

On Factoring a Class of Complex Symmetric Matrices Without Pivoting

By Steven M. Serbin*

Abstract. Let $A = B + iC$ be a complex, symmetric $n \times n$ matrix with B and C each real, symmetric and positive definite. We show that the LINPACK diagonal pivoting decomposition $U^{-1}A(U^{-1})^T = D$ proceeds without the necessity for pivoting. In particular, when B and C are band matrices, bandwidth is preserved.

In this brief note we wish to investigate the problem of solving the linear system

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

where A is a complex, symmetric ($A = A^T$), $n \times n$ matrix and $z, b \in \mathbb{C}^n$. In general, such systems do not admit simple methods like Cholesky factorization without pivoting, as would be the case for A positive definite real symmetric or complex Hermitian. The best currently available scheme is a version of Bunch's diagonal pivoting method, (for symmetric indefinite systems), which may be found in LINPACK (cf. Chapter 5 of LINPACK guide [3] for details). The essential step at each stage of the elimination is the choice of either 1 by 1 or 2 by 2 blocks for the purpose of pivoting for numerical stability. Bunch et al. [1], [2] discuss various complete and partial pivoting strategies.

For certain applications, in particular one arising in the use of the (2, 2) Padé method for the solution of certain systems of ordinary differential equations [4], A has additional structure which we would like to exploit to expedite the solution of (1). Namely, we will assume that

$$(2) \quad A = B + iC,$$

where B and C are both real, symmetric, positive definite. Further, for the particular application in mind, B and C are also band matrices.

The point in question here is the necessity of row/column interchanges in the pursuance of the diagonal pivoting method. Bunch and Kaufman [1] discuss the situation for band matrices; they report that when the bandwidth is greater than five, if 2 by 2 pivots are required, then the bandwidth is not preserved. On the other hand, if symmetry is ignored, then bandwidth may be preserved, but only at the expense of doubling storage. We show here that in our case, no pivoting is required in the LINPACK algorithm, hence bandwidth will be preserved. Our analysis can be modified easily to treat the usual LDL^T decomposition.

Received November 26, 1979.

1980 *Mathematics Subject Classification*. Primary 65F05.

* Research supported by USARO Grant DAAG29-78-C-0024.

The LINPACK algorithm seeks to produce a factorization

$$(3) \quad U^{-1}A(U^{-1})^T = \mathcal{D},$$

where U is the product of upper triangular block "elementary eliminators" and permutation matrices and \mathcal{D} is a symmetric block diagonal matrix with 1 by 1 or 2 by 2 diagonal blocks. The elimination works from the last column to the first, rather than as in a standard Gaussian elimination. At an intermediate step, the working matrix is of the form

$$(4) \quad \left[\begin{array}{c|c} A_k & \mathbf{0} \\ \hline \mathbf{0} & D_k \end{array} \right],$$

A_k being symmetric and (k by k), D_k symmetric, block diagonal ($n - k$) by ($n - k$). The decision on how (and whether) to pivot first rests on a comparison of $\lambda_k \equiv \max_{i \leq j \leq k-1} |(A_k)_{jk}|$, the largest off-diagonal element in the last column of A_k , to the diagonal element. For complex z , the algorithm uses $|z| \equiv |\text{real}(z)| + |\text{imag}(z)|$. Specifically, if

$$(5) \quad |(A_k)_{kk}| \geq \alpha \lambda_k \quad (\alpha \doteq .6404),$$

then *no interchanges are required and a 1 × 1 block elimination is used.*

Since B and C are positive definite, at the first step ($k = n$), clearly the condition (5) is met (it is true with $\alpha = 1$). Hence the elimination proceeds without interchange. If we can show that the ($n - 1$) by ($n - 1$) matrix A_{n-1} , obtained from the elimination step, inherits the property that its real and imaginary parts are positive definite, then an induction shows that the entire algorithm may proceed without *any* interchanges. Since we shall only investigate one elimination step, the subscripts may be abandoned. The technique modifies an argument found in Wendroff [5].

Since no pivoting is required at the first step, we can write

$$(6) \quad A = \begin{bmatrix} B & \mathbf{a} \\ \mathbf{a}^T & d \end{bmatrix},$$

with

$$(7) \quad d = \sigma + i\tau, \quad \sigma, \tau > 0,$$

and

$$(8) \quad \mathbf{a} = \alpha + i\beta.$$

The elimination is performed by letting

$$(9) \quad \mathbf{m} = -d^{-1}\mathbf{a},$$

and

$$(10) \quad U = \begin{bmatrix} I & \mathbf{m} \\ \mathbf{0} & 1 \end{bmatrix},$$

so that

$$(11) \quad UAU^T = \begin{bmatrix} B - \mathbf{m}\mathbf{d}\mathbf{m}^T & \mathbf{0} \\ \mathbf{0} & d \end{bmatrix}.$$

