

An empirical estimate of the relative error of the computed solution \bar{x} of $Ax = b$

By G. Loizou*

An empirical method is described for providing, with a minimum of effort, useful estimates for the relative errors in individual components of the computed solution \bar{x} of the system of linear equations $Ax = b$. The estimate for the error is in no sense a rigorous upper bound, but in view of the economy of effort it is felt that extensions of the method described here provide a valuable addition to known methods of *a posteriori* error analyses.

(First received June 1967, and in revised form, September 1967)

1. Introduction

The problem to be discussed in this paper is that of estimating empirically the accuracy of the computed vector \bar{x} , obtained in the solution of

$$Ax = b, \quad (1)$$

where A will be assumed to be an $n \times n$ non-singular matrix. The computed solution can be effected by either a direct or an iterative method. Here attention is given to the direct method of Gaussian elimination which has been studied extensively by Wilkinson (1963 and 1965). As is often the case the matrix A is very or fairly ill-conditioned and the computed solution \bar{x} is a poor approximation to the true solution. The reason for this is attributed to the accumulation of rounding errors, cancellation and the ill-conditioning of the matrix A . The accumulation of rounding errors has been studied (see, for example, Wilkinson, 1963 and 1965) by using the elegant technique of backward error analysis in floating or fixed-point arithmetic. However, once a computed solution \bar{x} is obtained one would like to know approximately the accuracy. In this paper a practical method is developed which gives component-wise estimates of the accuracy of the computed solution \bar{x} . It must be emphasised that the experimental relative error obtained by the empirical method described in this paper is not an upper bound of the actual relative error in the computed solution; it does, however, give a good approximation to the order of magnitude of the actual relative error. The method, as it will be seen, was applied to very ill-conditioned matrices and the results have been fairly good. Tables are given for various matrices and, as is shown from the ALGOL procedure of the method, almost no extra computing time is required to ascertain the accuracy of \bar{x} , since most of the quantities needed to estimate it are obtained progressively in the various stages of the elimination method for solving (1).

2. Description of the method

If A is a non-singular matrix then, in general, it can be factorised in the form $A = LU$, where L is unit lower

triangular and U is upper triangular (Gaussian elimination). The factorisation, when it exists, is unique.

Denote the initial set of equations by $A_0x = b^0$. The factorisation which is performed in $n - 1$ major steps is equivalent to premultiplication of A_0 by elementary matrices M_i , $i = 1(1)n - 1$, where the i th step is given by

$$A_i = M_i M_{i-1} \dots M_2 M_1 A_0 = P_i A_0,$$

$$i = 1(1)n - 1, P_{n-1} = L^{-1} \text{ and } A_{n-1} = U.$$

The forward substitution, which gives $Ux = L^{-1}b^0$, is similarly equivalent at the i th step to

$$b^i = M_i M_{i-1} \dots M_2 M_1 b^0 = P_i b^0,$$

$$i = 1(1)n - 1, \text{ and } b^{n-1} = L^{-1}b^0.$$

The backward substitution gives progressively the computed solution vector \bar{x} , whose components \bar{x}_r , $r = n(-1)1$, are given by

$$\bar{x}_r = \frac{u \cdot y}{u_{rr}}, \quad r = n(-1)1, \quad \text{where}$$

$$u = (-u_{r+1}, -u_{r+2}, \dots, -u_{rn}, b_r^{n-1})^T,$$

$$y = (\bar{x}_{r+1}, \bar{x}_{r+2}, \dots, \bar{x}_n, 1)^T,$$

and u^T denotes the transpose of u .

Having considered the above structure of the solution of (1), an empirical method is now devised for computing the component-wise relative loss of accuracy of the computed solution vector x . Define

$$p_r = \frac{|u_{rr}|}{\max_{i=0(1)r-1} |a_{rr}^{(i)}|}, \quad r = 1(1)n,$$

$$q_r = \frac{\|(u_{r+1}, u_{r+2}, \dots, u_{rn}, L_r^{-1}b^0)\|}{\max_{i=0(1)r-1} \|(a_{r+1}^{(i)}, a_{r+2}^{(i)}, \dots, a_{rn}^{(i)}, b_r^{(i)})\|}, \quad r = 1(1)n,$$

where $a_{n+1}^{(n-1)} = u_{n+1} = 0$, L_r^{-1} denotes the r throw of the matrix L^{-1} , the underlying norm $\|\cdot\|$ being the euclidean vector norm, and

$$s_r = \frac{|u \cdot y|}{\|u\| \|y\|}, \quad r = 1(1)n.$$

REMARK. The quantities p_r, q_r estimate respectively the loss of relative accuracy in the r th pivot and the remaining part of the r th pivotal row.

