

## Unsymmetric Positive Definite Linear Systems

Gene H. Golub  
*Department of Computer Science*  
*Stanford University*  
*Stanford, California 94305*

and

Charles Van Loan  
*Department of Computer Science*  
*Cornell University*  
*Ithaca, New York 14853*

Dedicated to Alston S. Householder  
on the occasion of his seventy-fifth birthday.

Submitted by G. W. Stewart

---

### ABSTRACT

Is it necessary to pivot when solving an unsymmetric positive definite linear system  $Ax = b$ ? Define  $T = (A + A^T)/2$  and  $S = (A - A^T)/2$ . It is shown that pivoting is unnecessary if the quantity  $\|ST^{-1}S\|_2/\|A\|_2$  is suitably small with respect to the working machine precision.

---

### 1. INTRODUCTION

An  $n \times n$  real matrix  $A$  is positive definite if  $x^T Ax > 0$  for all nonzero  $x$  in  $R^n$ . Setting  $A = T + S$ , where

$$T = (A + A^T)/2, \quad (1.1)$$

“the symmetric part of  $A$ ,” and

$$S = (A - A^T)/2 \quad (1.2)$$

“the skew-symmetric part of  $A$ ,” we see that  $A$  is positive definite if and only if  $T$  is positive definite. Such matrices have a number of important properties

[4]. In particular, if  $A$  is positive definite, then each of its principal submatrices is nonsingular and therefore the factorization

$$\begin{aligned} A &= LDM^T, \\ L, M &\text{ unit lower triangular,} \\ D &= \text{diag}(d_1, \dots, d_n) \end{aligned} \tag{1.3}$$

exists.

The  $LDM^T$  factorization is a slight rearrangement of the more familiar  $LU$  factorization ( $U = DM^T$ ) and can be computed in  $n^3/3$  operations using either Gaussian elimination or Crout reduction techniques [6, p. 131 ff.]. Denote computed quantities with the "hat" notation " $\hat{\cdot}$ ," and suppose  $\hat{L}$ ,  $\hat{D}$ , and  $\hat{M}$  are used to solve  $Ax = b$ . Using the general roundoff error analysis of deBoor and Pinkus [3], it can be shown that the computed solution  $\hat{x}$  satisfies

$$(A + E)\hat{x} = b, \tag{1.4}$$

where

$$|E| \leq uc_n |\hat{L}| |\hat{D}| |\hat{M}^T|. \tag{1.5}$$

Here,  $u$  is the machine precision (i.e.,  $\beta^{1-t}$ , where  $t$  digit, base  $\beta$  arithmetic is used), and  $c_n$  is a constant which depends linearly on  $n$  and whose exact value depends upon the details of the algorithm. The matrix absolute value  $|\cdot|$  is defined by  $|Z| = (|z_{ij}|)$ , and the relation  $|Z| \leq |Y|$  implies  $|z_{ij}| \leq |y_{ij}|$  for all  $i$  and  $j$ .

For general matrices it is customary to pivot and thereby determine the  $LDM^T$  factorization of a row permuted version of  $A$ . The resulting matrix  $|\hat{L}| |\hat{D}| |\hat{M}^T|$  can then be suitably bounded, implying from (1.4) and (1.5) that  $\hat{x}$  satisfies a "nearby" linear system.

For several important classes of matrices, however, pivoting is not a numerical necessity. This is always a welcome state of affairs from the standpoint of program simplicity and efficiency. For example, if  $A$  is both positive definite and symmetric, then (a)  $M = L$  in (1.3), (b) computational requirements are reduced to  $n^3/6$ , and (c) the matrix  $E$  in (1.4) can be shown to satisfy

$$\|E\|_\infty \leq uq_n \|A\|_\infty, \tag{1.6}$$

where  $q_n$  is a constant quadratic in  $n$  and  $\|A\|_\infty = \max_i \sum_j |a_{ij}|$ . See [3] for details. Since rounding errors of order  $u\|A\|_\infty$  are usually present in  $A$  to

begin with, it can be concluded from (1.4) and (1.6) that pivoting is unnecessary when solving *symmetric* positive definite linear systems.

