

NUREG/CR-1943  
LA-8705-MS

TERA

The Incomplete Lower-Upper  
Conjugate Gradient (ILUCG) Method for  
Solving Large, Sparse Linear Systems



University of California



LOS ALAMOS SCIENTIFIC LABORATORY

Post Office Box 1663 Los Alamos, New Mexico 87545

8108260022

An Affirmative Action/Equal Opportunity Employer

This report was not edited by the Technical Information staff.

NOTICE

This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor any agency thereof, or any of their employees, makes any warranty, expressed or implied, or assumes any legal liability or responsibility for any third party's use, or the results of such use, of any information, apparatus, product or process disclosed in this report, or represents that its use by such third party would not infringe privately owned rights.

# **The Incomplete Lower-Upper Conjugate Gradient (ILUCG) Method for Solving Large, Sparse Linear Systems**

Charles H. Neil\*

Manuscript submitted: January 1981

Date published: February 1981

Prepared for  
Division of Reactor Safety Research  
Office of Nuclear Regulatory Research  
US Nuclear Regulatory Commission  
Washington, DC 20555

NRC FIN No. A7016

\*Graduate Research Assistant. 3209-30th Street, Lubbock, TX 79410.



CONTENTS

ABSTRACT - - - - -	1
I. INTRODUCTION - - - - -	1
II. THE METHOD OF CONJUGATE GRADIENTS - - - - -	2
III. THE INCOMPLETE L-U FACTORIZATION -- CONJUGATE GRADIENT METHOD- - -	7
IV. NUMERICAL RESULTS - - - - -	8
V. CONCLUSION - - - - -	9
REFERENCES - - - - -	12

TABLE

I ILUGG METHOD VS GAUSS-SEIDEL METHOD FOR VALUES OF $\delta$ - - - - -	10
--	----

THE INCOMPLETE LOWER-UPPER CONJUGATE GRADIENT (ILUCG)  
METHOD FOR SOLVING LARGE, SPARSE LINEAR SYSTEMS

by  
Charles H. Neil

ABSTRACT

The method of Conjugate Gradients is combined with an incomplete factorization of the coefficient matrix to produce an iterative method for approximate solution of systems of linear equations. The method is suited for systems where the coefficient matrix is large, sparse, and nonsymmetric. A comparison is presented in this document of the performance of the method vs the Gauss-Seidel method for test problems arising in the TRAC Code for reactor hydrodynamics.

I. INTRODUCTION

The discretization of partial differential equations in mathematical physics often results in the problem of solving systems of linear equations, the coefficients of which form a large, sparse matrix. The method of Conjugate Gradients (Refs. 1, 2, and 3) was originally proposed as a general iterative technique for solving linear systems but is not competitive with other iterative methods when applied to the full matrix. More recent work (Refs. 4, 5, and 6) has been concerned with applying the Conjugate Gradient method to a system modified by an incomplete decomposition of the coefficient matrix. An accurate approximate solution can thereby be obtained in relatively few iterations.

Much of the work to date deals with symmetric coefficient matrices. In the solution of hydrodynamic equations in the Transient Reactor Analysis Code<sup>7</sup> (TRAC) developed at the Los Alamos National Laboratory, however, the resulting seven-stripe pressure matrix is nonsymmetric. This paper presents a generalized Conjugate Gradient (CG) method applied to an incomplete factorization for such nonsymmetric matrices, and gives results of numerical tests of the method.

## II. THE METHOD OF CONJUGATE GRADIENTS

We will be concerned in what follows with solving the system

$$Mx = k \quad (1)$$

Here,  $M$  is an  $N$ -dimensional nonsingular matrix, not necessarily symmetric; let us denote the exact solution by  $h$ . Suppose we have available a linearly independent set of vectors  $\{\pi_0, \pi_1, \dots\}$ . We will be interested in studying an iterative method of the form

$$x_{i+1} = x_i + a_i \pi_i, \quad i = 0, 1, \dots; \quad x_0 \text{ arbitrary}, \quad (2)$$

which produces the exact solution after a finite number of steps.

