

On approximating the inverse of a matrix

ION PĂVĂLOIU

ABSTRACT. In this note we deal with two problems: the first regards the efficiency in approximating the inverse of a matrix by the Shulz-type methods, and the second is the problem of evaluating the errors in the approximation of the inverses of the perturbed matrices.

1. INTRODUCTION

In this note we deal with two problems: the first regards the efficiency in approximating the inverse of a matrix by the Shulz-type methods, and the second is the problem of evaluating the errors in the approximation of the inverses of the perturbed matrices.

As it is well known, given a nonsingular matrix $A \in \mathbb{R}^{m \times m}$ and a matrix $D_0 \in \mathbb{R}^{m \times m}$ such that

$$\|I - AD_0\| \leq q < 1 \quad (1)$$

with $q \in \mathbb{R}$ and I the m -th order unit matrix, then, for $k \in \mathbb{N}, k \geq 2$ fixed, the sequence of matrices $(D_n)_{n \geq 0}$ given by

$$F_n = I - AD_n \quad (2)$$

$$D_{n+1} = D_n \left(I + F_n + F_n^2 + \cdots + F_n^{k-1} \right), \quad n = 0, 1, \dots$$

is convergent and $\lim_{n \rightarrow \infty} D_n = A^{-1}$. Moreover, $(F_n)_{n \geq 0}$ verifies

$$F_{n+1} = F_n^k, \quad n = 0, 1, \dots \quad (3)$$

The methods of type (2) represent generalizations of the well known Shulz method. Relation (3) shows that the convergence order of sequence $(D_n)_{n \geq 0}$ is $k, k \geq 2$.

We introduce the notion of efficiency index of method (2). We notice that at each iteration step, the number of the matrix sums required is equal to the number of matrix products which appear in (2). Moreover, for computing the sum $I + F_n + \cdots + F_n^k$ we may use a method similar to the Horner scheme, i.e.

$$I + F_n + F_n^2 + \cdots + F_n^{k-1} = \{ \{ [(F_n + I) F_n + I] F_n + I \} + \cdots \}. \quad (4)$$

Received: 28.02.2003; In revised form: 30.10.2003

Key words and phrases. *Inverse of a matrix, perturbed matrix, efficiency in approximating the inverse of a matrix.*

In this way the matrix sums required reduce to sums in which one term is the identity matrix. This remark is also valid for the term $F_n = I - AD_n$. The operation consisting of one matrix product and one matrix sum (regardless of their order) we call it computing unit.

Definition 1. *The efficiency index of method (2) is given by*

$$E_k = k^{1/s}, \quad (5)$$

where $s \in \mathbb{N}$ represents the number of computing units required at each iteration step of method (2).

This definition is given by analogy to the efficiency index introduced by A.M. Ostrowski in [2]. The definition may also be motivated by the following reasoning.

From (3) it follows

$$\|F_{n+1}\| \leq \|F_n\|^k, \quad n = 0, 1, \dots \quad (6)$$

whence

$$\|F_{n+1}\| \leq \|F_0\|^{k^{n+1}}, \quad n = 0, 1, \dots \quad (7)$$

The above inequalities lead to the following error bounds:

$$\|A^{-1} - D_n\| \leq \|A^{-1}\| \|F_0\|^{k^n}, \quad n = 0, 1, \dots \quad (8)$$

Consider now two methods of type (2), having the convergence orders k_1 and k_2 respectively. Assume that, for achieving the same precision, these methods require n_1 respectively n_2 iteration steps. Then (8) implies

$$k_1^{n_1} = k_2^{n_2}. \quad (9)$$

The total number of computing units is $n_1 s_1$ in the first case and $n_2 s_2$ in the second case.

It is clear now that the method with convergence order k_1 is more efficient than the other if

$$n_1 s_1 < n_2 s_2. \quad (10)$$

Relations (9) and (10) lead us to

$$k_1^{1/s_1} > k_2^{1/s_2}. \quad (11)$$

Taking into account Definition 1, it follows that among the methods of type (2) for different values of k , the most efficient is given by the one with high efficiency index.

