

Backward Error Analysis for Gaussian Elimination

A. K. Cline

Department of Computer Sciences

University of Texas at Austin

Austin, Texas 78712

The purpose of this note is to show that the Gaussian elimination algorithm when performed on a matrix A with floating point arithmetic can be shown to result in the form

$$LU = A + E$$

where L is unit, lower triangular, U is upper triangular, and E is a matrix of perturbations. Formally written, the algorithm is

$$a_{i,j}^1 = a_{i,j}, \text{ for } i, j = 1, \dots, n \quad (1)$$

and for $k = 1, \dots, n-1$

$$u_{k,j} = a_{k,j}^k, \text{ for } j = k, \dots, n. \quad (2)$$

$$l_{i,k} = a_{i,k}^k / u_{k,k}, \quad (3)$$

and

$$a_{i,j}^{k+1} = a_{i,j}^k - l_{i,k} u_{k,j}. \quad (4)$$

Since floating point arithmetic is being employed, equations (3) and (4) are actually

$$l_{i,k} = \text{fl}(a_{i,k}^k / u_{k,k}), \quad (5)$$

and

$$a_{i,j}^{k+1} = \text{fl}(a_{i,j}^k - l_{i,k} u_{k,j}). \quad (6)$$

(Notice that the floating point arithmetic has no affect on equation (2) .) Using the model of errors introduced by this arithmetic, these equations could be recast as

$$l_{i,k} = a_{i,k}^k / u_{k,k} (1 + \mathbf{d}_1), \quad (7)$$

and

$$a_{i,j}^{k+1} = (a_{i,j}^k - l_{i,k} u_{k,j} (1 + \mathbf{d}_2)) (1 + \mathbf{d}_3). \quad (8)$$

where $|\mathbf{d}_1|, |\mathbf{d}_2|, \text{ and } |\mathbf{d}_3| \leq \mathbf{e}_0$, the machine unit floating point precision. (The dependency of the \mathbf{d} s upon $i, j, \text{ and } k$ has been omitted from the notation for clarity.)

Consider the error in “backward fashion” (i.e., the operations were performed correctly on perturbed data), equations (7) and (8) could be stated as

$$l_{i,k} = (a_{i,k}^k + \mathbf{e}_{i,k}^k) / u_{k,k}, \quad (9)$$

and

$$a_{i,j}^{k+1} = a_{i,j}^k + \mathbf{e}_{i,j}^k - l_{i,k} u_{k,j}. \quad (10)$$

These equations actually define the perturbations $\mathbf{e}_{i,k}^k$ and $\mathbf{e}_{i,j}^k$, thus

$$\mathbf{e}_{i,k}^k = l_{i,k} u_{k,k} - a_{i,k}^k, \quad (11)$$

and

$$\mathbf{e}_{i,j}^k = l_{i,k} u_{k,j} - a_{i,j}^k + a_{i,j}^{k+1}. \quad (12)$$

Equating (7) and (9), we obtain

$$\mathbf{e}_{i,k}^k = \mathbf{d}_1 a_{i,k}^k, \quad (13)$$

and we may conclude that

$$|\mathbf{e}_{i,k}^k| \leq \mathbf{e}_0 |a_{i,k}^k|. \quad (14)$$

Notice from (8) that

$$l_{i,k} u_{k,j} = (a_{i,j}^k - a_{i,j}^{k+1}) / (1 + \mathbf{d}_3) / (1 + \mathbf{d}_2), \quad (15)$$

so (12) could be rewritten as

$$\mathbf{e}_{i,j}^k = a_{i,j}^k (1 / (1 + \mathbf{d}_2) - 1) - a_{i,j}^{k+1} (1 / ((1 + \mathbf{d}_3)(1 + \mathbf{d}_2)) - 1). \quad (16)$$

Ignoring terms of size \mathbf{e}_0^2 , this is

$$\mathbf{e}_{i,j}^k = a_{i,j}^k (-\mathbf{d}_2) - a_{i,j}^{k+1} (-\mathbf{d}_2 - \mathbf{d}_3), \quad (17)$$

so

$$|\mathbf{e}_{i,j}^k| \leq \mathbf{e}_0 |a_{i,j}^k| + 2\mathbf{e}_0 |a_{i,j}^{k+1}| \leq 3\mathbf{e}_0 \max\{|a_{i,j}^k|, |a_{i,j}^{k+1}|\}. \quad (18)$$

If we sum equation (10) over values of k from 1 through $i-1$ and assume $i \leq j$, we obtain

$$a_{i,j}^2 + \dots + a_{i,j}^i = a_{i,j}^1 + \dots + a_{i,j}^{i-1} + (\mathbf{e}_{i,j}^1 + \dots + \mathbf{e}_{i,j}^{i-1}) - (l_{i,1} u_{1,j} + \dots + l_{i,i-1} u_{i-1,j}), \quad (19)$$

and after canceling identical terms

$$a_{i,j}^i = a_{i,j}^1 + (\mathbf{e}_{i,j}^1 + \dots + \mathbf{e}_{i,j}^{i-1}) - (l_{i,1} u_{1,j} + \dots + l_{i,i-1} u_{i-1,j}), \quad (20)$$