Our main result is stated as

THEOREM. *Under the conditions stated on A, the matrix $B - \mathbf{m}\mathbf{d}\mathbf{m}^T$ in (11) has real and imaginary parts both positive definite.*

Proof. Write $B = S + iT$. S and T themselves are positive definite, being principal submatrices of positive definite matrices. By direct computation,

$$(12) \quad \begin{aligned} \mathbf{m}\mathbf{d}\mathbf{m}^T = & \{[\sigma(\alpha\alpha^T - \beta\beta^T) + \tau(\beta\alpha^T + \alpha\beta)] \\ & + i[\sigma(\beta\alpha^T + \alpha\beta^T) + \tau(\beta\beta^T - \alpha\alpha^T)]\}/(\sigma^2 + \tau^2). \end{aligned}$$

It suffices (due to the symmetries $S \leftrightarrow T, \alpha \leftrightarrow \beta, \sigma \leftrightarrow \tau$) for us just to investigate the real part of $B - \mathbf{m}\mathbf{d}\mathbf{m}^T$. Let $\mathbf{x} \in \mathbf{R}^{n-1}$ be arbitrary ($\mathbf{x} \neq \mathbf{0}$). Now, for any $\xi \in \mathbf{R}$,

$$(13) \quad \begin{bmatrix} \mathbf{x}^T & \xi \end{bmatrix} \begin{bmatrix} S & \alpha \\ \alpha^T & \sigma \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \xi \end{bmatrix} > 0,$$

or in scalar form

$$(14) \quad 0 < \sum_{i,j=1}^{n-1} s_{ij}x_i x_j + 2\xi \sum_{i=1}^{n-1} \alpha_i x_i + \sigma\xi^2,$$

or

$$(15) \quad -2\xi \sum_{i=1}^{n-1} \alpha_i x_i - \sigma\xi^2 < \sum_{i,j=1}^{n-1} s_{ij}x_i x_j.$$

Now, from (12), if $Q = \text{Re}\{B - \mathbf{m}\mathbf{d}\mathbf{m}^T\}$, then

$$(16) \quad \begin{aligned} (\sigma^2 + \tau^2)\mathbf{x}^T Q \mathbf{x} &= \left\{ \sum_{i,j=1}^{n-1} [(\sigma^2 + \tau^2)s_{ij} - \sigma(\alpha_i\alpha_j - \beta_i\beta_j) - \tau(\beta_i\alpha_j + \alpha_i\beta_j)] x_i x_j \right\} \\ &> (\sigma^2 + \tau^2) \left[-2\xi \sum_{i=1}^{n-1} \alpha_i x_i - \sigma\xi^2 \right] \\ &\quad - \sum_{i,j=1}^{n-1} [\sigma(\alpha_i\alpha_j - \beta_i\beta_j) + \tau(\beta_i\alpha_j + \alpha_i\beta_j)] x_i x_j, \quad \text{by (15).} \end{aligned}$$

Now, if we abbreviate $\sum_{i=1}^{n-1} \alpha_i x_i = \mu, \sum_{i=1}^{n-1} \beta_i x_i = \nu$, we have

$$(17) \quad \begin{aligned} (\sigma^2 + \tau^2)\mathbf{x}^T Q \mathbf{x} &> (\sigma^2 + \tau^2) [-2\xi\mu - \sigma\xi^2] - [\sigma(\mu^2 - \nu^2) + 2\tau\mu\nu] \\ &= -(\sigma^2 + \tau^2)\sigma \left(\xi + \frac{\mu}{\sigma} \right)^2 + \frac{1}{\sigma}(\tau\mu - \sigma\nu)^2. \end{aligned}$$

Choosing $\xi = -\mu/\sigma$, we have that for all $\mathbf{x} \neq \mathbf{0}$,

$$(\sigma^2 + \tau^2)\mathbf{x}^T Q \mathbf{x} > \frac{1}{\sigma}(\tau\mu - \sigma\nu)^2 \geq 0,$$

which completes the proof.

Department of Mathematics
University of Tennessee
Knoxville, Tennessee 37916

1. J. R. BUNCH & L. KAUFMAN, "Some stable methods for calculating inertia and solving symmetric linear systems," *Math Comp.*, v. 31, 1977, pp. 163–179.
2. J. R. BUNCH, L. KAUFMAN & B. N. PARLETT, "Decomposition of a symmetric matrix," *Numer. Math.*, v. 27, 1976, pp. 95–109.
3. J. J. DONGARRA, C. B. MOLER, J. R. BUNCH & G. W. STEWART, *LINPACK User's Guide*, SIAM, Philadelphia, Pa., 1979.
4. G. FAIRWEATHER, "A note on the efficient implementation of certain Padé methods for linear parabolic problems," *BIT*, v. 18, 1978, pp. 106–109.
5. B. WENDROFF, *Theoretical Numerical Analysis*, Academic Press, New York, 1966.