

General Disclaimer

One or more of the Following Statements may affect this Document

- This document has been reproduced from the best copy furnished by the organizational source. It is being released in the interest of making available as much information as possible.
- This document may contain data, which exceeds the sheet parameters. It was furnished in this condition by the organizational source and is the best copy available.
- This document may contain tone-on-tone or color graphs, charts and/or pictures, which have been reproduced in black and white.
- This document is paginated as submitted by the original source.
- Portions of this document are not fully legible due to the historical nature of some of the material. However, it is the best reproduction available from the original submission.

NATIONAL AERONAUTICS AND SPACE ADMINISTRATION

Technical Memorandum 33-627

*Sparse Matrix Methods Based on Orthogonality
and Conjugacy*

C. L. Lawson

(NASA-CR-133231) SPARSE MATRIX METHODS
BASED ON ORTHOGONALITY AND CONJUGACY (Jet
Propulsion Lab.) 66 p HC \$5.50 CSCL 12A

N73-26617

Unclas
G3/19 08275



**JET PROPULSION LABORATORY
CALIFORNIA INSTITUTE OF TECHNOLOGY
PASADENA, CALIFORNIA**

June 15, 1973

NATIONAL AERONAUTICS AND SPACE ADMINISTRATION

Technical Memorandum 33-627

*Sparse Matrix Methods Based on Orthogonality
and Conjugacy*

C. L. Lawson

**JET PROPULSION LABORATORY
CALIFORNIA INSTITUTE OF TECHNOLOGY
PASADENA, CALIFORNIA**

June 15, 1973

PREFACE

The work described in this report was performed by the Data Systems Division of the Jet Propulsion Laboratory.

CONTENTS

Part I.	Introduction	1
	Chapter 1. Introduction	1
	Chapter 2. Mathematical Background	6
Part II.	The Conjugate Gradient Algorithm and Variations	8
	Chapter 3. Solving a Consistent System, $Ax=b$, where A is Symmetric and Non- negative Definite	8
	Chapter 4. Solving a Least Squares Problem, $Ax \approx b$	19
	Chapter 5. Solving a Consistent System $Ax=b$	22
Part III.	Other Algorithms	26
	Chapter 6. Solving a Least Squares Problem, $Ax \approx b$	26
	Chapter 7. The Theoretical Equivalence of Algorithms CGLS and ITLS	39
	Chapter 8. Solving a Consistent System $Ax=b$	44
	Chapter 9. The Theoretical Equivalence of Algorithms CGC and ITC	53
	Chapter 10. Solving a Consistent System $Ax=b$ where A is Symmetric	55
	References	66

FIGURE

3.1	Vectors and subspaces related to Algorithm CG	13
-----	---------------------------------------------------------	----

ABSTRACT

A matrix having a high percentage of zero elements is called sparse. In the solution of systems of linear equations or linear least squares problems involving large sparse matrices significant saving of computer cost can be achieved by taking advantage of the sparsity. This Memorandum derives and describes the well known conjugate gradient algorithm and a set of related algorithms which are applicable to such problems.

Control of accuracy is a serious problem with this class of methods. We plan to devote a subsequent study to methods of controlling algorithms of this class.

PART I
Introduction

Chapter 1 Introduction

A matrix A is called sparse if a large proportion of its elements are zero. Significant savings of execution time and computer storage can be realized in solving systems of linear equations, $Ax=b$, or least squares problems, $Ax \cong b$, if the matrix A is sparse and if special solution methods are used which take advantage of the sparseness.

Most direct solution methods such as Gaussian elimination or Householder orthogonal triangularization perform transformations on the given matrix A which significantly increase the number of nonzero elements. The two main ideas in developing sparse matrix methods have been

1. Reorganize direct elimination methods to reduce the growth of the number of nonzero elements.

and

2. Use iterative methods so that the original sparse matrix A is used throughout the computation.

In Reid (1971) it is pointed out that the conjugate gradient (CG) method for solving a system $Ax=b$, with A symmetric and positive definite, is well suited for use when A is sparse. This method shares with iterative methods the feature of continually using the initial matrix A but it shares with direct methods the property of theoretically terminating after not more than n iterations.

More specifically the major computational cost in each iteration of the CG method arises from the multiplication of A times a vector. If A is an $n \times n$

matrix and the proportion of nonzero elements in A is ρ then the multiplication of A times an n-vector requires ρn^2 multiplications and additions if multiplication by zero elements of A are skipped. This algorithm theoretically reaches the solution vector in at most n iterations so the number of multiplications and additions would theoretically be at most ρn^3 .