Noticing that $a_{i,j}^i = u_{i,j} = l_{i,i} u_{i,j}$ (since the diagonal elements of L are one) and that $a_{i,j}^1 = a_{i,j}$, we have

$$\sum_{k=1}^i l_{i,k} u_{k,j} = a_{i,j} + e_{i,j}, \quad (21)$$

where $e_{i,j} = \sum_{k=1}^{i-1} \mathbf{e}_{i,j}^k$. Similarly for $j < i$, if we sum (10) for values of k from 1 through $j-1$, we obtain

$$a_{i,j}^2 + \dots + a_{i,j}^j = a_{i,j}^1 + \dots + a_{i,j}^{j-1} + (\mathbf{e}_{i,j}^1 + \dots + \mathbf{e}_{i,j}^{j-1}) - (l_{i,1} u_{1,j} + \dots + l_{i,j-1} u_{j-1,j}). \quad (22)$$

Again, after canceling identical terms, we have

$$a_{i,j}^j = a_{i,j}^1 + (\mathbf{e}_{i,j}^1 + \dots + \mathbf{e}_{i,j}^{j-1}) - (l_{i,1} u_{1,j} + \dots + l_{i,j-1} u_{j-1,j}). \quad (23)$$

From (9) we have $l_{i,j} u_{j,j} = a_{i,j}^j + \mathbf{e}_{i,j}^j$, and thus

$$\sum_{k=1}^j l_{i,k} u_{k,j} = a_{i,j} + e_{i,j}, \quad (24)$$

where $e_{i,j} = \sum_{k=1}^j \mathbf{e}_{i,j}^k$.

Notice that for $i \leq j$, from (18)

$$|e_{i,j}| \leq \sum_{k=1}^{i-1} |\mathbf{e}_{i,j}^k| \leq 3(i-1)\mathbf{e}_0 \max\{|a_{i,j}^1|, \dots, |a_{i,j}^i|\}, \quad (25)$$

and for $j < i$, from (18) and (14)

$$|e_{i,j}| \leq \sum_{k=1}^j |e_{i,j}^k| \leq 3(j-1)\mathbf{e}_0 \max\{|a_{i,j}^1|, \dots, |a_{i,j}^j|\} + \mathbf{e}_0 |a_{i,j}^j|. \quad (26)$$

Thus, in general

$$|e_{i,j}| \leq 3\mathbf{e}_0 \min\{i-1, j\} \max\{|a_{i,j}^1|, \dots, |a_{i,j}^{\min\{i,j\}}|\}, \quad (27)$$

for $i, j = 1, \dots, n$.

A slightly different approach results in the production of the factorization

$$L^{n-1} \dots L^2 L^1 (A + E) = U$$

where A, E , and U are exactly as before and $L^{n-1} \dots L^2 L^1$ is inverse of L . This variant of the algorithm is more efficient when row interchanges are introduced since only portions of rows are interchanged here whereas entire rows must be interchanged with the LU factorization. (As before however the actual permutations will be ignored to avoid the great complexity they add to the notation.)

Equation (9) could be also be stated as

$$l_{i,k} a_{k,k}^k = a_{i,k}^k + \mathbf{e}_{i,k}^k \quad (28)$$

which is

$$-l_{i,k} a_{k,k}^k + 1(a_{i,k}^k + \mathbf{e}_{i,k}^k) = 0 \quad (29)$$

and equation (10) could be stated as

$$-l_{i,k} a_{k,j}^k + 1(a_{i,j}^k + \mathbf{e}_{i,j}^k) = a_{i,j}^{k+1}. \quad (30)$$

By defining

$$a_{i,j}^{k+1} = a_{i,j}^k \quad (31)$$

for $i = 1, \dots, k$ and $j = 1, \dots, n$ and for $i = 1, \dots, n$ and $j = 1, \dots, k$, and

$$\mathbf{e}_{i,j}^k = 0 \quad (32)$$

for $i = 1, \dots, k-1$ and $j = 1, \dots, n$ and for $i = 1, \dots, n$ and $j = 1, \dots, k$, and constructing the lower triangular matrix L^k with unit diagonal, $-l_{i,k}$ in the i,k position, and zeroes elsewhere, we have

$$L^k (A^k + E^k) = A^{k+1}, \quad (33)$$

for $k = 1, \dots, n-1$. Using the fact that

$$L^{k-1} \dots L^2 L^1 E^k = E^k \quad (34)$$

(which follows from the structure of the L^i matrices) a simple induction argument shows that

$$L^k \dots L^2 L^1 (A^1 + E^1 + \dots + E^k) = A^{k+1} \quad (35)$$

for $k = 1, \dots, n-1$. By defining $E = E^1 + \dots + E^{n-1}$ (which is the same E as in the previous analysis), we have

$$L^{n-1} \dots L^2 L^1 (A + E) = A^n, \quad (36)$$

but from (2) and (30) we see that $U = A^n$ and this is the factorization we sought.