One way to accomplish this is to replace each vector,  $\pi_i$ , by a vector

$$p_i = \sum_{j=0}^i c_{ij} \pi_j,$$

such that

$$\langle p_i, p_j \rangle = 0, \quad i \neq j$$

(that is, orthogonalize the set  $\{\pi_i\}$ ), and then select the constants  $a_i$  so that the error in the  $i+1$  <sup>st</sup> iterate,  $|h - x_{i+1}|^2$ , is minimized. This will occur when  $h - x_{i+1}$  is orthogonal to the set  $\{p_0, p_1, \dots, p_i\}$  and  $a_i$  is given by



$$a_i = \frac{\langle h - x_0, p_i \rangle}{\langle p_i, p_i \rangle} . \quad (3)$$

If it ever occurs that the initial error,  $e_0 = h - x_0$ , lies in the subspace  $S_k = \text{span} \{ p_0, \dots, p_k \}$ , then we have

$$h - x_{k+1} = h - x_0 - \sum_{j=0}^k a_j p_j \in S_k .$$

As  $h - x_{k+1}$  also was chosen orthogonal to  $S_k$ , it must be the zero vector, and therefore the method terminates.

We are thus motivated to try to select the vectors  $\{ p_i \}$  and  $\{ \pi_i \}$  in such a way that  $e_0 \in S_k$  for  $k$  is as small as possible. Furthermore, the exact solution,  $h$ , is unavailable a priori and so we must choose the vectors without explicit knowledge of  $e_0$ . The following allows us to find a way of doing this.

Lemma: Suppose  $T$  is a symmetric, positive definite matrix and  $u$  is some chosen vector. Then  $u \in \text{span} \{ Tu, T^2u, \dots, T^{k-\ell}u \}$ , where

$k$  = number of distinct eigenvalues of  $T$  and

$\ell$  = number of vanishing components corresponding to distinct eigenvalues in the eigenvector decomposition of  $u$ .

Proof Since  $T$  is symmetric and positive definite, there exists a matrix  $R$  (of eigenvectors of  $T$ ) such that  $R^{-1}TR = \Lambda$ , where  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{k-1}, \lambda_k, \dots, \lambda_k)$ , and each  $\lambda_j$  is distinct.

Let  $d$  be the vector of coefficients in the eigenvector decomposition of  $u$ ; thus  $u = Rd$ . We then have

$$T^j u = R \Lambda^j d, \text{ for } j = 0, 1, \dots .$$

We wish to be able to write for some constants  $c_j$ ,

$$u = c_1 T u + \dots + c_{k-l} T^{k-l} u,$$

or

$$R d = c_1 R \Lambda d + \dots + c_{k-l} R \Lambda^{k-l} d.$$

Equivalently, we want

$$d = c_1 \Lambda d + \dots + c_{k-l} \Lambda^{k-l} d. \quad (4)$$

Among the first  $k$  components of  $d$ , there are only  $k-l$  nonzero entries, say

$d_{j_1}, d_{j_2}, \dots, d_{j_{k-l}}$ , are nonzero. After dividing both sides of these

equations by the components of  $d$ , we can re-write the nonzero equations of (4) as

$$Lc = \underline{1}, \text{ where } L = \begin{bmatrix} \lambda_{j_1} & \lambda_{j_1}^2 & \dots & \lambda_{j_1}^{k-l} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \lambda_{j_{k-l}} & \lambda_{j_{k-l}}^2 & \dots & \lambda_{j_{k-l}}^{k-l} \end{bmatrix}.$$

$c = (c_1, \dots, c_{k-l})^T$ , and  $\underline{1}$  is a  $k-l$  dimensional vector of 1's. As the  $\lambda$ 's are all distinct,  $L$  is non-singular:  $c$  exists and the proof is complete.

Returning now to the problem of choosing the spanning set of vectors

$\{\pi_0, \pi_1, \dots\}$ , consider the set

$$\{M^T M e_0, (M^T M)^2 e_0, \dots\}.$$



We note that  $Me_0 = M(h - x_0) = k - Mx_0$ , so that once  $x_0$  is selected, all these vectors are available. Furthermore, in view of the lemma, if  $M^T M$  has many repeated eigenvalues (or in the fortunate event that  $x_0$  is chosen so that  $e_0$  is defective in many eigendirections of  $M^T M$ ), then the method (1) terminates in relatively few steps.

In creating the set  $\{p_i\}$ , orthogonalization of each  $(M^T M)^j e_0$  with respect to  $(M^T M)^i e_0$ ,  $i = 0, 1, \dots, j-1$ , is too expensive. We produce an equivalent spanning set  $\{p_0, p_1, \dots\}$  as follows.

Denote by  $r_i$  the  $i^{\text{th}}$  residual:

$$r_i = k - Mx_i, \quad i = 0, 1, \dots \quad (5)$$

Having produced  $\{p_0, p_1, \dots, p_i\}$ ,  $i \geq 0$ , obtain  $p_{i+1}$  by orthogonalizing  $M^T M e_{i+1} = M^T r_{i+1}$  with respect to  $p_i$  (instead of orthogonalizing  $(M^T M)^{i+1} e_0 = (M^T M)^i M^T r_0$  with respect to  $p_i$ ). The resulting set  $\{p_j\}$  has the same span as does  $\{(M^T M)^j e_0\}$ . The advantage of this procedure is that if  $b_i$  is chosen so that  $p_{i+1} = M^T r_{i+1} + b_i p_i$  is orthogonal to  $p_i$ , then automatically  $p_{i+1}$  is orthogonal to  $p_0, p_1, \dots, p_{i-1}$ .