Is pivoting necessary when  $A$  is positive definite but unsymmetric? Buckley [1,2] considered this question for positive definite systems of the form

$$(I+S)x = b, \quad S^T = -S. \quad (1.7)$$

These systems arise from the discretization of certain differential operators. It turns out that if

$$(I+S) = LU$$

is the  $LU$  factorization of  $(I+S)$ , then the pivots  $u_{kk}$  satisfy

$$1 \leq u_{kk} \leq 1 + \|S\|_2^2, \quad (1.8)$$

where the 2-norm of a matrix  $Z$  is defined by

$$\|Z\|_2^2 = \text{the largest eigenvalue of } Z^T Z.$$

Since the pivots are nicely bounded away from zero, Buckley concluded that it is safe not to pivot when applying Gaussian elimination to (1.7).

This conclusion has to be qualified for general positive definite systems, as the example

$$A = \begin{bmatrix} \epsilon & 1 \\ -1 & \epsilon \end{bmatrix} = \begin{bmatrix} \epsilon & 0 \\ 0 & \epsilon \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \quad (1 \gg \epsilon > 0)$$

indicates. This matrix has a perfect condition number in the 2-norm,

$$\|A\|_2 \|A^{-1}\|_2 = (\epsilon^2 + 1)^{1/2} (\epsilon^2 + 1)^{-1/2} = 1,$$

and therefore, if Gaussian elimination with partial pivoting is used to solve  $Ax = b$ , the solution will be correct to nearly full working precision. However, if the solution is obtained via the factorization

$$\begin{bmatrix} \epsilon & 1 \\ -1 & \epsilon \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ -1/\epsilon & 1 \end{bmatrix} \begin{bmatrix} \epsilon & 0 \\ 0 & \epsilon - 1/\epsilon \end{bmatrix} \begin{bmatrix} 1 & 1/\epsilon \\ 0 & 1 \end{bmatrix},$$

then rounding errors of order  $u/\epsilon$  can be expected to contaminate the result.

The problem, of course, is the presence of large entries in  $L$ ,  $D$ , and  $M$ . In the next section we show that this can only occur if  $\|ST^{-1}S\|_2$  is large, where the matrices  $T$  and  $S$  are given by (1.1) and (1.2) respectively. Notice that in the above example, this quantity has the value of  $1/\epsilon$ .

## 2. THEORETICAL BOUNDS

In this section we put finite precision arithmetic aside and examine how big the factors  $L$ ,  $D$ , and  $M$  can get when  $A = LDM^T$  is positive definite. Ignoring the distinction between exact and computed quantities, we conclude from (1.5) that a bound on  $|L||D||M^T|$  is what we need. For this purpose we shall use the Frobenius norm:

$$\|Z\|_F^2 = \sum_i \sum_j |z_{ij}|^2.$$

**THEOREM.** *Let  $A$  be an  $n \times n$  positive definite matrix having the factorization*

$$A = LDM^T,$$

where  $L$  and  $M$  are unit lower triangular and  $D = \text{diag}(d_1, \dots, d_n)$ . If

$$T = (A + A^T)/2 \quad \text{and} \quad S = (A - A^T)/2,$$

then

$$\||L||D||M^T|\|_F \leq n[\|T\|_2 + \|ST^{-1}S\|_2]. \quad (2.1)$$

*Proof.* Let the Cholesky decomposition of the symmetric positive definite matrix  $T$  be given by

$$T = GG^T, \quad G \text{ lower triangular,}$$

and partition  $L^{-1}$  and  $M^{-1}$  into their respective rows as follows:

$$L^{-1} = \begin{bmatrix} y_1^T \\ \vdots \\ y_n^T \end{bmatrix}, \quad M^{-1} = \begin{bmatrix} z_1^T \\ \vdots \\ z_n^T \end{bmatrix}.$$

