

## A Stability Analysis of Incomplete LU Factorizations

By Howard C. Elman\*

**Abstract.** The combination of iterative methods with preconditionings based on incomplete LU factorizations constitutes an effective class of methods for solving the sparse linear systems arising from the discretization of elliptic partial differential equations. In this paper, we show that there are some settings in which the incomplete LU preconditioners are not effective, and we demonstrate that their poor performance is due to numerical instability. Our analysis consists of an analytic and numerical study of a sample two-dimensional non-self-adjoint elliptic problem discretized by several finite-difference schemes.

**1. Introduction.** The preconditioned conjugate gradient method [3], [4] and preconditioned iterative methods for nonsymmetric linear systems [1], [6], [18], [23] are effective methods for solving the large sparse linear systems arising from the discretization of elliptic partial differential equations. Two good preconditioners are the incomplete LU factorization (ILU) [16] and the modified incomplete LU factorization (MILU) [5], [11], each of which makes use of an approximate factorization of the coefficient matrix into the product of a sparse lower triangular matrix,  $L$ , and a sparse upper triangular matrix,  $U$ . For the symmetric positive-definite systems derived from the finite-difference discretization of self-adjoint elliptic problems on a uniform  $n \times n$  grid, it is known that the MILU preconditioning produces a reduction of the condition number from  $O(n^2)$  to  $O(n)$ . The ILU preconditioning does not improve the conditioning in this way, but it has been observed empirically to generate a linear system most of whose eigenvalues are clustered near one. The effectiveness of both techniques for the nonsymmetric linear systems derived from non-self-adjoint elliptic problems has been demonstrated in many numerical experiments [2], [7], [8], [21], although the analysis from the symmetric case has not been generalized.

Let  $A$  denote the coefficient matrix, let  $Q = LU$  denote the approximate factorization of  $A$ , and let  $R = Q - A$ . Loosely speaking, the analysis for symmetric discretized elliptic equations examines the effect of  $R$  on vectors  $u$  whose values come from a smooth function evaluated at the mesh points. In particular, a heuristic explanation of the difference between the MILU and ILU factorizations is that the individual entries of the vector  $Ru$  satisfy  $[Ru]_j = O(1/n)$  for the MILU factorization, whereas  $[Ru]_j = O(1)$  for the ILU factorization [11]. The MILU factorization

---

Received April 18, 1985.

1980 *Mathematics Subject Classification.* Primary 65F10, 65N20, 15A06; Secondary 65N05.

\*The work presented in this paper was supported by the U. S. Office of Naval Research under contract N00014-82-K-0814 and by the U. S. Army Research Office under contract DAAG-83-0177. The work was performed while the author was at the Department of Computer Science, Yale University.

*Current address:* Department of Computer Science, University of Maryland, College Park, Maryland.

©1986 American Mathematical Society  
0025-5718/86 \$1.00 + \$.25 per page

has a higher order of accuracy as an approximation to  $A$  than the ILU factorization (see also [19]). In this sense, the analysis is reminiscent of the notion of “consistency” of difference schemes for ordinary differential equations [10]. The notion of order of accuracy also extends to the nonsymmetric case (see, e.g., [14]). In this regime, the MILU factorization is also of higher order of accuracy, and it has been demonstrated to be more effective in many numerical experiments [7], [8].

In this paper, we show that *stability* can also play a role in the performance of the incomplete LU preconditioners when they are applied to discretized non-self-adjoint elliptic equations. We show, using a model problem, that there are nonsymmetric linear systems that can cause difficulty for either the ILU preconditioning or the MILU preconditioning, and that the source of this difficulty is instability of the computations involving the triangular factors. Our analysis is similar to stability analysis of methods for ordinary differential equations [10]. It shows that the performance of incomplete factorizations is sensitive to both the values of the coefficients of the elliptic operator and the choice of difference scheme used to discretize the problem. In Section 2, we present the model problem and give examples of numerical difficulties exhibited by both preconditioners. In Section 3, we construct constant-coefficient approximations to the two factorizations based on an asymptotic analysis of the values of their coefficients, and in Section 4, we present the stability analysis for these simplified factorizations. In Section 5, we demonstrate with numerical experiments that the presence of instability correlates with a degradation of the performance of the preconditioners.

**2. The Model Problem and Some Numerical Examples.** In this section, we present a model problem and briefly discuss some numerical experiments that demonstrate some difficulties encountered by the ILU and MILU preconditioners. Consider the constant-coefficient elliptic equation

$$(2.1) \quad -\Delta u + 2P_1u_x + 2P_2u_y = f$$

on the unit square  $\Omega = \{(x,y) | 0 \leq x,y \leq 1\}$ , with Dirichlet boundary conditions  $u = \hat{f}$  on  $\partial\Omega$ . Discretizing (2.1) on an  $n \times n$  grid gives rise to a sparse linear system of equations

$$(2.2) \quad Au = g,$$

of order  $N = n^2$ . We consider two difference schemes for approximating the first derivatives in (2.1). In the first scheme, we use second-order centered differences

$$u_x \approx \frac{u_{i+1,j} - u_{i-1,j}}{2h}, \quad u_y \approx \frac{u_{i,j+1} - u_{i,j-1}}{2h},$$

where  $h = 1/(n+1)$ . In the second scheme, we use first-order upwind differences, i.e., backward differences if the coefficient is positive, forward differences if the coefficient is negative. For the upwind scheme, we restrict our attention to the case  $P_1, P_2 \geq 0$ , so that the differences are given by

$$u_x \approx \frac{u_{ij} - u_{i-1,j}}{h}, \quad u_y \approx \frac{u_{ij} - u_{i,j-1}}{h}.$$

For both schemes, we use standard second-order centered differences for the Laplacian [9],

$$\Delta u \approx \frac{u_{i+1,j} - 2u_{ij} + u_{i-1,j}}{h^2} + \frac{u_{i,j+1} - 2u_{ij} + u_{i,j-1}}{h^2}.$$



TABLE 1  
*Number of iterations of Orthomin(1) to convergence*

	ILU	MILU
Problem 1.	21	18
Problem 2.	$\infty$	7
Problem 3.	32	$\infty$

TABLE 2  
*Two leftmost and two rightmost eigenvalues of the symmetric parts of the coefficient matrices*

	ILU		MILU	
	Leftmost	Rightmost	Leftmost	Rightmost
Problem 1.	0.172, 0.196	1.31, 1.32	0.240, 0.247	2.589, 2.593
Problem 2.	-49.3, -12.6	18.2, 56.1	0.681, 0.681	1.02, 1.03
Problem 3.	0.043, 0.075	1.56, 1.57	-2.90E10, -4.36E8	4.36E8, 2.90E10

We discretize on a  $31 \times 31$  grid using centered differences for the first derivatives. The preconditioning is applied from the right, i.e., given the incomplete factorization  $Q = LU$ , the preconditioned problem to be solved is

$$(2.5) \quad A Q^{-1} \tilde{u} = g, \quad u = Q^{-1} \tilde{u}.$$

Using Orthomin(1) [6] as the basic iterative method, we attempt to solve each of the three problems with both the ILU and the MILU preconditioning. Table 1 contains the number of iterations of Orthomin(1) needed to satisfy the stopping criterion of  $\|r_i\|/\|r_0\| \leq 10^{-6}$  (where  $\|v\|$  denotes the Euclidean norm  $(v, v)^{1/2}$ ). The symbol  $\infty$  indicates that the residual norms  $\{\|r_i\|\}$  stopped decreasing at some point of the iteration, i.e., that the iteration failed to converge to the solution.

Note that if the symmetric part,  $(A Q^{-1} + (A Q^{-1})^T)/2$ , of the coefficient matrix in (2.5) is positive-definite, then Orthomin(1) is guaranteed to generate a sequence of iterates whose residual norms are strictly decreasing [6], [7]. The symmetric part is indefinite if and only if it has both nonpositive and nonnegative eigenvalues. In Table 2, we list the two leftmost and two rightmost eigenvalues of the symmetric parts of the six coefficient matrices tested. This data confirms that the problems in which the failures occurred have indefinite symmetric parts. Taken together, the results from the two tables also indicate that, at least when centered differences are used for the discretization of (2.1), the performance of the incomplete factorization preconditionings is very sensitive to the values of the coefficients of the first derivatives. (We remark that the symmetric part of the original matrix  $A$  is the discrete negative Laplacian, which is positive-definite [9].)

**3. The Incomplete Factorizations.** In this section, we define the ILU and MILU factorizations and construct simplified constant-coefficient approximations to each of them that will lend themselves to a stability analysis. We consider incomplete factorizations in which the lower and upper triangular factors,  $L$  and  $U$ , have the same sparsity structure as the lower and upper parts of  $A$ , respectively, and  $U$  is unit



recurrence of (3.1). For the MILU factorization, all but  $4n - 4$  diagonals satisfy the last recurrence of (3.2).\*\*\* To facilitate an analysis, we construct constant-coefficient ILU and MILU factorizations based on an asymptotic analysis of these two recurrences. We make use of a relationship between the recurrences for  $\{\alpha_j\}$  and continued fractions.

As motivation, consider the first block (i.e., the first  $n$  rows) of the ILU factorization, whose diagonal entries satisfy

$$(3.3) \quad \alpha_1 = a, \quad \alpha_j = a - bd/\alpha_{j-1}, \quad 2 \leq j \leq n.$$

Expanding for the  $j$ th value gives

$$(3.4) \quad \alpha_j = a - \frac{bd}{a - \frac{bd}{a - \frac{bd}{a - \dots}}} \left. \vphantom{\frac{bd}{a - \frac{bd}{a - \frac{bd}{a - \dots}}}} \right\} j - 1 \text{ divides.}$$

In the language of continued fractions,  $\alpha_j$  is the  $(j - 1)$ st approximant of the continued fraction [22, p. 14]

$$a - \frac{bd}{a - \frac{bd}{a - \dots}}$$

The convergence of continued fractions (i.e., of the approximants) of this form is well understood. We state without proof the result we need, which is taken essentially verbatim from [22, Theorem 8.2, p. 39].

**THEOREM 3.1.** *The continued fraction*

$$(3.5) \quad 1 + \frac{\xi}{1 + \frac{\xi}{1 + \dots}}$$

converges for any (complex) number  $\xi$ , except when  $\xi$  is real and satisfies  $\xi < -\frac{1}{4}$ . For real  $\xi \geq -\frac{1}{4}$ , the limiting value is  $(1 + \sqrt{1 + 4\xi})/2$ .

