m})}{\sum_{m=0}^k \beta_m} \right) \tilde{E}_n - \\ & - \frac{\sum_{m=1}^k m(\beta_m - \beta_{k-m})}{\sum_{m=0}^k \beta_m} D^{\frac{1}{2}} \tilde{E}_n D^{\frac{1}{2}} + O(\varepsilon^2). \end{aligned}$$

Now we show that

$$(2.12) \quad \frac{\sum_{m=1}^k m(\beta_m - \beta_{k-m})}{\sum_{m=0}^k \beta_m} = \frac{1}{2}$$

The equality (2.12) is equivalent to

$$\sum_{i=0}^k (2k+1) \binom{2k+1}{2i} = \sum_{i=1}^k 4i \binom{2k+1}{2i}.$$

Using the facts that

$$\sum_{i=0}^k \binom{2k+1}{2i} = 2^{2k}, \quad \sum_{i=1}^k \binom{2k}{2i-1} = 2^{2k-1},$$

$$2i \binom{2k+1}{2i} = (2k+1) \binom{2k}{2i-1} \quad \text{we have}$$

$$\sum_{i=0}^k (2k+1) \binom{2k+1}{2i} = (2k+1) 2^{2k}$$

$$\sum_{i=1}^k 4i \binom{2k+1}{2i} = 2 \sum_{i=1}^k (2k+1) \binom{2k}{2i-1} = (2k+1) 2^{2k}.$$

So, it holds (2.12).

Finally,

$$\tilde{E}_{n+1} = \frac{1}{2} (\tilde{E}_n - D^{\frac{1}{2}} \tilde{E}_n D^{-\frac{1}{2}}) + O(\varepsilon^2),$$

$$E_{n+1} = \frac{1}{2} (E_n - A^{\frac{1}{2}} E_n A^{-\frac{1}{2}}) + O(\varepsilon^2).$$

Assuming that n is large enough and that ε is small enough, we have

$$E_{n+1} \approx \frac{1}{2} (E_n - A^{\frac{1}{2}} E_n A^{-\frac{1}{2}}).$$

The above equality is equivalent to

$$\text{Vec } E_{n+1} \approx \text{Vec } E_n \quad \text{where} \quad G \in C^{n^2, n^2},$$

$$G = I - A^{\frac{1}{2}} * A^{-\frac{1}{2}} \quad (A * B \text{ denote the Kronecker product}).$$

Hence for $j = 1, 2, \dots$, $\text{Vec } E_{n+j} \approx G^j \text{Vec } E_n$.

Consequently the errors in the succeeding steps will not be magnified if $\rho(G) < 1$. The spectrum of G , $\sigma(G)$ can be shown to be [5]

$$\sigma(G) = \left(g_{i,j} : g_{i,j} = 0.5 \left(1 - \sqrt{\frac{a_i}{a_j}} \right), 1 \leq i, j \leq n \right).$$

So $\rho(G) < 1$ is equivalent to

$$(2.13) \quad \left| 1 - \sqrt{\frac{a_i}{a_j}} \right| < 2 .$$

The above inequality gives a sufficient condition for the method (I) to be locally stable.

3. An Alternative Method

We propose an alternative method to compute $A^{-1/2}$. This method will be shown to be locally stable. Let $A \in C^{n,n}$ be nonsingular, we define the sequences (T_n) and (S_n) as follows

$$(II) \quad \begin{cases} T_{n+1} = T_n \left(\sum_{m=0}^k \alpha_m S_n^m \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1}, & T_0 = I \\ S_{n+1} = S_n \left(\left(\sum_{m=0}^k \alpha_m S_n^m \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \right)^2, & S_0 = A \end{cases}$$

where α_m , β_m and k are as in (2.1).