From the equation  $LDM^T = GG^T + S$  we find

$$DM^TL^{-T} = (G^TL^{-T})^T(G^TL^{-T}) + L^{-1}SL^{-T},$$

where  $L^{-T}$  denotes the inverse of the transpose of  $L$ . Since  $M^TL^{-T}$  is unit upper triangular and since the diagonal of the skew symmetric matrix  $L^{-1}SL^{-T}$  is zero, it follows by comparing  $(k, k)$  entries that

$$d_k = \|G^T y_k\|_2^2, \quad k = 1, \dots, n. \quad (2.2)$$

Thus, if  $D^{1/2} = \text{diag}(\sqrt{d_1}, \dots, \sqrt{d_n})$  we obtain

$$\|G^TL^{-T}D^{-1/2}\|_F^2 = \sum_{k=1}^n \left[ \frac{\|G^T y_k\|_2}{\sqrt{d_k}} \right]^2 = n.$$

Similarly,

$$M^{-1}LD = (G^TM^{-T})^T(G^TM^{-T}) + M^{-1}SM^{-T}$$

implies

$$\|G^TM^{-T}D^{-1/2}\|_F^2 = n.$$

Now since  $(LD^{1/2})(D^{1/2}M^T) = GG^T + S$ , we have the identities

$$LD^{1/2} = (G + SG^{-T})(G^TM^{-T}D^{-1/2})$$

and

$$D^{1/2}M^T = (G^TL^{-T}D^{-1/2})^T(G^T + G^{-1}S).$$

Taking norms and using the inequalities  $\|YZ\|_F \leq \|Y\|_F \|Z\|_2$  and  $\|YZ\|_F \leq \|Y\|_2 \|Z\|_F$ , we have

$$\|LD^{1/2}\|_F \leq \|G + SG^{-T}\|_2 \|G^TM^{-T}D^{-1/2}\|_F$$

$$\leq \sqrt{n} \|G + SG^{-T}\|_2$$

and

$$\begin{aligned} \|D^{1/2}M^T\|_F &\leq \|(G^T L^{-T} D^{-1/2})^T\|_F \|G^T + G^{-1}S\|_2 \\ &\leq \sqrt{n} \|G^T + G^{-1}S\|_2, \end{aligned}$$

which imply

$$\begin{aligned} \| |L| |D| |M^T| \|_F &= \| |LD^{1/2}| |D^{1/2}M^T| \|_F \\ &\leq \|LD^{1/2}\|_F \|D^{1/2}M^T\|_F \\ &\leq n \|G + SG^{-T}\|_2 \|G^T + G^{-1}S\|_2. \end{aligned}$$

The theorem now follows by using the identities  $\|Z\|_2^2 = \|ZZ^T\|_2 = \|Z^T Z\|_2$  to prove

$$\begin{aligned} \|G + SG^{-T}\|_2^2 &= \|(G + SG^{-T})(G + SG^{-T})^T\|_2 \\ &= \|(G + SG^{-T})(G^T - G^{-1}S)\|_2 \\ &= \|GG^T - SG^{-T}G^{-1}S\|_2 \\ &= \|T - ST^{-1}S\|_2 \\ &\leq \|T\|_2 + \|ST^{-1}S\|_2 \end{aligned} \tag{2.3}$$

and similarly

$$\|G^T + G^{-1}S\|_2^2 \leq \|T\|_2 + \|ST^{-1}S\|_2. \quad \blacksquare \tag{2.4}$$

What this theorem implies for positive definite linear system solvers will be discussed in the next section. However, as affirmed by our earlier  $2 \times 2$  example, it is clear that a large  $\|ST^{-1}S\|_2$  may imply unacceptably large  $|L||D||M^T|$  and a consequent need for pivoting.

We conclude this section by generalizing Buckley's result (1.8) concerning the pivots  $d_k$  (his  $u_{kk}$ ).