* Computer Centre, Queen Mary College, London, E.1.

Present address: Mathematics Department, Birkbeck College, London, W.C.1.

A quantity t_r is defined progressively (as the elements of \bar{x} , \bar{x}_r , $r = n(-1)1$, are computed in the backward substitution) in the following way:

$$t_r = \frac{[1 + (\text{rel. loss of acc. } (\bar{x}_n))^2 + \dots + (\text{rel. loss of acc. } (\bar{x}_{r+1}))^2]^{1/2}}{n + 1 - r},$$

where $\bar{x}_{n+1} = 0$, $r = n(-1)1$.

The relative loss of accuracy of $(u.y)$ is defined by

$$w_r = q_r \times s_r \times t_r, \quad r = 1(1)n,$$

whilst that of \bar{x}_r is defined by the quantity h_r , where

$$h_r = p_r \times w_r, \quad r = 1(1)n,$$

provided $p_r \times w_r > 10 \times 2^{-t}$, 2^{-t} representing the least significant digit in the word of a binary digital computer. If, however, the above side condition is not satisfied then $h_r = p_r$, if $p_r < w_r$; otherwise $h_r = w_r$. In such a case it seems more appropriate that $h_r = 0.1 \times p_r$ or $h_r = 0.1 \times w_r$, since in either case w_r or p_r must have some bearing on the evaluation of h_r .

Finally, the quantity h_r , namely the measure of the relative loss of accuracy of \bar{x}_r , should satisfy

$$h_r \times (\text{relative error } (\bar{x}_r)) \div 2^{-t}, \quad r = 1(1)n. \quad (2)$$

In the ALGOL procedure that follows no pivoting has been used and consequently the process can break down even when A is very well-conditioned: for example

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

Even more serious (because it happens far more frequently) it may be numerically unstable when A is well-conditioned: for example

$$A = \begin{bmatrix} \epsilon & 1 \\ 1 & \epsilon \end{bmatrix},$$

where $\epsilon = 0(10^{-10})$. In fact the process is unnecessarily unstable whenever a principal submatrix of A is much more ill-conditioned than A itself. These 'unnecessary' failures and instability may, usually, be avoided by partial pivoting (see, for example, Wilkinson, 1963 and 1965).

Nevertheless, the emphasis here is not on the solution of (1), but on the empirical computation of the component-wise relative loss of accuracy of \bar{x} . Certainly one can easily introduce partial pivoting into the procedure following closely the algorithm given by Bowdler, Martin, Peters and Wilkinson (1966).

3. ALGOL procedure

Formal parameter list:

n order of the matrix A .

a elements of the matrix A stored as $n \times n$ array.

b elements of b stored as $n \times 1$ array.

h elements of the computed relative loss of accuracy of \bar{x} stored as $n \times 1$ array.