For, we have

$$\begin{aligned} \langle p_{i+1}, p_{i-j} \rangle &= \langle M^T r_{i+1} + b_i p_i, p_{i-j} \rangle \\ &= \langle M^T r_{i+1}, p_{i-j} \rangle + b_i \langle p_i, p_{i-j} \rangle \\ &= \langle M^T M (h - x_{i+1}), p_{i-j} \rangle \\ &= \langle h - x_{i+1}, M^T M p_{i-j} \rangle, \end{aligned}$$

and  $h - x_{i+1}$  was chosen orthogonal to the span  $\{M^T r_0, (M^T M)M^T r_0, \dots, (M^T M)^{i-j} M^T r_0\}$ , which contains  $(M^T M)^j p_{i-j}$  for  $j \geq 1$ .

Thus,  $\langle p_{i+1}, p_{i-j} \rangle = 0, j \geq 1$ .

Therefore, to carry out the iterations we need only save the vectors  $x_i$ ,  $r_i$ , and  $p_i$ .

Thus we write the Conjugate Gradient Algorithm as follows. Choose  $x_0$  arbitrarily; set  $r_0 = k - Mx_0$ , and  $p_0 = M^T r_0$ .

For  $i = 0, 1, \dots$ , set

$$a_i = \frac{\langle r_i, r_i \rangle}{\langle p_i, p_i \rangle}$$

$$x_{i+1} = x_i + a_i p_i$$

$$r_{i+1} = r_i - a_i M p_i$$

$$b_i = \frac{\langle r_{i+1}, r_{i+1} \rangle}{\langle r_i, r_i \rangle}$$

$$p_{i+1} = M^T r_{i+1} + b_i p_i.$$

It should be noted that these formulae for  $a_i$  and  $r_{i+1}$  are equivalent to Eqs. (3) and (5), respectively (see Refs. 2 and 3).

As previously observed, the CG method terminates in relatively few steps if the matrix  $M$  has large clustering of eigenvalues. Furthermore, at each step of the iteration the vector  $x_i$  minimizes  $\|h - x\|^2$  within a certain subspace. With our choice of vectors  $\{p_i\}$ , each one being a linear combination of  $\{(M^T M)^j e_0\}_{j=0}^{i-1}$ , much of the change in  $x_i$  occurs in eigendirections corresponding to the largest eigenvalues of  $M$ . Accordingly, a large reduction in

error can be anticipated within the first few iterations. It is therefore possible that  $x_i$  may be an accurate approximation to  $h$  for  $i \ll N$ , and it is this situation we will try to produce.

### III. The INCOMPLETE L-U FACTORIZATION CONJUGATE GRADIENT METHOD

Let us now turn to the problem of solving the system

$$Ax = b, \tag{6}$$

where  $A$  is a large, sparse, nonsymmetric matrix. Suppose we can carry out the usual L-U decomposition of  $A$  (that is, where  $L$  is unit lower triangular and  $U$  is upper triangular), with the exception that elements in  $L$  or  $U$  corresponding to zero elements in  $A$  are considered zero and are neither computed nor stored. The resulting incomplete L-U decomposition of  $A$  will require no more storage than does  $A$  (and will not be as costly as the complete L-U decomposition wherein  $L$  and  $U$  are generally no longer sparse).

We have  $LU = A + E$ , where the matrix  $E$  contains errors because of the incomplete decomposition. Hence,

$$L^{-1}AU^{-1} = I - L^{-1}EU^{-1}.$$

If the matrix  $E$  is in some sense small, then the matrix  $L^{-1}AU^{-1}$  is an approximate identity -- in particular,  $L^{-1}AU^{-1}$  should tend to have many eigenvalues close to unity, and by previous remarks, the conjugate gradient method, applied to a system with  $L^{-1}AU^{-1}$  as the coefficient matrix, should converge rapidly.

We therefore re-write Eq. (6) as

$$L^{-1}AU^{-1}ux = L^{-1}b.$$

Letting  $M = L^{-1}AU^{-1}$ , and writing the algorithm in terms of  $x$  (instead of  $Ux$ ), we can write the ILUCG algorithm as follows.

