"EST IMATES OF THE ROUNDOFF ERROR IN THE SOLUTION OF A SYSTEM OF CONDITIONAL EQUATIONS" BY V. I. GORDONOVA

TRANSLATED BY LINDA KAUFMAN



いていたのない、日本のないないないないないないないないないないであるである

١

X

STAN-CS-70-164 JUNE 1970

This doment			on oved
for put' te .		4	
els mondet	···· 1		

COMPUTER SCIENCE DEPARTMENT School of Humanities and Sciences STANFORD UNIVERSITY



CLEARINGHOUSE for Federal Scientific & Technical Information Springfield Va. 22151

ESTIMATES OF THE ROUNDOFF ERROR IN THE SOLUTION OF A SYSTEM OF CONDITIONAL EQUATIONS

by V. I. Gordonova

> Translated by Linda Kaufman

Reproduction in whole or in part is permitted for any purpose of the United States Government.

The preparation of this report was sponsored by the Office of Naval Research under grant number N0013-67-A-0112-0029, the National Science Foundation under grant number NSF GJ 408 and the Atomic Energy Commission under grant number AT (04-3)-326, PA 30.

- Aler

ESTIMATES OF THE ROUNDOFF ERROR IN THE SOLUTION OF A SYSTEM OF CONDITIONAL EQUATIONS

Ъу

V. I. Gordonova Zh. Vychislitel. Matem. i Matem. Fiziki, Vol. 9, No. 4 July-August 1969, pp. 775-83

In the present work we will examine estimates of the equivalent perturbation of roundoff errors in the solution of a system of conditional equations by the method of least squares (Method A) and by a method which was proposed by D. K. Faddeev, V. N. Faddeeva, and V. N. Kublanovskaya in a joint report at a conference on numerical methods in Kiev in 1966 (Method B).

Let us examine the system of conditional equations:

 $Ax = f \tag{1}$

with a rectangular matrix A having N rows and n columns, where generally N >> n. Method A leads to the system of normal equations

$$\mathbf{A}^{\mathrm{T}}\mathbf{A}\mathbf{x} = \mathbf{A}^{\mathrm{T}}\mathbf{f} \tag{2}$$

with a symmetric positive definite matrix $A^{T}A$ of rank n. We will assume that the solution of (2) is found by the method of square roots, always taking advantage of the accumulation of scalar products, independently of how one computes the elements of system (2).

Method B leads to a left orthogonal transformation of (1) into

 $\mathbf{Px} = \mathbf{L} \tag{2'}$

The term "equivalent perturbation" seems to refer to inverse roundoff analysis.

where P = QA, l = Qf, matrix P has non-null elements only in the right upper triangle \tilde{P} of rank n. Let \tilde{l} be the vector whose components are the first n components of the vector Qf. The triangular system

$$\tilde{\mathbf{P}}\mathbf{x} = \hat{\boldsymbol{\ell}}$$
(3)

is equivalent to system (2).

The total error in both methods is composed of the roundoff error in reading in the coefficients and the right-hand terms of (2) and (3) and the roundoff error during the solving of these systems. Since triangular systems may be solved very exactly ([1, Chapter 4]), we can neglect the roundoff error in the solution of (3) and in the backsolution part of the method of square roots in the solution of (2).

Because of the equivalence of (2) and (3) it does not matter whether one calculates the equivalent perturbation of roundoff errors of Methods A and B in terms of (2) or (3). We will do the calculations in terms of system (2) since this is more convenient. Everywhere below, if it is not specifically stated, we will use the symbols adopted in [1] and the Euclidean norm of the matrices and vectors.

1) Let us examine in the first place the errors of Method A. Because of the roundoff in the calculation of the scalar products, the elements of the matrix $A^{T}A$ and the vector $A^{T}f$ will be obtained with a certain error; i.e., instead of (2) we obtain

$$Bx = k \tag{4}$$

where $B = A^{T}A + \triangle(A^{T}A), k = A^{T}f + \triangle(A^{T}f).$

The norms of the error matrix $\triangle(A^TA)$ and the error vector $\triangle(A^Tf)$ essentially depend on the method of calculating of scalar products in the machine.

