



THE UNIVERSITY
of MANCHESTER



**Stability of the
Diagonal Pivoting Method
with Partial Pivoting**

N. J. Higham

Numerical Analysis Report No. 265

July 1995

Manchester Centre for Computational Mathematics
Numerical Analysis Reports

DEPARTMENTS OF MATHEMATICS

Reports available from: And over the World-Wide Web from URLs
Department of Mathematics <http://www.ma.man.ac.uk/MCCM/MCCM.html>
University of Manchester <ftp://vtx.ma.man.ac.uk/pub/narep>
Manchester M13 9PL
England

Stability of the Diagonal Pivoting Method with Partial Pivoting

Nicholas J. Higham*

July 16, 1995

Abstract

LAPACK and LINPACK both solve symmetric indefinite linear systems using the diagonal pivoting method with the partial pivoting strategy of Bunch and Kaufman (1977). No proof of the stability of this method has appeared in the literature. It is tempting to argue that the diagonal pivoting method is stable for a given pivoting strategy if the growth factor is small. We show that this argument is false in general, and give a sufficient condition for stability. This condition is not satisfied by the partial pivoting strategy, because the multipliers are unbounded. Nevertheless, using a more specific approach we are able to prove the stability of partial pivoting, thereby filling a gap in the body of theory supporting LAPACK and LINPACK.

Key words. symmetric indefinite matrix, diagonal pivoting method, LDL^T factorization, partial pivoting, growth factor, numerical stability, rounding error analysis, LAPACK, LINPACK.

AMS subject classifications. primary 65F05, 65G05

1 Introduction

LAPACK is renowned for the numerical reliability of the algorithms it employs. The *LAPACK Users' Guide* [1] states that “almost all the algorithms in LAPACK (as well as LINPACK and EISPACK) are [normwise backward] stable” [1, p. 74], and the algorithms not covered by this statement are known to be stable in appropriately weakened senses. The analyses to back up these claims of stability are spread

*Department of Mathematics, University of Manchester, Manchester, M13 9PL, England (na.nhigham@na-net.ornl.gov). This work was supported by Engineering and Physical Sciences Research Council grants GR/H5213 and GR/H/94528.

that searches only two columns at each stage and so requires only $O(n^2)$ comparisons. The LAPACK driver routines `xSYSV` (simple) and `xSYSVX` (expert) and the LINPACK routines `xSIFA/xSISL` all use the diagonal pivoting method with partial pivoting to solve a linear system with a symmetric (indefinite) coefficient matrix.

To describe the partial pivoting strategy it suffices to define the pivot choice for the first stage of the factorization. Recall that s denotes the size of the pivot block.

Algorithm 1 (Bunch–Kaufman Partial Pivoting Strategy) This algorithm determines the pivot for the first stage of the diagonal pivoting method with partial pivoting applied to a symmetric matrix $A \in \mathbb{R}^{n \times n}$.

- $$\alpha := (1 + \sqrt{17})/8 \ (\approx 0.64)$$
- $$\lambda := \|A(2:n, 1)\|_\infty$$
- If $\lambda = 0$ there is nothing to do on this stage of the elimination.
- $$r := \min\{i \geq 2: |a_{i1}| = \lambda\}$$
- if $|a_{11}| \geq \alpha\lambda$
- (1) $s = 1, \Pi = I$
- else
- $$\sigma := \left\| \begin{bmatrix} A(1:r-1, r) \\ A(r+1:n, r) \end{bmatrix} \right\|_\infty$$
- if $|a_{11}|\sigma \geq \alpha\lambda^2$
- (2) $s = 1, \Pi = I$
- else if $|a_{rr}| \geq \alpha\sigma$
- (3) $s = 1$ and choose Π to swap rows and columns 1 and r .
- else
- (4) $s = 2$ and choose Π to swap rows and columns 2 and r ,
so that $|(\Pi A \Pi^T)_{21}| = \lambda$.
- end
- end

To understand the partial pivoting strategy it helps to consider the matrix

$$\begin{bmatrix} a_{11} & \dots & \lambda & \dots & \dots & \dots \\ \vdots & & \vdots & & & \\ \lambda & \dots & a_{rr} & \dots & \sigma & \dots \\ \vdots & & \vdots & & & \\ \vdots & & \sigma & & & \\ \vdots & & \vdots & & & \end{bmatrix},$$

and to note that the pivot is one of a_{11} , a_{rr} and $\begin{bmatrix} a_{11} & \lambda \\ \lambda & a_{rr} \end{bmatrix}$ (or, rather, since $\lambda = |a_{r1}|$, this matrix with λ replaced by a_{r1}).

The value of the constant $\alpha = (1 + \sqrt{17})/8$ is determined by regarding α as a free parameter and equating a bound for the element growth over two $s = 1$ stages to a bound for the element growth over one $s = 2$ stage; see [5] or [14] for the details.

A growth factor can be defined for the diagonal pivoting method in just the same way as for Gaussian elimination:

$$\rho_n = \frac{\max_{i,j,k} |a_{ij}^{(k)}|}{\max_{i,j} |a_{ij}|},$$

where the $a_{ij}^{(k)}$ are the elements of the Schur complements arising in the course of the factorization. From the derivation of the constant α it is easy to show that $\rho_n \leq (1 + 1/\alpha)^{n-1} = (2.57)^{n-1}$ for partial pivoting, which is larger than the bound 2^{n-1} for Gaussian elimination with partial pivoting (GEPP). But, it seems that as for GEPP, large element growth is rare in practice [5], [9].