Theorem 3.1. Let $A \in C^{n,n}$ be nonsingular and diagonalizable. Assume that A has no negative real eigenvalues. Then $\lim_{n \rightarrow \infty} T_n = A^{-1/2}$, $\lim_{n \rightarrow \infty} S_n = I$, where $A^{-1/2}$ is the principal inverse square root of A . Further

$$\|T_n - A^{-1/2}\| = O(\|T_n - A^{-1/2}\|^{2k+1}).$$

Proof. Let

$$(3.1) \quad L_n = V^{-1} T_n V, \quad H_n = V^{-1} S_n V$$

Now (II) reduces to

$$(3.2) \quad \begin{cases} L_{n+1} = L_n \left(\sum_{m=0}^k \alpha_m H_n^m \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1}, & L_0 = I \\ H_{n+1} = H_n \left(\left(\sum_{m=0}^k \alpha_m H_n^m \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \right)^2, & H_0 = D \end{cases}$$

From equations (3.2) it follows that L_j and H_j are diagonal matrices. Let

$$(3.3) \quad L_i = \text{diag} (l_i^{(j)}, \dots, l_n^{(j)}), \quad H_j = \text{diag} (h_i^{(j)}, \dots, h_n^{(j)})$$

Equation (3.2) is equivalent to n scalar sequences

$$(3.4) \quad \begin{cases} l_i^{(0)} = 1, & h_i^{(0)} = a_i, \\ l_i^{(n+1)} = l_i^{(n)} \frac{\sum_{m=0}^k \alpha_m (h_i^{(n)})^m}{\sum_{m=0}^k \beta_m (h_i^{(n)})^m}, & h_i^{(n+1)} = h_i^{(n)} \left(\frac{\sum_{m=0}^k \alpha_m (h_i^{(n)})^m}{\sum_{m=0}^k \beta_m (h_i^{(n)})^m} \right)^2 \end{cases}$$

Similarly as in Lemma One can show that

$$(3.5) \quad l_i^{(n+1)} - \frac{1}{\sqrt{a_i}} = \frac{2 \left(l_i^{(n)} - \frac{1}{\sqrt{a_i}} \right)^{2k+1}}{\sqrt{a_i} \left(l_i^{(n)} + \frac{1}{\sqrt{a_i}} \right)^{2k+1} + \left(-l_i^{(n)} + \frac{1}{\sqrt{a_i}} \right)^{2k+1}}$$

and

$$(3.6) \quad \lim_{n \rightarrow \infty} l_i^{(n)} = \frac{1}{\sqrt{a_i}}$$

Since the sequences $(l_i^{(n)})$ converges and it holds (3.4), hence

$$\lim_{n \rightarrow \infty} h_i^{(n)} = 1.$$

From (3.1), (3.3), (3.5), (3.6) and (3.7) it follows

$$\lim_{n \rightarrow \infty} T_n = A^{-1/2}, \quad \lim_{n \rightarrow \infty} S_n = 1 \quad \text{and}$$

$$T_{n+1} - A^{-1/2} = 2A^{-1/2}((A^{-1/2} + T_n)^{2k+1} + (A^{-1/2} - T_n)^{2k+1})^{-1}(T_n - A^{-1/2})^{2k+1}.$$

Taking the norm of the above equation, the bound in the theorem is established.

Stability Analysis

Assume that at the n -th step errors P_n and Q_n are introduced in T_n and S_n respectively. Let \tilde{T}_n and \tilde{S}_n be the computed matrices of this step. Now $\tilde{T}_n = T_n + P_n$, $\tilde{S}_n = S_n + Q_n$ where $P_n = O(\varepsilon)$ and $Q_n = O(\varepsilon)$.

We define computations give

$$\begin{aligned} P_{n+1} &= T_n \left(\left(\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} S_n^i Q_n S_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m S_n^m \right) - \left(\sum_{m=0}^k \alpha_m S_n^m \right) \right. \\ &\quad \times \left. \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \left(\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} S_n^i Q_n S_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \right) \\ &\quad + P_n \left(\sum_{m=0}^k \alpha_m S_n^m \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} + O(\varepsilon^2), \\ \tilde{P}_{n+1} &= L_n \left(\left(\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} H_n^i \tilde{Q}_n H_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} - \left(\sum_{m=0}^k \alpha_m H_n^m \right) \right. \\ (3.8) \quad &\quad \times \left. \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \left(\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} H_n^i \tilde{Q}_n H_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \right) \\ &\quad + \tilde{P}_n \left(\sum_{m=0}^k \alpha_m H_n^m \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} + O(\varepsilon^2) \\ Q_{n+1} &= S_n \left(\left(\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} S_n^i Q_n S_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \left(\sum_{m=0}^k \alpha_m S_n^m \right) \right. \end{aligned}$$