The storage required is just that needed for the sparse matrix A plus four n-vectors. Since A is symmetric it has only approximately $\rho n^2/2$ potentially distinct nonzero elements. Such a matrix can be stored in ρn^2 locations or even in less space if the nonzero elements are located in some known regular pattern.

In comparison direct solution of this problem using an efficient stable method such as Cholesky factorization would in general require about $n^2/2$ storage locations and $n^3/6$ multiplications and additions. Thus the CG method will require less storage than the Cholesky method if $\rho < 1/2$ and will require fewer multiplications and additions if $\rho < 1/6$.

These cut-off values for ρ should be taken only as very rough guidelines. Storage management for the sparse matrix method may increase its execution time.

More serious is the lack of numerical stability in the CG method. If the matrix A has a large condition number the intermediate vectors computed by the algorithm which are theoretically orthogonal or A-conjugate (defined in Chapter 2) may not even come close to having these properties. I feel that this means that a reliable subroutine for the CG method must include some monitoring of the error generated in intermediate quantities.

The purpose of this report is to collect in one place, and with some consistency of notation, the statements and theoretical justifications of the conjugate gradient algorithm and a number of other algorithms having very

similar characteristics with regard to mathematical theory, operation counts, and storage requirements. We subsequently plan to produce Fortran subroutines for some of these algorithms and study particularly the effectiveness and reliability of various techniques for monitoring accuracy and testing for termination in these subroutines.

The CG algorithm was invented independently and simultaneously (~ 1951) by M. R. Hestenes and E. Stiefel [see Hestenes and Stiefel (1952)]. The paper by Craig (1955) which discusses methods of this type also references related work by Fox, Huskey, Wilkinson, Lanczos, Forsythe, and Rosser in the 1948-1952 period.

After providing some mathematical background in Chapter 2 the CG algorithm is presented in Chapter 3. In Chapters 4 and 5 algorithms are given which result directly from replacing the matrix A in the CG algorithm by $A^T A$ or AA^T respectively. The use of $A^T A$ covers the case of a least squares problem whereas the use of AA^T allows solution of consistent systems, $Ax=b$, in which A is not necessarily symmetric or positive definite (or even square).

This replacement of A by $A^T A$ and by AA^T occurs in Craig (1955) and is also treated in Faddeev and Faddeeva (1963).

More recently Reid (1971) made a strong case for the usefulness of the CG method in large sparse problems. Interest in sparse problems also stimulated Paige (1972) to derive two algorithms of similar character. We will refer to these algorithms as PAIGE-I and PAIGE-II. Paige's derivation of these algorithms is based on a bidiagonalization method given by Golub and Kahan (1965) which has its mathematical roots in a method of Lanczos (1950) for tridiagonalizing a symmetric matrix. We present Paige's least squares algorithm, PAIGE-II, in Chapter 6, where we call it ITLS to distinguish our particular statement of the algorithm. In Chapter 7 we show that ITLS theoretically generates the same

sequence of approximate solution vectors as the algorithm of Chapter 4. The intermediate steps are sufficiently different however to make it of interest to investigate the numerical performance of each of these two algorithms.

In Chapter 8 we present an algorithm ITC which is a subset of Paige's algorithm PAIGE-I. The algorithm PAIGE-I includes provision for handling a least squares problem. It appears to me that for least squares problems this algorithm is not competitive with PAIGE-II or the least squares algorithm of Chapter 4 and thus I have presented only the subset of PAIGE-I which handles a consistent system of linear equations. Paige noted that this subset of PAIGE-I is equivalent to the algorithms given by Faddeev and Faddeeva (1963) and by Craig (1955) which are described in Chapter 5 of this report. This equivalence is verified in Chapter 9.

Finally in Chapter 10 we give an algorithm due to C. C. Paige and M. A. Saunders (personal correspondence, 1972) for the symmetric consistent problem, $Ax = b$. This algorithm is called SYMMLQ by Paige and Saunders. Note that whereas the CG algorithm requires only one matrix-vector multiplication per iteration the other algorithms discussed in Chapters 4-9 each require two matrix-vector multiplications per iteration or else the preliminary computation of $A^T A$ or AA^T . The algorithm, SYMMLQ, requires only one matrix-vector multiplication per iteration and thus is nominally the most economical method described in this report for the indefinite symmetric consistent problem. This algorithm is however notably more complicated than the other algorithms in this report.

Saunders has also developed a modification to Paige's least squares algorithm, PAIGE-II, (our Chapter 6) however we have omitted this as it is more complicated than Paige's algorithm and in a few test cases which we ran its sequence of approximate solution vectors was approximately one step behind the sequence generated by Paige's algorithm.