In the carrying out of all operations in a machine with a t-digit accuracy, the elements of $\triangle(A^TA)$ and $\triangle(A^Tf)$, which we will designate respectively by $\triangle b_{ij}$ and $\triangle k_i$, may be estimated on the basis of [1, Chapter 3] as

$$|\Delta b_{ij}| \le N2^{-t_1} \|a_i\| \|a_j\|$$
, $|\Delta k_j| \le N2^{-t_1} \|a_i\| \|f\|$,

if the calculations are executed with floating point (fl). Here and later $t_1 = t - 0.08406$, and a_1 is the i-th column of the matrix A. Hence, we obtain

$$\begin{split} \|\triangle(A^{T}A)\| &\leq N2^{-t_{1}} (\sum_{i,j} \|a_{i}\|^{2} \|a_{j}\|^{2})^{1/2} \\ &= N2^{-t_{1}} (\sum_{i} \|a_{i}\|^{2} \sum_{j} \|a_{j}\|^{2})^{1/2} = N2^{-t_{1}} \|A\|^{2} , \\ \|\triangle(A^{T}f\| \leq N2^{-t_{1}} \|A\| \|f\|. \end{split}$$
(5)

If the calculations are executed in fixed point (fi), we get correspondingly

$$\|\Delta(A^{T}A)\| \leq Nn2^{-t-1}, \|\Delta(A^{T}f)\| \leq Nn^{1/2} 2^{-t-1}.$$
 (6)

Here it is assumed $||a_i|| \le 1-N2^{-t-1}$, $||f|| \le 1-N2^{-t-1}$, which guarantees the possibility of calculating in fixed point.

If the scalar products are calculated with double precision, then the estimate under consideration is practically independent of N. In particular, in the case of floating point (fl_2) , according to [1, Chapter 3],

$$|\Delta b_{ij}| \le 2^{-t} (a_i^T a_j) + \frac{3}{2} N 2^{-2t+0.08406} ||a_i|| ||a_j||$$

Assuming $\frac{3}{2} N2^{-t} < 0.1$, we obtain

$$|\Delta b_{ij}| \le 2^{-t} (a_i^T a_j) + 0.11 \cdot 2^{-t} ||a_i|| ||a_j||$$

Using the relation $|(a_i^T a_j)| \leq ||a_j|| ||a_j||$, we find

$$||\Delta(A^{T}A)| \leq 1.11 \cdot 2^{-t} ||A||^{2}$$
 (7)

In the same way,

$$\|\Delta(A^{T}f)\| \leq 1.11 \cdot 2^{-t} \|A\| \|f\|$$
 (8)

In the case of fixed point (fi₂), we have

$$\|\Delta(A^{T}A)\| \le n2^{-t-1}, \|\Delta(A^{T}f)\| \le n^{1/2} 2^{-t-1}$$
 (9)

with the assumption that $||a_{i}|| \leq 1-2^{-t-1}$, $||f|| \leq 1-2^{-t-1}$.

Let us now estimate the equivalent perturbation due to the roundoff error in the application of the forward step in the method of square roots, i.e., in the decomposition of the matrix of system (4) into the product of two triangular matrices. It is known that the triang.lar factors S and S^{T} of the matrix B that are really obtained in the machine are the exact factors of a certain matrix B+E, that is

$$B + E = SS^{T}.$$
 (10)