COROLLARY.

$$\frac{1}{\|T^{-1}\|_2} \leq d_k \leq \|T\|_2 + \|ST^{-1}S\|_2 \quad (k=1, \dots, n). \tag{2.5}$$

*Proof.* Since  $L^{-T}$  is a unit upper triangular matrix, its  $k$ th column,  $y_k$ , has a 2-norm of at least unity. Hence, from (2.2) we obtain

$$\begin{aligned} d_k &= \|G^T y_k\|_2^2 = y_k^T T y_k \\ &> \frac{\|y_k\|_2^2}{\|T^{-1}\|_2} > \frac{1}{\|T^{-1}\|_2}. \end{aligned}$$

To establish the upper bound in (2.5), we once again use the identity

$$\begin{aligned} LD^{1/2} &= (GG^T + S)M^{-T}D^{-1/2} \\ &= (G + SG^{-T})(G^T M^{-T} D^{-1/2}). \end{aligned}$$

Let  $e_k$  denote the  $k$ th column of the identity. Using (2.2) and (2.3), we find

$$\begin{aligned} d_k &\leq \|(LD^{1/2})e_k\|_2^2 = \|(G + SG^{-T})(G^T M^{-T} D^{-1/2} e_k)\|_2^2 \\ &\leq \|G + SG^{-T}\|_2^2 \|(G^T M^{-T})e_k\|_2^2 / d_k \\ &\leq \|T\|_2 + \|ST^{-1}S\|_2. \quad \blacksquare \end{aligned}$$

### 3. PRACTICAL OBSERVATIONS

If we assume that

$$\| |\hat{L}| |\hat{D}| |\hat{M}^T| \|_F \doteq \| |L| |D| |M^T| \|_F, \quad (3.1)$$

then (1.5), (2.1), and the inequality  $\|T\|_2 \leq \|A\|_2$  imply that the computed solution  $\hat{x}$  satisfies

$$(A + E)\hat{x} = b, \quad (3.2)$$

where

$$\|E\|_F \leq \text{unc}_n[\|A\|_2 + \|ST^{-1}S\|_2] \quad (3.3)$$

The symbols “ $\doteq$ ” and “ $\leq$ ” remind us that these are heuristic results. By



setting  $S=0$  in (3.3) we “essentially” obtain the same bound that Wilkinson [7] derived in his classical roundoff error analysis of the Cholesky factorization  $A = GG^T$ .

Now from the usual sensitivity analysis of linear systems [6, p. 194 ff.] we know that if  $Ax = b$  and  $(A + F)y = b$ , then roughly speaking

$$\frac{\|x - y\|_2}{\|x\|_2} < \frac{\|F\|_2}{\|A\|_2} \kappa_2(A)$$

where

$$\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2$$

is the condition of  $A$ . Hence, from (3.2) and (3.3) we obtain the key result

$$\frac{\|x - \hat{x}\|_2}{\|x\|_2} < unc_n \left[ 1 + \frac{\|ST^{-1}S\|_2}{\|A\|_2} \right] \kappa_2(A). \quad (3.4)$$

On the other hand, the corresponding result for  $\hat{x}_{\text{piv}}$ , the computed solution obtained via Gaussian elimination with pivoting, is roughly of the form

$$\frac{\|x - \hat{x}_{\text{piv}}\|_2}{\|x\|_2} < u\rho n^3 \kappa_2(A), \quad (3.5)$$

where  $\rho$  is the “growth factor.” Assuming  $\rho \neq 1$  and ignoring the unimportant factors  $nc_n$  and  $n^3$ , we see that the main distinction between (3.4) and (3.5) is the factor

$$\Omega = \frac{\|ST^{-1}S\|_2}{\|A\|_2}.$$

Based on dozens of “randomly generated” examples, our experience has been that it is safe not to pivot in problems having a small  $\Omega$ . In other words, heuristic (3.4) can be trusted in practice. However, if  $\Omega$  is large, then it is more difficult to draw conclusions from (3.4). A large  $\Omega$  may certainly imply a need for pivoting, as the  $2 \times 2$  example of Section 1 indicates. But to show that an accurate  $x$  can be obtained despite a large  $\Omega$ , we consider the linear system