We wish to thank John Reid and Gene Golub for kindling our awareness of the value of this class of methods for sparse problems. We thank Chris Paige and Michael Saunders for sharing their partially completed current research work with us and David Saunders for supplying Fortran implementations of Michael Saunders' two algorithms.

Chapter 2 Mathematical Background

Two real n -vectors x and y are mutually orthogonal if their inner product, denoted by $x^T y$ or $y^T x$, is zero. They are orthonormal if they are mutually orthogonal and are each of unit euclidean length, i. e. $\|x\| \equiv (x^T x)^{1/2} = 1$ and $\|y\| \equiv (y^T y)^{1/2} = 1$.

A generalization of orthogonality is the notion of A-conjugacy where A is a symmetric matrix. Two vectors x and y are A-conjugate if $x^T A y$ (or equivalently $y^T A x$) is zero. This notion of A-conjugacy is most commonly defined only for a positive definite symmetric matrix A since then one has the convenient property that $x^T A x > 0$ for all $x \neq 0$.

We will use the notion of A-conjugacy under the weaker assumption that A is nonnegative definite symmetric matrix but limit the vectors being considered to those lying in the row space of A . For such vectors the property $x^T A x > 0$ for $x \neq 0$ still holds.

If the set of vectors $\gamma_i = \{v^{(1)}, \dots, v^{(i)}\}$ are mutually orthogonal and a vector $y^{(i+1)}$ is linearly independent of the set γ_i then a vector $v^{(i+1)}$ orthogonal to the set γ_i can be defined by

$$(2.1) \quad v^{(i+1)} = y^{(i+1)} - \frac{v^{(1)T} y^{(i+1)}}{v^{(1)T} v^{(1)}} v^{(1)} - \dots - \frac{v^{(i)T} y^{(i+1)}}{v^{(i)T} v^{(i)}} v^{(i)}$$

This is the formula of Gram-Schmidt orthogonalization.

A similar formula exists for extending a mutually A-conjugate set of vectors. Thus if the set of vectors $\mathcal{U}_i = \{u^{(1)}, \dots, u^{(i)}\}$ are mutually A-conjugate and a vector $z^{(i+1)}$ is linearly independent of the set \mathcal{U}_i then a vector $u^{(i+1)}$, A-conjugate to \mathcal{U}_i , is defined by

$$(2.2) \quad u^{(i+1)} = z^{(i+1)} - \frac{u^{(1)T} A z^{(i+1)}}{u^{(1)T} A u^{(1)}} u^{(1)} - \dots - \frac{u^{(i)T} A z^{(i+1)}}{u^{(i)T} A u^{(i)}} u^{(i)}$$

if all of the denominators are nonzero. Nonzero denominators are assured if A is positive definite symmetric or if A is nonnegative definite symmetric and all of the vectors of the set \mathcal{U}_i lie in the row space of A.

The algorithms to be described in this report have the common feature that in constructing sequences of mutually orthogonal (or mutually A-conjugate) vectors the new linearly independent vector $y^{(i+1)}$ (or $z^{(i+1)}$) will be constructed in such a way that it is already orthogonal (or A-conjugate) to all but one or two of the vectors in the set \mathcal{V}_i (or \mathcal{U}_i). This permits economy in storage and in computation time since only the most recent one or two of the vectors in the set \mathcal{V}_i (or \mathcal{U}_i) need to be retained in storage and only the terms involving these one or two retained vectors need to be computed in Equations (2.1) or (2.2).

PART II

The Conjugate Gradient Algorithm and Variations

Chapter 3 Solving a Consistent System, $Ax=b$, where A is Symmetric and Nonnegative Definite

Let A be an $n \times n$ symmetric nonnegative definite matrix and let b be an n -vector in the column space (range space) of A . We wish to find an n -vector x satisfying

$$(3.1) \quad Ax = b$$

If $\text{Rank}(A) < n$ the solution of Problem (3.1) is nonunique. In this case there is a unique solution vector, \hat{x} , in the row space of A . This vector \hat{x} is the minimal length solution vector for the problem and is the solution vector which the algorithm to be described constructs.

The algorithm to be described is the conjugate gradient method due to Hestenes and Stiefel (1952). Presentations of this method appear in Faddeev and Fadeeva (1963, pp. 392-405), Beckman (1960), Reid (1971), Fox (1965, pp. 208-214), and Householder (1964, pp. 139-141). This method is usually described under the assumption that the matrix A is positive definite although, as will be seen from the discussion to follow, the theory of the method is also valid for a nonnegative definite matrix if it is assumed that b is in the column space of A .