2 Stability of the Diagonal Pivoting Method

Since the growth factor for the diagonal pivoting method with partial pivoting is bounded, and is usually small in practice, does it not follow that the method is stable in the same sense as for GEPP? This is a tempting argument, and one that is neither used nor warned against in the existing literature. However, it is easy to show that the argument is false, by exhibiting an example where the diagonal pivoting method has a small growth factor but is unstable. An example (not produced by partial pivoting) is, with $n = 3$ and with a 2×2 pivot followed by a 1×1 pivot,

$$\begin{aligned} A &= \begin{bmatrix} 1 & -(1 + \epsilon^2) & -\epsilon \\ -(1 + \epsilon^2) & 1 & -\epsilon \\ -\epsilon & -\epsilon & -1 \end{bmatrix} \\ &= \begin{bmatrix} 1 & & \\ 0 & 1 & \\ \epsilon^{-1} & \epsilon^{-1} & 1 \end{bmatrix} \begin{bmatrix} 1 & -(1 + \epsilon^2) & \\ -(1 + \epsilon^2) & 1 & \\ & & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & \epsilon^{-1} \\ & 1 & \epsilon^{-1} \\ & & 1 \end{bmatrix} = LDL^T, \end{aligned} \quad (2.1)$$

where $\epsilon > 0$. The growth factor ρ_n is 1, yet $\|L\|_\infty / \|A\|_\infty$ is unbounded as $\epsilon \rightarrow 0$, which suggests that the factorization, however it is computed, may not provide a stable way to solve linear systems $Ax = b$ in finite precision arithmetic. The instability is confirmed by a MATLAB experiment, in which the unit roundoff $u = 2^{-53} \approx 1.1 \times 10^{-16}$. We solved a linear system $Ax = b$, where $b = A[1 \ 2 \ 3]^T$, in two different ways. First, we computed the factorization in (2.1) using the diagonal pivoting method, as specified in (1.1) (with $H = I$), taking a 2×2 pivot on the first step and using GEPP to solve linear systems involving this pivot. For comparison, we evaluated the explicit formulae for the LDL^T factors in (2.1), and used the explicit inverse of $D(1:2, 1:2)$ when solving the linear system involving D . Table 2.1 shows the normwise relative backward error of the computed solution \hat{x} ,

$$\begin{aligned} \eta_\infty(\hat{x}) &:= \min\{\epsilon : (A + \Delta A)\hat{x} = b + \Delta b, \quad \|\Delta A\|_\infty \leq \epsilon \|A\|_\infty, \quad \|\Delta b\|_\infty \leq \epsilon \|b\|_\infty\} \\ &= \frac{\|b - A\hat{x}\|_\infty}{\|A\|_\infty \|\hat{x}\|_\infty + \|b\|_\infty} \end{aligned}$$

ϵ	Diagonal pivoting	Explicit factors
10^{-1}	9e-17	6e-16
10^{-2}	5e-17	2e-14
10^{-3}	3e-15	5e-11
10^{-4}	7e-14	4e-9
10^{-5}	6e-13	6e-8
10^{-6}	1e-13	1e-6
10^{-7}	4e-11	1e-7

Table 2.1: Backward error for computed solution of indefinite system of order 3.

(see [16] or [14, Th. 7.1] for a proof of the latter equality), which would be of order u for a stable solution method. As ϵ decreases the computations become unstable. We note that stability is obtained if, in (1.1), we take the natural 1×1 pivot a_{11} instead of the ill conditioned 2×2 pivot $A(1:2, 1:2)$; interestingly, though, the 2×2 pivot shares with those chosen by the Bunch–Kaufman partial pivoting strategy the property that it is indefinite. Partial pivoting is stable on this example.

We conclude that a small growth factor is not, by itself, enough to guarantee stability of the diagonal pivoting method. A sufficient condition for stability can be obtained by regarding the block LDL^T factorization computed by the diagonal pivoting method as a special case of a block LU factorization. Error analysis for block LU factorization is given by Demmel, Higham and Schreiber [8], and a suitable modification of this analysis gives the following result: if linear systems involving 2×2 pivots are solved in a normwise backward stable fashion then the condition

$$\|L\|_\infty \|D\|_\infty \|L^T\|_\infty \leq c_n \|A\|_\infty, \quad (2.2)$$

for a modest constant c_n , is sufficient to ensure that the diagonal pivoting method produces a factorization with a small relative residual and provides computed solutions to linear systems that have a small backward error. Unfortunately, condition (2.2) does not hold for the partial pivoting strategy of Bunch and Kaufman, as is shown by the following example. For $\epsilon > 0$, the diagonal pivoting method with partial pivoting produces the factorization, with $P = I$,

$$A = \begin{bmatrix} 0 & \epsilon & 0 \\ \epsilon & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & & \\ 0 & 1 & \\ 1/\epsilon & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & \epsilon & \\ \epsilon & 0 & \\ & & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1/\epsilon \\ & 1 & 0 \\ & & 1 \end{bmatrix} = LDL^T.$$

As $\epsilon \rightarrow 0$, $\|L\|_\infty \|D\|_\infty \|L^T\|_\infty / \|A\|_\infty \rightarrow \infty$, and indeed the multipliers are unbounded. Even 1×1 pivots can lead to arbitrarily large elements in L , as the following example

with $0 < \epsilon < \alpha$ shows (again, partial pivoting selects $P = I$):

$$A = \begin{bmatrix} \epsilon^2 & \epsilon & \epsilon \\ \epsilon & 0 & 1 \\ \epsilon & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & & \\ 1/\epsilon & 1 & \\ 1/\epsilon & 0 & 1 \end{bmatrix} \begin{bmatrix} \epsilon^2 & & \\ & -1 & \\ & & -1 \end{bmatrix} \begin{bmatrix} 1 & 1/\epsilon & 1/\epsilon \\ & 1 & 0 \\ & & 1 \end{bmatrix} = LDL^T.$$

It is worth emphasizing that large elements in a factor of a matrix do not necessarily imply that the factorization is unstable. For example, in the (point) LDL^T factorization of a symmetric positive definite matrix A with $D = \text{diag}(d_{ii})$, $d_{ii} > 0$, the ratio $\|L\|_\infty/\|A\|_\infty$ can be arbitrarily large, yet the factorization is guaranteed to be stable. One such example is, with $\epsilon > 0$,

$$A = \begin{bmatrix} \epsilon^2 & \epsilon \\ \epsilon & 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \epsilon^{-1} & 1 \end{bmatrix} \begin{bmatrix} \epsilon^2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & \epsilon^{-1} \\ 0 & 1 \end{bmatrix}.$$

Our conclusion is that existing results for LU factorization and block LU factorization do not directly imply the stability of the diagonal pivoting method with partial pivoting. Any proof of stability must make use of the particular properties of the partial pivoting strategy.