$$\begin{aligned}
& \times \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} - \left(\sum_{m=0}^k \alpha_m S_n^m \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \left(\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} S_n^i Q_n S_n^{m-i-1} \right) \\
& \times \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \left(\sum_{m=0}^k \alpha_m S_n^m \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} + \left(\sum_{m=0}^k \alpha_m S_n^m \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \\
& \times \left(\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} S_n^i Q_n S_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} - \left(\left(\sum_{m=0}^k \alpha_m S_n^m \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \right)^2 \\
& \times \left(\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} S_n^i Q_n S_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \\
& + Q_n \left(\left(\sum_{m=0}^k \alpha_m S_n^m \right) \left(\sum_{m=0}^k \beta_m S_n^m \right)^{-1} \right)^2 + O(\varepsilon^2),
\end{aligned}$$

$$\begin{aligned}
\tilde{Q}_{n+1} &= H_n \left(\left(\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} H_n^i \tilde{Q}_n H_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \left(\sum_{m=0}^k \alpha_m H_n^m \right) \right. \\
& \times \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} - \left(\sum_{m=0}^k \alpha_m H_n^m \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \left(\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} H_n^i \tilde{Q}_n H_n^{m-i-1} \right) \\
& \times \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \left(\sum_{m=0}^k \alpha_m H_n^m \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} + \left(\sum_{m=0}^k \alpha_m H_n^m \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \\
& \times \left(\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} H_n^i \tilde{Q}_n H_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} - \left(\left(\sum_{m=0}^k \alpha_m H_n^m \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \right)^2 \\
(3.9) \quad & \times \left(\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} H_n^i \tilde{Q}_n H_n^{m-i-1} \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \\
& \left. + Q_n \left(\left(\sum_{m=0}^k \alpha_m H_n^m \right) \left(\sum_{m=0}^k \beta_m H_n^m \right)^{-1} \right)^2 + O(\varepsilon^2). \right.
\end{aligned}$$

Writing (3.8) and (3.9) element-wise we have

$$\tilde{q}_{rs}^{(n+1)} = b_{rs}^{(n)} \tilde{q}_{rs}^{(n)},$$

$$\tilde{p}_{rs}^{(n+1)} = v_{rs}^{(n)} \tilde{q}_{rs}^{(n)} + c_{rs}^{(n)} \tilde{p}_{rs}^{(n)},$$

$$v_{rs}^{(n)} = l_r^{(n)} \left(\frac{\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} (h_r^{(n)})^i (h_s^{(n)})^{m-i-1}}{\sum_{m=0}^k \beta_m (h_s^{(n)})^m} \right. \\ \left. \frac{\left(\sum_{m=0}^k \alpha_m (h_r^{(n)})^m \right) \left(\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} (h_r^{(n)})^i (h_s^{(n)})^{m-i-1} \right)}{\left(\sum_{m=0}^k \beta_m (h_r^{(n)})^m \right) \left(\sum_{m=0}^k \beta_m (h_s^{(n)})^m \right)} \right),$$

$$b_{rs}^{(n)} = b_r^{(n)} \left(\frac{\left(\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} (h_r^{(n)})^i (h_s^{(n)})^{m-i-1} \right) \left(\sum_{m=0}^k \alpha_m (h_s^{(n)})^m \right)}{\left(\sum_{m=0}^k \beta_m (h_s^{(n)})^m \right)^2} \right)$$

$$\frac{\left(\sum_{m=0}^k \alpha_m (h_r^{(n)})^m \right) \left(\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} (h_r^{(n)})^i (h_s^{(n)})^{m-i-1} \right) \left(\sum_{m=0}^k \alpha_m (h_s^{(n)})^m \right)}{\left(\sum_{m=0}^k \beta_m (h_r^{(n)})^m \right) \left(\sum_{m=0}^k \beta_m (h_s^{(n)})^m \right)^2}$$