procedure solve $G(a,b,h,n)$; **value** n ; **array** a,b,h ; **integer** n ; **comment** *Solves* $Ax=b$. *The unsymmetric matrix, A , is stored in the $n \times n$ array $a[i,j]$, $i=1(1)n$, $j=1(1)n$, b is stored in the array $b[i]$, $i=1(1)n$. The decomposition $A=LU$, where L is a unit lower triangular matrix and U an upper triangular matrix, is performed and overwritten on A , omitting the unit diagonal of L . The solution vector is overwritten on b , whilst the computed relative loss of accuracy is stored in the array $h[i]$, $i=1(1)n$. The method will fail if A , modified by the rounding errors, is singular or if the pivot is zero or almost zero, since no pivoting (whether it be partial or total) is used;*

```
begin array  $e[0:n]$ ,  $k,l,u,w,d[1:n]$ ; real  $m$ ; integer  $i,j,r$ ;
for  $i := 1$  step 1 until  $n$  do
begin  $u[i] := \text{abs}(a[i,i])$ ;  $m := 0$ ;
    for  $j := i+1$  step 1 until  $n$  do  $m := m + a[i,j] \times a[i,j]$ ;
     $k[i] := \text{sqrt}(m + b[i] \times b[i])$ 
end;
for  $r := 1$  step 1 until  $n-1$  do
for  $i := r+1$  step 1 until  $n$  do
begin  $a[i,r] := a[i,r]/a[r,r]$ ;  $b[i] := b[i] - a[i,r] \times b[r]$ ;
     $m := 0$ ;
    for  $j := i+1$  step 1 until  $n$  do
begin  $a[i,j] := a[i,j] - a[i,r] \times a[r,j]$ ;
        if  $j > i$  then  $m := m + a[i,j] \times a[i,j]$ 
    end;
     $m := l[i] := \text{sqrt}(m + b[i] \times b[i])$ ;
    if  $k[i] < m$  then  $k[i] := m$ ;
     $m := \text{abs}(a[i,i])$ ;
    if  $u[i] < m$  then  $u[i] := m$ 
end;  $l[1] := k[1]$ ;
for  $i := 1$  step 1 until  $n$  do
begin  $m := \text{abs}(a[i,i])$ ;
     $u[i] := \text{if } m < u[i] \text{ then } m/u[i] \text{ else } 1$ ;
     $k[i] := \text{if } l[i] < k[i] \text{ then } l[i]/k[i] \text{ else } 1$ 
end;
for  $i := n$  step  $-1$  until 1 do
begin  $m := b[i]$ ; if  $i \neq n$  then
    for  $j := i+1$  step 1 until  $n$  do  $m := m - a[i,j] \times b[j]$ ;
     $d[i] := \text{abs}(m)$ ;  $b[i] := m/a[i,i]$ ;
     $m := 1$ ;
    for  $r := i+1$  step 1 until  $n$  do  $m := m + b[r] \times b[r]$ ;
     $d[i] := d[i]/(l[i] \times \text{sqrt}(m))$ 
end;
 $e[n] := 1$ ;  $m := 1$ ;
for  $i := n$  step  $-1$  until 1 do begin
     $w[i] := k[i] \times d[i] \times e[i]$ ; if  $w[i] \times u[i] < 0.72_{10} - 11$  then begin
if  $u[i] > w[i]$  then  $h[i] := w[i] \times 0.1$  else  $h[i] := u[i] \times 0.1$  end
    else  $h[i] := u[i] \times w[i]$ ;  $m := m + h[i] \times h[i]$ ;
     $e[i-1] := \text{sqrt}(m)/(n-i+2)$  end
end;
```

4. Results and discussion

In what follows $\text{cond}(A)$ denotes the condition number of the matrix A , namely

$$\text{cond}(A) = \text{lub}(A) \text{lub}(A^{-1}) = \|A\| \|A^{-1}\|$$

Table 1

$$A = \begin{bmatrix} 5.000 & 2.100 & 0.050 & 0.910 & 0.871 & 0.030 & 0.059 & 0.006 & 0.871 & 0.061 & 0.005 \\ 0.560 & 6.200 & 0.030 & -0.820 & 0.090 & 0.073 & 0.525 & 0.305 & -3.009 & 0.561 & 0.003 \\ -2.000 & 0.050 & 6.800 & -0.005 & 0.008 & 0.012 & 0.053 & 0.007 & 0.623 & -0.192 & -1.520 \\ 3.100 & 0.003 & 0.059 & 7.300 & 0.053 & -0.092 & -0.059 & 0.001 & 0.007 & -0.003 & 2.250 \\ 4.200 & 0.058 & 0.067 & 0.003 & 8.900 & 0.562 & 0.923 & 0.005 & 0.010 & 0.050 & -0.007 \\ 3.200 & 1.500 & 0.331 & -0.523 & -0.008 & 9.700 & 0.004 & 0.001 & -0.009 & 1.500 & -0.098 \\ -2.000 & 0.050 & 0.060 & 0.071 & 0.075 & -0.502 & 10.500 & 0.002 & 0.053 & -0.100 & 0.005 \\ 1.500 & 1.600 & 1.700 & -1.800 & 0.062 & -0.065 & -0.068 & 11.700 & -2.200 & 0.050 & 0.001 \\ 0.050 & 0.003 & 0.007 & 0.600 & -0.700 & -0.500 & 0.009 & 0.012 & 12.300 & -0.300 & 5.000 \\ 1.200 & 1.350 & -0.059 & 0.350 & 6.030 & -0.570 & -0.029 & 0.635 & 0.920 & 13.900 & 1.000 \\ -4.300 & -0.600 & 0.052 & 0.900 & 0.300 & 0.800 & 0.825 & -0.920 & -5.300 & -2.000 & 14.900 \end{bmatrix}$$

$$b = (13.2450, -2.8925, 2.9930, 28.2870, 29.7075, 38.5690, 34.4875, 38.1775, 79.4180, 97.5050, 49.5855)^T.$$