To use Theorem 3.1 to examine the convergence of (3.3) or (3.4), consider  $\hat{\alpha}_j = \alpha_j/a$ . These quantities satisfy

$$\hat{\alpha}_1 = 1, \quad \hat{\alpha}_j = 1 - \frac{bd/a^2}{\hat{\alpha}_{j-1}}, \quad 2 \leq j \leq n.$$

With  $\xi = -bd/a^2$ , the approximant for  $\hat{\alpha}_j$  analogous to (3.4) is the  $(j - 1)$ st approximant of a continued fraction of the form (3.5). The convergence result requires that  $\xi \geq -\frac{1}{4}$ , which follows for both the centered- and upwind-difference schemes by direct computation. Hence,  $\{\hat{\alpha}_j\}$  converges to

$$\hat{\alpha} = \frac{1 + \sqrt{1 - 4bd/a^2}}{2},$$

and  $\{\alpha_j\}$  converges to

$$(3.6) \quad \alpha^{(1)} = a\hat{\alpha} = \frac{a + \sqrt{a^2 - 4bd}}{2}.$$

---

\*\*\*Actually, with the convention that the off-diagonal entries are zero in the appropriate indices, all the diagonal entries satisfy the last recurrences (3.1) or (3.2). We specify the exceptions explicitly to emphasize that there are exceptions.

For the centered-difference scheme, the limiting value is  $2 + \sqrt{3 + p_1^2}$ , and for the upwind scheme, it is  $2 + p_1 + p_2 + \sqrt{(2 + p_1 + p_2)^2 - (1 + 2p_1)}$ . We remark that the limiting value (3.6) is equal to the larger root obtained by formally substituting a constant value  $\alpha^{(1)}$  in place of  $\alpha_j$  and  $\alpha_{j-1}$  in (3.3) and solving the resulting quadratic equation for  $\alpha^{(1)}$ .

This argument establishes the convergence of the sequence defined by (3.3) as  $j \rightarrow \infty$ , although it does not prove that the first  $n$  values are near the limiting value. Moreover, most of the diagonal entries of  $L$  outside of the first block satisfy the more complicated recurrence given by the last of the defining equations for  $\{\alpha_j\}$ , which for the ILU factorization can be written as

$$(3.7) \quad \alpha_j = a - bd/\alpha_{j-1} - ce/\alpha_{j-n}.$$

We attempt to gain some insight into this recurrence by examining a simplified version of it. By analogy with the result for the first block, let

$$(3.8) \quad \alpha = \frac{a + \sqrt{a^2 - 4(bd + ce)}}{2},$$

which is the larger root obtained from formally substituting the constant  $\alpha$  in place of  $\alpha_j$ ,  $\alpha_{j-1}$  and  $\alpha_{j-n}$  in (3.7) and solving for  $\alpha$ . Consider the alternative recurrence

$$(3.9) \quad \alpha_1 = a - ce/\alpha, \quad \alpha_j = a - ce/\alpha - bd/\alpha_{j-1}, \quad j \geq 2,$$

which is a simplification of (3.7) in which  $\alpha_{j-n}$  is replaced by  $\alpha$  of (3.8).

**THEOREM 3.2.** *Let  $\alpha$  be given by (3.8). Then for the values of  $a, b, c, d$  and  $e$  of either difference scheme, the sequence defined by (3.9) is convergent with limit  $\alpha$ . The limiting values are*

$$\begin{aligned} \alpha &= 2 + \sqrt{2 + p_1^2 + p_2^2} && \text{for centered differences,} \\ \alpha &= 2 + p_1 + p_2 + \sqrt{1 + (1 + p_1 + p_2)^2} && \text{for upwind differences.} \end{aligned}$$

*Proof.* The values of  $\alpha$  for the two schemes are obtained by substituting into (3.8) the values of  $a, b, c, d$  and  $e$  from Section 2. The simplified recurrence (3.9) has the same form as the recurrence of (3.3), with  $a$  replaced by  $a - ce/\alpha$ . For both schemes, a straightforward computation shows that  $a - ce/\alpha$  is greater than zero and hence nonzero. We apply Theorem 3.1 in the same manner as above: if  $\xi = -bd/(a - ce/\alpha)^2$  is greater than or equal to  $-\frac{1}{4}$ , then the sequence  $\{\alpha_j/(a - ce/\alpha)\}$  converges to

$$\frac{1 + \sqrt{1 - 4bd/(a - ce/\alpha)^2}}{2},$$

and  $\{\alpha_j\}$  converges to

$$(3.10) \quad \bar{\alpha} = \frac{(a - ce/\alpha) + \sqrt{(a - ce/\alpha)^2 - 4bd}}{2}.$$

The condition  $\xi \geq -\frac{1}{4}$  is equivalent to  $bd \leq -\frac{1}{4}(a - ce/\alpha)^2$ . To see that the latter inequality holds for both schemes, first note that  $\alpha \geq 2 + \sqrt{2} > 0$ . For the centered scheme, we wish to show that

$$(3.11) \quad 1 - p_1^2 \leq 4 - \frac{2(1 - p_2^2)}{\alpha} + \frac{1}{4} \left( \frac{1 - p_2^2}{\alpha} \right)^2.$$

If  $p_2^2 \geq 1$ , then (3.11) holds trivially. If  $p_2^2 < 1$ , then  $0 < 1 - p_2^2 < 1$ , and the right-hand side of (3.11) is greater than  $4 - 2/(2 + \sqrt{2}) > 3$ , from which the inequality follows. For the upwind scheme, the condition  $\xi \geq -\frac{1}{4}$  is equivalent to

$$(3.12) \quad 1 + 2p_1 \leq \frac{1}{4} \left( 4 + 2(p_1 + p_2) - \frac{1 + 2p_2}{\alpha} \right)^2.$$

But  $(1 + 2p_2)/\alpha \leq 2/(2 + \sqrt{2}) + p_2$ . Inequality (3.12) then follows from direct computation.

It remains to show that  $\tilde{\alpha} = \alpha$ , i.e., that

$$\alpha = \frac{(a - ce/\alpha) + \sqrt{(a - ce/\alpha)^2 - 4bd}}{2},$$

or, equivalently, that

$$(3.13) \quad \sqrt{(a - ce/\alpha)^2 - 4bd} = 2\alpha - (a - ce/\alpha).$$

Rewriting the right-hand side of (3.13) as  $(2\alpha^2 - a\alpha + ce)/\alpha$ , it can be shown by direct computation that this quantity is positive for both schemes. Hence, it suffices to show that the square of both sides of (3.13) are equal, which simplifies to

$$\alpha^2 - a\alpha + (bd + ce) = 0.$$

The solution to this equation is  $\alpha$  of (3.8), which completes the proof.  $\square$

Note that we made use of three inequalities in this proof:

1.  $a - ce/\alpha > 0$ , which ensures that  $\xi$  is well defined, and figures in the choice of the positive square root of  $(a - ce/\alpha)^2$  in (3.10);
2.  $2\alpha - (a - ce/\alpha) \geq 0$ , which allows us to square both sides of (3.13);
3.  $\xi \geq -\frac{1}{4}$ , so that Theorem 3.1 can be applied.

For the MILU factorization, we rewrite the last recurrence of (3.2) as

$$(3.14) \quad \alpha_j = a - b(d + e)/\alpha_{j-1} - c(d + e)/\alpha_{j-n}.$$

The roots of the quadratic equation obtained by substituting  $\alpha$  for  $\alpha_j$ ,  $\alpha_{j-1}$  and  $\alpha_{j-n}$  in (3.14) are

$$(3.15) \quad \frac{a \pm \sqrt{a^2 - 4(b + c)(d + e)}}{2}.$$

Let  $\alpha$  denote the root with larger modulus. As above, we examine the simpler recurrence derived by replacing  $\{\alpha_{j-n}\}$  with  $\alpha$ :

$$(3.16) \quad \begin{aligned} \alpha_1 &= a - c(d + e)/\alpha, \\ \alpha_j &= (a - c(d + e)/\alpha) - b(d + e)/\alpha_{j-1}, \quad j \geq 2. \end{aligned}$$

We have the following convergence result:

**THEOREM 3.3.** *For either difference scheme, let  $\alpha$  denote the larger root given by (3.15) and let  $\xi = -b(d + e)/(a - c(d + e)/\alpha)^2$ . For the values of  $a, b, c, d$  and  $e$  of the upwind scheme, the sequence (3.16) is convergent with limit  $\alpha$ . For the values of the centered-difference scheme, the sequence (3.16) is convergent to  $\alpha$  provided that  $\xi \geq -\frac{1}{4}$ ,  $a - c(d + e)/\alpha > 0$  and  $2\alpha - a + c(d + e)/\alpha \geq 0$ . The limiting values are*

$$\begin{aligned} \alpha &= 2 + |p_1 + p_2| && \text{for centered differences,} \\ \alpha &= 2(1 + p_1 + p_2) && \text{for upwind differences.} \end{aligned}$$

*Proof.* The values of  $\alpha$  for the two schemes are obtained by substituting into (3.15) the values of  $a, b, c, d$  and  $e$  from Section 2. For the upwind scheme, as in the proof of Theorem 3.2,  $a - c(d + e)/\alpha > 0$  follows from direct substitution. The condition  $\xi \geq -\frac{1}{4}$  is equivalent to

$$2(1 + 2p_1) \leq \frac{1}{4} \left( 4 + 2p_1 + 2p_2 - \frac{2(1 + 2p_2)}{\alpha} \right)^2 = \left( 2 + p_1 + p_2 - \frac{1 + 2p_2}{\alpha} \right)^2.$$

But  $\alpha \geq 2$ , so that

$$2 + p_1 + p_2 - \frac{1 + 2p_2}{\alpha} \geq \frac{3}{2} + p_1.$$

Consequently, it suffices to show that  $(\frac{3}{2} + p_1)^2 \geq 2(1 + 2p_1)$ , which follows directly. For the centered scheme, we are only concerned with values of  $p_1$  and  $p_2$  for which  $a - c(d + e)/\alpha > 0$  and  $\xi \geq -\frac{1}{4}$ .

Applying Theorem 3.1 to both difference schemes, the sequence defined by (3.16) is convergent with limit

$$\tilde{\alpha} = \frac{(a - c(d + e)/\alpha) + \sqrt{(a - c(d + e)/\alpha)^2 - 4b(d + e)}}{2}.$$

The proof that  $\tilde{\alpha} = \alpha$  is identical to the analogous proof of Theorem 3.2. The condition  $2\alpha - a + c(d + e)/\alpha \geq 0$  again follows for the upwind scheme from direct computation.  $\square$

The extra assumptions made for the centered-difference scheme in Theorem 3.3 are not valid for all real  $p_1$  and  $p_2$ . The following result gives sufficient conditions for the inequalities to hold. We defer a proof to the Appendix. An illustration of these values is shown in Figure 1. Note that for most values where Lemma 3.1 does not hold, either  $p_1$  or  $p_2$  is large, and these values are of somewhat limited practical interest [20]. See the Appendix for some other comments on these values.