We shall determine in the following section the optimal method, i.e., having the high efficiency index, when $k \in \mathbb{N}$, $k \geq 2$.

2. OPTIMAL EFFICIENCY INDEX

Assume that we use (4) at each iteration step in (2). It can be easily seen that for the sum in (4) there are needed $k - 2$ matrix products. Relation (2) shows that 2 more matrix products are required at each iteration step, so in total we need k matrix products.

Taking into account (5), it follows that the *efficiency index* of method (2) is given by

$$\bar{E}_k = k^{1/k}. \quad (12)$$

Considering the function $f : (0, +\infty) \rightarrow \mathbb{R}$, $f(x) = x^{\frac{1}{x}}$, it can be easily seen that this function attains a maximum value at $x = e$. Since f is increasing on $(0, e)$ and decreasing on $(e, +\infty)$, it follows that \bar{E}_k is the largest for $k = 3$.

We have proved the following result.

Theorem 2. *Among the methods (2) for $k = \mathbb{N}$, $k \geq 2$, the method with highest efficiency index is given by:*

$$\begin{cases} F_n = I - AD_n \\ D_{n+1} = D_n (I + F_n + F_n^2), \quad n = 0, 1, \dots \end{cases} \quad (13)$$

with D_0 verifying $\|I - AD_0\| \leq q < 1$.

By (4), the above method may be written as

$$\begin{aligned} F_n &= I - AD_n \\ D_{n+1} &= D_n [(F_n + I) F_n + I], \quad n = 0, 1, \dots \end{aligned} \quad (14)$$

In this case, (7) becomes

$$\|F_{n+1}\| \leq \|F_0\|^{3^{n+1}}, \quad n = 0, 1, \dots \quad (15)$$

and for the error bounds one has

$$\|A^{-1} - D_n\| \leq \|A^{-1}\| \|F_n\| \leq \|A^{-1}\| \|F_0\|^{3^n}, \quad n = 0, 1, \dots \quad (16)$$

It can be easily seen that under (1), one has the inequality

$$\|A^{-1}\| \leq \frac{\|D_0\|}{1 - \|F_0\|} \quad (17)$$

whence

$$\|A^{-1} - D_n\| \leq \|D_0\| \frac{\|F_0\|^{3^n}}{1 - \|F_0\|}, \quad n = 0, 1, \dots \quad (18)$$

Analogously, for any method of type (2) one may deduce the evaluation

$$\|A^{-1} - D_n\| \leq \|D_0\| \frac{\|F_0\|^{k^n}}{1 - \|F_0\|}, \quad n = 0, 1, \dots \quad (19)$$

3. ERROR BOUNDS IN CASE OF PERTURBED MATRICES

In practice, the elements of the matrix A are usually obtained as results of certain experiments, measurements, approximations etc. Therefore their values are altered by errors. Consequently we replace A by the approximation \tilde{A} . For a rigorous interpretation of the results, it is necessary to know an error bound $\varepsilon > 0$ for which

$$\|A - \tilde{A}\| \leq \varepsilon. \quad (20)$$

Instead of sequence $(D_n)_{n \geq 0}$ we consider $(\tilde{D}_n)_{n \geq 0}$, generated by

$$\begin{aligned} \tilde{F}_n &= I - \tilde{A}\tilde{D}_n; \\ \tilde{D}_{n+1} &= \tilde{D}_n \left(I + \tilde{F}_n + \tilde{F}_n^2 + \dots + \tilde{F}_n^{k-1} \right), \quad n = 0, 1, \dots \end{aligned} \quad (21)$$

We assume that the matrices \tilde{A} and \tilde{D}_0 above obey

$$\|I - \tilde{A}\tilde{D}_0\| \leq \bar{q} < 1. \quad (22)$$

It follows that \tilde{A} is invertible: $\exists \tilde{A}^{-1}$ and by (18) we get

$$\|\tilde{A}^{-1} - \tilde{D}_n\| \leq \|\tilde{D}_0\| \frac{\|\tilde{F}_0\|^{k^n}}{1 - \|\tilde{F}_0\|}, \quad n = 0, 1, \dots \quad (23)$$