COMPUTED SOLUTION	EXACT SOLUTION	EXP. REL. ERROR	ACTUAL REL. ERROR
0.5000000000	0.5	(0.34)10 ⁻⁸	0
1.1000000000	1.0	(0.39)10 ⁻⁹	0
1.5000000000	1.5	(0.44)10 ⁻⁹	0
2.0000000000	2.0	(0.11)10 ⁻⁸	0
2.5000000000	2.5	(0.61)10 ⁻⁹	0
3.0000000000	3.0	(0.42)10 ⁻⁹	0
3.5000000000	3.5	(0.29)10 ⁻⁹	0
4.0000000000	4.0	(0.16)10 ⁻⁹	0
4.5000000000	4.5	(0.17)10 ⁻⁹	0
5.0000000000	5.0	(0.81)10 ⁻¹⁰	0
5.5000000000	5.5	(0.73)10 ⁻¹¹	0

Table 2

$A = (a_{ij}) = (i+j-1)^{-1}$, $i, j=1(1)5$ (Hilbert matrix of order 5).
 $b = (1, 1, 1, 1, 1)^T$.

COMPUTED SOLUTION	EXACT SOLUTION	EXP. REL. ERROR	ACTUAL REL. ERROR
5.000000730	5	(0.13)10 ⁻⁷	(0.15)10 ⁻⁶
-120.0000125	-120	(0.18)10 ⁻⁷	(0.10)10 ⁻⁶
630.0000512	630	(0.29)10 ⁻⁶	(0.81)10 ⁻⁷
-1120.000075	-1120	(0.92)10 ⁻⁵	(0.67)10 ⁻⁷
630.0000359	630	(0.25)10 ⁻⁵	(0.57)10 ⁻⁷

in the case of symmetric norms. In particular, if the underlying vector norm for the lub norm is the euclidean (2-norm) norm, then $\text{cond}(A) = \sigma_1/\sigma_n$ or $\text{cond}(A) = |\lambda_1|/|\lambda_n|$, if A is Hermitian (symmetric for the real case), where

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0 \text{ and } |\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n| > 0$$

are the singular values and eigenvalues of A respectively. $|A|$ denotes the matrix whose elements are the absolute values of the elements of the matrix A . Also, whenever 0 appears in the column 'Actual Rel. Error', it is inferred that the actual relative error is bounded by 2^{-t} , and for the Atlas computer on which the calculations were performed, $2^{-t} = 2^{-37} \doteq (7.275957614)10^{-12}$, taking into

Table 3

$A = L^{-1}HL$, $H = (h_{ij}) = (i+j-1)^{-1}$, $i, j=1(1)6$.
 $b = (38.43333333, -6.933333334, 2.307142857, -1.069841270, 0.6126984129, -0.4038961034)^T$.

COMPUTED SOLUTION	EXACT SOLUTION	EXP. REL. ERROR	ACTUAL REL. ERROR
0.9999999575	1	(0.66)10 ⁻⁸	(0.42)10 ⁻⁷
2.000001487	2	(0.34)10 ⁻⁸	(0.74)10 ⁻⁶
2.999986025	3	(0.35)10 ⁻⁶	(0.47)10 ⁻⁵
4.000069004	4	(0.13)10 ⁻³	(0.17)10 ⁻⁴
4.999764560	5	(0.78)10 ⁻⁴	(0.47)10 ⁻⁴
6.000635130	6	(0.10)10 ⁻²	(0.11)10 ⁻³

consideration the 'noise-level', although the mantissa on Atlas has 40 binary digits including the sign digit.

Table 1 is for a well-conditioned matrix, whilst Tables 2 and 3 are for very ill-conditioned matrices. In particular, Table 3 is for the Hilbert matrix of order 6 pre- and post-multiplied by

$$L^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 & 0 \\ 1 & -2 & 1 & 0 & 0 & 0 \\ -1 & 3 & -3 & 1 & 0 & 0 \\ 1 & -4 & 6 & -4 & 1 & 0 \\ -1 & 5 & -10 & 10 & -5 & 1 \end{bmatrix}$$

and

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 2 & 1 & 0 & 0 & 0 \\ 1 & 3 & 3 & 1 & 0 & 0 \\ 1 & 4 & 6 & 4 & 1 & 0 \\ 1 & 5 & 10 & 10 & 5 & 1 \end{bmatrix}$$

respectively. Table 3 illustrates the well-known result that a small residual vector, in this case $(0.0000000000, -0.0000000001, -0.0000000000, 0.0000000000, -0.0000000000, 0.0000000000)$, does not necessarily imply a highly accurate computed solution vector \bar{x} . On the other hand Table 4, again for a very ill-conditioned matrix, illustrates the fact that a large residual vector, in this case $(0.0000000000, 0.0000000000, 0.0000610352, 0.0312500000)$, does not necessarily imply an inaccurate computed solution vector \bar{x} .