$$+ \frac{\left(\sum_{m=0}^k \alpha_m (h_r^{(n)})^m \right) \left(\sum_{m=1}^k \alpha_m \sum_{i=0}^{m-1} (h_r^{(n)})^i (h_s^{(n)})^{m-i-1} \right)}{\left(\sum_{m=0}^k \beta_m (h_r^{(n)})^m \right) \left(\sum_{m=0}^k \beta_m (h_s^{(n)})^m \right)}$$

$$\left. \frac{\left(\sum_{m=0}^k \alpha_m (h_r^{(n)})^m \right) \left(\sum_{m=1}^k \beta_m \sum_{i=0}^{m-1} (h_r^{(n)})^i (h_s^{(n)})^{m-i-1} \right)}{\left(\sum_{m=0}^k \beta_m (h_r^{(n)})^m \right)^2 \left(\sum_{m=0}^k \beta_m (h_s^{(n)})^m \right)} \right\} + \\
 + \left(\frac{\sum_{m=0}^k \alpha_m (h_s^{(n)})^m}{\sum_{m=0}^k \beta_m (h_s^{(n)})^m} \right)^2, \\
 c_{rs}^{(n)} = \frac{\sum_{m=0}^k \alpha_m (h_s^{(n)})^m}{\sum_{m=0}^k \beta_m (h_s^{(n)})^m}.$$

Let

$$e_{rs}^{(n)} = \begin{bmatrix} \tilde{q}_{rs}^{(n)} \\ \tilde{p}_{rs}^{(n)} \end{bmatrix}, \quad w_{rs}^{(n)} = \begin{bmatrix} b_{rs}^{(n)} & 0 \\ v_{rs}^{(n)} & c_{rs}^{(n)} \end{bmatrix}.$$

Now we have $e_{rs}^{(n)} = w_{rs}^{(n)} e_{rs}^{(n)} + O(\varepsilon^2)$.

From (3.6) and (3.7) we can write $w_{rs}^{(n)}$ as

$$w_{rs}^{(n)} = w_{rs} + O(\varepsilon^{(n)}), \quad w_{rs} = \begin{bmatrix} 0 & 0 \\ -1/2\sqrt{a_i} & 1 \end{bmatrix}$$

where $\varepsilon^{(n)}$ is sufficiently small for large n . The matrix w_{rs} has the eigenvalues 0 and 1, let y_0 and y_1 be the corresponding eigenvectors, so

$$e_{rs}^{(n)} = u_0^{(n)} y_0 + u_1^{(n)} y_1.$$

For sufficiently small ε and large n we have

$$e_{rs}^{(n+k)} \approx w_{rs}^k e_{rs}^{(n)} = u_1^{(n)} y_1, \quad k = 1, 2, \dots.$$

Consequently $|e_{rs}^{(n+k)}| = |e_{rs}^{(n+1)}|$ the method (II) is locally stable.

It is easy to show that for iterations (I) and (II), the operation counts for one stage of each iterations, measured in flops are as follows

flops per stage	general A	symmetric positive definite A
method (I)	$(k+4)n^3$	$(k+4)n^3/2$
method (II)	$(k+4)n^3$	$(k+4)n^3/2$

An alternative approach to compute $A^{-1/2}$ is via computing the eigenvectors of the matrix A . Let $A \in C^{n,n}$ be nonsingular, where A is diagonalizable by V and has eigenvalue matrix D . The inverse square root of A can be defined as the matrix

$$A^{-1/2} = V D^{-1/2} V^{-1} \quad (3.11)$$

where $D^{-1/2} = \text{diag}(1/\sqrt{a_1}, \dots, 1/\sqrt{a_n})$ and the matrix V is eigenvector matrix [5].