We are interested in conditions which ensure that \tilde{A} is nonsingular. We consider the identity

$$I - \tilde{A}^{-1}A = \tilde{A}^{-1}(\tilde{A} - A)$$

which implies

$$\|I - \tilde{A}^{-1}A\| \leq \|\tilde{A}^{-1}\| \varepsilon,$$

whence, by (17) we get

$$\|I - \tilde{A}^{-1}A\| \leq \frac{\varepsilon \|\tilde{D}_0\|}{1 - \|\tilde{F}_0\|} \quad (24)$$

This relation shows that for the existence of the inverse for $\tilde{A}^{-1}A$ it suffices that

$$r = \frac{\varepsilon \|\tilde{D}_0\|}{1 - \|\tilde{F}_0\|} < 1 \quad (25)$$

whence for ε we get the condition

$$\varepsilon < \frac{1 - \|F_0\|}{\|\tilde{D}_0\|}. \quad (26)$$

Further,

$$A^{-1} = \left(\tilde{A}^{-1}A\right)^{-1} \tilde{A}^{-1}$$

whence

$$\|A^{-1}\| \leq \|\tilde{A}^{-1}\| \|\tilde{A}^{-1}A\| \leq \frac{\|\tilde{D}_0\|}{1 - \|\tilde{F}_0\| - \varepsilon \|\tilde{D}_0\|}$$

and (26) attracts $1 - \|\tilde{F}_0\| - \varepsilon \|\tilde{D}_0\| > 0$.

The following inequality can be easily proved

$$\|A^{-1} - \tilde{A}^{-1}\| \leq \frac{\|\tilde{D}_0\|^2 \varepsilon}{\left(1 - \|\tilde{F}_0\|\right) \left(1 - \|\tilde{F}_0\| - \varepsilon \|\tilde{D}_0\|\right)} \quad (27)$$

which, together with (23) leads to

$$\|A^{-1} - \tilde{D}_n\| \leq \frac{\|\tilde{D}_0\|}{1 - \|\tilde{F}_0\|} \left[\frac{\|\tilde{D}_0\| \varepsilon}{1 - \|\tilde{F}_0\| - \varepsilon \|\tilde{D}_0\|} + \|F_0\|^{t^n} \right], \quad n = 0, 1, \dots \quad (28)$$

This inequality provides a priori evaluations for the error. If we want to stop the iterations at a certain step \bar{n} such that $\|\tilde{F}_{\bar{n}}\| \leq \varepsilon_1$, $\varepsilon_1 > 0$ given, then by (7) and (17) it follows

$$\|\tilde{A}^{-1} - \tilde{D}_{\bar{n}}\| \leq \frac{\|\tilde{D}_0\|}{1 - \|\tilde{F}_0\|} \varepsilon_1,$$

which, together with (27) lead to

$$\|A^{-1} - \tilde{D}_{\bar{n}}\| \leq \frac{\|\tilde{D}_0\|}{1 - \|\tilde{F}_0\|} \left[\varepsilon_1 + \frac{\varepsilon \|\tilde{D}_0\|}{1 - \|\tilde{F}_0\| - \varepsilon \|\tilde{D}_0\|} \right]$$

which is an a posteriori error bound.

REFERENCES

- [1] Herzberger J., *Explizite Shulz Verfahren höherer Ordnung zur Approximation der re-
versen Matrix*, Z. Angew Math. und Mech. 1988, Bd. 68, No. 5, pp. 494-496
- [2] Ostrowski M.A., *Solution of equations in euclidian and Banach spaces*, Academic Press.
New York and London (1975)
- [3] Stickel E., *On a class of high order methods for investing matrices*, Z. Angew Math.
und Mech. 1987, Bd. 67, No. 7, pp. 334-336

"T. POPOVICIU" INSTITUTE OF NUMERICAL ANALYSIS
STR. FÂNTÂNELE, NR.57, BLOC B7, SC.II, ETAJ 5, AP. 67-68
CLUJ-NAPOCA, ROMANIA
E-mail address: `pavaloiu@ictp.acad.ro`