The only claims of stability that we have found in the literature are in the paper by Bunch, Kaufman and Parlett [6] and in the *LINPACK Users' Guide* [9, p. 5.19]; in both cases, residual bounds of the form $\|A - \widehat{L}\widehat{D}\widehat{L}^T\|_\infty \leq p(n)\rho_n\|A\|_\infty u$ are stated without proof, where p is a polynomial; we prove a result of this form and, in Theorem 4.2, a backward error result for the computed solution of $Ax = b$. We note that much of Bunch's analysis of the diagonal pivoting method in [3] is specific to complete pivoting, so his analysis does not readily yield results for partial pivoting.

In the rest of the paper we present a new analysis to show that partial pivoting is indeed a stable pivoting strategy for the diagonal pivoting method.

3 Background Results from Error Analysis

We collect in this section some standard error analysis results that will be needed later. For our model of floating point arithmetic we take

$$fl(x \text{ op } y) = (x \text{ op } y)(1 + \delta), \quad |\delta| \leq u, \quad \text{op} = +, -, *, /, \quad (3.1)$$

where u is the unit roundoff. All the results we quote remain true under a weaker model that accommodates machines without a guard digit [14, §2.4], provided some of the constants are increased slightly.

We introduce the constant

$$\gamma_n = \frac{nu}{1 - nu},$$

which carries with it the implicit assumption that $nu < 1$. Useful properties are (a) $\gamma_m + \gamma_n + \gamma_m\gamma_n \leq \gamma_{m+n}$ and (b) if $c \geq 1$ then $c\gamma_n \leq \gamma_{cn}$.

Proofs of the following results can be found in [14]. First, for matrix multiplication,

$$fl(AB) = AB + \Delta, \quad |\Delta| \leq \gamma_n |A| |B|, \quad A \in \mathbb{R}^{m \times n}, \quad B \in \mathbb{R}^{n \times p}.$$

Second, if $T \in \mathbb{R}^{n \times n}$ is a nonsingular triangular matrix and the system $Tx = b$ is solved by substitution then

$$(T + \Delta T)\hat{x} = b, \quad |\Delta T| \leq \gamma_n |T|. \quad (3.2)$$

Third, if a linear system $Ax = b$, where $A \in \mathbb{R}^{n \times n}$, is solved without breakdown by Gaussian elimination without pivoting, then the computed solution satisfies

$$(A + \Delta A)\hat{x} = b, \quad |\Delta A| \leq 2\gamma_n |\hat{L}||\hat{U}|, \quad (3.3)$$

where \hat{L} and \hat{U} are the computed LU factors.

We will use the norm defined by

$$\|A\|_M = \max_{i,j} |a_{ij}|$$

(for which $\|AB\|_M \leq n\|A\|_M\|B\|_M$ is the best bound of this form that holds for all $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$).

4 Error Analysis

4.1 2×2 Linear Systems

Crucial to the error analysis that follows is a backward error result for the solution of linear systems involving 2×2 pivots. Note that, in the notation of Algorithm 1, the pivot is

$$E = \begin{bmatrix} a_{11} & a_{r1} \\ a_{r1} & a_{rr} \end{bmatrix}, \quad |a_{r1}| = \lambda.$$

For this subsection and the later analysis, it is convenient to tabulate the conditions that must hold for a 2×2 pivot to be selected:

$$|a_{11}| < \alpha \lambda, \quad (4.1a)$$

$$|a_{11}|\sigma < \alpha \lambda^2, \quad (4.1b)$$

$$|a_{rr}| < \alpha \sigma, \quad (4.1c)$$

$$|a_{11}||a_{rr}| < \alpha^2 \lambda^2, \quad (4.1d)$$

where the fourth inequality is a consequence of the previous two (note that (4.1c) implies $\sigma \neq 0$).

Suppose, first, that linear systems $Ex = b$ are solved by GEPP. By (4.1a), $|a_{11}| < \alpha|a_{r1}| < |a_{r1}|$, so GEPP interchanges rows 1 and 2 of E and computes the LU factorization

$$PE = \begin{bmatrix} a_{r1} & a_{rr} \\ a_{11} & a_{r1} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \frac{a_{11}}{a_{r1}} & 1 \end{bmatrix} \begin{bmatrix} a_{r1} & a_{rr} \\ 0 & a_{r1} - \frac{a_{11}a_{rr}}{a_{r1}} \end{bmatrix} = LU.$$

From (3.3), we have the backward error result

$$(PE + \Delta E)\hat{x} = Pb, \quad |\Delta E| \leq 2\gamma_2|\widehat{L}||\widehat{U}|.$$

Now

$$|L||U| \leq \begin{bmatrix} |a_{r1}| & \left| \frac{a_{11}a_{rr}}{a_{r1}} \right| + \left| a_{r1} - \frac{a_{11}a_{rr}}{a_{r1}} \right| \\ |a_{11}| & (2\alpha^2 + 1)|a_{r1}| \end{bmatrix},$$

using (4.1d). It follows that

$$(E + \widetilde{\Delta E})\hat{x} = b, \quad |\widetilde{\Delta E}| \leq 2\gamma_2 \begin{bmatrix} |a_{11}| & 2|a_{r1}| \\ |a_{r1}| & |a_{rr}| \end{bmatrix} \leq 4\gamma_2|E|, \quad (4.2)$$

using the numerical value of α specified in Algorithm 1. Strictly, we should append “ $+O(u^2)$ ” to this bound, to account for replacing $|\widehat{L}||\widehat{U}|$ by a bound for $|L||U|$; we omit the second order term for the moment and reinstate it later. Note that the result (4.2) holds trivially for a 1×1 pivot E .

