

Petr Liebl; Markéta Nováková

A method for dealing with ill-conditioned symmetric linear systems

*Aplikace matematiky*, Vol. 15 (1970), No. 6, 407--412

Persistent URL: <http://dml.cz/dmlcz/103314>

## Terms of use:

© Institute of Mathematics AS CR, 1970

Institute of Mathematics of the Academy of Sciences of the Czech Republic provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This paper has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

A METHOD FOR DEALING WITH ILL-CONDITIONED SYMMETRIC LINEAR SYSTEMS

PETR LIEBL and MARKÉTA NOVÁKOVÁ

(Received October 21, 1969)

1. Let us first introduce the notation used here. Let  $\mathbf{A}$  be a regular real hermitian matrix of order  $n$  with elements  $a_{ik}$ . Let  $\mathbf{A}^{-1} = \mathbf{B}$  have elements  $b_{ik}$ . Denote the eigenvalues of  $\mathbf{A}$ , arranged by ascending value, by  $\lambda_1, \lambda_2, \dots, \lambda_n$ , so that

$$(1) \quad 0 < |\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_n|.$$

Let the corresponding eigenvectors, written for convenience as row-vectors, be  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ , with coordinates  $\mathbf{v}_i = (v_{i1}, v_{i2}, \dots, v_{in})$ , so that

$$(2) \quad \mathbf{v}_i \mathbf{A} = \lambda_i \mathbf{v}_i.$$

Assume the eigenvectors normalized, so that

$$(3) \quad \sum_{k=1}^n v_{ik}^2 = 1, \quad i = 1, 2, \dots, n,$$

and orthogonal to each other, which is always possible. Introduce the matrices

$$\mathbf{A} = \begin{pmatrix} \lambda_1, & 0, & 0 \\ 0, & \lambda_2, & 0 \\ 0, & 0, & \lambda_n \end{pmatrix} \quad \text{and} \quad \mathbf{V} = \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_n \end{pmatrix} = \begin{pmatrix} v_{11}, v_{12}, \dots, v_{1n} \\ v_{21}, v_{22}, \dots, v_{2n} \\ \dots \dots \dots \\ v_{n1}, v_{n2}, \dots, v_{nn} \end{pmatrix}.$$

Then it is well known that  $\mathbf{V}^T = \mathbf{V}^{-1}$  and hence

$$(4) \quad \sum_{i=1}^n v_{ik}^2 = 1, \quad k = 1, 2, \dots, n.$$

Further it holds that  $\mathbf{A} = \mathbf{V}^T \mathbf{A} \mathbf{V}$  and  $\mathbf{B} = \mathbf{V}^T \mathbf{A}^{-1} \mathbf{V}$ . These equalities can be written in elements as follows

$$(5) \quad a_{ik} = \sum_j v_{ji} \lambda_j v_{jk}$$

$$(6) \quad b_{ik} = \sum_j v_{ji} \lambda_j^{-1} v_{jk}.$$

The summations, here as elsewhere, will be assumed to extend from 1 to  $n$ , unless explicitly stated otherwise. By a matrix norm, a mapping of the set of matrices into the interval  $\langle 0, \infty \rangle$  with certain properties is meant. Of the properties of matrix norms, we shall need the fact that any norm of the square matrix  $\mathbf{M}$  is not less than the value of any of the eigenvalues of  $\mathbf{M}$ . In this paper, we shall use throughout the “row sum” norm

$$(7) \quad \|\mathbf{M}\| = \max_i \sum_k |m_{ik}|$$

where  $m_{ik}$  is the element in row  $i$  and column  $k$  of  $\mathbf{M}$ . There is common agreement that the “ill-conditionedness” of a system of linear algebraic equations can be described, to a certain degree, by the value of a “condition number” of the coefficient matrix. A matrix norm given, the corresponding condition number of a regular matrix  $\mathbf{M}$  is defined as the product of the norm of  $\mathbf{M}$  itself and the norm of  $\mathbf{M}^{-1}$ . The condition number of a singular matrix is usually put equal to  $\infty$ . From properties of norms it follows that a condition number of any matrix is never less than 1. Using the norm defined by (7), introduce the condition number for any regular matrix  $\mathbf{M}$