Assume the existence of an integer k ($1 \leq k \leq n$), matrices $V_{n \times k}$ and $D_{k \times k}$ and a k -vector \hat{p} such that

$$(3.2) \quad V = [v^{(1)}, \dots, v^{(k)}]$$

$$(3.3) \quad D = \text{Diag}\{d_1, \dots, d_k\} \quad d_i > 0, \quad i=1, \dots, k$$

$$(3.4) \quad V^T A V = D$$

and

$$(3.5) \quad \hat{x} = V \hat{p}$$

Note that Equation (3.4) implies that the vectors $v^{(i)}$ are mutually A-conjugate.

If such matrices V and D are available Problem (3.1) can be attacked as follows: Left multiply Equation (3.1) by V^T obtaining

$$(3.6) \quad V^T A x = g$$

where g is defined by

$$(3.7) \quad g = V^T b$$

Introduce the change of variables

$$(3.8) \quad x = V p$$

in Equation (3.6) obtaining

$$(3.9) \quad V^T A V p = g$$

which due to Equation (3.4) may also be written as

$$(3.10) \quad Dp = g$$

Thus Problem (3.1) could be solved by computing $g = V^T b$, solving the diagonal system of equations $Dp = g$ for its solution vector \hat{p} , then computing the solution vector \hat{x} for Problem (3.1) as $\hat{x} = V\hat{p}$.

The algorithm to be described constructs the A-conjugate vectors $v^{(j)}$ one at a time and as each such vector is produced its contribution to g , p , x , and the residual vector r is determined. Thus only one of the vectors $v^{(j)}$ needs to be maintained in storage at any one time.

It will be convenient to define auxiliary vectors

$$(3.11) \quad w^{(i)} = Av^{(i)} \quad i=1, \dots, k$$

successive approximations to the solution vector

$$(3.12) \quad x^{(0)} = 0$$

$$(3.13) \quad x^{(i)} = \sum_{j=1}^i v^{(j)} p_j = x^{(i-1)} + v^{(i)} p_i \quad i=1, \dots, k$$

and corresponding residual vectors

$$(3.14) \quad r^{(0)} = b$$

$$(3.15) \quad r^{(i)} = b - Ax^{(i)} = b - \sum_{j=1}^i Av^{(j)} p_j$$

$$= b - \sum_{j=1}^i w^{(j)} p_j = r^{(i-1)} - w^{(i)} p_i \quad i=1, \dots, k$$

It happens that the residual vectors $r^{(0)}, r^{(1)}, \dots$ occurring in this algorithm are mutually orthogonal. The algorithm alternates between producing a vector in the orthogonal sequence $r^{(0)}, r^{(1)}, \dots$ and a vector in the A-conjugate sequence $v^{(1)}, v^{(2)}, \dots$. The various relations which exist between these two sequences permits the algorithm to be remarkably concise.

(3.16) Algorithm CG The Conjugate Gradient Algorithm for Solving a Consistent System $Ax=b$ where A is Symmetric and Nonnegative Definite
 [Due to Hestenes and Stiefel (1952)]

<u>Step</u>	<u>Description</u>
1	$x^{(0)} := 0, r^{(0)} := b, v^{(1)} := b$
2	If $b=0$ set $i:=0$ and go to Step 13
3	$i:=1$
4	$w^{(i)} := Av^{(i)}$
5	$p_i := (r^{(i-1)} T_r^{(i-1)}) / (v^{(i)} T_w^{(i)})$
6	$x^{(i)} := x^{(i-1)} + v^{(i)} p_i$
7	$r^{(i)} := r^{(i-1)} - w^{(i)} p_i$
8	Theoretical termination test: If $r^{(i)}=0$ go to Step 13 Practical termination test: If $\ r^{(i)}\ $ is sufficiently small go to Step 13
9	$\beta_i := (r^{(i)} T_r^{(i)}) / (r^{(i-1)} T_r^{(i-1)})$

<u>Step</u>	<u>Description</u>
10	$v^{(i+1)} := r^{(i)} + v^{(i)} \beta_i$
11	$i := i + 1$
12	Go to Step 4
13	$k := i$
14	Stop

Figure (3.1) is provided as an aid to understanding Algorithm CG. All vectors in the first column of Figure (3.1) lie in the same one-dimensional subspace, \mathcal{S}_1 . For general $l > 1$ if the vectors b, Ab, \dots, Ab^{l-1} are linearly independent let \mathcal{S}_l denote the l -dimensional subspace spanned by these vectors. Then the first k vectors in each row of Figure (3.1) are also linearly independent and span the same subspace \mathcal{S}_l .