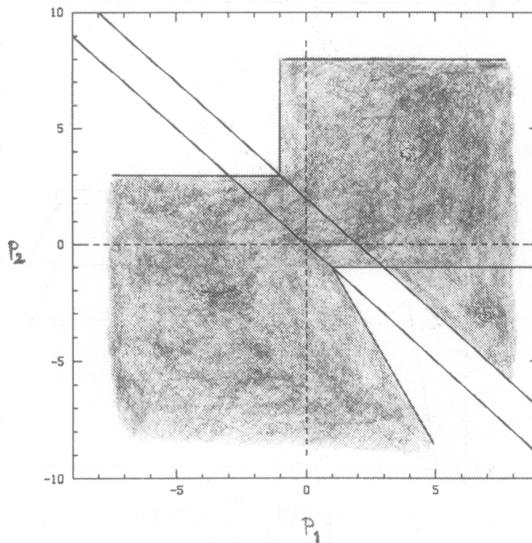


FIGURE 1  
 Values of  $p_1$  and  $p_2$  where Theorem 3.3 holds for centered differences and MILU.

LEMMA 3.1. For the values of  $a, b, c, d$  and  $e$  of the centered-difference scheme, the inequalities  $\xi \geq -\frac{1}{4}$ ,  $a - c(d + e)/\alpha > 0$ , and  $2\alpha - a + c(d + e)/\alpha \geq 0$  hold for all  $p_1$  and  $p_2$  satisfying

1.  $3 \leq p_2 \leq 8$  and  $p_1 \geq -1$ ; or
2.  $-1 \leq p_2 \leq 3$  and  $p_1$  arbitrary; or
3.  $p_2 \leq -1$  and either  $p_1 \leq \frac{1}{2}(1 - p_2)$  or  $p_1 \geq 2 - p_2$ .

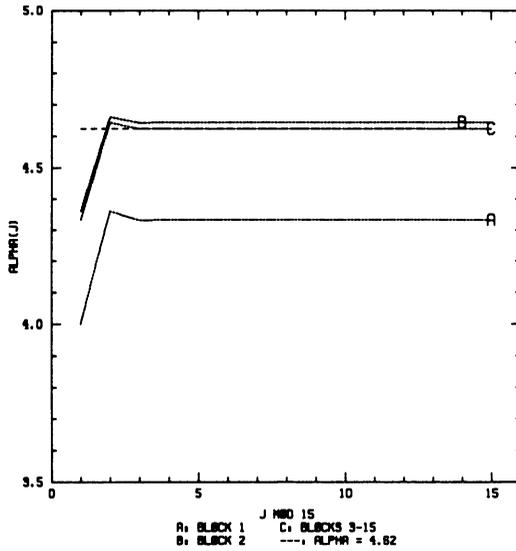


FIGURE 2  
 Convergence of diagonal values,  $P_1 = P_2 = 25$ ,  
 $h = \frac{1}{16}$ , centered differences, ILU factorization.

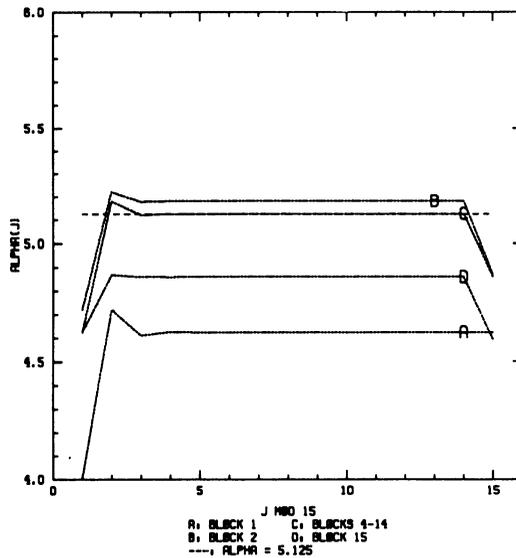


FIGURE 3  
 Convergence of diagonal values,  $P_1 = P_2 = 25$ ,  
 $h = \frac{1}{16}$ , centered differences, MILU factorization.

We define the constant-coefficient ILU and MILU factorizations as in (3.1), except we take the diagonal entries  $\{\alpha_j\}$  to have the limiting values  $\{\alpha\}$  of Theorems 3.2 and 3.3. We denote the off-diagonal entries  $\beta_j, \gamma_j, \delta_j$  and  $\eta_j$  of the constant-coefficient factorizations by  $\beta, \gamma, \delta$  and  $\eta$ , respectively.

The results of the theorems do not rigorously prove that the true factors converge in any sense to the constant-coefficient factors. The simplifying assumptions that the diagonal values from the previous blocks are constant and that they take on the values derived from solving the specified quadratic equations are *not* true. (For example, the limiting values for the first block are, in general, different from the quadratic roots used in the theorems.) Moreover, the convergence results do not say anything about how close the first  $n$  values are to the limits, which is the only useful information in this context. Nevertheless, our numerical experience supports the idea that the constant-coefficient factors are reasonable approximations. In Figure 2, we graph the computed values of the ILU diagonals  $\{\alpha_j\}$  and the limiting value  $\alpha$  of

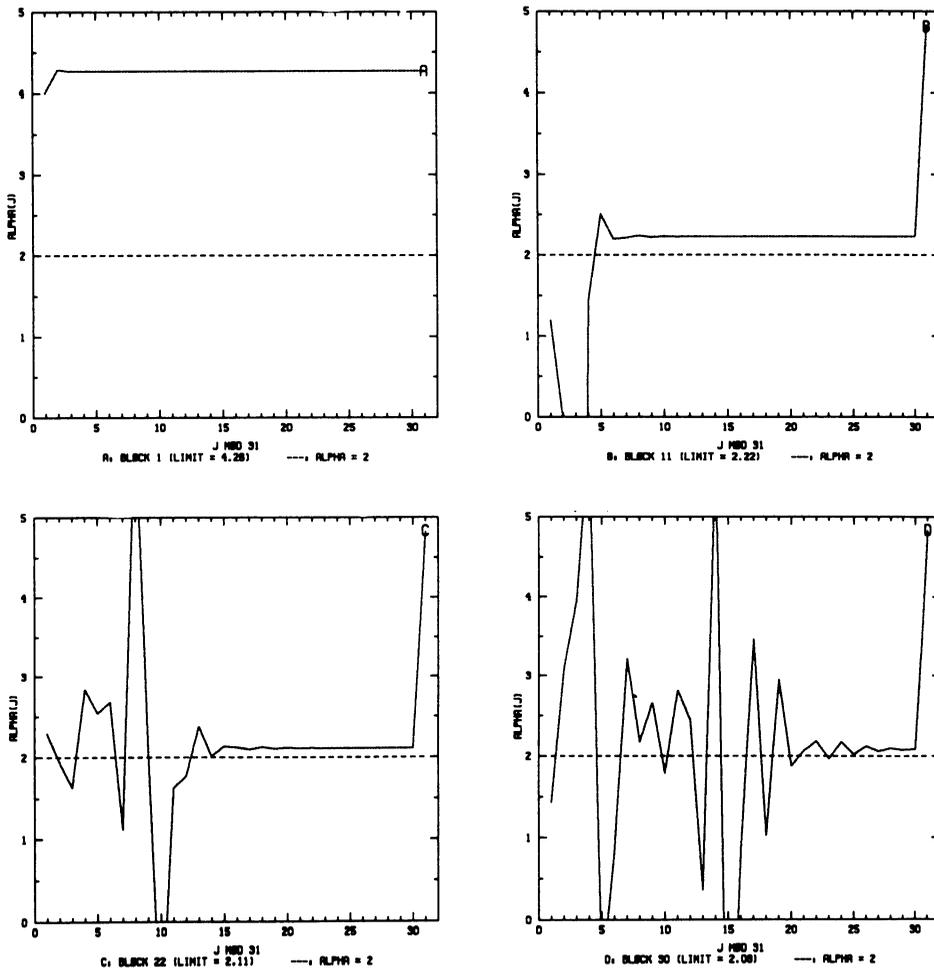


FIGURE 4

Convergence of diagonal values,  $P_1 = -50, P_2 = 50,$   
 $h = \frac{1}{32},$  centered differences, MILU factorization.

Theorem 3.2, for a  $15 \times 15$  grid ( $h = \frac{1}{16}$ ) and  $P_1 = P_2 = 25$  (so that  $p_1 = p_2 = 1.5625$ ). As expected, the limiting value for the first block is different from  $\alpha$ , but starting from the third block, most of the subsequent diagonal values are virtually indistinguishable from  $\alpha$ . Figure 3 graphs the analogous data for the same problem and difference scheme with the MILU factorization. In this case, the values from the last block and the last values within the other blocks differ from the limit because they satisfy a different recurrence (see (3.2)). For the same problem with the upwind schemes, both factorizations show the same qualitative behavior. Similarly, for centered differences on a finer grid with the same values of  $p_1$  and  $p_2$  ( $h = \frac{1}{32}$  and  $P_1 = P_2 = 50$ ), the qualitative behavior is identical.

We have not seen any examples in which the sequences (3.9) and (3.16) are convergent, but the true diagonal sequences  $\{\alpha_j\}$  are not convergent. However, we do have cases where convergence of the diagonal sequence is slower than in the examples above. For example, Figure 4 graphs the computed and limiting values for  $P_1 = -50$ ,  $P_2 = 50$ ,  $h = \frac{1}{32}$ , centered differences with the MILU factorization, from four blocks of the factorization. Within each block, the diagonal values approach a limit, and as the factorization proceeds, the limit for each block gets closer to the limit of Theorem 3.3. However, the initial values within each block oscillate, and these oscillations are both larger in magnitude and occur at higher indices (mod  $n$ ) within the block at the later stages of the factorization.

**4. Stability Analysis.** In this section, we examine the numerical stability of the lower triangular and upper triangular solves performed with the constant-coefficient factorizations introduced in Section 3. An alternative analysis using Fourier techniques in a slightly different setting is given in [15]. Given an incomplete factorization  $LU$ , the preconditioning operation consists of a pair of triangular solves of the form

$$(4.1) \quad Lv = w, \quad Uv = w.$$

Typically these operations are performed once each per iteration of the basic iterative method.