The following estimates are verifiable for the elements e, of matrix E:

$$|\mathbf{e}_{ji}| \leq \begin{cases} |\mathbf{s}_{ij}\mathbf{s}_{jj}|^{2^{-t}}, i > j \\ |\mathbf{s}_{ji}\mathbf{s}_{ii}|^{2^{-t}}, i < j \\ \mathbf{s}_{ii}^{2} 2^{-t}, i = j \end{cases}$$
(11)

with an accuracy up to terms of $O(2^{-2t})$ in calculations with floating point and

$$|\mathbf{e}_{ij}| \leq \begin{cases} 0.5s_{ii} \ 2^{-t} &, \ i > j \\ 0.5s_{jj} \ 2^{-t} &, \ i < j \\ 1.00001s_{ii} \ 2^{-t} &, \ i = j \end{cases}$$
(12)

in calculations with fixed point. In the latter case, if $|b_{ij}| \leq 1-1.00001 \cdot 2^{-t}$ for all i,j and if matrix B is not very badly conditioned, then $|s_{ij}| < 1$ for all i,j.

Considering (4) and (10), we get that the numerical decomposition is exactly the decomposition of a perturbed matrix, i.e., $A^{T}A + C = SS^{T}$, where

$$C = \Delta(A^{T}A) + E.$$
 (13)

The norm of C is indeed of interest to us as the norm of the total error in the coefficients of system (2), while the norm of the vector $\triangle(A^Tf)$ is the norm of the error in the right-hand side of the system.

From (11) in calculations with floating point, neglecting terms of order 2^{-2t} , we have

$$\sum_{i,j=1}^{n} e_{ij}^{2} \leq 2^{-2t} \sum_{\substack{i=1 \ j=1}}^{n} (\sum_{j=1}^{i-1} s_{ij}^{2} s_{jj}^{2} + 4s_{ii}^{4} + \sum_{j=i+1}^{n} s_{ji}^{2} s_{ii}^{2})$$

$$\leq 2 \cdot 2^{-2t} (\sum_{\substack{i,j=1 \ i,j=1}}^{n} s_{ij}^{2} s_{jj}^{2} + \sum_{\substack{i,j=1 \ i,j=1}}^{n} s_{ji}^{2} s_{ii}^{2})$$

$$= 4 \cdot 2^{-2t} \sum_{\substack{i,j=1 \ i,j=1}}^{n} s_{ij}^{2} s_{jj}^{2}$$

$$\leq 4 \cdot 2^{-2t} \max_{j} s_{jj}^{2} \sum_{\substack{i,j=1}}^{n} s_{ij}^{2} = 4 \cdot 2^{-2t} \|S\|^{4}.$$

Considering that

$$\sum_{i=1}^{n} \sum_{j=1}^{2} = \sum_{i=1}^{i} \sum_{j=1}^{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{2} \sum_{i=1}^{n} \sum_{j=1}^{2} \sum_{i=1}^{2} \sum$$

we obtain $||S|| = ||A||[1 + O(N2^{-t})]$, where

$$\|\mathbf{E}\| \leq 2 \cdot 2^{-t} \|\mathbf{A}\|^{2} (1 + 0(N2^{-t})).$$
 (14)

As V. V. Voyevodin observed, these considerations permit us to obtain an estimate of the equivalent perturbation for the method of square roots which is n times better than that suggested in [2], without the assumption of accumulation.

Actually from the above explanation it follows that with an accuracy up to quantities of order $O(2^{-2t})$

$$\|\mathbf{E}\| \leq 2 \cdot 2^{-t} \left[\sum_{j \in \mathbf{i}} (\sum_{j \in \mathbf{i}} s_{ij}^2)^2 \right]^{1/2} \leq 2 \cdot 2^{-t} \|\mathbf{s}^{\mathrm{T}}\mathbf{s}\|$$
$$= 2 \cdot 2^{-t} \|\mathbf{s}\mathbf{s}^{\mathrm{T}}\| = 2 \cdot 2^{-t} \|\mathbf{s}\|.$$

Passing from the Euclidean norm to the spectral norm, we obtain

$$\|\mathbf{E}\| \le 2 \cdot 2^{-t} (\operatorname{Sp} B)^{1/2} \max_{i} s_{ii} \le 2 \cdot 2^{-t} (\operatorname{n} \max_{i} \lambda_{i}^{2})^{1/2}$$
$$= 2 \cdot 2^{-t} \operatorname{n}^{1/2} \|\mathbf{B}\|_{2}.$$