Choose  $x_0$  arbitrarily; set  $r_0 = b - Ax_0$  and

$$p_0 = (U^T U)^{-1} A^T (LL^T)^{-1} r_0 .$$

For  $i = 0, 1, \dots$ , until satisfied, set

$$a_i = \frac{\langle r_i, (LL^T)^{-1} r_i \rangle}{\langle p_i, U^T U p_i \rangle}$$

$$x_{i+1} = x_i + a_i p_i$$

$$r_{i+1} = r_i - a_i A p_i$$

$$b_i = \frac{\langle r_{i+1}, (LL^T)^{-1} r_{i+1} \rangle}{\langle r_i, (LL^T)^{-1} r_i \rangle}$$

$$p_{i+1} = (U^T U)^{-1} A^T (LL^T)^{-1} r_{i+1} + b_i p_i .$$

It should be noted that the inversion and vector multiplication are carried out quickly because of the triangular and sparse nature of  $L$  and  $U$ .

#### IV. NUMERICAL RESULTS

This algorithm was applied to several test cases designed to simulate

numerical problems in the TRAC code for reactor hydrodynamics. In this application, the matrix A was a nonsymmetric 440 x 440 real matrix with unit main diagonal and with negative off-diagonal entries occurring in a pattern following a seven-point discretization in cylindrical coordinates of the Laplace operator.

The matrix was loaded uniformly across rows such that  $\sum_j = 1 - \rho_j$ , where  $\sum_j$  denotes the sum of the absolute values of the off-diagonal entries in the  $i^{\text{th}}$  row, and  $\rho_j$  is a pseudo-random number in the interval  $[0, \delta]$ , with  $\delta$  an input parameter. By allowing  $\delta \rightarrow 0^+$ , the degree of diagonal dominance of the matrix A is reduced (and its condition is worsened).

Table I contains results of a comparison of the performance of the ILUCG method and the Gauss-Seidel method on the same linear system for values of  $\delta$  ranging from 0.500 to 0.001. As shown in the table, below a certain level of diagonal dominance, the Gauss-Seidel method failed to converge to the desired residual norm within 1000 iterations, whereas in all cases the ILUCG method converged. Because of the relative expense of the ILUCG method per iteration (it is roughly 8.2 times as costly per iteration as Gauss-Seidel for matrices of this type) it is important to observe the column labeled "Ratio," which is the ratio of the number of iterations of Gauss-Seidel (GS) to ILUCG; only when this figure is greater than 8.2 will ILUCG be the preferred method.

Figure 1 illustrates the norm of the residual of the two methods, plotted vs number of multiplications, for the case  $\delta = 0.060$ . As this figure shows, the error decreases faster for ILUCG than for GS, although it is initially greater for ILUCG than for GS. The plot is typical of all cases, with the divergence of the curves more pronounced for smaller values of  $\delta$ .

Finally, Fig. 2 is a plot of the 440 eigenvalues of the original matrix A and of the matrix used in the ILUCG algorithm,  $L^{-1}AU^{-1}$ , for  $\delta = 0.060$ . The eigenvalues of A are uniformly spread from 1.97 to 0.029; but those of  $L^{-1}AU^{-1}$  are, for the most part, quite close to 1.0, indicating a reason for the effectiveness of the ILUCG method.

## V. CONCLUSION

The reported results indicate that the ILUCG method warrants consideration in situations requiring high accuracy or when the coefficient matrix is poorly

conditioned (or both). However, the cost per iteration of ILUCG in its present form makes it not a preferred method in less severe instances.

TABLE I

ILUCG METHOD vs GAUSS-SEIDEL METHOD FOR VALUES OF  $\delta$

$\delta$	Number of iterations		Ratio (Break-even = 8.2)
	GS	ILUCG	
0.500	49	16	3.1
0.400	61	19	3.2
0.300	81	22	3.7
0.200	122	26	4.7
0.100	241	34	7.1
0.080	301	36	8.4
0.060	401	39	10.3
0.040	601	42	14.3
0.020	a	47	-
0.010	a	49	-
0.008	a	50	-
0.004	a	52	-
0.001	a	55	-

<sup>a</sup>Failed to converge within 1000 iterations.