Consider the lower triangular solve. The typical computation for the  $j$ th entry of  $v$  has the form

$$v_j = \frac{1}{\alpha} (w_j - \beta v_{j-1} - \gamma v_{j-n}),$$

where  $v_{j-1}$  and  $v_{j-n}$  have been previously computed. Equivalently, most entries of  $v$  satisfy the  $n$ th order inhomogeneous recurrence relation

$$(4.2) \quad \alpha v_j + \beta v_{j-1} + \gamma v_{j-n} = w_j.$$

Our stability analysis is based on an analysis of this recurrence, whose solution is uniquely defined once  $n$  initial values, say  $v_1, \dots, v_n$ , are given. (We note that as in computation of the factors, not every step of the backsolve satisfies this recurrence, since there are  $2n - 1$  cases corresponding to grid points next to the boundary where the computation of  $v_j$  is simpler.) We will define the stability of the lower triangular solve in terms of properties of the characteristic polynomial associated with (4.2), which is given by

$$\lambda_n(z) = \alpha z^n + \beta z^{n-1} + \gamma.$$

Consider the homogeneous case,  $w = 0$ . If  $\lambda_n$  has  $n$  distinct roots  $\{z_1, \dots, z_n\}$ , then it is well known (see, e.g., [13]) that for  $j > n$  the homogeneous solution of (4.2) is

$$(4.3) \quad v_j = c_1 z_1^j + \dots + c_n z_n^j,$$

where  $c_1, \dots, c_n$  are determined by the initial values  $v_1, \dots, v_n$ . If some root  $z_s$  of  $\lambda_n$  has modulus greater than 1, then  $z_s^j$  will grow as  $j$  increases, and any error in  $c_s$  (caused, say, by an error in the initial conditions) will result in increasing errors in the solution (4.3). If a root  $z_s$  has multiplicity  $m$ , then the homogeneous solution also contains linear combinations of

$$(4.4) \quad j(j-1)\dots(j-t+1)z_s^{j-t}, \quad t = 1, \dots, m-1,$$

as components; if  $z_s$  has modulus greater than or equal to 1, then any errors in the coefficients of (4.4) will also be enhanced with increasing  $j$ . Therefore, we say that the lower triangular solve is *stable* if all the roots of its characteristic polynomial are less than or equal to 1 in modulus and no root with unit modulus has multiplicity greater than 1. It is *unstable* otherwise.

The typical computation in the upper triangular solve is

$$v_j + \delta v_{j+1} + \eta v_{j+n} = w_j,$$

where  $v_{j+1}$  and  $v_{j+n}$  are given. To fit this computation into the setting of recurrences, we renumber the unknowns so that they are computed in order of increasing indices, i.e., let  $\hat{v}_j = v_{N-j}$ . Then, the upper triangular computation has the form

$$\hat{v}_j + \delta \hat{v}_{j-1} + \eta \hat{v}_{j-n} = w_{N-j},$$

and the associated characteristic polynomial is

$$(4.5) \quad \mu_n(z) = z^n + \delta z^{n-1} + \eta = 0.$$

We say that the upper triangular solve is stable if all the roots of its characteristic polynomial have modulus less than or equal to 1, and no root with unit modulus has multiplicity greater than 1. It is unstable otherwise.

We first note that for the ILU factorization and for the MILU factorization with the upwind scheme, there are no multiple roots of unit modulus. Any multiple root of  $\lambda_n(z)$  must also be a root of  $\lambda'_n(z) = z^{n-2}(n\alpha z + (n-1)\beta)$ . Thus, the only possibility for a nonzero multiple root is  $z^* = -(n-1)\beta/n\alpha$ . For ILU and the centered-difference scheme,

$$|z^*| < \frac{|\beta|}{\alpha} = \frac{|1 + p_1|}{2 + \sqrt{2 + p_1^2 + p_2^2}} < 1.$$

The same argument works for the upwind scheme, for both ILU and MILU. For  $\mu_n$ , the only possibility for a nonzero multiple root is  $z^{**} = -(n-1)\delta/n$ , and the identical argument shows that  $|z^{**}| < 1$ . For MILU and centered differences, the same reasoning shows that there are no multiple roots if  $p_1$  and  $p_2$  have the same sign, or if either  $|p_1| \leq 1$  or  $|p_2| \leq 1$ . There *can* be multiple roots otherwise. However, we will show at the end of this section that such roots play no role, so in the following we will ignore the possibility of multiple roots.

It is sufficient, then, to examine the moduli of the largest roots of the two characteristic polynomials. For  $\lambda_n$ , we distinguish among four cases:

1.  $\beta \leq 0, \gamma \leq 0;$       3.  $\beta \geq 0, \gamma \leq 0;$
2.  $\beta \geq 0, \gamma \geq 0;$       4.  $\beta \leq 0, \gamma \geq 0.$

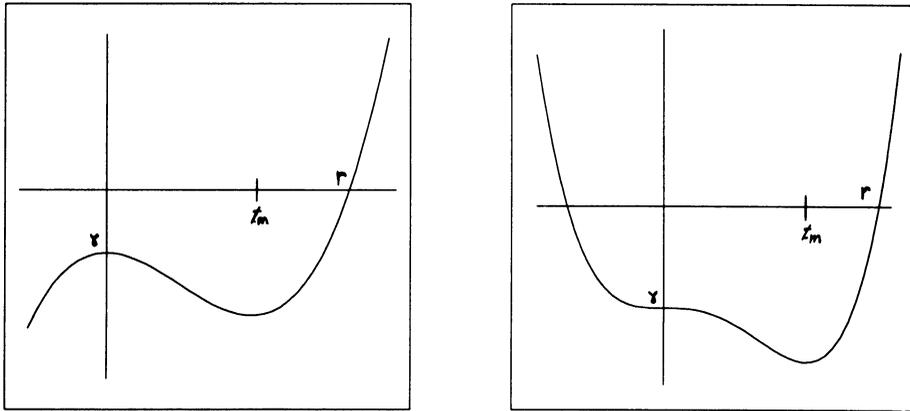


FIGURE 5

Shape of characteristic polynomials,  $\beta < 0, \gamma < 0$ .  
 Left side:  $n$  odd. Right side:  $n$  even.

For Case 1 and nonzero  $\beta$  and  $\gamma$ , it can be shown using elementary calculus that the graph of the real values of  $\lambda_n$  has one of the two shapes given in Figure 5, depending on whether  $n$  is odd or even. In particular,  $\lambda_n(0) = \gamma < 0$ ,  $\lambda_n(t)$  has a local minimum at  $t_m = -((n - 1)/n)(\beta/\alpha) > 0$ , and  $\lambda_n(t)$  is strictly decreasing for  $t$  between 0 and  $t_m$  and strictly increasing for  $t > t_m$ . Therefore,  $\lambda_n$  has precisely one positive real root,  $r$ , and  $\lambda_n(t) > 0$  for all  $t > r$ . The same argument works for  $\beta < 0, \gamma = 0$ . If  $\beta = 0$  and  $\gamma < 0$ , then  $\lambda_n(t)$  is strictly increasing for  $t > 0$  and again has precisely one positive real root,  $r$ . (When  $\beta = \gamma = 0$ , 0 is an  $n$ -fold root of  $\lambda_n$  and the solve is stable. In this trivial case, the lower triangular matrix is actually a diagonal matrix and we will not discuss it further.)

We claim that the largest positive root  $r$  is a root with largest modulus. For if  $\hat{r}e^{i\theta}$  is any root with largest modulus, then by definition  $\hat{r} \geq r$  and

$$\alpha \hat{r}^n e^{in\theta} = -\beta \hat{r}^{n-1} e^{i(n-1)\theta} - \gamma.$$

Taking the modulus of both sides and applying the triangle inequality,

$$\alpha \hat{r}^n = |-\beta \hat{r}^{n-1} e^{i(n-1)\theta} - \gamma| \leq |\beta| \hat{r}^{n-1} + |\gamma|,$$

i.e.,  $\lambda_n(\hat{r}) \leq 0$ . Since  $\hat{r} \geq 0$ , it follows that  $\hat{r} \leq r$ . Consequently,  $r = \hat{r}$ . Whether or not  $r$  is greater than 1 can be determined from the sign of  $\lambda_n(1)$ . For if  $\lambda_n(1) < 0$ , then since 1 is a positive number, it follows that 1 must lie to the left of  $r$ , i.e.,  $r > 1$ . Conversely, if  $\lambda_n(1) \geq 0$ , then  $r \leq 1$ . Hence, for Case 1, the lower triangular solve is stable if and only if  $\lambda_n(1) \geq 0$ .

For Cases 2–4, we can only give a partial analysis, which characterizes stability only for certain parities of  $n$ . Consider Case 2 and  $n$  odd. Then,

$$(4.6) \quad \lambda_n(-z) = -(\alpha z^n - \beta z^{n-1} - \gamma) = -(\alpha z^n + \hat{\beta} z^{n-1} + \hat{\gamma}),$$

where  $\hat{\beta}$  and  $\hat{\gamma}$  are less than or equal to 0. Up to a sign, the polynomial on the right of (4.6) has the form considered for Case 1. Hence, the analysis for Case 1 implies that for odd  $n$ , the lower triangular solve is stable if and only if  $\lambda_n(-1) \leq 0$ .

Similarly, for Case 3, when  $n$  is even,

$$\lambda_n(-z) = \alpha z^n + (-\beta) z^{n-1} + \gamma,$$

which again has the form studied for Case 1. Thus, in Case 3 with  $n$  even, the lower triangular solve is stable if and only if  $\lambda_n(-1) \geq 0$ .

Finally, for Case 4, we consider a recurrence analogous to (4.2) corresponding to a reordering of unknowns. Recall that in Section 2 we assumed that the unknowns are ordered by lines along the  $x$ -direction. For the lower triangular solve of (4.1), the computations for the unknowns  $v_j$  are unchanged (except for the order in which they are done) if the unknowns are computed along the lines in the  $y$ -direction. If we renumber the unknowns for the lower triangular solve according to this ordering, then the typical recurrence in the new ordering is

$$\alpha v_j + \gamma v_{j-1} + \beta v_{j-n} = w_j,$$

and the corresponding characteristic polynomial is

$$\hat{\lambda}_n(z) = \alpha z^n + \gamma z^{n-1} + \beta.$$

Hence, with the alternative ordering, Case 4 is equivalent to Case 3 with the roles of  $\beta$  and  $\gamma$  interchanged. The lower triangular solve is stable for even  $n$  if and only if  $\hat{\lambda}_n(-1) \geq 0$ . (The proof that  $\hat{\lambda}_n$  also has no multiple roots is the same as that for  $\lambda_n$ ; we omit the details.)

We summarize these observations as follows:

**THEOREM 4.1.** *Necessary and sufficient conditions for the lower triangular solve to be stable are:*

1. for  $\beta \leq 0, \gamma \leq 0$ :  $\lambda_n(1) = \alpha + \beta + \gamma \geq 0$ ;
2. for  $\beta \geq 0, \gamma \geq 0$  and  $n$  odd:  $\lambda_n(-1) = -\alpha + \beta + \gamma \leq 0$ ;
3. for  $\beta \geq 0, \gamma \leq 0$  and  $n$  even:  $\lambda_n(-1) = \alpha - \beta + \gamma \geq 0$ ;
4. for  $\beta \leq 0, \gamma \geq 0$  and  $n$  even:  $\hat{\lambda}_n(-1) = \alpha - \gamma + \beta \geq 0$ .