This estimate is n times better than the one obtained in [2], for example. For fixed point, an estimate analogous to (14), derived from (12) with the assumption that $|s_{ij}| < 1$, has the form

$$\|\mathbf{E}\| \le n2^{-t-1}(1 + O(\frac{1}{n})) .$$
 (15)

Using the relations (5)-(15), we obtain finally

$$\|C\| \leq N2^{-t_1} \|A\|^2 (1 + O(\frac{1}{N})), \|\Delta(A^T f)\| \leq N2^{-t_1} \|A\| \|f\|;$$
 (fe)

$$\|\Im_{\|} \leq 2.71 \cdot 2^{-t} \|A\|^{2}, \|\Delta(A^{T_{f}})\| \leq 1.11 \cdot 2^{-t} \|A\| \|f\|; \qquad (fl_{2})$$

$$\|C\| \le Nn2^{-t-1}(1 + O(\frac{1}{N})), \|\Delta(A^{T}f)\| \le Nn^{1/2} 2^{-t-1}; \quad (fi)$$

$$\|C\| \le n2^{-t}(1 + O(\frac{1}{n})), \|\Delta(\mathbf{A}^{T}\mathbf{f})\| \le n^{1/2} 2^{-t-1}, \qquad (fi_{2})$$

respectively, for the calculation of the elements of $A^{T}A$ and $A^{T}f$ in the cases of fl, fl_{2} , fi, fi₂.

7

F

2) We will now estimate the equivalent perturbation for the errors in Method B, which is equivalently an estimate of the errors in the elements of the system

$$P^{T}Px = P^{T}\ell$$
 (16)

which were obtained because of the inaccurate calculation of P and ℓ . Let us denote by ΔP and $\Delta \ell$, respectively, the matrix and the vector error. Because of these errors, instead of (16) we obtain the perturbed system $(P + \Delta P)^{T}(P + \Delta P)x = (P + \Delta P)^{T}(\ell + \Delta \ell)$. Neglecting the products $\Delta P^{T} \Delta P$ and $(\Delta P)^{T} \Delta \ell$, we obtain for the perturbations the approximate equalities

$$\triangle (\mathbf{P}^{\mathrm{T}}\mathbf{P}) = \mathbf{P}^{\mathrm{T}} \triangle \mathbf{P} + (\triangle \mathbf{P})^{\mathrm{T}}\mathbf{P}, \quad \triangle (\mathbf{P}^{\mathrm{T}}\boldsymbol{\ell}) = \mathbf{P}^{\mathrm{T}} \triangle \boldsymbol{\ell} + (\triangle \mathbf{P})^{\mathrm{T}}\boldsymbol{\ell},$$

from which

$$\|\triangle(\mathbf{P}^{\mathbf{T}}\mathbf{P})\| \leq 2\|\mathbf{P}\| \| |\Delta \mathbf{P}\| , \|\triangle(\mathbf{P}^{\mathbf{T}}\boldsymbol{\ell})\| \leq \|\mathbf{P}\| \| |\Delta \boldsymbol{\ell}\| + \|\Delta \mathbf{P}\| \| \| \boldsymbol{\ell}\| .$$

Because of the orthogonality of the matrix of transformation Q, we have

||P|| = ||QA|| = ||A|| and ||l|| = ||Qf|| = ||f||,

whence

$$\|\triangle(\mathbf{P}^{\mathbf{T}}\mathbf{P})\| \leq 2\|\mathbf{A}\| \|\Delta\mathbf{P}\| , \|\triangle(\mathbf{P}^{\mathbf{T}}\boldsymbol{\ell})\| \leq \|\mathbf{A}\| \|\triangle\boldsymbol{\ell}\| + \|\boldsymbol{\ell}\| \|\Delta\mathbf{P}\| .$$
(17)

In order to obtain final results it is necessary to estimate the norms of $\triangle P$ and $\triangle l$. These estimates essentially depend on the actual method of obtaining P, i.e., the method of transforming the system of simultaneous equations into system (2'). To obtain the matrix `P we will eliminate

the elements a_{ij} of matrix A for which i > j. We will perform the elimination with the help of a matrix of rotation or reflection [3]. Moreover, we will designate by $\alpha_1, \alpha_2, \ldots$ constants, which depend on the actual method of rounding in the machine. According to the assumptions of [1], these constants are not more than a few units or 1-2 tens.