# RESIDUALS

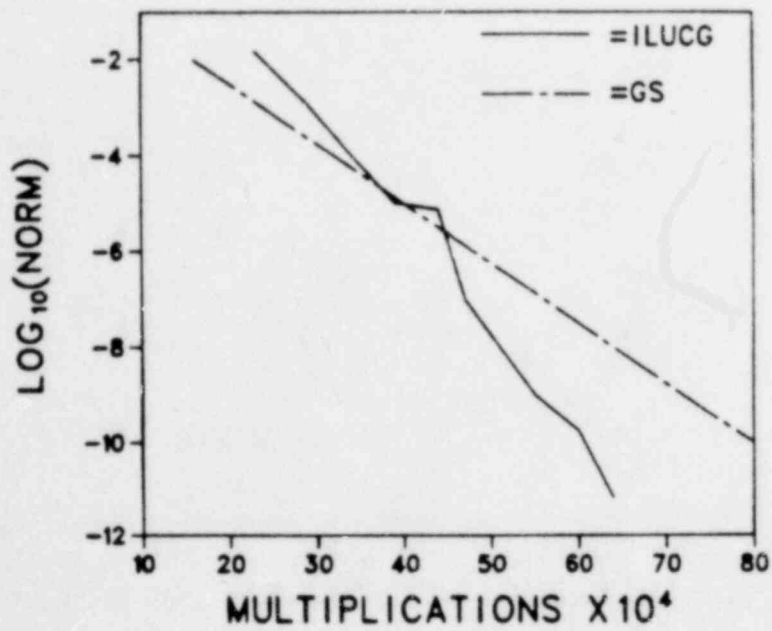


Fig. 1. Residual error.

# EIGENVALUES

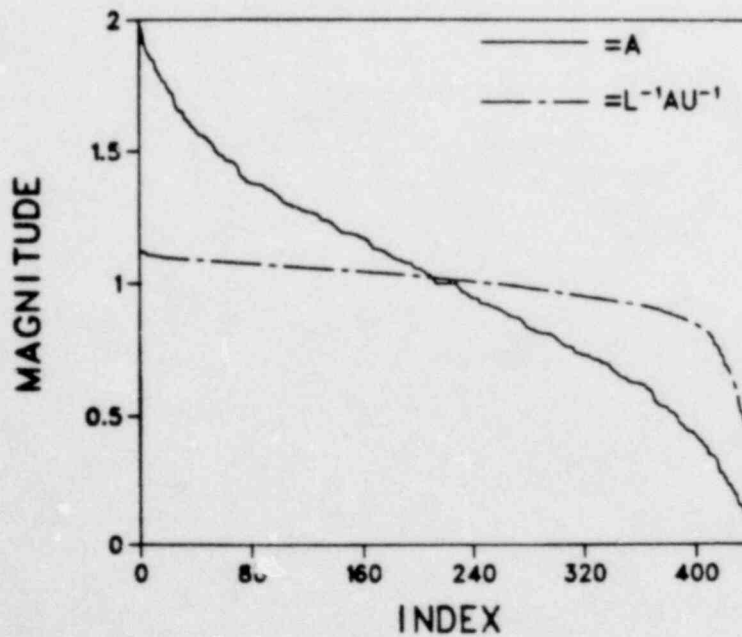


Fig. 2. Matrix eigenvalues.

## REFERENCES

1. M. R. Hestenes and E. Steifel, "Methods of Conjugate Gradients for Solving Linear Systems," J. of Research of the National Bureau of Standards 59, 409-436(December, 1952).
2. M. R. Hestenes, "The Conjugate-Gradient Method for Solving Linear Systems," Proc. Sympos. Appl. Math., Numerical Analysis, (McGraw-Hill, New York, 1956) Vol. VI, pp. 83-102.
3. J. K. Reid, "On the Method of Conjugate Gradients for the Solution of Large Sparse Systems of Linear Equations," Proc. Conf. on Large Sparse Systems of Linear Equations (Academic Press, New York, 1971).
4. J. A. Meijerink and H. A. van der Vorst, "An Iterative Solution Method for Linear Systems of Which the Coefficient Matrix is a Symmetric M-Matrix," Math. Comp. 31, 148-162 (January, 1977).
5. P. Concus, G. Golub, and D. O'Leary, "A Generalized Conjugate Gradient Method for the Numerical Solution of Elliptic Partial Differential Equations," Lawrence Berkeley Laboratory report LBL-4604 (1975).
6. D. S. Kershaw, "The Incomplete Cholesky-Conjugate Gradient Method for the Iterative Solution of Systems of Linear Equations," J. Comp. Phys 26, 43-65 (1978).
7. "TRAC-PIA, An Advanced Best Estimate Computer Program for PWR LOCA Analysis," Los Alamos Scientific Laboratory report LA-7777-MS (May 1979).

DISTRIBUTION

	<u>Copies</u>
Nuclear Regulatory Commission, R4, Bethesda, Maryland	388
Technical Information Center, Oak Ridge, Tennessee	2
Los Alamos National Laboratory, Los Alamos, New Mexico	<u>50</u>
	440