The identical analysis can be used to determine the stability of the upper triangular solve, based on the largest characteristic root of  $\mu_n$  in (4.5). We again distinguish among four cases, depending on the signs of  $\delta$  and  $\eta$ .

**THEOREM 4.2.** *Necessary and sufficient conditions for the upper triangular solve to be stable are:*

1. for  $\delta \leq 0, \eta \leq 0$ :  $\mu_n(1) = 1 + \delta + \eta \geq 0$ ;
2. for  $\delta \geq 0, \eta \geq 0$  and  $n$  odd:  $\mu_n(-1) = -1 + \delta + \eta \leq 0$ ;
3. for  $\delta \geq 0, \eta \leq 0$  and  $n$  even:  $\mu_n(-1) = 1 - \delta + \eta \geq 0$ ;
4. for  $\delta \leq 0, \eta \geq 0$  and  $n$  even:  $\hat{\mu}_n(-1) = 1 - \eta + \delta \geq 0$ , where  $\hat{\mu}_n(z) = z^n + \eta z^{n-1} + \delta$ .

Given values of  $p_1$  and  $p_2$  (determined by  $P_1, P_2$  and  $h$ ) and choice of difference scheme, Theorems 4.1 and 4.2 can be used to determine whether the ILU or MILU factorization results in a stable preconditioner. We now characterize some particular classes of problems and factorizations. As we will show in Section 5, the conclusions of the theorems also appear to hold for the parities of  $n$  not covered by the analysis. Hence, in the following we do not limit our conclusions to particular parities. Recall

that for centered-difference discretization,  $\beta = -(1 + p_1)$ ,  $\gamma = -(1 + p_2)$ ,  $\delta = (-1 + p_1)/\alpha$  and  $\eta = (-1 + p_2)/\alpha$ .

1. *Centered Differences*,  $p_1, p_2 \geq 0$ , *ILU*. By Theorem 3.2,  $\alpha = 2 + \sqrt{2 + p_1^2 + p_2^2}$ . For the lower triangle,  $\beta \leq 0$  and  $\gamma \leq 0$ , so that Case 1 of Theorem 4.1 applies. The lower triangular solve is stable if and only if

$$(4.7) \quad \lambda_n(1) = \sqrt{2 + p_1^2 + p_2^2} - (p_1 + p_2) \geq 0.$$

After simplifying, (4.7) reduces to  $p_1 p_2 \leq 1$ , which is the necessary and sufficient condition for stability of the lower triangular solve.

For the upper triangle, all four cases of Theorem 4.2 can occur, with  $\delta \leq 0$  if and only if  $p_1 \leq 1$  and  $\eta \leq 0$  if and only if  $p_2 \leq 1$ . We examine Case 2 in detail, noting without proof that the upper triangular solve is stable for the other three cases. Case 2 corresponds to  $\{(p_1, p_2) \mid p_1, p_2 > 1\}$ . By Theorem 4.2, the solve is stable for odd  $n$  when  $\mu_n(-1) \leq 0$ . After scaling by  $\alpha$ , this is equivalent to

$$(4.8) \quad p_1 + p_2 - 4 \leq \sqrt{2 + p_1^2 + p_2^2}.$$

There are now three main subcases. For Case 2a, if  $p_1 + p_2 - 4 \leq 0$ , then (4.8) is trivially true. Otherwise, squaring both sides and simplifying gives  $p_2(p_1 - 4) \leq 4p_1 - 7$  as the condition for stability. This is trivially true when  $p_1 = 4$ . Cases 2b and 2c are determined by the two branches of the hyperbola  $p_2 = (4p_1 - 7)/(p_1 - 4)$ . Case 2b consists of the set of points  $(p_1, p_2)$  in the upper right quadrant of  $\mathbf{R}^2$  for which  $p_1 < 4$  and  $p_2 \geq (4p_1 - 7)/(p_1 - 4)$ . Case 2c consists of the points in the upper right quadrant for which  $p_1 > 4$  and  $p_2 \leq (4p_1 - 7)/(p_1 - 4)$ . A diagram of the stability region in the  $p_1 - p_2$  plane for the upper solve, with labels for the various cases and subcases, is given in Figure 6. (The hyperbola branch for Case 2b is not shown.)

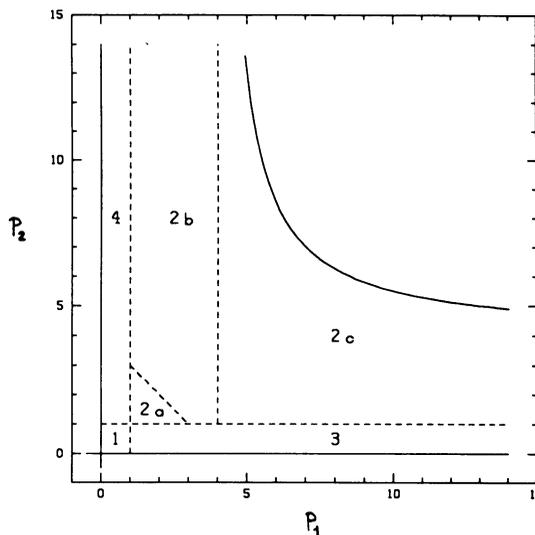


FIGURE 6

Labeled stability regions,  $p_1, p_2 \geq 0$ , upper triangle, *ILU*.

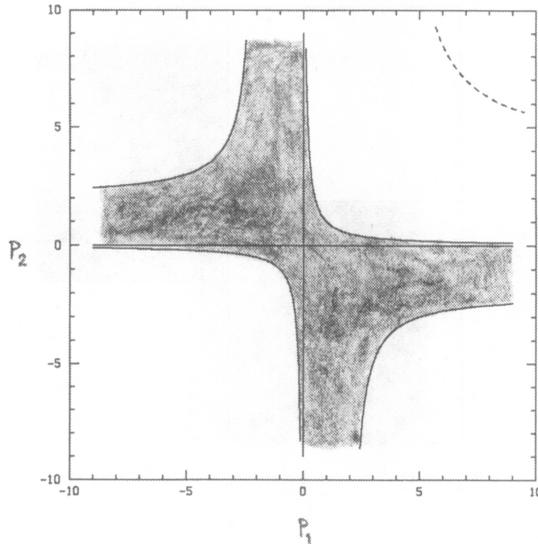


FIGURE 7  
Stability regions for the ILU factorization.

The stability region for this problem class is the intersection of the regions for the lower and upper triangular solves. This region is shown in the first quadrant in Figure 7. The dashed curve is the boundary of the stability region for the upper triangle, showing that the lower solve has more stringent stability restrictions than the upper solve. In the simple case of  $p_1 = p_2 = p$ , the stability bound is  $p \leq 1$ .

2. *Centered Differences,  $p_1 \leq 0, p_2 \geq 0, ILU$ .* The analysis for both the lower triangle and the upper triangle is essentially the same as that given for the upper triangle in the previous example. For both solves, the stability region consists of the set of points  $(p_1, p_2)$  in the second quadrant of  $\mathbf{R}^2$  satisfying

$$p_2 \leq \frac{2p_1 + 1}{p_1 + 2}.$$

This region is shown in the second quadrant of Figure 7. In the case of  $-p_1 = p_2 = p$ , the stability bound is  $p \leq 2 + \sqrt{3}$ .

3. *Centered Differences,  $p_1, p_2 \geq 0, MILU$ .* For the lower triangle, Case 1 of Theorem 4.1 applies, and  $\lambda_n(1) = 0$ . Hence, the lower triangular solve is always stable. For the upper triangle, all four cases can occur, and the solve is stable in each case. The analysis for all four cases is straightforward. We present Case 3 as representative:  $\delta \geq 0$  ( $p_1 \geq 1$ ) and  $\eta \leq 0$  ( $0 \leq p_2 \leq 1$ ). The solve is stable for even  $n$  if and only if  $\mu_n(-1) \geq 0$ . But

$$\alpha\mu_n(-1) = 2(1 + p_2) \geq 2 > 0.$$

4. *Centered Differences,  $p_1 \leq 0, p_2 \geq 0, MILU$ .* For this class of problems,  $\gamma \leq 0$ , and either  $\beta \leq 0$  for  $-1 \leq p_1 \leq 0$  and Case 1 of Theorem 4.1 applies; or  $\beta \geq 0$  for  $p_1 \leq -1$  and Case 3 applies. In the first instance,

$$\lambda_n(1) = \alpha + \beta + \gamma = \begin{cases} 0 & \text{if } p_2 \geq -p_1, \\ -2(p_1 + p_2) & \text{if } p_2 \leq -p_1. \end{cases}$$

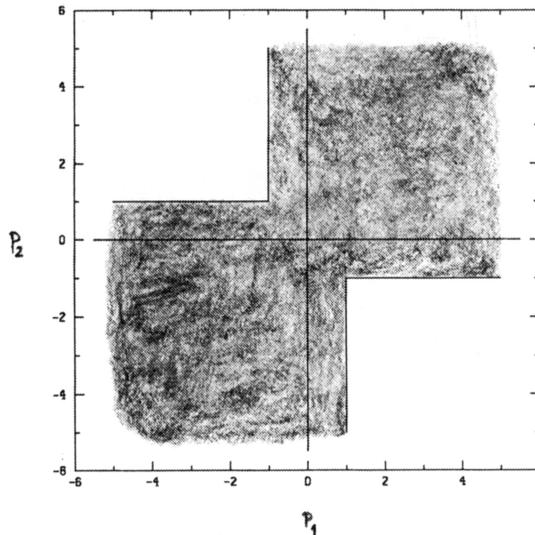


FIGURE 8

*Stability regions for the MILU factorization.*

Both these expressions are nonnegative, so that the lower solve is stable for  $-1 \leq p_1 \leq 0$ . In the second instance,

$$\lambda_n(-1) = \alpha - \beta + \gamma = \begin{cases} 2(1 + p_1) & \text{if } p_2 \geq -p_1, \\ 2(1 - p_2) & \text{if } p_2 \leq -p_1. \end{cases}$$

The first expression is negative for  $p_1 \leq -1$ , and the second expression is nonnegative for  $0 \leq p_2 \leq 1$ . Thus, the lower solve is stable for  $-1 \leq p_1 \leq 0$  or  $0 \leq p_2 \leq 1$ . The analysis for the upper triangle is similar and gives the same stability region.