(1) The transformation of matrix A is accomplished with the help of a succession of elementary rotation matrices T_{ij} in a cyclic order (Method B_{l}). Each of these rotations eliminates the element standing in the (i,j)-th position.

The roundoff error during the corresponding process of eliminating the subdiagonal elements of the square matrix was investigated in [1, Chapter 3], where elimination by columns was examined. In our case it is more convenient and necessary to eliminate elements by rows, i.e., in the order $(2,1), (3,1), (3,2), (4,1), \ldots, (n,n-1), (n+1,1), \ldots, (n+1,n), \ldots,$ (N,n). It can be shown that the roundoff error in the elimination of elements by rows and columns is the same.

Without stating the calculations, which are like those examined in [1, Chapter 3], but which are even more cumbersome, let us write the final result for the i-th column Δ_i of the error matrix ΔP :

$$\|_{\Delta_{1}}\| \leq \alpha_{1}^{2^{-t}} [n(N-n) + \frac{n(n-1)}{2}]^{1/2} (N+n-2)^{1/t} (1+6\cdot 2^{-t})^{(N+n-3)} \|_{a_{1}} \|$$
(18)

in the case of computing with floating point. In the same way an estimate, with the substitution of $\|\mathbf{f}\|$ for $\|\mathbf{a}_{\mathbf{j}}\|$, is verifiable for the error of transforming the column of the right-hand side. Here the calculation of scalar products with double precision has not been assumed. This cannot

essentially change the estimate since, in the process under consideration, we do not encounter the calculation of scalar products of a vector of more than the second order.

In computing with fixed point

$$\|\Delta_{1}\| \leq \alpha_{2} 2^{-t} [n(N-n) + \frac{n(n-1)}{2}]; \qquad (19)$$

moreover, for it to be possible to compute with fixed point it is sufficient that

$$\|a_{1}\| \leq 1 - \alpha_{2} 2^{-t} [n(N-n) + \frac{n(n+1)}{2}]$$

The same estimate is correct for the error of rotating the right-hand side.

The estimate obtained is exactly like that given in [1, Chapter 3], where actually the fact that the transformed matrix is square is not used.

Considering that $\|\Delta P\| = (\sum_{i=1}^{n} \|\Delta_i\|^2)^{1/2}$, we obtain from (18)

$$\|\Delta \mathbf{P}\| \leq \alpha_1 \operatorname{Nn}^{1/2} 2^{-t} \|A\|$$
, $\|\Delta \boldsymbol{\ell}\| \leq \alpha_1 \operatorname{Nn}^{1/2} 2^{-t} \|\mathbf{f}\|$

for floating point. In the same way from (19) we obtain

$$\|\Delta \mathbf{P}\| \leq \alpha_2 \operatorname{Nn}^{1/2} 2^{-t} \|\mathbf{A}\|, \quad \|\Delta \mathbf{\ell}\| \leq \alpha_2 \operatorname{Nn}^{2^{-t}}$$

for fixed point.

(2) Errors can be reduced essentially if one uses rotation matrices with the order of elimination of the unknowns that is suggested in [4] (Method B_2).