To verify that the algorithm CG is mathematically correct we must show that the denominator in Step 5 is positive for $i \leq k$, that Step 5 defines components p_i satisfying Equation (3.10), and that the vectors $v^{(i)}$ produced at Step 10 are mutually A -conjugate. It will also be seen that the residual vectors $r^{(i)}$ are mutually orthogonal.

Assume Algorithm CG has been executed for $i=1, \dots, l-1$, that the set of vectors $\{r^{(0)}, \dots, r^{(l-1)}\}$ are mutually orthogonal and the set of vectors $\{v^{(1)}, \dots, v^{(l)}\}$ are mutually A -conjugate. With this assumption of A -conjugacy we include the assumption that $v^{(i)T} A v^{(i)} > 0, i=1, \dots, l$.

We also assume that

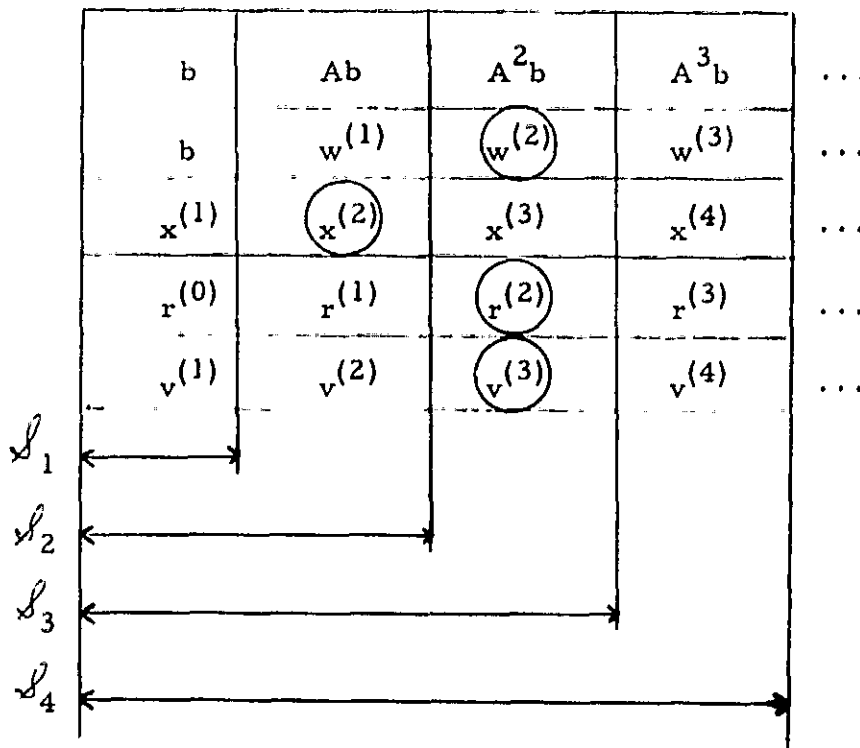


FIGURE (3.1)

Vectors and subspaces related to Algorithm CG. The four circles identify vectors which would be constructed during the second iteration of the algorithm.

$$\begin{aligned}
\mathcal{V}_i &= \text{Span}\{b, w^{(1)}, \dots, w^{(i-1)}\} \\
&= \text{Span}\{r^{(0)}, \dots, r^{(i-1)}\} \\
&= \text{Span}\{v^{(1)}, \dots, v^{(i)}\} \quad i=1, \dots, \ell
\end{aligned}$$

The reader may find it convenient to refer to Figure (3.1) and consider $\ell=2$ for definiteness. The quantities to the left of the circles result from earlier iterations and the circled quantities will be computed during the second iteration.

Consider now the quantities computed during the ℓ^{th} iteration. The denominator in Step 5 is $v^{(\ell)T} A v^{(\ell)}$ which by assumption is positive.

It must be verified that p_ℓ , computed at Step 5 satisfies Equation (3.10), i.e. that

$$(3.17) \quad p_\ell = g_\ell / d_\ell = (v^{(\ell)T} b) / (v^{(\ell)T} A v^{(\ell)})$$

The denominator in Equation (3.17) is clearly identical with the denominator in Step 5. As to the numerators we have

$$\begin{aligned}
v^{(\ell)T} b &= v^{(\ell)T} (r^{(\ell-1)} + \sum_{j=1}^{\ell-1} w^{(j)} p_j) \\
&= v^{(\ell)T} r^{(\ell-1)} \\
&= (r^{(\ell-1)} + v^{(\ell-1)} \beta_{\ell-1})^T r^{(\ell-1)} \\
&= r^{(\ell-1)T} r^{(\ell-1)}
\end{aligned}$$