The stability regions for problem classes 3 and 4 are shown in the first and second quadrants of Figure 8. For both ILU and MILU, the stability regions in the third quadrant are the reflections over the diagonal line  $p_2 = -p_1$  of the regions in the first quadrant, and the stability regions in the fourth quadrant are the reflections over the line  $p_2 = p_1$  of the regions of the second quadrant. This can be seen by replacing  $p_1$  and  $p_2$  by  $-q_1$  and  $-q_2$ , respectively, in the analysis of the four examples above. Figures 7 and 8 show the full stability regions.

We return to the question of multiple roots for MILU and centered differences. Recall that the characteristic polynomials may have multiple roots in the second or fourth quadrants for some choices of  $p_1$  and  $p_2$  for which both have modulus greater than one. However, as we have just shown, the requirement that the maximal root be bounded by one is also violated for these values of  $p_1$  and  $p_2$ , so that the presence of multiple roots does not affect the stability analysis.

5. *Upwind Differences.* By assumption, we are examining the upwind scheme only for  $p_1 \geq 0$ ,  $p_2 \geq 0$ . For the lower triangle of both the ILU and MILU factorizations,  $\beta = -(1 + 2p_1)$ ,  $\gamma = -(1 + 2p_2)$ . Both these quantities are negative, so that Case 1 of Theorem 4.1 applies. For ILU,  $\alpha = 1 + \xi + \sqrt{1 + \xi^2}$ , where  $\xi = 1 + p_1 + p_2 \geq 1$ , and  $\lambda_n(1) = 1 - \xi + \sqrt{1 + \xi^2} > 1$ . Hence, the lower triangular solve is always stable. For MILU,  $\alpha = 2\xi$ , so that  $\lambda_n(1) = 0$  and the lower triangular solve is also always stable. For the upper triangle of both factorizations,  $\delta = -1/\alpha$ ,

$\eta = -1/\alpha$ , so that Case 1 of Theorem 4.2 applies. But

$$\mu_n(1) = 1 - \frac{2}{\alpha} = \begin{cases} 1 - \frac{2}{1 + \xi + \sqrt{1 + \xi^2}} \geq 1 - \frac{2}{2 + \sqrt{2}} > 0 & \text{for ILU,} \\ 1 - \frac{2}{2\xi} \geq 0 & \text{for MILU.} \end{cases}$$

Hence, for the upwind scheme, both upper triangular solves are always stable.

**5. Correlation of Numerical Performance and Stability.** In this section, we show that there is a correlation between the numerical properties of the true preconditioned operators and the stability of the constant-coefficient preconditioning solves. Let  $\hat{A}$  denote the preconditioned operator,  $AQ^{-1}$ . For various values of the parameters  $p_1$  and  $p_2$  (determined by  $P_1$ ,  $P_2$  and  $h$ ), we examine four properties of the preconditioned matrix and linear system:

1. The extreme eigenvalues of the symmetric part  $(\hat{A} + \hat{A}^T)/2$ ;
2. The extreme real parts of the eigenvalues of  $\hat{A}$ ;
3. The performance of Orthomin(1) with preconditioning by the incomplete LU factorizations;
4. The performance of GMRES(20) [18] with preconditioning by the incomplete LU factorizations.

All computations were performed on a VAX11-780 in double precision (55-bit mantissa).

The eigenvalues were computed by Arnoldi's method with Chebyshev acceleration, which can compute eigenvalues with algebraically largest or smallest real parts efficiently if these eigenvalues are well separated from the interior eigenvalues [17]. This method repeatedly computes matrix-vector products of the form  $AQ^{-1}v$  and, for the symmetric part,  $(AQ^{-1})^T v$ , so that the preconditioning triangular solves figure prominently in the computations. (The transpose operations  $L^{-T}v$  and  $U^{-T}v$  have the same stability properties as  $L^{-1}v$  and  $U^{-1}v$ , respectively.) The stopping criterion for the eigenpair estimates  $(\lambda, v)$  is  $\|\hat{A}v - \lambda v\| \leq 10^{-6}$ , where  $\|v\| = 1$ . (For entries in the tables below marked with an asterisk (\*), convergence of the eigenvalue computations was slow and the stopping criterion was not satisfied. In these cases, the extreme eigenvalues are not well separated, but the values shown give a reasonable idea of the approximate values of the set of extreme eigenvalues [17].)

For positive integer  $k$ , Orthomin( $k$ ) and GMRES( $k$ ) generate a sequence of approximate solutions  $\{x_j\}$  to (2.2) that minimize  $\|g - Ax_j\|$  over a space of dimension at most  $k + 1$  [6], [18]. Recall that Orthomin( $k$ ) is known to converge only when the coefficient matrix has positive-definite symmetric part [6]. In our experience, Orthomin( $k$ ) is more robust when more directions are used; we use  $k = 1$  to try to identify when a preconditioning is weak. GMRES( $k$ ) will solve arbitrary nonsingular problems for large enough  $k$ , although for any given value of  $k$  it is only guaranteed to compute a sequence of iterates with nonincreasing residuals. For testing these methods, the right-hand side of (2.1) was chosen so that (2.4) is the continuous solution, and the initial guess for the discrete solution was zero, as in Section 2. In the examples below, we allowed these methods to run for at most 100 iterations, and  $\infty$  indicates that the residual norms  $\{\|r_i\|\}$  stopped decreasing at some point during the run.

TABLE 3

*Eigenvalues and performance of iterative solvers, centered differences,*  
 $P_1 = P_2 = P$ ,  $h = \frac{1}{32}$ , ILU. Predicted stability bound  $p = 1$ .

$P$	$p = Ph$	Eigenvalues of SymmetricPart		Real Parts of Eigenvalues of Operator		Iterations of Orthomin(1)	Iterations of GMRES(20)
		Leftmost	Rightmost	Leftmost	Rightmost		
20	.625	.159E0	.111E1	.694E0	.110E1	19	11
30	.9375	.696E0	.102E1	.949E0	.102E1	6	6
40	1.25	-.148E1	.548E1	.939E0*	.119E1	17	8
50	1.5625	-.493E2	.561E2	.866E0*	.138E1	$\infty$	11
60	1.875	-.393E3	.401E3	.800E0*	.152E1	$\infty$	13
100	2.1875	-.393E5	.393E5	.621E0*	.219E1	$\infty$	27
150	4.6875	-.520E6	.520E6	.466E0*	.702E1	$\infty$	74
175	5.4688	-.115E7	.115E7	-.989E2	.445E1	$\infty$	> 100
200	6.25	-.223E7	.222E7	-.775E4	.553E1	$\infty$	> 100
225	7.0313	-.495E7	.477E7	-.177E6	Overflow	$\infty$	$\infty$

We first consider examples from the four classes of problems derived from centered differences at the end of Section 4. In Table 3, we treat the first class: centered differences,  $p_1, p_2 \geq 0$ , ILU preconditioning. For fixed  $h = \frac{1}{32}$ , we examine various values of  $P = P_1 = P_2 > 0$ , with  $p = Ph$ . The stability boundary for this set of problems is  $p = 1$ . The most dramatic correlation between stability and performance is in the eigenvalues of the symmetric part. As  $p$  increases through 1, the leftmost computed eigenvalues change from positive to negative, and both extreme eigenvalues grow rapidly as  $p$  increases. Failure of Orthomin(1) to converge to the correct solution coincides almost exactly with the set of values  $p$  giving negative eigenvalues for the symmetric part. The eigenvalue computations for the matrix  $AQ^{-1}$  itself appear to be less sensitive to stability, although for  $p \gg 1$ , the leftmost real parts also become large and negative. Similarly, the performance of GMRES(20) degrades for large values of  $p$ , but it is less sensitive than that of Orthomin(1).

In Table 4, we consider examples from the second class of problems of Section 4: centered differences,  $p_1 \leq 0, p_2 \geq 0$ , ILU preconditioning. Again, we fix  $h = \frac{1}{32}$  and vary  $P = -P_1 = P_2 > 0$ . The stability bound for this set of problems is  $p = Ph = 2 + \sqrt{3} \approx 3.732$ . The results are similar to those of the previous example. As  $p$  increases through the stability bound, the leftmost eigenvalues of the symmetric part change in sign, and the performance of Orthomin(1) degrades. The diminished effectiveness of GMRES(20) coincides more closely with instability of the preconditioning than in the previous example, although the eigenvalues of  $\hat{A}$  again appear to be less sensitive than those of the symmetric part.<sup>†</sup>

Tables 5 and 6 show the results for problems from the third and fourth classes, respectively: centered differences, MILU preconditioning, and either  $p_1 = p_2 = p > 0$  or  $p = -p_1 = p_2 > 0$ . For all the values considered, Theorem 3.3 holds. The analysis of Section 4 predicts that there are no stability restrictions for the problems of Table 5. The numerical results are consistent with this: all of the extreme eigenvalues vary smoothly with  $p$ , and neither iterative method has difficulty solving

<sup>†</sup>Of course, the real parts are not the only indicators of large eigenvalues. Although the Chebyshev-Arnoldi method is not specifically designed to find eigenvalues with large imaginary parts, the Arnoldi computation does compute a set of estimates whose real parts lie between the extreme real parts. We did not observe any such eigenvalue estimates with imaginary parts that were significantly larger than their extreme real parts for either of the first two examples.