Let us denote by M the number of cycles required for the transformation of A into triangular form. The estimate computed in [4] for our case takes on the form

$$\|\Delta \mathbf{P}\| \leq \alpha_3 2^{-t} \mathbf{M} (1 + 6 \cdot 2^{-t})^{\mathbf{M} - 1} \|\mathbf{A}\| , \|\Delta \mathbf{e}\| \leq \alpha_3 2_{\mathbf{M}}^{-t} (1 + 6 \cdot 2^{-t})^{\mathbf{M} - 1} \|\mathbf{f}\|$$
(20)

for floating point, and

$$\|\Delta P\| \leq \alpha_{4} 2^{-t} M^{1/2} n^{1/2} [n(N-n) + \frac{n(n-1)}{2}]^{1/2} ,$$

$$\|\Delta e\| \leq \alpha_{4} 2^{-t} M^{1/2} [n(N-n) + \frac{n(n-1)}{2}]^{1/2}$$
(21)

for fixed point.

For an estimate of the value of M let us note that the number of cycles is independent of the actual realization of the process suggested in [4] if one does not consider zero elements of the initial matrix or any elements accidentally zeroed in one elementary transformation. For the elimination of the m-l elements of the matrix consisting of m rows and one column, $[log_2(m-1)] + l$ cycles are required. Here the square brackets denote the integer part.

Let the matrix have N rows and n columns. For the elimination of all the elements of the first column except the first element, one requires $[\log_2(N-1)] + 1$ cycles. With these it could happen that some of the elements of other columns are eliminated. However, even if one disregards the last situation for the elimination of elements of the second column, $[\log_2(N-2)] + 1$ cycles are required, etc.

Finally, we obtain

$$M \le \sum_{k=1}^{n} [\log_2(N-k) + n \le n[[\log_2(N-1)] + 1] .$$

This estimate is a little excessive, but not by more than 4-5 times for N < 100000.

Using this estimate for M, we find from (20) and (21)

$$\|\Delta P\| \leq \alpha_3 n \log_2 N \cdot 2^{-t} \|A\| , \|\Delta \ell\| \leq \alpha_3 n \log_2 N \cdot 2^{-t} \|f\|$$

for floating point and

$$\|\Delta P\| \le \alpha_{4} n^{3/2} (N \log_{2} N)^{1/2}, \|\Delta \ell\| \le \alpha_{4} n (N \log_{2} N)^{1/2}$$

for fixed point.

(3) Using a matrix of rotation (Method B_3) for the elimination of the clements of A appears most expedient in that case where the scalar products are calculated with double precision. Moreover, the estimates for ΔP and Δl are practically independent of N. Let us assume here that the calculation is carried out in floating point. The results obtained in [1, Chapter 3] go for rectangular matrices A and give

$$\|\Delta P\| \leq \alpha_5(n-1)2^{-t} \|A\|, \|\Delta \ell\| \leq \alpha_5(n-1)2^{-t} \|f\|.$$

Having substituted the estimate received for ΔP and $\Delta \ell$ into (17), we obtain a final estimate of the norm of the error matrix $\Delta(P^{T}P)$ and the error vector $\Delta(P^{T}\ell)$; namely,

for method B₁:

$$|\Delta(\mathbf{P}^{\mathbf{T}}\mathbf{P})|| \leq \alpha_1 \operatorname{Nn}^{1/2} 2^{-t} ||\mathbf{A}||^2, \ ||\Delta(\mathbf{P}^{\mathbf{T}}\boldsymbol{\ell})|| \leq \alpha_1 \operatorname{Nn}^{1/2} 2^{-t} ||\mathbf{A}|| \ ||\mathbf{f}|| \tag{f}\boldsymbol{\ell}$$

$$\|\Delta(\mathbf{P}^{\mathrm{T}}\mathbf{P})\| \leq \alpha_{2} \mathrm{Nn}^{2} 2^{-t}, \quad \|\Delta(\mathbf{P}^{\mathrm{T}}\boldsymbol{\ell})\| \leq \alpha_{2} \mathrm{Nn}^{3/2} 2^{-t}; \quad (fi)$$

for method B₂:

$$\|\Delta(\mathbf{P}^{\mathbf{T}}\mathbf{P})\| \leq \alpha_{3} \ln \log_{2} \mathbb{N} \cdot 2^{-t} \|\mathbf{A}\|^{2}, \ \|\Delta(\mathbf{P}^{\mathbf{T}}\boldsymbol{\ell})\| \leq \alpha_{3} \ln \log_{2} \mathbb{N} \cdot 2^{-t} \|\mathbf{A}\| \ \|\mathbf{f}\| \qquad (f\boldsymbol{\ell})$$

$$\|\Delta(\mathbf{P}^{T}\mathbf{P})\| \leq \alpha_{\mu} n^{2} (N \log_{2} N)^{1/2}, \|\Delta(\mathbf{P}^{T}\boldsymbol{\ell})\| \leq \alpha_{\mu} n^{3/2} (N \log_{2} N)^{1/2}; \quad (fi)$$

for method B₂:

$$\|\Delta(\mathbf{P}^{\mathbf{T}}\mathbf{P})\| \leq \alpha_{5}(n-1)2^{-t} \|A\|^{2}, \ \|\Delta(\mathbf{P}^{\mathbf{T}}\boldsymbol{\ell})\| \leq \alpha_{5}(n-1)2^{-t} \|A\| \|\|\mathbf{f}\| . \qquad (\mathbf{f}\boldsymbol{\ell}_{2}).$$

Comparing the obtained results, we see that the estimates of the equivalent perturbations for the matrix of system 2 have the form $2^{-t}\varphi(N,n)||A||^2$ and $2^{-t}\psi(N,n)$, respectively, for the different methods of calculation. In the table the order of magnitude of the functions $\varphi(N,n)$ and $\psi(N,n)$ are set forth $(N \gg n)$.



In this table it is seen that a comparison of Methods A and B, in the sense of majorizing the estimate, goes as a rule in favor of Method A. Method B_2 is the elimination method.

Let us go now from the equivalent perturbations to the error in the solution of the system. It is not difficult to construct an example in which with Method B_1 one obtains an order of the norm of the error in the solution which is equal to the largest estimate of Method A without accumulation. Let us examine, for example, the system with a matrix of coefficients and a right-hand vector, respectively, of the form

$$a_{ij} = \begin{cases} 0.5 & i = j, \\ 0 & i \neq j, i \le n, \\ \epsilon_1 & i > n, \end{cases} \quad \begin{array}{c} 1/n & i \le n; \\ f_i = \\ 0 & i > n; \\ \end{array}$$

where $\epsilon \ll 1$, so that $n(N-n) \epsilon < 1$.

Let us consider that computations are carried out with fixed point, and that the elementary matrix rotations are computed exactly. Assume that multiplication by these matrices is equally exact. After each multiplication by an elementary matrix of rotation, one rounds off the elements obtained up to a t digit number with fixed point, which gives an error of 2^{-t-1} . It is possible to assume that in this situation the elements of ΔP , which stand on the main diagonal and above, have the form $(N-n)2^{-t-1} + O(n(N-n) \in 2^{-t})$. Also, the components of the vector $\Delta \ell$ have this form with numbers which are not larger than n.

Let us designate by Δx the vector of the error of the solution. When $(\tilde{P} + \Delta \tilde{P})(x + \Delta x) = \tilde{l} + \Delta \tilde{l}$, then, neglecting the product $\Delta \tilde{P}\Delta x$, we obtain $\Delta x = \tilde{P}^{-1}(\Delta \tilde{l} - \Delta \tilde{P}x)$. Having computed \tilde{P}^{-1} and x, we obtain $\|\Delta x\| = O((N-n)^{1/2}2^{-t})$. The same order for $\|\Delta x\|$ is obtained in Method A if one uses the identity $\Delta x = (A^T A)^{-1}(\Delta (A^T f) - Cx)$ and the maximum estimates for $\Delta (A^T l)$ and C. In conclusion, let us take note of a fact which is connected to the practical application of the methods under consideration. The application of Methods B_2 and B_3 requires storage in memory of all the elements of the matrix A, while the application of Methods A and B_1 permits a row-by-row introduction of the information. The latter allows one a practically limit-less way to increase the values of N. In the row-by-row introduction of information in Method A with accumulation of scalar products, one demands in addition n^2 + n work cells for the storage of intermediate values during the calculation of the elements of A^TA and A^Tf . Actually, in this case the coefficients (and the right-hand side) of system (2) can be considered in a parallel fashion and each of these intermediate values, written down in 2t digits, can be stored in 2 cells of memory.