In the general case the above approach to compute $A^{-1/2}$ is computationally expensive (for large n) and this is a major drawback when using (3.11) to compute $A^{-1/2}$. Namely, computation of V requires $n^4/3$ flps (we need solved n systems of linear equations of the form $(A - a_i C)x_i = 0$ by method Gaussian elimination).

The above approach to compute $A^{-1/2}$ can be effective, as we will see in the next section, if A is hermitian matrix. In this case we can use the Jacobi iterative method to compute the eigenvalues and the eigenvectors of the matrix A . If the matrices D and V are computed after k iterations, then computation of $A^{-1/2}$ requires approximately

$$(3.12) \quad 3kn^3/2 \text{ flops.}$$

4. Numerical Examples

In this section we will use the Frobenius matrix norm

$$|A|_F = \sqrt{\sum_{i,j} |a_{i,j}|^2} \quad \text{and the error} \quad e_n = \|AX_n^2 - 1\|_F.$$

Example 1.

$$A = \begin{bmatrix} 5 & 4 & 1 & 1 \\ 4 & 5 & 1 & 1 \\ 1 & 1 & 4 & 2 \\ 1 & 1 & 2 & 4 \end{bmatrix}.$$

For $k = 1$ the method (I) converges in 3 iterations and $e_3 = 5.62E - 7$.

For $k = 2$ the method (I) converges in 2 iterations and $e_2 = 1.12E - 6$.

For $k = 3$ the method (I) converges in 2 iterations and $e_2 = 1.38E - 5$.

For $k = 4$ the method (I) converges in 2 iterations and $e_2 = 2.65E - 5$.

For $k = 5$ the method (I) converges in 2 iterations and $e_3 = 5.9E - 5$.

For $k = 6$ the method (I) converges in 1 iterations and $e_1 = 1.99E - 3$.

The eigenvalues of A are 1, 2, 5, 10. This example illustrates that condition (2.13) is not a necessary condition for numerical stability of method (I).

This example illustrates the corespondence between the parameter k and the computational complexity of the method (I).

Example 2.

$$A = \begin{bmatrix} 4 & 1 & 1 \\ 2 & 4 & 1 \\ 0 & 1 & 4 \end{bmatrix}.$$

The matrix A is not diagonalizable. For $k = 4$ the method (I) converges in 2 iterations and $e_2 = 4.13E - 7$.

Example 3.

$$A = \begin{bmatrix} 0 & 0.07 & 0.27 & -0.33 \\ 1.31 & -0.36 & 1.21 & 0.41 \\ 1.06 & 2.86 & 1.49 & -1.43 \\ -2.64 & -1.84 & -0.24 & -2.01 \end{bmatrix}.$$

The eigenvalues of A are 0.03, 3.03, $-1.97 \pm i$. Method (I) diverges as predicated by stability analysis.

For $k = 1$ the method (II) converges in 4 iterations and $e_4 = 1.12E - 4$.

For $k = 2$ the method (II) converges in 3 iterations and $e_3 = 9.64E - 3$.

For $k = 3$ the method (II) converges in 2 iterations and $e_2 = 7.92E - 4$.

This example illustrates the correspondence between the parameter k and the computational complexity of the method (II).

In [2] the following quadratically convergent method was proposed to find $A^{-1/2}$.

$$(III) \quad \begin{cases} P_0 = A, & Q_0 = I \\ P_{n+1} = 0.5(P_n + Q_n^{-1}) \\ Q_{n+1} = 0.5(Q_n + P_n^{-1}) \\ \lim_{n \rightarrow \infty} Q_n = A^{-1/2} \end{cases}$$

In [1] the following quadratically convergent methods were proposed to find $A^{-1/2}$.

$$(IV) \quad \begin{cases} X_0 = I \\ X_{n+1} = 2X_n(I + AX_n^2)^{-1} \\ \lim_{n \rightarrow \infty} X_n = A^{-1/2} \end{cases}$$

$$(V) \quad \begin{cases} T_{n+1} = T_n(I + S_n), & T_0 = I \\ S_{n+1} = S_n^2(2I - S_n^2)^{-1} & S_0 = (I - A)(I + A)^{-1} \\ \lim_{n \rightarrow \infty} T_n = A^{-1/2} \end{cases}$$

In [3] the following quadratically convergent method was proposed to compute $A^{-1/2}$ for the special case when A is hermitian positive definite.

$$(VI) \begin{cases} X_0 = \nu I \\ X_{n+1} = X_n + \frac{1}{4} X_n (I - X_n A X_n) + \frac{1}{4} (I - X_n A X_n) X_n \end{cases}$$

where $0 < \nu < \sqrt{3/\rho(A)}$.