Steps 6 and 7 are clearly consistent with Equations (3.13) and (3.15). If $r^{(\ell)} = 0$ the algorithm terminates, setting $k=\ell$. Otherwise with $r^{(\ell)} \neq 0$ we proceed to verify that $r^{(\ell)}$ is orthogonal to the subspace \mathcal{S}_ℓ . Using the basis $\{r^{(0)}, \dots, r^{(\ell-1)}\}$ of \mathcal{S}_ℓ it suffices to verify that $r^{(i)T} r^{(\ell)} = 0$ for $i=0, \dots, \ell-1$.

$$(3.18) \quad r^{(i)T} r^{(\ell)} = r^{(i)T} r^{(\ell-1)} - r^{(i)T} w^{(\ell)} p_\ell$$

Both right side terms are zero for $i=0, \dots, \ell-2$. For $i=\ell-1$ substitute the definition of p_ℓ from Step 5 into Equation (3.18) obtaining

$$\begin{aligned} r^{(\ell-1)T} r^{(\ell)} &= r^{(\ell-1)T} r^{(\ell-1)} - \frac{[r^{(\ell-1)T} w^{(\ell)}] \cdot [r^{(\ell-1)T} r^{(\ell-1)}]}{v^{(\ell)T} w^{(\ell)}} \\ &= 0 \end{aligned}$$

since

$$v^{(\ell)T} w^{(\ell)} = [r^{(\ell-1)} + v^{(\ell-1)} \beta_{\ell-1}]^T w^{(\ell)} = r^{(\ell-1)T} w^{(\ell)}$$

Next it must be shown that $v^{(\ell+1)}$ defined at Step 10 is A-conjugate to the subspace \mathcal{S}_ℓ . We use the basis $\{v^{(1)}, \dots, v^{(\ell)}\}$ for \mathcal{S}_ℓ and verify that $v^{(i)T} A v^{(\ell+1)} = 0$ for $i=1, \dots, \ell$.

$$(3.19) \quad v^{(i)T} A v^{(\ell+1)} = v^{(i)T} A r^{(\ell)} + v^{(i)T} A v^{(\ell)} \beta_\ell$$

For $i=1, \dots, \ell-1$ both right side terms are zero; the first because $v^{(i)T}A$ is in \mathcal{A}_ℓ and thus is orthogonal to $r^{(\ell)}$ and the second because $v^{(i)}$ is in $\mathcal{S}_{\ell-1}$ and thus is A-conjugate to $v^{(\ell)}$. For $i=\ell$ substitute the definition of β_ℓ from Step 9 into Equation (3.19) obtaining

$$\begin{aligned}
 (3.20) \quad v^{(\ell)T}A v^{(\ell+1)} &= v^{(\ell)T}A r^{(\ell)} + \frac{(v^{(\ell)T}A v^{(\ell)})(r^{(\ell)T}r^{(\ell)})}{r^{(\ell-1)T}r^{(\ell-1)}} \\
 &= w^{(\ell)T}r^{(\ell)} + p_\ell^{-1}r^{(\ell)T}r^{(\ell)} \\
 &= [w^{(\ell)} + p_\ell^{-1}r^{(\ell)}]^T r^{(\ell)} \\
 &= p_\ell^{-1}r^{(\ell-1)T}r^{(\ell)} = 0
 \end{aligned}$$

Finally we verify that $v^{(\ell+1)T}A v^{(\ell+1)} > 0$. Let h denote the rank of A .

Since A is nonnegative definite there exists an $h \times n$ matrix R of rank h such that

$$A = R^T R$$

and clearly the row space of R is identical with the row (and column) space of A .

Let \mathcal{A} denote the row (and column) space of A . From Steps 1, 7, and 10 all of the vectors $r^{(i)}$ and $v^{(i)}$ produced by the algorithm lie in the subspace \mathcal{A} . Since $v^{(\ell+1)}$ is constructed at Step 10 as the sum of the nonzero vector $r^{(\ell)}$ and a vector $v^{(\ell)}\beta_i$ orthogonal to $r^{(\ell)}$ it follows that $v^{(\ell+1)} \neq 0$. Since $v^{(\ell+1)}$ is a nonzero vector

in the subspace \mathcal{L} the vector $Rv^{(\ell+1)}$ must also be nonzero. It follows that $[Rv^{(\ell+1)}]^T Rv^{(\ell+1)} > 0$ or equivalently $v^{(\ell+1)T} A v^{(\ell+1)} > 0$, as was to be shown.