The author wishes to thank V. V. Voyevodin for posing the problem and for guidance.

REFERENCES

- J. H. Wilkinson. <u>The Algebraic Eigenvalue Problem</u>. Oxford, Clarendon Press, 1965.
- [2] J. H. Wilkinson. Apriori error analysis of algebraic processes. <u>Reports of the Proceedings of the International Congress of</u> <u>Mathematics</u>, Moscow, 1966.
- [3] D. K. Faddeev, V. N. Faddeeva. <u>Computational Methods of Linear</u> <u>Algebra</u>. Translated by R. C. Williams. San Francisco: W. H. Freeman, 1963.
- [4] V. V. Voevodin (Voyevodin). On the order of elimination of unknowns. <u>USSR Computational Mathematics and Mathematical</u> <u>Physics</u>, 1966, vol. 6, no. 4, pp. 203-06. R. C. Glass, translation editor. Oxford, Pergamon Press, Ltd.

		_			
DOCUMENT C	CONTROL DATA - R &	D			
(Security classification of title, body of abstract and ind BIGINA TING ACTIVITY (Compare) authori	exing ennotation must be ent	ered when the	ECURITY CLASSIFICATION		
Computer Science Department		Unclassified			
Stanford University	2	20. GROUP			
Stanford, California 94305					
EPORT TITLE					
ESTIMATES OF THE ROUNDOFF ERROR IN ' A SYSTEM OF CONDITIONAL EQUATIONS	THE SOLUTION OF	. <u></u>			
ESCRIPTIVE NOT ES (Type of report and inclusive dates) Manuscrint for Publication (Mechnic	el Report)				
UTHOR(3) (First name, middle initial, last name)					
V. I. Gordonova (Translated by Linda	a Kaufman)		· · · · · · · · · · · · · · · · · · ·		
June 1070	TE. TOTAL NO. OF		10. NO. OF REFS		
CONTRACT OR GRANT NO	10 Se. ORIGINATOR'S P	EPORT NUM	4 BER(\$)		
N00014-67-A-0112-0029					
PROJECT NO.	5	STAN-CS-	70-164		
NR 044-211					
	this report)	NO(S) (MHy C	mer numbers mat may be satighe		
		non	e		
DISTRIBUTION STATEMENT					
SUPPLEMENTARY NOTES	12. SPONSORING MIL	12. SPONSORING MIL, TARY ACTIVITY			
	Office	Office of Naval Research			
ABSTRACT					
Using backward error analysis, this least-squares solution of a system of different methods. The first one er A ^T Ax=A ^T f and the second is one prop in 1966. This latter method involve	paper compares the of conditional equi- ntails solving the posed by Faddeev, es multiplying the into upper trians	ne round nations normal Faddeeva system gular for	off error in the Ax=f by two equations a, and Kublanovskaya by othogonal rm.		
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					
matrices to transform the matrix A					

4

22

18.00

C. States

Security Classification							
14. KEY WORDS		LINKA		LINK		LINKC	
	ROLE		ROLE		HOLE	WT	
Roundoff Error							
Backward Error Analysia							
Backward Error Analysis							
Least Squares Conditional Equations							
	1	1		i			
						1	
	1						
	1						
					Í		
	1 1				1		
					1		
	i ľ						
DD		INCT	ASSTET	ED			
(PAGE 2)		Security	Classific				

UNCLASSIFIED

.