The main alternative to using GEPP to solve the systems $Ex = b$ is to use the explicit inverse of E , as is done in the implementations of the diagonal pivoting method with partial pivoting in LAPACK and LINPACK (see the auxiliary routine `xLASYF` in LAPACK and `xSIFA` in LINPACK). In both LAPACK and LINPACK, $Ex = b$ is solved by evaluating

$$x = \frac{1}{a_{r1} \left(\frac{a_{11}}{a_{r1}} \cdot \frac{a_{rr}}{a_{r1}} - 1 \right)} \begin{bmatrix} \frac{a_{rr}}{a_{r1}} & -1 \\ -1 & \frac{a_{11}}{a_{r1}} \end{bmatrix} b, \quad (4.3)$$

which corresponds to using an explicit formula for the inverse of a 2×2 matrix (or, equivalently, Cramer’s rule), with scaling to avoid overflow. The term

$$\mu = \frac{a_{11}}{a_{r1}} \cdot \frac{a_{rr}}{a_{r1}} - 1$$

appears to be a potential source of instability, since for arbitrary a_{11} , a_{r1} and a_{rr} the relative error in the computed $\hat{\mu}$ is unbounded. However, by exploiting the condition (4.1d) for a 2×2 pivot, which we rewrite as

$$\frac{|a_{11}||a_{rr}|}{a_{r1}^2} \leq \alpha^2,$$

we can obtain a very satisfactory error bound for $\hat{\mu}$. Using the model (3.1) we have

$$\hat{\mu} = \left(\frac{a_{11}}{a_{r1}} \cdot \frac{a_{rr}}{a_{r1}} (1 + \delta_1)(1 + \delta_2)(1 + \delta_3) - 1 \right) (1 + \delta_4),$$

where $|\delta_i| \leq u$, $i = 1:4$, which implies [14, Lemma 3.1]

$$\hat{\mu} = \frac{a_{11}}{a_{r1}} \cdot \frac{a_{rr}}{a_{r1}} (1 + \theta_4) - (1 + \delta_4), \quad |\theta_4| \leq \gamma_4.$$

Hence

$$\begin{aligned} |\mu - \hat{\mu}| &\leq \gamma_4 \left(\frac{|a_{11}a_{rr}|}{a_{r1}^2} + 1 \right) \leq \gamma_4(\alpha^2 + 1) \\ &\leq \gamma_4 \left(\frac{1 + \alpha^2}{1 - \alpha^2} \right) |\mu| < 3\gamma_4|\mu|. \end{aligned}$$

It is then straightforward to show that, denoting the matrix in (4.3) by Z ,

$$\hat{x} = (a_{r1}\mu)^{-1}(Z + \Delta Z)b, \quad |\Delta Z| \leq \gamma_{30}|Z|.$$

Thus $b - E\hat{x} = -E((a_{r1}\mu)^{-1}\Delta Z)b$, so that

$$\begin{aligned} |b - E\hat{x}| &\leq \gamma_{30}|E||E^{-1}||b| \\ &\leq \gamma_{30}|E||E^{-1}||E||x| \\ &\leq \gamma_{180}|E||x|, \end{aligned} \tag{4.4}$$

using (A.3). The Oettli–Prager theorem [15], [14, Th. 7.3] then implies that

$$(E + \Delta E)\hat{x} = b, \quad |\Delta E| \leq \gamma_{180}|E|.$$

Again, strictly a second order term should be added to the bound, this time to account for the fact that $|x|$ rather than $|\hat{x}|$ appears on the right-hand side of (4.4).

The conclusion is that whether the linear system $Ex = b$ involving the 2×2 pivot is solved by GEPP or by using the explicit inverse, we have

$$(E + \Delta E)\hat{x} = b, \quad |\Delta E| \leq \gamma_c|E|, \tag{4.5}$$

for an integer constant c . It is worth stressing that such a result does not hold for an arbitrary 2×2 (symmetric) matrix E —we have fully exploited the pivoting conditions in the derivation.

4.2 Componentwise Backward Error Analysis

Now we carry out a componentwise backward error analysis of the diagonal pivoting method. We make only one assumption about the pivoting strategy: that (4.5) holds for the 2×2 pivots. For convenience, we assume, without loss of generality, that no interchanges are needed, which amounts to redefining $A := PAP^T$ in (1.2).

To begin, we consider the first stage of the factorization, using the notation of (1.1). The submatrix $L_{21} = CE^{-1} \in \mathbb{R}^{(n-s) \times s}$ satisfies $L_{21}E = C$ or $EL_{21}^T = C^T$. If l_j is the j th column of L_{21}^T and c_j is the j th column of C^T , then, from (4.5),

$$(E + \Delta E_j)\widehat{l}_j = c_j, \quad |\Delta E_j| \leq \gamma_c |E|.$$

Hence, overall,

$$\widehat{L}_{21}E = C + \Delta C, \quad |\Delta C| \leq \gamma_c |\widehat{L}_{21}| |E|. \quad (4.6)$$

We assume that the Schur complement is computed as $S = B - L_{21}C^T$, so that²

$$\widehat{S} = B - \widehat{L}_{21}C^T + \Delta S, \quad |\Delta S| \leq \gamma_{s+1} (|B| + |\widehat{L}_{21}| |C^T|). \quad (4.7)$$

The remaining stages of the diagonal pivoting method factorize the Schur complement as $S = L_S D_S L_S^T$, and we assume, inductively, that the computed factors satisfy

$$\widehat{L}_S \widehat{D}_S \widehat{L}_S^T = \widehat{S} + \Delta_S, \quad |\Delta_S| \leq d(n-s, u) (|\widehat{S}| + |\widehat{L}_S| |\widehat{D}_S| |\widehat{L}_S^T|),$$

where $d(n-s, u)$ is a constant depending on $n-s$ and u . We therefore have computed factors \widehat{L} and \widehat{D} of A that satisfy

$$\begin{aligned} \widehat{L} \widehat{D} \widehat{L}^T &:= \begin{bmatrix} I & 0 \\ \widehat{L}_{21} & \widehat{L}_S \end{bmatrix} \begin{bmatrix} E & 0 \\ 0 & \widehat{D}_S \end{bmatrix} \begin{bmatrix} I & \widehat{L}_{21}^T \\ 0 & \widehat{L}_S^T \end{bmatrix} \\ &= \begin{bmatrix} E & E \widehat{L}_{21}^T \\ \widehat{L}_{21} E & \widehat{L}_{21} E \widehat{L}_{21}^T + \widehat{L}_S \widehat{D}_S \widehat{L}_S^T \end{bmatrix} \\ &= \begin{bmatrix} E & (C + \Delta C)^T \\ C + \Delta C & \widehat{L}_{21} E \widehat{L}_{21}^T + \widehat{S} + \Delta_S \end{bmatrix} \\ &= \begin{bmatrix} E & (C + \Delta C)^T \\ C + \Delta C & B + (\widehat{L}_{21} E \widehat{L}_{21}^T - \widehat{L}_{21} C^T) + \Delta_S + \Delta_S \end{bmatrix}. \end{aligned}$$