TABLE 4

*Eigenvalues and performance of iterative solvers, centered differences,  $-P_1 = P_2 = P$ ,  $h = \frac{1}{32}$ , ILU. Predicted stability bound  $p = 3.732$ .*

$P$	$p = Ph$	Eigenvalues of Symmetric Part		Real Parts of Eigenvalues of Operator		Iterations of Orthomin(1)	Iterations of GMRES(20)
		Leftmost	Rightmost	Leftmost	Rightmost		
50	1.5625	.433E - 1	.156E1	.863E0	.135E1*	32	19
60	1.875	.415E - 1	.177E1	.807E0*	.129E1*	32	19
100	2.1875	.340E - 1	.380E1	.562E0*	.219E1	43	31
110	3.4375	.320E - 1	.512E1	.575E0*	.264E1	14	13
120	3.75	.289E - 1	.744E1	.527E0*	.327E1	63	55
130	4.0625	-.176E1	.121E2	.547E0*	.418E1	83	76
140	4.375	-.999E1	.242E2	.554E0*	.548E1	$\infty$	98
150	4.6875	-.481E2	.681E2	.399E0*	.702E1	$\infty$	> 100
175	5.4688	-.261E4	.251E4	-.989E2	.445E1	$\infty$	> 100
200	6.25	-.764E5	.687E5	-.775E4	.553E1	$\infty$	> 100
225	7.0313	-.125E7	.107E7	-.177E6	Overflow	$\infty$	> 100

TABLE 5

*Eigenvalues and performance of iterative solvers, centered differences,  $P_1 = P_2 = P$ ,  $h = \frac{1}{32}$ , MILU. No stability bound.*

$P$	$p = Ph$	Eigenvalues of Symmetric Part		Real Parts of Eigenvalues of Operator		Iterations of Orthomin(1)	Iterations of GMRES(20)
		Leftmost	Rightmost	Leftmost	Rightmost		
30	.9375	.996E0	.106E1	.100E1	.104E1	4	4
50	1.5625	.681E0	.103E1	.780E0	.100E1	7	7
100	2.1875	.370E0	.121E1	.479E0*	.100E1	12	12
150	4.6875	.141E0	.136E1	.352E0	.100E1*	16	15
200	6.25	-.100E - 1	.146E1	.251E0	.100E0	19	18
225	7.0313	-.663E - 1	.150E1	.215E0	.978E0*	20	19

TABLE 6

*Eigenvalues and performance of iterative solvers, centered differences,  $-P_1 = P_2 = P$ ,  $h = \frac{1}{32}$ , MILU. Predicted stability bound  $p = 1$ .*

$P$	$p = Ph$	Eigenvalues of Symmetric Part		Real Parts of Eigenvalues of Operator		Iterations of Orthomin(1)	Iterations of GMRES(20)
		Leftmost	Rightmost	Leftmost	Rightmost		
30	.9375	.845E0	.149E2	.100E1	.323E1	52	35
31	.9688	.834E0*	.138E2	.100E1	.364E1	52	35
32	1	.819E0*	.134E2	.100E1	.395E1	51	36
33	1.0313	-.134E4	.159E3	-.708E0	.874E1	$\infty$	55
34	1.0625	-.138E4	.138E4	-.112E3	.675E1	$\infty$	> 100
36	1.125	-.287E3	.299E3	-.112E2	.442E1	$\infty$	> 100
50	1.5625	-.290E11	.290E11	-.343E2	.567E3	$\infty$	$\infty$

the preconditioned problem. (We will comment below on the change in sign of the leftmost eigenvalue of the symmetric part at  $p = 6.25$ .) The stability bound for the problems of Table 6 is  $p = 1$ . For this set of problems, all four performance indicators change dramatically when  $p$  increases through 1. (Note in particular that the values of  $p$  are significantly higher in Table 5 than in Table 6.)

As we noted in Section 4, some of the stability analysis of Theorems 4.1 and 4.2 applies only for certain parities of  $n$ , but the numerical performance seems independent of parity. As an example, consider the case of ILU preconditioning with centered differences in the second quadrant (see Figure 7). The analysis at the stability boundary makes use of Case 3 of Theorem 4.1 and Case 4 of Theorem 4.2, both of which require  $n$  to be even. The performance of ILU shown in Table 4 is for  $n = 31$ , i.e.,  $n$  odd, suggesting that parity plays no role. Further evidence is given in Table 7, which shows that for values of  $p$  near the stability boundary the performance for  $n = 32$  ( $h = \frac{1}{33}$ ) is essentially the same.

TABLE 7

Check on effect of parity of  $n$ . Eigenvalues and performance of iterative solvers, centered differences,  $-P_1 = P_2 = P$ ,  $h = \frac{1}{33}$ , ILU.

$P$	$p = Ph$	Eigenvalues of Symmetric Part		Real Parts of Eigenvalues of Operator		Iterations of Orthomin(1)	Iterations of GMRES(20)
		Leftmost	Rightmost	Leftmost	Rightmost		
110	3.33	.309E - 1	.468E1	.582E0*	.249E1	48	38
120	3.64	.287E - 1	.659E1	.555E0*	.305E1	58	55
130	3.94	-.506E1	.102E2	.590E0*	.385E1	86	75
140	4.24	-.591E1	.187E2	.631E0*	.500E1	> 100	81

TABLE 8

Performance when only one first derivative term is present.

Eigenvalues of the symmetric part, centered differences,  $P_1 = 0$ ,  $h = \frac{1}{32}$ .

$P_2$	$p_2 = P_2 h$	ILU		MILU	
		Leftmost	Rightmost	Leftmost	Rightmost
30	.9375	.858E - 1	.123E0	.242E0	.305E1
50	1.5625	.172E0	.132E1	.240E0	.259E1
100	2.1875	.386E1	.147E1*	.322E0	.213E1*
150	4.6875	.437E0	.152E1*	.398E0	.192E1*
200	6.25	.477E0	.153E1*	.455E0*	.178E1
225	7.0313	.494E0*	.152E1	.477E0	.173E1*

TABLE 9

Performance for upwind differencing. Eigenvalues of the symmetric part, upwind differences,  $P_1 = P_2 = P > 0$ ,  $h = \frac{1}{32}$ .

$P$	$p = Ph$	ILU		MILU	
		Leftmost	Rightmost	Leftmost	Rightmost
30	.9375	.553E - 1	.118E1	.629E0	.254E1
50	1.5625	.803E - 1	.115E1	.846E0	.198E1
100	2.1875	.150E0	.111E1	.951E0	.151E1
150	4.6875	.218E0	.109E1	.973E0	.135E1
200	6.25	.281E0	.107E1	.980E0	.126E1
225	7.0313	.309E0	.106E1	.983E0	.124E1

Our previous experience with these preconditionings has been on problems with just one first derivative term present, for which the coefficient is positive. Both preconditioners have been very successful in solving such problems [7], [8]. The stability analysis of Section 4 suggests that neither preconditioning suffers from instability in these cases, or for problems where upwind differences are used. Tables 8 and 9 show that the extreme eigenvalues of the symmetric parts for some problems of these types are well behaved, as expected.

Note that we are not addressing the question of *accuracy* of the incomplete factorization, and we cannot be certain that it is instability rather than inaccuracy that is causing the difficulties exhibited in Tables 3, 4, and 6. However, only in cases in which at least one of the triangular solves is unstable in the sense of Section 4 do the computed eigenvalues increase dramatically as they do in these three tables. We have encountered problems with unfavorable eigenvalue distributions that appear to be due to inaccuracy of the incomplete factorization. One such case is the example of Table 5: the eigenvalues of the symmetric part turn negative when  $p \geq 6.25$ , but they do not change dramatically in magnitude with small changes in  $p$ . Another

TABLE 10

*Inaccuracy vs. instability. Eigenvalues of the symmetric part, centered differences,  $P_2 = 100$ ,  $h = \frac{1}{64}$ , MILU. Stability bound  $p_1 = -1$ .*

$P_1$	$p_1 = P_1 h$	Leftmost	Rightmost
-10	-.1563	-.633E0	.402E1
-20	-.3125	-.992E0	.509E1
-40	-.625	-.196E1	.886E1
-60	-.9375	-.955E1	.221E2
-64	-1.0	-.282E2	.546E2
-68	-1.0625	-.463E3	.529E9
-72	-1.125	-.715E6	.716E6
-76	-1.1875	-.152E9	.152E9

example is as follows: we let  $h = \frac{1}{64}$  and  $P_2 = 100$  be fixed so that  $p_2 = 1.5625$ , we vary  $P_1 \leq 0$ , and we use centered differences and the MILU factorization. (This problem is a member of the fourth class of problems analyzed at the end of Section 4.) The extreme eigenvalues of the symmetric part are shown in Table 10. The leftmost eigenvalues are negative for all values of  $p_1$  considered, but they are fairly well behaved until the stability bound of  $p_1 = -1$  is reached, after which they quickly diverge.

We conclude with two observations that we have been unable to explain. First, in Tables 3 and 4, the extreme real parts of the eigenvalues of the operator  $AQ^{-1}$  are the same for large  $p$ , although the signs of  $p_1$  are opposite. Indeed, we have observed in some other tests with centered differences and ILU that for large  $p_1$  and  $p_2$ , the extreme eigenvalues of  $AQ^{-1}$  appear to be independent of the signs of  $p_1$  and  $p_2$ .<sup>‡</sup> Second, for very large values of the parameters  $p_1$  and  $p_2$  in Tables 3, 6, and 10, the computed extreme eigenvalues of the symmetric part are opposite in sign but nearly equal in magnitude. We do not have good explanations for these phenomena.

**6. Appendix.** In this section we prove Lemma 3.1. To simplify notation, in the proof we use the symbols “ $x$ ” and “ $y$ ” instead of “ $p_1$ ” and “ $p_2$ ”. The three inequalities to be considered are then:

$$(6.1) \quad \xi \geq -\frac{1}{4}: \quad (1+x)(2-(x+y)) \leq \frac{1}{4} \left( 4 - \frac{(1+y)(2-(x+y))}{2+|x+y|} \right)^2;$$

$$(6.2) \quad a - c(d+e)/\alpha > 0: \quad 4 - \frac{(1+y)(2-(x+y))}{2+|x+y|} > 0;$$

$$(6.3) \quad 2\alpha - a + c(d+e)/\alpha \geq 0: \quad 2|x+y| + \frac{(1+y)(2-(x+y))}{2+|x+y|} \geq 0.$$

We partition the plane into five regions:

1.  $x + y \leq 0$ ;
2.  $0 \leq x + y \leq 2$  and  $1 + y \geq 0$ ;
3.  $0 \leq x + y \leq 2$  and  $1 + y \leq 0$ ;

<sup>‡</sup>We can show that if  $|p_1| = |p'_1|$ ,  $|p_2| = |p'_2|$ , then the error matrices  $R = Q - A$  and  $R' = Q' - A'$  for the constant-coefficient factorizations are similar under a diagonal similarity transformation. This result holds for both the ILU and MILU factorizations, but we have not been able to make use of this observation.

$$4. 2 \leq x + y \text{ and } 1 + y \geq 0;$$

$$5. 2 \leq x + y \text{ and } 1 + y \leq 0.$$

These regions are depicted in Figures 9–11. Note that in Region 1,

$$(2 - (x + y))/(2 + |x + y|) = 1,$$

so that all three inequalities are simpler in this case.

In Region 1, (6.1) is equivalent to

$$\phi(x, y) \equiv \frac{1}{4} - x - \frac{1}{2}x^2 + xy + \frac{1}{4}y^2 \geq 0.$$

It is straightforward to show that  $\phi = 0$  and  $\nabla\phi = 0$  on the line  $y = 1 - 2x$ , and  $\nabla^2\phi$  is positive semidefinite everywhere. But  $\phi$  is quadratic, so for any  $X = (x, y)^T$  and  $X_0 = (x_0, y_0)^T$  such that  $y_0 = 1 - 2x_0$ ,

$$\phi(X) = \frac{1}{2}(X - X_0)^T(\nabla^2\phi)(X - X_0) \geq 0.$$

In Region 2, the right-hand side of (6.1) is greater than

$$4 - (1 + y)(2 - (x + y)),$$

which can be seen to be greater than or equal to  $(1 + x)(2 - (x + y))$  by direct computation. In Region 4, the left-hand side of (6.1) is negative when  $x \geq -1$ , so that the inequality holds trivially. The same argument works for (all of) Region 5. We have not been able to establish (6.1) for Region 3 or for Region 4 with  $x < -1$ , although numerical tests suggest that it also holds in these regions. The shaded area in Figure 9 shows the set of points in the plane for which we have demonstrated that (6.1) holds.