In conclusion, the method is obviously more empirical rather than mathematical in the strict analytic sense. It must be emphasised again that the experimental relative error is not an upper bound of the actual relative error in \bar{x} ; it does, however, give a good approximation to the order of magnitude of the actual relative error. The tables given are only a small cross-section of many more examples that have been tested and yielded, in general, good results for the computed (experimental) component-wise relative error. The results appear to be more useful than those obtained from the usually quoted formulae (see Bauer, 1963)

$$\frac{\|\bar{x} - A^{-1}b\|}{\|A^{-1}b\|} \leq \text{cond}(A) \frac{\|b - A\bar{x}\|}{\|b\|}, \quad (3)$$

and (see, for example, Khabaza, 1964, p. 91)

$$\frac{\|\bar{x} - A^{-1}b\|}{\|A^{-1}b\|} \leq \text{cond}(A) \frac{\|\delta A\|}{\|A\|}, \quad (4)$$

References

- BAUER, F. L. (1966). Genauigkeitsfragen bei der Lösung linearer Gleichungssysteme, *Zeitschrift für Angewandte Mathematik und Mechanik*, Vol. 7, p. 409.
- BAUER, F. L. (1963). Optimally Scaled Matrices, *Numerische Mathematik*, Vol. 5, p. 73.
- BOWDLER, H. J., MARTIN, R. S., PETERS, G., and WILKINSON, J. H. (1966). Solution of Real and Complex Systems of Linear Equations, *Numerische Mathematik*, Vol. 8, p. 217.
- CHARTRES, B. A., and GEUDER, J. C. (1967). Computable Error Bounds for Direct Solution of Linear Equations, *J. Assoc. Comp. Mach.*, Vol. 14, p. 63.
- FORSYTHE, G. E. (1960). Crout with Pivoting, *Commun. Ass. Comp. Mach.*, Vol. 3, p. 507.
- INTERNATIONAL COMPUTERS AND TABULATORS LTD. (1965). The I.C.T. Atlas 1 Computer Programming Manual for Atlas Basic Language (ABL). I.C.T. publication No. CS 348A.
- KHABAZA, I. M. (1964). *A Course of Lectures on Matrix Computations*, London: Institute of Computer Science.
- WILKINSON, J. H. (1963). *Rounding Errors in Algebraic Processes*, London: Her Majesty's Stationery Office; New Jersey: Prentice Hall.
- WILKINSON, J. H. (1965). *The Algebraic Eigenvalue Problem*, London: Oxford Clarendon Press, pp. 189-264.

Table 4

$$A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 9.9 & 99.8 & 999.7 \\ 1 & 99.8 & 9999.9 & 999999.8 \\ 1 & 999.7 & 999999.8 & 999999999.9 \end{bmatrix}$$

$$b = (10, 4319, 4030199.5, 4003001999.4)^T.$$

COMPUTED SOLUTION	EXACT SOLUTION	EXP. REL. ERROR	ACTUAL REL. ERROR
0.9999999353	1	$(0.13)10^{-8}$	$(0.65)10^{-7}$
2.000000072	2	$(0.24)10^{-7}$	$(0.35)10^{-7}$
2.999999993	3	$(0.82)10^{-8}$	$(0.23)10^{-8}$
4.000000000	4	$(0.92)10^{-11}$	0

provided $\|A^{-1}\| \|\delta A\| \ll 1$, which give upper bounds for the relative error in terms of the condition number and the relative residuum for (3) and in terms of the condition number and the relative perturbation for (4). More important, the method does not require either the computation of the condition number or the residual vector or upper bounds for $|\delta A|$.

Acknowledgements

The author wishes to thank Mr. I. M. Khabaza, Director of the Computer Centre of Queen Mary College, London, for helpful discussions during the preparation of this work, and the Science Research Council for a maintenance grant. He also thanks the referee for his most helpful criticisms. Part of this work is from a Ph.D. Thesis approved by the University of London.