Now, from (4.6) we have the inequalities

$$|\widehat{L}_{21} E \widehat{L}_{21}^T - \widehat{L}_{21} C^T| \leq \gamma_c |\widehat{L}_{21}| |E| |\widehat{L}_{21}^T|$$

and

$$|\widehat{L}_{21}| |C^T| \leq (1 + \gamma_c) |\widehat{L}_{21}| |E| |\widehat{L}_{21}^T|. \quad (4.8)$$

²If the Schur complement is computed as $S = B - L_{21}EL_{21}^T$ then the same bound (4.9) ensues.

Using (4.7) and (4.8) we have

$$|\widehat{S}| \leq (1 + \gamma_{s+1})(|B| + (1 + \gamma_c)|\widehat{L}_{21}||E||\widehat{L}_{21}^T|).$$

Overall, then, we have

$$\widehat{L}\widehat{D}\widehat{L}^T = A + \Delta A,$$

where $\Delta A_{11} = 0$, $|\Delta A_{21}| \leq \gamma_c|\widehat{L}_{21}||E|$, and

$$\begin{aligned} |\Delta A_{22}| &\leq \gamma_c|\widehat{L}_{21}||E||\widehat{L}_{21}^T| + \gamma_{s+1}(|B| + (1 + \gamma_c)|\widehat{L}_{21}||E||\widehat{L}_{21}^T|) \\ &\quad + d(n-s, u)((1 + \gamma_{s+1})(|B| + (1 + \gamma_c)|\widehat{L}_{21}||E||\widehat{L}_{21}^T|) + |\widehat{L}_S||\widehat{D}_S||\widehat{L}_S^T|) \\ &\leq (\gamma_c + d(n-s, u)(1 + \gamma_c))|B| + (\gamma_c(2 + \gamma_c) + d(n-s, u)(1 + \gamma_c)^2)|\widehat{L}_{21}||E||\widehat{L}_{21}^T| \\ &\quad + d(n-s, u)|\widehat{L}_S||\widehat{D}_S||\widehat{L}_S^T| \\ &\leq (\gamma_c(2 + \gamma_c) + d(n-s, u)(1 + \gamma_c)^2)(|B| + |\widehat{L}_{21}||E||\widehat{L}_{21}^T| + |\widehat{L}_S||\widehat{D}_S||\widehat{L}_S^T|). \end{aligned}$$

Hence

$$\widehat{L}\widehat{D}\widehat{L}^T = A + \Delta A, \quad |\Delta A| \leq d(n, u)(|A| + |\widehat{L}||\widehat{D}||\widehat{L}^T|), \quad (4.9)$$

where $d(n, u)$ is clearly of the form $p(n)u + O(u^2)$, where p is a linear polynomial.

Now we analyse the substitution stages when the LDL^T factorization is used to solve a linear system $Ax = b$. From (3.2) and (4.5), the computed solutions to the three systems $Ly_1 = b$, $Dy_2 = y_1$, $L^T x = y_2$ satisfy

$$\begin{aligned} (\widehat{L} + \Delta L_1)\widehat{y}_1 &= b, & |\Delta L_1| &\leq \gamma_n|\widehat{L}|, \\ (\widehat{D} + \Delta D)\widehat{y}_2 &= \widehat{y}_1, & |\Delta D| &\leq \gamma_c|\widehat{D}|, \\ (\widehat{L} + \Delta L_2)^T \widehat{x} &= \widehat{y}_2. \end{aligned}$$

Thus

$$b = (\widehat{L} + \Delta L_1)(\widehat{D} + \Delta D)(\widehat{L} + \Delta L_2)^T \widehat{x} = (A + \Delta A + \Delta A_2)\widehat{x},$$

where $|\Delta A|$ is bounded in (4.9) and

$$|\Delta A_2| \leq \gamma_{2n+c}|\widehat{L}||\widehat{D}||\widehat{L}^T| + O(u^2).$$

On bringing back into account the row and column interchanges, we obtain the following result.

Theorem 4.1 *Let $A \in \mathbb{R}^{n \times n}$ be symmetric and let \widehat{x} be a computed solution to the linear system $Ax = b$ produced by the diagonal pivoting method with any pivoting strategy. If for all linear systems involving 2×2 pivots (4.5) holds, then*

$$(A + \Delta A)\widehat{x} = b, \quad |\Delta A| \leq p(n)u(|A| + P^T|\widehat{L}||\widehat{D}||\widehat{L}^T|P) + O(u^2), \quad (4.10)$$

where p is a linear polynomial and $PAP^T \approx \widehat{L}\widehat{D}\widehat{L}^T$ is the factorization computed by the diagonal pivoting method.

The bound in (4.10) is analogous to the bound in (3.3) that holds for Gaussian elimination. We have already seen that the assumption (4.5) in Theorem 4.1 holds for the partial pivoting strategy of Bunch and Kaufman, provided linear systems $Ex = b$ are solved by GEPP or by using the explicit inverse. It is easy to show that this assumption also holds for the complete pivoting strategy of Bunch and Parlett [7] under the same conditions (interestingly, for the 2×2 pivots E that arise with the Bunch–Parlett strategy, GEPP applied to a $Ex = b$ is identical to Gaussian elimination with complete pivoting).