$$Ax = (GG^T + S)x = b, \quad (3.6)$$



where

$$G = (g_{ij}), \quad g_{ij} = \begin{cases} 0, & i < j, \\ 1, & i = j, \\ -1, & i > j, \end{cases}$$

$$S = (s_{ij}), \quad s_{ij} = \begin{cases} -10 \times (-1)^{j+i(i-1)/2}, & i < j, \\ 0, & i = j, \\ 10 \times (-1)^{j+i(i-1)/2}, & i > j, \end{cases}$$

and  $b$  is chosen such that the exact solution is given by  $x = (1, 1, \dots, 1)^T$ . The system has the property that

$$\Omega = O(4^n),$$

and so by increasing the dimension  $n$  we can arbitrarily increase the bound in (3.4). Using IBM/370 short precision arithmetic ( $u = 10^{-6}$ ), we solved (3.6) for  $n = 5, 10, 15, 20,$  and  $25$ . Table 1 summarizes our findings.

TABLE 1

(a)	(b)	(c)	(d)	(e)	(f)
$n$	$\frac{\ x - \hat{x}\ _2}{\ x\ _2}$	$u\Omega\kappa_2(A)$	$u\kappa_2(A)$	$\   \hat{L}   \hat{D}   \hat{M}^T  \ _F$	Range of $\hat{d}_k$
5	$4 \times 10^{-6}$	$1 \times 10^{-2}$	$2 \times 10^{-5}$	$8 \times 10^2$	$1-10^2$
10	$1 \times 10^{-5}$	$2 \times 10^1$	$3 \times 10^{-5}$	$2 \times 10^5$	$1-10^2$
15	$2 \times 10^{-5}$	$3 \times 10^4$	$5 \times 10^{-5}$	$3 \times 10^5$	$1-10^2$
20	$2 \times 10^{-5}$	$4 \times 10^7$	$7 \times 10^{-5}$	$5 \times 10^5$	$1-10^2$
25	$4 \times 10^{-5}$	$6 \times 10^{10}$	$1 \times 10^{-4}$	$8 \times 10^5$	$1-10^2$

We have tabulated  $u\Omega\kappa_2(A)$  and  $u\kappa_2(A)$  because these are the key factors in the upper bounds of (3.4) and (3.5) respectively. It is clear from columns (b) and (c) in the table that the upper bound in (3.4) is extremely pessimistic. For all the values of  $n$  selected, the matrix  $A$  is fairly well conditioned, as evidenced by column (d), which indicates the error which can be expected when Gaussian elimination with partial pivoting is applied to (3.6). There is no way we know to explain why such good solutions were obtained without pivoting, in view of the largeness of  $\Omega$  and the numbers in column (e). Perhaps a more refined analysis would reveal that the size of the pivots has a

critical bearing on the relative error, for as we see in these examples, the  $d_k$  are rather nicely behaved.

As a final example, we demonstrate how our analysis can be applied to the linear system which arises from the discretization of the boundary value problem

$$y''(x) = p(x)y'(x) + q(x)y(x) + r(x) \quad (-1 \leq x \leq 1),$$

$$y(-1) = \alpha, \quad y(1) = \beta, \quad q(x) \geq 0.$$

Although this is a simple problem, it illustrates how an estimate of  $\Omega$  can be obtained from information relating to the underlying differential equation.

Set  $h = 2/(n+1)$  and  $x_i = a + ih$ ,  $i = 0, \dots, n+1$ . Following Ortega [5, p. 96 ff.], if we replace the above derivatives with appropriate divided differences, then we are led to a linear system  $Ay = d$  in the unknowns  $y_i \doteq y(x_i)$  ( $i = 1, \dots, n$ ), where

$$A = \begin{bmatrix} a_1 & -c_1 & & & \circ \\ -b_2 & a_2 & -c_2 & & \\ & \ddots & \ddots & \ddots & \\ \circ & & -b_{n-1} & a_{n-1} & -c_{n-1} \\ & & & b_n & a_n \end{bmatrix}$$