$$(8) \quad C(\mathbf{M}) = \|\mathbf{M}\| \|\mathbf{M}^{-1}\|.$$

2. Consider the system of  $n$  linear algebraic equations

$$(9) \quad \mathbf{A}\mathbf{x}^T = \mathbf{y}^T$$

where  $\mathbf{A}$  is the matrix introduced above,  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  is a given row-vector, and  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  is the row-vector of the unknowns. Assume now the, say,  $q$ -th eigenvector  $\mathbf{v}_q$  and the corresponding eigenvalue  $\lambda_q$  to be known. We may write, using (2) and (9)

$$\mathbf{v}_q \mathbf{x}^T = \lambda_q^{-1} \mathbf{v}_q \mathbf{A} \mathbf{x}^T = \lambda_q^{-1} (\mathbf{v}_q \mathbf{y}^T)$$

where the leftmost member is a linear combination of the unknowns with known coefficients and the rightmost member a known number. So, we have derived a new equation for the unknowns  $x_k$

$$(10) \quad K v_{q1} x_1 + K v_{q2} x_2 + \dots + K v_{qn} x_n = \lambda_q^{-1} K \sum_k v_{qk} y_k$$

where  $K \neq 0$  is an arbitrary parameter. The idea now is to replace one, the  $p$ -th, say, of the original equations in (9) by the equation (10), obtaining the system

$$(9a) \quad \mathbf{A}^{(p,q)} \mathbf{x}^T = \mathbf{y}^{(p,q)T}$$

where

$$\mathbf{A}^{(p,q)} = \begin{bmatrix} a_{11}, & a_{12}, & \dots, & a_{1n} \\ a_{21}, & a_{22}, & \dots, & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{p-1,1}, & a_{p-1,2}, & \dots, & a_{p-1,n} \\ Kv_{q1}, & Kv_{q2}, & \dots, & Kv_{qn} \\ a_{p+1,1}, & a_{p+1,2}, & \dots, & a_{p+1,n} \\ \dots & \dots & \dots & \dots \\ a_{n1}, & a_{n2}, & \dots, & a_{nn} \end{bmatrix}$$

and

$$\mathbf{y}^{(p,q)} = (y_1, y_2, \dots, y_{p-1}, \lambda_q^{-1} K \sum_k v_{qk} y_k, y_{p+1}, \dots, y_n).$$

Obviously, the solution  $\mathbf{x}$  of (9) is a solution of (9a) as well. When the system (9) is ill-conditioned this is often caused by  $\lambda_1$  being in value much smaller than the other ones. In this situation it is however rather easy to find  $\lambda_1$  and  $\mathbf{v}_1$ . It is expected, from heuristic considerations, that by replacing one equation in (9) by (10) for  $q = 1$ , the "conditionedness" of the system is substantially improved. A theorem showing that this is indeed the case is proved in sec. 4.

3. Denote by  $\mathbf{A}^{(p,1)}$  the matrix obtained from matrix  $\mathbf{A}$  by replacing its  $p$ -th row by the vector  $K\mathbf{v}_1$ .

$\mathbf{A}$  being given,  $\mathbf{A}^{(p,1)}$  depends on  $p$  and  $K$ . It is easily verified that  $\mathbf{A}^{(p,1)} = \mathbf{T}\mathbf{A}$  where

$$\mathbf{T} = \begin{bmatrix} 1, & 0, & \dots, & 0 \\ 0, & 1, & \dots, & 0 \\ \frac{K}{\lambda_1} v_{11}, & \frac{K}{\lambda_1} v_{12}, & \dots, & \frac{K}{\lambda_1} v_{1n} \\ 0, & 0, & \dots, & 1 \end{bmatrix} \leftarrow p$$

If now

$$(12) \quad v_{1p} \neq 0,$$

$\mathbf{T}^{-1}$  exists and is found to be

$$(13) \quad \mathbf{T}^{-1} = \begin{bmatrix} 1, & 0, & \dots, & 0 \\ 0, & 1, & \dots, & 0 \\ \dots & \dots & \dots & \dots \\ -\frac{v_{11}}{v_{1p}}, & -\frac{v_{12}}{v_{1p}}, & \dots, & \frac{\lambda_1}{Kv_{1p}}, \dots, & -\frac{v_{1n}}{v_{1p}} \\ \dots & \dots & \dots & \dots & \dots \\ 0, & 0, & \dots, & 1 \end{bmatrix} \leftarrow p$$

and

$$(14) \quad \mathbf{A}^{(p,1)^{-1}} = \mathbf{A}^{-1}\mathbf{T}^{-1} = \mathbf{B}\mathbf{T}^{-1}.$$