4.3 Normwise Analysis for Partial Pivoting

To show that the diagonal pivoting method is stable for a particular pivoting strategy, we need to show that the matrix $|\widehat{L}||\widehat{D}||\widehat{L}^T|$ is suitably bounded. We now specialise to partial pivoting. For partial pivoting, \widehat{L} can be arbitrarily large, so stability is not an immediate consequence of Theorem 4.1. We therefore need to look closely at the elements of the matrix $|\widehat{L}||\widehat{D}||\widehat{L}^T|$. For simplicity, we bound the matrix $|L||D||L^T|$ containing the exact factors, which makes only a second order change to the overall bounds, since $|\widehat{L}||\widehat{D}||\widehat{L}^T| = |L||D||L^T|$.

Initially, we examine the contribution from the blocks of L and D produced by the first stage of the factorization. For this more delicate part of the analysis we take full account of the interchanges in our notation. Note that

$$\begin{aligned} |L||D||L^T| &= \begin{bmatrix} I & \\ |L_{21}| & |L_S| \end{bmatrix} \begin{bmatrix} |E| & \\ & |D_S| \end{bmatrix} \begin{bmatrix} I & |L_{21}^T| \\ & |L_S^T| \end{bmatrix} \\ &= \begin{bmatrix} |E| & |E||L_{21}^T| \\ |L_{21}||E| & |L_{21}||E||L_{21}^T| + |L_S||D_S||L_S^T| \end{bmatrix}. \end{aligned} \quad (4.11)$$

We first bound

$$F := |L_{21}||E| = |CE^{-1}||E| \in \mathbb{R}^{(n-s) \times s}.$$

For a 1×1 pivot, F is a vector with elements $|c_i e_{11}^{-1}| |e_{11}|$, each of which is trivially bounded by $\max_{i,j} |a_{ij}|$.

Now consider a 2×2 pivot. Algorithm 1 dictates that Π in (1.1) swaps rows and columns 2 and r so that, as noted earlier,

$$E = \begin{bmatrix} a_{11} & a_{r1} \\ a_{r1} & a_{rr} \end{bmatrix}, \quad |a_{r1}| = \lambda.$$

Using (A.1) and (4.1a), we have

$$e_i^T F \leq (e_i^T |C|) |E^{-1}| |E|$$

$$\begin{aligned}
&\leq \frac{1}{1-\alpha^2} [\lambda \quad \sigma] \begin{bmatrix} 1+\alpha^2 & \frac{2|a_{rr}|}{\lambda} \\ \frac{2|a_{11}|}{\lambda} & 1+\alpha^2 \end{bmatrix} \\
&\leq \frac{1}{1-\alpha^2} [(1+\alpha^2)\lambda + 2\alpha\sigma \quad 2|a_{rr}| + (1+\alpha^2)\sigma] \\
&\leq \frac{\max_{i,j} |a_{ij}|}{1-\alpha^2} [\alpha^2 + 2\alpha + 1 \quad \alpha^2 + 3] \\
&\leq \max_{i,j} |a_{ij}| [5 \quad 6]. \tag{4.12}
\end{aligned}$$

Next, we need to bound

$$G := |L_{21}| |E| |L_{21}^T| = |CE^{-1}| |E| |E^{-1}C^T|.$$

First, consider a 1×1 pivot. In cases (1) and (2) of Algorithm 1 we have

$$g_{ij} = |c_i e_{11}^{-1}| |e_{11}| |e_{11}^{-1} c_j| = \frac{|a_{i+1,1}| |a_{j+1,1}|}{|a_{11}|} \leq \frac{\lambda^2}{|a_{11}|} \leq \begin{cases} \frac{\lambda}{\alpha}, & \text{case (1),} \\ \frac{\sigma}{\alpha}, & \text{case (2).} \end{cases}$$

In case (3),

$$\begin{aligned}
|g_{ij}| &= \frac{|a_{lr}| |a_{mr}|}{|a_{rr}|} \quad (l, m \neq r) \\
&\leq \frac{\sigma^2}{|a_{rr}|} \leq \frac{\sigma}{\alpha}.
\end{aligned}$$

For a 1×1 pivot, then, $|g_{ij}| \leq \alpha^{-1} \max_{i,j} |a_{ij}| < 2 \max_{i,j} |a_{ij}|$.

For a 2×2 pivot (case (4) of Algorithm 1), using (A.2) we have

$$\begin{aligned}
|g_{ij}| &\leq (e_i^T |C|) (|E^{-1}| |E| |E^{-1}|) |C^T| e_j \\
&\leq \frac{3+\alpha^2}{(1-\alpha^2)^2 \lambda^2} [\lambda \quad \sigma] \begin{bmatrix} |a_{rr}| & \lambda \\ \lambda & |a_{11}| \end{bmatrix} \begin{bmatrix} \lambda \\ \sigma \end{bmatrix} \\
&= \frac{3+\alpha^2}{(1-\alpha^2)^2 \lambda^2} (\lambda^2 (|a_{rr}| + \sigma) + \sigma (\lambda^2 + |a_{11}| \sigma)) \\
&= \frac{3+\alpha^2}{(1-\alpha^2)^2} \left(|a_{rr}| + 2\sigma + \frac{\sigma^2 |a_{11}|}{\lambda^2} \right) \\
&\leq \frac{3+\alpha^2}{(1-\alpha^2)^2} (3+\alpha) \max_{i,j} |a_{ij}| \quad (\text{using (4.1b)}) \\
&= 36 \max_{i,j} |a_{ij}|. \tag{4.13}
\end{aligned}$$

The remaining blocks of $|L||D||L^T|$ are composed of blocks of L and D that make up LDL^T factors of Schur complements of A . But every Schur complement satisfies

$$\|S\|_M \leq \rho_n \|A\|_M,$$

where ρ_n is the growth factor. Hence, applying the bounds above recursively to the $(2, 2)$ block in (4.11), we deduce the (pessimistic) bound

$$\| |L||D||L^T| \|_M \leq 36n\rho_n \|A\|_M. \quad (4.14)$$

We mention in passing that in early drafts of this paper we had a weaker version of (4.5) in which $|E|$ in the bound was replaced by $|E| + |a_{r1}|e_2e_2^T$. We were still able to obtain a satisfactory bound for $\| |L||D||L^T| \|_M$, indicating that partial pivoting is somewhat more tolerant of how the 2×2 systems are solved than might be thought from the analysis above.