Inequality (6.2) holds immediately in Regions 3 and 4, since the quantity subtracted is negative. The inequality reduces to  $y < 3$  in Region 1, and for Region 2, using the fact that  $(2 - (x + y))/(2 + |x + y|) < 1$ , the same condition is sufficient.

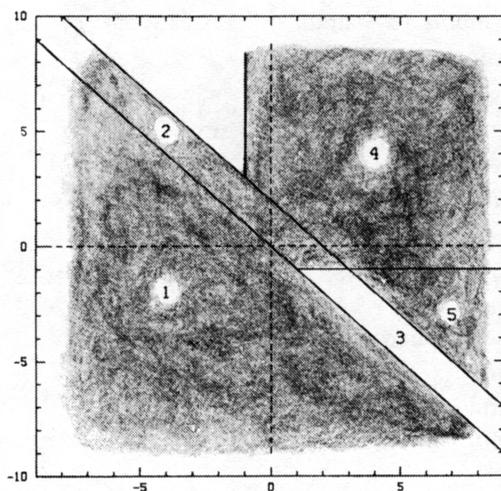


FIGURE 9

Values of  $x$  and  $y$  where inequality (6.1) holds.

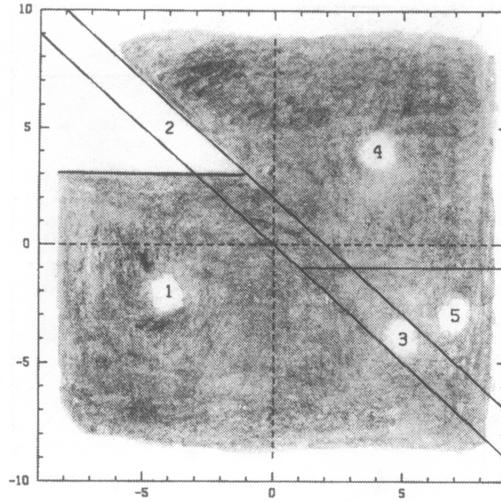


FIGURE 10  
 Values of  $x$  and  $y$  where inequality (6.2) holds.

For Region 5, after being scaled by  $2 + x + y$ , the right-hand side of (6.2) can be shown to be greater than

$$8 + 3x + y = 8 + x + y + 2x \geq 10 + 2x \geq 13,$$

since  $x + y \geq 2$ ,  $-y \geq 1$  and  $x \geq 2 - y \geq 3$ . Figure 10 shows the set of points where (6.2) holds.

Inequality (6.3) holds trivially in Regions 2 and 5. In Region 1, the inequality reduces to  $y \leq 1 - 2x$ . In Region 3, the fact that  $(2 - (x + y))/(2 + |x + y|) < 1$  makes the right-hand side greater than  $2x + 3y + 1$ , which is greater than 0 for  $y > -\frac{1}{3}(1 + 2x)$ . This determines a small subset of the region. Finally, in Region 4,  $2 + |x + y| \geq 4$ , so that if  $-1 \leq y \leq 7$ , then

$$\frac{(1 + y)}{2 + |x + y|} \leq 2.$$

Consequently, the right-hand side of (6.3) is greater than

$$2(x + y) + 2(2 - (x + y)) = 4 > 0.$$

Hence, (6.3) holds in the shaded region of Figure 11.

The region specified by Lemma 3.1 and Figure 1 is the intersection of the regions of Figures 9–11. Recall that the inequalities of the lemma give sufficient conditions for the convergence result in Theorem 3.3 for MILU and centered differences. Numerical evidence suggests that conditions (6.2) and (6.3) are also necessary for the diagonal sequences in the (true) MILU factorization to be convergent. For example, the values  $p_1 = -4$  and  $p_2 = 4$  violate condition (6.2), and with these values (from  $h = \frac{1}{32}$ ,  $P_1 = -128$  and  $P_2 = 128$ ), the MILU diagonal sequence appears not to be convergent. Similarly, the values  $p_1 = 2$  and  $p_2 = -2$  violate condition (6.3), and the MILU sequence is also not convergent for them (from  $h = \frac{1}{32}$ ,  $P_1 = 64$  and  $P_2 = -64$ ).

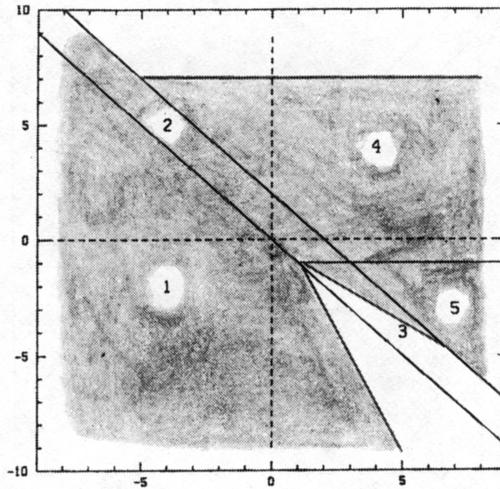


FIGURE 11

Values of  $x$  and  $y$  where inequality (6.3) holds.

**Acknowledgments.** The author would like to thank Stan Eisenstat for a part of the proof of Case 1 of Theorem 4.1, and Youcef Saad for his help in using his Chebyshev-Arnoldi eigenvalue routines for computing eigenvalues and for pointing out reference [22].

Department of Computer Science  
Yale University  
New Haven, Connecticut 06520

1. OWE AXELSSON, "Conjugate gradient type methods for unsymmetric and inconsistent systems of linear equations," *Linear Algebra Appl.*, v. 29, 1980, pp. 1-16.

2. OWE AXELSSON & I. GUSTAFSSON, "A modified upwind scheme for convective transport equations and the use of a conjugate gradient method for the solution of non-symmetric systems of equations," *J. Inst. Math. Appl.*, v. 23, 1979, pp. 321-337.

3. R. CHANDRA, *Conjugate Gradient Methods for Partial Differential Equations*, Ph.D. Thesis, Department of Computer Science, Yale University, 1978. Also available as Technical Report 129.

4. P. CONCUS, G. H. GOLUB & DIANNE P. O'LEARY, "A generalized conjugate gradient method for the numerical solution of elliptic partial differential equations," in *Sparse Matrix Computations* (James R. Bunch and Donald J. Rose, eds.), Academic Press, New York, 1976, pp. 309-332.

5. TODD DUPONT, RICHARD P. KENDALL & H. H. RACHFORD, JR., "An approximate factorization procedure for solving self-adjoint elliptic difference equations," *SIAM J. Numer. Anal.*, v. 5, 1968, pp. 559-573.

6. S. C. EISENSTAT, H. C. ELMAN & M. H. SCHULTZ, "Variational iterative methods for nonsymmetric systems of linear equations," *SIAM J. Numer. Anal.*, v. 20, 1983, pp. 345-357.

7. H. C. ELMAN, *Iterative Methods for Large, Sparse, Nonsymmetric Systems of Linear Equations*, Ph.D. Thesis, Department of Computer Science, Yale University, 1982. Also available as Technical Report 229.

8. H. C. ELMAN, "Preconditioned conjugate gradient methods for nonsymmetric systems of linear equations," in *Advances in Computer Methods for Partial Differential Equations, IV* (R. Vichnevetsky and R. S. Stepleman, eds.), International Association for Mathematics and Computers in Simulation, IMACS, 1981, pp. 409-417.

9. GEORGE E. FORSYTHE & WOLFGANG R. WASOW, *Finite-Difference Methods for Partial Differential Equations*, Wiley, New York, 1960.

10. C. WILLIAM GEAR, *Numerical Initial Value Problems in Ordinary Differential Equations*, Prentice-Hall, Englewood Cliffs, N.J., 1971.

11. IVAR GUSTAFSSON, "A class of first order factorizations," *BIT*, v. 18, 1978, pp. 142–156.
12. IVAR GUSTAFSSON, *Stability and Rate of Convergence of Modified Incomplete Cholesky Factorization Methods*, Ph.D. Thesis, Department of Computer Sciences, Chalmers University of Technology and the University of Göteborg, 1978. Also available as Technical Report 77.04R.
13. PETER HENRICI, *Elements of Numerical Analysis*, Wiley, New York, 1964.
14. J. M. HYMAN & T. A. MANTEUFFEL, "High-order sparse factorization methods for elliptic boundary value problems," in *Advances in Computer Methods for Partial Differential Equations*, V (R. Vichnevetsky and R. S. Stepleman, eds.), IMACS, 1984, pp. 551–555.
15. W. LINIGER, "On factored discretizations of the Laplacian for the fast solution of Poisson's equation on general regions," *BIT*, v. 24, 1984, pp. 592–608.
16. J. A. MEIJERINK & H. A. VAN DER VORST, "An iterative solution method for linear systems of which the coefficient matrix is a symmetric  $M$ -matrix," *Math. Comp.*, v. 31, 1977, pp. 148–162.
17. YOUSSEF SAAD, "Chebyshev acceleration techniques for solving nonsymmetric eigenvalue problems," *Math. Comp.*, v. 42, 1984, pp. 567–588.
18. Y. SAAD & M. H. SCHULTZ, *GMRES: A Generalized Minimal Residual Algorithm for Solving Nonsymmetric Linear Systems*, Technical Report 254, Department of Computer Science, Yale University, 1983.
19. P. E. SAYLOR, "Second order strongly implicit symmetric factorization methods for the solution of elliptic difference equations," *SIAM J. Numer. Anal.*, v. 11, 1974, pp. 894–908.
20. A. SEGAL, "Aspects of numerical methods for elliptic singular perturbation problems," *SIAM J. Sci. Statist. Comput.*, v. 3, 1982, pp. 327–349.
21. HENK A. VAN DER VORST, "Iterative solution methods for certain sparse linear systems with a nonsymmetric matrix arising from PDE-problems," *J. Comput. Phys.*, v. 44, 1981, pp. 1–19.
22. H. S. WALL, *Analytic Theory of Continued Fractions*, Van Nostrand, New York, 1948.
23. D. M. YOUNG & KANG C. JEA, "Generalized conjugate gradient acceleration of nonsymmetrizable iterative methods," *Linear Algebra Appl.*, v. 34, 1980, pp. 159–194.