Denoting the elements of  $\mathbf{A}^{(p,1)^{-1}}$  by  $\tilde{b}_{ik}$ , we have from (13) and (14)

$$\begin{aligned} \tilde{b}_{ik} &= b_{ik} - v_{1k}v_{1p}^{-1}b_{ip} \quad \text{for } k \neq p \\ \tilde{b}_{ip} &= \lambda_1 K^{-1} v_{1p}^{-1} b_{ip} \end{aligned}$$

and using further (6)

$$(15) \quad \tilde{b}_{ik} = \sum_j v_{ji} \lambda_j^{-1} (v_{jk} - v_{1k} v_{1p}^{-1} v_{jp}) \quad \text{for } k \neq p$$

$$(16) \quad \tilde{b}_{ip} = \lambda_1 K^{-1} v_{1p}^{-1} \sum_j v_{ji} \lambda_j^{-1} v_{jp}.$$

Note that in (15), the term for  $j = 1$  vanishes. We may therefore in (15) assume the summation to run from 2 to  $n$  and write

$$(17) \quad \tilde{b}_{ik} = \sum_{j \neq 1} \lambda_j^{-1} v_{ji} v_{jk} - v_{1k} v_{1p}^{-1} \sum_{j \neq 1} \lambda_j^{-1} v_{ji} v_{jp}.$$

**4.** Our further aim is to find upper bounds for the values of the elements of  $(\mathbf{A}^{(p,1)})^{-1}$  and hence an upper bound for  $\|(\mathbf{A}^{(p,1)})^{-1}\|$ . First we observe, using (3), that for any  $r, s$

$$(18) \quad \sum_j |v_{jr}| |v_{js}| \leq \frac{1}{2} \sum_j (v_{jr}^2 + v_{js}^2) = 1.$$

Now we start from (17), use (1) and (18) and obtain

$$\begin{aligned} (19) \quad |\tilde{b}_{ik}| &\leq \sum_{j \neq 1} |\lambda_j^{-1}| |v_{ji}| |v_{jk}| + |v_{1k}| |v_{1p}^{-1}| \sum_{j \neq 1} |\lambda_j^{-1}| |v_{ji}| |v_{jp}| \leq \\ &\leq |\lambda_2^{-1}| (\sum_j |v_{ji}| |v_{jk}| + |v_{1k}| |v_{1p}^{-1}| \sum_j |v_{ji}| |v_{jp}|) \leq \\ &\leq |\lambda_2^{-1}| (1 + |v_{1k}| |v_{1p}^{-1}|) \quad \text{for } k \neq p. \end{aligned}$$

From (16), using (1) and (18) we have

$$(20) \quad |\tilde{b}_{ip}| \leq |\lambda_1| |K^{-1}| |v_{1p}^{-1}| |\lambda_1^{-1}| \sum_j |v_{ji}| |v_{jp}| \leq |K^{-1}| |v_{1p}^{-1}|.$$

Putting together (19) and (20) we have

$$(21) \quad \sum_k |\tilde{b}_{ik}| \leq (n-1) |\lambda_2^{-1}| + |v_{1p}^{-1}| |\lambda_2^{-1}| \sum_{k \neq p} |v_{1k}| + |v_{1p}^{-1}| |K^{-1}|.$$

Note that the right hand side of the inequality is independent of  $i$  and hence is an upper bound for  $\|(\mathbf{A}^{(p,1)})^{-1}\|$  itself. Now everything is ready to prove the

**Theorem. Putting**

$$(22) \quad K = \left( \sum_k |v_{1k}| \right)^{-1} \|\mathbf{A}\|$$

and choosing  $p$  in such a way that

$$(23) \quad |v_{1p}| \geq |v_{1k}| \quad \text{for all } k,$$

it holds that

$$(24) \quad C(\mathbf{A}^{(p,1)}) < 3n|\lambda_1| |\lambda_2^{-1}| C(\mathbf{A}).$$