Using the bound (4.14) in Theorem 4.1 we obtain the following normwise backward stability result for partial pivoting.

Theorem 4.2 *Let $A \in \mathbb{R}^{n \times n}$ be symmetric and let \hat{x} be a computed solution to the linear system $Ax = b$ produced by the diagonal pivoting method with the partial pivoting strategy of Bunch and Kaufman, where linear systems involving 2×2 pivots are solved by GEPP or by use of the explicit inverse. Then*

$$(A + \Delta A)\hat{x} = b, \quad \|\Delta A\|_M \leq p(n)\rho_n u \|A\|_M + O(u^2), \quad (4.15)$$

where p is a quadratic.

Theorem 4.2 has the same form as Wilkinson's result for GEPP applied to a nonsymmetric system (see, e.g., [14, §9.2]), though of course the numerical value of ρ_n is usually different for the two methods.

5 Discussion

The backward error matrix ΔA in (4.9) is necessarily symmetric, but that in (4.15) is not, in general. However, we can take ΔA in (4.15) to be symmetric, at the cost of increasing the bound by a factor n , because of the following result of Bunch, Demmel and Van Loan [4]: if $(A + G)y = b$ then there exists $H = H^T$ such that $(A + H)y = b$ with $\|H\|_2 \leq \|G\|_2$ and $\|H\|_F \leq \sqrt{2}\|G\|_F$.

Sorensen and Van Loan [10, §5.3.2] modify the Bunch–Kaufman partial pivoting strategy by redefining, in Algorithm 1,

$$\sigma = \|A(:, r)\|_\infty$$

This small change has the pleasing effect of ensuring that for a positive definite matrix no interchanges are done (and that, as for the Bunch–Kaufman strategy, only 1×1

pivots are used in this case). At the same time it leaves the growth factor bound unchanged, and all our analysis remains valid for this variant.

For sparse symmetric matrices, Duff, Reid and co-authors compute the block LDL^T factorization using a pivoting strategy very different from that of Bunch and Kaufman [11], [12], [13]. We describe the strategy in [13] as it applies to the first stage of the factorization: a_{11} is defined to be an acceptable 1×1 pivot, from the point of view of numerical stability, if

$$|a_{11}| \geq \theta \max_{i>1} |a_{i1}|, \quad (5.1)$$

where $\theta \in (0, 1/2]$ is a tolerance; the matrix

$$D_1 = \begin{bmatrix} a_{11} & a_{r1} \\ a_{r1} & a_{rr} \end{bmatrix}$$

is an acceptable 2×2 pivot if

$$\|D_k^{-1}\|_\infty \max\{|a_{ij}| : i \neq 1, r; j = 1, r\} \leq \theta^{-1}. \quad (5.2)$$

From among the acceptable pivots one is chosen that best preserves sparsity, according to some particular sparsity criterion. The conditions (5.1) and (5.2) ensure that $\|L\|_\infty$ is bounded by a multiple of θ^{-1} , which then implies bounds on the growth factor, and hence on $\|D\|_\infty$. The stability of this pivoting strategy is therefore immediate, since (2.2) is satisfied. An interesting contrast is that the Bunch–Kaufman strategy involves a fixed amount of searching for a pivot, and the reasons for its stability are subtle, whereas the Duff et al. strategy more directly forces stability by bounding the multipliers, but gives up the fixed amount of searching of the Bunch–Kaufman strategy.

Finally, we emphasize that the aim of this work was to obtain a rigorous backward error bound for the diagonal pivoting method with partial pivoting. The actual performance of the method is affected by the size of the growth factor. More work is needed to investigate the behaviour of the growth factor, about which less is known than the growth factor for Gaussian elimination with partial pivoting. Although the unboundedness of $\|L\|_\infty$ does not preclude backward stability, it does have implications for the practical behaviour of the method; see Ashcraft, Grimes and Lewis [2] for a thorough study for both dense and sparse matrices.

A Appendix

In this appendix we bound three matrix expressions involving a 2×2 pivot from partial pivoting,

$$E = \begin{bmatrix} a_{11} & a_{r1} \\ a_{r1} & a_{rr} \end{bmatrix} \quad |a_{r1}| = \lambda.$$

First, we note that

$$|\det(E)| = |a_{r1}^2 - a_{11}a_{rr}| \geq \lambda^2 - \alpha^2\lambda^2 = (1 - \alpha^2)\lambda^2,$$

using (4.1d). Hence

$$\begin{aligned} |E^{-1}||E| &\leq \frac{1}{(1 - \alpha^2)\lambda^2} \begin{bmatrix} |a_{rr}| & \lambda \\ \lambda & |a_{11}| \end{bmatrix} \begin{bmatrix} |a_{11}| & \lambda \\ \lambda & |a_{rr}| \end{bmatrix} \\ &= \frac{1}{1 - \alpha^2} \begin{bmatrix} \frac{|a_{11}||a_{rr}|}{\lambda^2} + 1 & \frac{2|a_{rr}|}{\lambda} \\ \frac{2|a_{11}|}{\lambda} & \frac{|a_{11}||a_{rr}|}{\lambda^2} + 1 \end{bmatrix} \\ &\leq \frac{1}{1 - \alpha^2} \begin{bmatrix} 1 + \alpha^2 & 2\frac{|a_{rr}|}{\lambda} \\ 2\frac{|a_{11}|}{\lambda} & 1 + \alpha^2 \end{bmatrix}, \end{aligned} \tag{A.1}$$