and

$$a_i = 2 + h^2 q_i, \quad q_i = q(x_i),$$

$$b_i = 1 + p_i h/2, \quad p_i = p(x_i),$$

$$c_i = 1 - p_i h/2.$$

The symmetric matrix  $T = (A + A^T)/2$  has the form  $T = T_1 + T_2$ , where

$$T_1 = \begin{bmatrix} a_1 & -1 & & & \circ \\ -1 & a_2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ \circ & & & -1 & a_n \end{bmatrix}$$

and

$$T_2 = \begin{bmatrix} 0 & w_2 & & & \circ \\ w_2 & 0 & w_3 & & \\ & \ddots & \ddots & \ddots & \\ & & & & w_n \\ \circ & & & w_n & 0 \end{bmatrix},$$

$$w_i = \frac{-h}{4}(p_i - p_{i-1}),$$

while the skew symmetric matrix  $S = (A - A^T)/2$  has the form

$$S = \begin{bmatrix} 0 & s_2 & & & \circ \\ -s_2 & 0 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & & 0 & s_n \\ \circ & & & -s_n & 0 \end{bmatrix},$$

$$s_i = (p_{i-1} + p_i)h/2.$$

Now if we assume that  $p'(x)$  is continuous on  $[-1, 1]$  and that

$$\max_{-1 < x < 1} |p'(x)| = p'_{\max} \leq 2, \tag{3.7}$$

then

$$|w_i| = \frac{h^2}{4} \left| \frac{p_i - p_{i-1}}{h} \right| \leq \frac{h^2}{4} p'_{\max} \leq \frac{h^2}{2},$$

and thus

$$\|T_2\|_2 \leq h^2.$$

Since the smallest eigenvalue of  $T_1$  is bounded below,

$$\lambda_{\min}(T_1) \geq 2 - 2\cos\left(\frac{\pi h}{2}\right) \geq 2h^2 \quad (h < 1),$$

we obtain the following lower bound for the smallest eigenvalue of  $T$ :

$$\lambda_{\min}(T) \geq \lambda_{\min}(T_1) - \|T_2\|_2 \geq 2h^2 - h^2 = h^2.$$

Thus,  $T = (A + A^T)/2$  is positive definite, and moreover,

$$\|T^{-1}\|_2 \leq 1/h^2. \quad (3.8)$$

Setting

$$p_{\max} = \max_{-1 < x < 1} |p(x)|,$$

it is easy to verify that

$$\|S\|_2 \leq 2hp_{\max},$$

and therefore

$$\begin{aligned} \Omega &\leq \frac{\|ST^{-1}S\|_2}{\|A\|_2} \leq \|S\|_2^2 \|T^{-1}\|_2 \\ &\leq (2p_{\max})^2. \end{aligned}$$

Thus, if the function  $p(x)$  is both smooth enough and small enough, the matrix  $A$  is positive definite and pivoting is unnecessary.

*The authors are grateful to Professors Jim Ortega and Olof Widlund for sharing their thoughts on the problem discussed in this paper. The work of G.H.G. was supported by the Department of Energy.*

#### REFERENCES

- 1 A. Buckley, A note on matrices  $A = I + H$ ,  $H$  skew-symmetric, *Z. Angew. Math. Mech.* 54 (2):125-126 (1974).
- 2 A. Buckley, On the solution of certain skew-symmetric linear systems, *SIAM J. Numer. Anal.* 14:566-570 (1977).
- 3 C. deBoor and A. Pinkus, Backward error analysis for totally positive linear systems, *Numer. Math.* 27:485-490 (1977).



- 4 K. Fan, On real matrices with positive definite symmetric component, *Linear and Multilinear Algebra* 1:1-4 (1973).
- 5 J. Ortega, *Numerical Analysis—A Second Course*, Academic, New York, 1972.
- 6 G. W. Stewart, *Introduction to Matrix Computations*, Academic, New York, 1973.
- 7 J. H. Wilkinson, A priori error analysis of algebraic processes, in *Proceedings of the International Congress of Mathematics (Moscow, 1968)*, Izdat. Mir., Moscow, 1968, pp. 629-639.

*Received 19 September 1978*