**Proof.** Note first that with (23), (12) is satisfied and our considerations are valid. Obviously, from (3) it follows that

$$\sum_{k \neq p} |v_{1k}| < \sum_k |v_{1k}| \leq \sqrt{n}.$$

Also, from (3) and (23) it follows that  $|v_{1p}^{-1}| \leq \sqrt{n}$ . According to our remark on norms,  $\|\mathbf{A}\| \geq |\lambda_2|$ . Using these three facts and (22) and (21), we have

$$\|(\mathbf{A}^{(p,1)})^{-1}\| \leq (3n-1) |\lambda_2^{-1}| < 3n |\lambda_2^{-1}|.$$

Now, from (22) it follows that

$$(25) \quad \|\mathbf{A}^{(p,1)}\| = \|\mathbf{A}\|.$$

Finally,

$$C(\mathbf{A}^{(p,1)}) = \|\mathbf{A}^{(p,1)}\| \|(\mathbf{A}^{(p,1)})^{-1}\| < \|\mathbf{A}\| \cdot 3n |\lambda_2^{-1}| = 3n |\lambda_2^{-1}| \|\mathbf{A}^{-1}\|^{-1} C(\mathbf{A}).$$

As the eigenvalues of  $\mathbf{A}^{-1}$  are  $\lambda_i^{-1}$ ,  $\|\mathbf{A}^{-1}\| \geq |\lambda_1^{-1}|$ , and (24) follows.

**5.** Based on the theorem and on what was said before, an algorithm may be suggested. Assume a system (9) to be given, with  $|\lambda_1|$  much smaller than  $|\lambda_2|$ . First, by backward iteration or any other method, a good approximation to  $\lambda_1$  and the corresponding eigenvector  $\mathbf{v}_1$  is found. Then,  $p$  is chosen with respect to (23), and  $K$  is found according to (22). Finally, in (9), the  $p$ -th equation is replaced by the equation (10). The resulting system (9a) has the same solution,  $\mathbf{x}$ , as (9), while the condition number of  $\mathbf{A}^{(p,1)}$  is at least  $|\lambda_1^{-1} \lambda_2| (3n)^{-1}$  times smaller than the condition number of  $\mathbf{A}$ . It is apparent that the smaller the ratio  $|\lambda_1| |\lambda_2|^{-1}$ , the easier the evaluation of  $\lambda_1$  and  $\mathbf{v}_1$  and the more substantial the improvement in the condition number. It should be noted, however, that one aspect of the illconditioned-ness of the system, the sensitivity of  $\mathbf{x}$  to small changes in  $\mathbf{y}$ , cannot be changed by any adjustments of the system. In the described method, it is hidden in the sensitivity of the right hand side of (10) to changes in  $\mathbf{y}$ . This suggests that that one number should be evaluated with special attention to accuracy. Also, it should be kept in mind that the accuracy of the result substantially depends on the accuracy in determining  $\lambda_1$  and  $\mathbf{v}_1$ .

### References

- [1] *J. W. Head, G. M. Oulton*: The Solution of "Ill-conditioned" Linear Simultaneous Equations. Aircraft Engineering, Vol. XXX, 1958, Math., 309—312.
- [2] *Ф. Р. Гаутмахер*: Теория матриц. Изд. Наука, Москва 1966.
- [3] *J. H. Wilkinson*: Sounding Errors in Algebraic Processes. London, H. M. Stationery Office, 1963.

### Souhrn

## METODA PRO ZACHÁZENÍ SE ŠPATNĚ PODMÍNĚNÝMI SYMETRICKÝMI SOUSTAVAMI LINEÁRNÍCH ALGEBRAICKÝCH ROVNIC

PETR LIEBL, MARKÉTA NOVÁKOVÁ

V [1] je naznačena metoda, kterou lze „zlepšit“ podmíněnost soustavy lineárních rovnic se symetrickou maticí. V matici koeficientů se nahradí jeden řádek vlastním vektorem, příslušejícím k vlastnímu číslu v absolutní hodnotě nejmenšímu.

V této práci je tato metoda popsána, proveden výběr parametrů ve formulích metody a dokázána věta o dosaženém zlepšení.

*Authors' addresses:* Petr Liebl, Matematický ústav ČSAV v Praze, Žitná 28, Praha 1; Markéta Nováková, CSc., Matematicko-fyzikální fakulta KU, Malostranské nám. 25, Praha 1.