using (4.1d) again. Next,

$$\begin{aligned} |E^{-1}||E||E^{-1}| &\leq \frac{1}{(1 - \alpha^2)^2\lambda^2} \begin{bmatrix} 1 + \alpha^2 & 2\frac{|a_{rr}|}{\lambda} \\ 2\frac{|a_{11}|}{\lambda} & 1 + \alpha^2 \end{bmatrix} \begin{bmatrix} |a_{rr}| & \lambda \\ \lambda & |a_{11}| \end{bmatrix} \\ &= \frac{1}{(1 - \alpha^2)^2\lambda^2} \begin{bmatrix} (3 + \alpha^2)|a_{rr}| & (1 + \alpha^2)\lambda + 2\frac{|a_{11}||a_{rr}|}{\lambda} \\ 2\frac{|a_{11}||a_{rr}|}{\lambda} + (1 + \alpha^2)\lambda & (3 + \alpha^2)|a_{11}| \end{bmatrix} \\ &\leq \frac{1}{(1 - \alpha^2)^2\lambda^2} \begin{bmatrix} (3 + \alpha^2)|a_{rr}| & (1 + 3\alpha^2)\lambda \\ (1 + 3\alpha^2)\lambda & (3 + \alpha^2)|a_{11}| \end{bmatrix} \\ &\leq \frac{3 + \alpha^2}{(1 - \alpha^2)^2\lambda^2} \begin{bmatrix} |a_{rr}| & \lambda \\ \lambda & |a_{11}| \end{bmatrix}. \end{aligned} \tag{A.2}$$

Finally,

$$\begin{aligned} |E||E^{-1}||E| &\leq \frac{1}{1 - \alpha^2} \begin{bmatrix} |a_{11}| & \lambda \\ \lambda & |a_{rr}| \end{bmatrix} \begin{bmatrix} 1 + \alpha^2 & 2\frac{|a_{rr}|}{\lambda} \\ 2\frac{|a_{11}|}{\lambda} & 1 + \alpha^2 \end{bmatrix} \\ &= \frac{1}{1 - \alpha^2} \begin{bmatrix} (3 + \alpha^2)|a_{11}| & 2\frac{|a_{11}||a_{rr}|}{\lambda} + (1 + \alpha^2)\lambda \\ (1 + \alpha^2)\lambda + 2\frac{|a_{11}||a_{rr}|}{\lambda} & (3 + \alpha^2)|a_{rr}| \end{bmatrix} \end{aligned}$$

$$\begin{aligned}
&\leq \frac{1}{1-\alpha^2} \begin{bmatrix} (3+\alpha^2)|a_{11}| & (1+3\alpha^2)\lambda \\ (1+3\alpha^2)\lambda & (3+\alpha^2)|a_{rr}| \end{bmatrix} \\
&\leq \left(\frac{3+\alpha^2}{1-\alpha^2}\right) |E| \leq 6|E|.
\end{aligned} \tag{A.3}$$

Acknowledgements

It is a pleasure to thank Philip Gill and Michael Saunders for valuable comments, particularly at early stages of this work. I also thank Jim Bunch, Des Higham and John Lewis for suggesting improvements to draft manuscripts.

References

- [1] E. Anderson, Z. Bai, C. H. Bischof, J. W. Demmel, J. J. Dongarra, J. J. Du Croz, A. Greenbaum, S. J. Hammarling, A. McKenney, S. Ostrouchov, and D. C. Sorensen. *LAPACK Users' Guide, Release 2.0*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, second edition, 1995. ISBN 0-89871-345-5. xix+325 pp.
- [2] C. Cleve Ashcraft, Roger G. Grimes, and John G. Lewis. Accurate symmetric indefinite linear equation solvers. Manuscript, May 1995.
- [3] James R. Bunch. Analysis of the diagonal pivoting method. *SIAM J. Numer. Anal.*, 8(4):656–680, 1971.
- [4] James R. Bunch, James W. Demmel, and Charles F. Van Loan. The strong stability of algorithms for solving symmetric linear systems. *SIAM J. Matrix Anal. Appl.*, 10(4):494–499, 1989.
- [5] James R. Bunch and Linda Kaufman. Some stable methods for calculating inertia and solving symmetric linear systems. *Math. Comp.*, 31(137):163–179, 1977.
- [6] James R. Bunch, Linda Kaufman, and Beresford N. Parlett. Decomposition of a symmetric matrix. *Numer. Math.*, 27:95–109, 1976.
- [7] James R. Bunch and Beresford N. Parlett. Direct methods for solving symmetric indefinite systems of linear equations. *SIAM J. Numer. Anal.*, 8(4):639–655, 1971.
- [8] James W. Demmel, Nicholas J. Higham, and Robert S. Schreiber. Stability of block LU factorization. *Numerical Linear Algebra with Applications*, 2(2):173–190, 1995.
- [9] J. J. Dongarra, J. R. Bunch, C. B. Moler, and G. W. Stewart. *LINPACK Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1979. ISBN 0-89871-172-X.
- [10] Jack J. Dongarra, Iain S. Duff, Danny C. Sorensen, and Henk A. van der Vorst. *Solving Linear Systems on Vector and Shared Memory Computers*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1991. ISBN 0-89871-270-X. x+256 pp.
- [11] I. S. Duff, N. I. M. Gould, J. K. Reid, J. A. Scott, and K. Turner. The factorization of sparse symmetric indefinite matrices. *IMA J. Numer. Anal.*, 11:181–204, 1991.

- [12] Iain S. Duff and John K. Reid. MA27—A set of Fortran subroutines for solving sparse symmetric sets of linear equations. Technical Report AERE R10533, AERE Harwell Laboratory, July 1982. Published by Her Majesty's Stationary Office, London.
- [13] Iain S. Duff, John K. Reid, Neils Munskgaard, and Hans B. Nielsen. Direct solution of sets of linear equations whose matrix is sparse, symmetric and indefinite. *J. Inst. Maths Applics*, 23:235–250, 1979.
- [14] Nicholas J. Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1995. ISBN 0-89871-355-2. Approx xxiv+690 pp. In press.
- [15] W. Oettli and W. Prager. Compatibility of approximate solution of linear equations with given error bounds for coefficients and right-hand sides. *Numer. Math.*, 6:405–409, 1964.
- [16] J. L. Rigal and J. Gaches. On the compatibility of a given solution with the data of a linear system. *J. Assoc. Comput. Mach.*, 14(3):543–548, 1967.