For the method (V) the cost per iteration by [3] is

$$(4.1) \quad 3n^3 \text{ flops.}$$

Example 4.

$$A = \begin{bmatrix} 0.2 & 100 & 150 & 50 \\ 0 & 0.4 & 50 & 50 \\ 0 & 0 & 0.4 & 100 \\ 0 & 0 & 0 & 0.4 \end{bmatrix}.$$

Both methods (I) and (I) converge in 2 iterations and $e_2 = 0$.

For the method (III) it holds $e_5 = e_6 = e_7 = \dots = 1.05$,

For the method (V) it holds $e_5 = e_6 = e_7 = \dots = 0.95$.

We see that in this example the methods (I) and (II) are more precisely than the methods (III) and (V).

Example 5.

$$A = \begin{bmatrix} 0.002 & 1 & 1.5 & 0.5 \\ 0 & 0.003 & 0.5 & 0.5 \\ 0 & 0 & 0.003 & 1 \\ 0 & 0 & 0 & 0.005 \end{bmatrix}.$$

For the method (IV) it holds $e_{10} = e_{11} = e_{12} = \dots = 1$.

For $k = 4$ the method (II) converges in 3 iterations, $e_3 = 8.72E - 3$.

We see that in this example the method (II) is more precisely than the method (IV).

Example 6. Let A be 10×10 symmetric positive definite matrix defined by

$$a_{ij} = \begin{cases} 100 + ij & \text{if } i = j \\ ij & \text{if } i \neq j \end{cases}$$

For $k = 3$ the method (I) converges in 3 iterations, $e_3 = 6.95E - 6$.

Computation of X_3 requires by (3.10) approximately 10500 flops (for 3 iterations).

For the method (VI) the best results are obtained in the following cases:

for $\nu = 0.05$ the method (VI) converges in 6 iterations, $e_6 = 5.5E - 7$

for $\nu = 0.06$ the method (VI) converges in 5 iterations, $e_5 = 6.03E - 7$

for $\nu = 0.07$ the method (VI) converges in 6 iterations, $e_6 = 7.32E - 7$.

For the method (VI), in the best case ($\nu = 0.06$), computation of X_5 requires by (4.1) approximately 15000 flops for 5 iteration. We see that in this example the method (I) converges approximately 1.5 times faster than the method (VI).

If we use the Jacobi method and (3.11) then after 9 iterations we have $e_9 = 8.83E - 7$, computation of X_9 requires by (3.12) approximately 13500 flops (for 9 iterations).

Single precision calculations were used for the 6 examples.

References

- [1] N. Sherif, On the computation of a matrix inverse square root, *Computing* 46(1991), 295-305.
- [2] N.J. Higham, Newtons method for matrix square root, *Math. Comp.* 46(1986), 537-549.
- [3] B. Philippe, An algorithm to improve nearly orthonormal set of vectors on a vector processor, *SIAM Alg. Disc. Math.* 8(3)(1987), 396.
- [4] H.J. Higham, Computing real square roots of a real matrix, *Linear Algebra Appl.* 88/89(1987), 405-430.

- [5] P. Lancaster, Theory of matrices, *Academic Press*, New York 1969.
- [6] KyFan, A.J. Hoffman, Some metric inequalities in the space of matrices, *Proc. Amer. Math. Soc.* 6,111(1955).
- [7] G.W. Stewart, Introduction to matrix computation, *Academic Press*, New York 1974.

(University of Novi Sad, Technical Faculty „M. Pupin”, 23000 Zrenjanin, Yugoslawia).
Received on 4.11.1994 and, in revised form, on 3.2.1995.