As is apparent from the different expressions derived for p_i and for β_i there are a number of different ways in which the CG algorithm could be implemented. There are also different trade-offs possible between storage used and counts of arithmetic operations. Reid (1971) discusses over a dozen such variations. The form in which we have stated Algorithm CG is the one preferred by Reid.

Theoretical termination of the CG algorithm occurs when $r^{(k)}=0$ with $r^{(i)} \neq 0$ for $i=0, \dots, k-1$. Referring to Figure (3.1) we see that this means that the subspaces $\mathcal{L}_1, \dots, \mathcal{L}_k$ are all different but that $\mathcal{L}_{k+1} = \mathcal{L}_k$. Equivalently this means that the vector $A^k b$ lies in the subspace \mathcal{L}_k spanned by $\{b, Ab, \dots, A^{k-1}b\}$.

This implies that b is representable as a linear combination of some set of k eigenvectors, say $\{f^{(1)}, \dots, f^{(k)}\}$, of the nonnegative definite symmetric matrix A and is not representable as a linear combination of any smaller set of eigenvectors of A . Furthermore the eigenvalues $\{\lambda_1, \dots, \lambda_k\}$ associated with these eigenvectors are all positive. Thus there exist nonzero numbers c_i such that

$$b = \sum_{i=1}^k f^{(i)} c_i$$

and the minimal length solution vector \hat{x} is representable as

$$\hat{x} = \sum_{i=1}^k f^{(i)} (c_i / \lambda_i)$$

Note also that $\text{Span}\{f^{(1)}, \dots, f^{(k)}\} = \mathcal{V}'_k$ and that \hat{x} lies in \mathcal{V}'_k but not in \mathcal{V}'_{k-1} .

Most commonly the value of k will be n . However k will be less than n if and only if b is orthogonal to some eigenvectors of A . In particular b will necessarily be orthogonal to some eigenvectors of A if A has any multiple eigenvalues.

Chapter 4 Solving a Least Squares Problem $Ax \cong b$

Let A be an $m \times n$ matrix and let b be an m -vector. We wish to find an n -vector x which minimizes $\|b - Ax\|$. We denote this least squares problem by the notation

$$(4.1) \quad Ax \cong b$$

We permit either $m \geq n$ or $m < n$. If $\text{Rank}(A) < n$ the algorithm to be described constructs the unique minimal length solution vector, \hat{x} . This solution vector is characterized as being the only solution vector lying in the row space of A .

A vector x minimizes $\|b - Ax\|$ if and only if it satisfies the "normal equations"

$$(4.2) \quad A^T Ax = A^T b$$

Normal equations are always consistent since the right side vector, $A^T b$, lies in the row space of A which is also the column space of the matrix $A^T A$. The ranks of $A^T A$ and A are equal and their row spaces are the same. Thus if $A^T A$ is singular the unique solution of Problem (4.2) lying in the row space of $A^T A$ is also the unique minimal length solution of Problem (4.1).

Since $A^T A$ is nonnegative definite and Problem (4.2) is consistent the conjugate gradient algorithm CG (3.16) is directly applicable to Problem (4.2). Denote the residual vector for Problem (4.2) by

$$\bar{h} = A^T b - A^T Ax = A^T (b - Ax)$$

and introduce bars on various other symbols in Algorithm CG to distinguish the application of the algorithm to Problem (4.2). Then Algorithm CG adapted to Problem (4.1) can be written briefly as

(4.3) Algorithm CGLS Conjugate Gradient Algorithm for the Least Squares Problem $Ax \cong b$

$$(4.4) \quad \bar{x}^{(0)} = 0, \quad \bar{h}^{(0)} = A^T b, \quad \bar{v}^{(1)} = A^T b$$

If $\bar{h}^{(0)} = 0$ terminate.

Do for $i=1, 2, \dots$, until $\bar{h}^{(i)} = 0$

$$(4.5) \quad \bar{w}^{(i)} = A^T A \bar{v}^{(i)}$$

$$(4.6) \quad \bar{p}_i = \|\bar{h}^{(i-1)}\|^2 / (\bar{v}^{(i)T} \bar{w}^{(i)})$$

$$(4.7) \quad \bar{x}^{(i)} = \bar{x}^{(i-1)} + \bar{v}^{(i)} \bar{p}_i$$

$$(4.8) \quad \bar{h}^{(i)} = \bar{h}^{(i-1)} - \bar{w}^{(i)} \bar{p}_i$$

$$(4.9) \quad \bar{\theta}_i = \|\bar{h}^{(i)}\|^2 / \|\bar{h}^{(i-1)}\|^2$$

$$(4.10) \quad \bar{v}^{(i+1)} = \bar{h}^{(i)} + \bar{v}^{(i)} \bar{\theta}_i$$

One may wish to have the residual vector of the least squares problem (4.1)

$$(4.11) \quad \bar{r}^{(i)} = b - A\bar{x}^{(i)}$$

available at each iteration, possibly for use in a supplementary termination test.

This may be accomplished by the following revision of the algorithm [Faddeev and Faddeeva (1963, p. 403)].

(4.12) Algorithm [Alternate Form of CGLS]

$$\bar{x}^{(0)} := 0, \quad \bar{r}^{(0)} := b, \quad \bar{h}^{(0)} := A^T b, \quad \bar{v}^{(1)} := A^T b$$

If $\bar{h}^{(0)} = 0$ terminate.

Do for $i := 1, 2, \dots$, until $\bar{h}^{(i)} = 0$

$$\bar{u}^{(i)} := A \bar{v}^{(i)}$$

$$\bar{p}_i := \|\bar{h}^{(i-1)}\|^2 / \|\bar{u}^{(i)}\|^2$$

$$\bar{x}^{(i)} := \bar{x}^{(i-1)} + \bar{v}^{(i)} \bar{p}_i$$

$$\bar{r}^{(i)} := \bar{r}^{(i-1)} - \bar{u}^{(i)} \bar{p}_i$$

$$\bar{h}^{(i)} := A^T \bar{r}^{(i)}$$

$$\bar{B}_i := \|\bar{h}^{(i)}\|^2 / \|\bar{h}^{(i-1)}\|^2$$

$$\bar{v}^{(i+1)} := \bar{h}^{(i)} + \bar{v}^{(i)} \bar{B}_i$$

This latter form of the algorithm requires more storage since one must maintain the current values of the two m -vectors $\bar{u}^{(i)}$ and $\bar{r}^{(i)}$.

Chapter 5 Solving a Consistent System $Ax=b$

Let A be an $m \times n$ matrix and let b be an m -vector contained in the column space (range space) of A . Consider the problem of finding an n -vector x satisfying

$$(5.1) \quad Ax = b$$

We will permit either $m \leq n$ or $m > n$. If $\text{Rank}(A) < n$ the solution to this problem is nonunique. In this case there is a unique solution vector \hat{x} lying in the row space of A . This vector \hat{x} is the unique minimal length solution vector for Problem (5.1). The algorithm to be described constructs this solution vector \hat{x} .

Since we seek a solution vector \hat{x} in the row space of A the solution vector \hat{x} will be representable in the form

$$(5.2) \quad \hat{x} = A^T \hat{y}$$

for some (not necessarily unique) m -vector \hat{y} . Making the change of variables $x = A^T y$ in Problem (5.1) we obtain the problem

$$(5.3) \quad AA^T y = b$$

This is a consistent problem with a symmetric nonnegative definite matrix AA^T . Thus the conjugate gradient algorithm CG (3.16) can be applied to solve Problem (5.3).

The resulting algorithm for solving a consistent system $Ax=b$ may be written as follows where the notation of Algorithm CG (3.16) is changed by writing \bar{p} , $\bar{\beta}$, \bar{w} , \bar{r} , \bar{q} , \bar{y} , and AA^T in place of p , β , w , r , v , x , and A respectively. Furthermore,

motivated by Equation (5.2) we introduce the sequence of approximate solution vectors

$$(5.4) \quad \bar{x}^{(i)} = A^T \bar{y}^{(i)}$$

(5.5) Algorithm

$$\bar{y}^{(0)} = 0, \bar{x}^{(0)} = 0, \bar{r}^{(0)} = b, \bar{q}^{(1)} = b$$

If $\bar{r}^{(0)} = 0$ terminate.

Do for $i=1, 2, \dots$, until $\bar{r}^{(i)} = 0$

$$\bar{w}^{(i)} = AA^T \bar{q}^{(i)}$$

$$\bar{p}_i = \|\bar{r}^{(i-1)}\|^2 / (\bar{q}^{(i)T} \bar{w}^{(i)})$$

$$\bar{y}^{(i)} = \bar{y}^{(i-1)} + \bar{q}^{(i)} \bar{p}_i$$

$$\bar{x}^{(i)} = \bar{x}^{(i-1)} + A^T \bar{q}^{(i)} \bar{p}_i$$

$$\bar{r}^{(i)} = \bar{r}^{(i-1)} - \bar{w}^{(i)} \bar{p}_i$$

$$\bar{\beta}_i = \|\bar{r}^{(i)}\|^2 / \|\bar{r}^{(i-1)}\|^2$$

$$\bar{q}^{(i+1)} = \bar{r}^{(i)} + \bar{q}^{(i)} \bar{\beta}_i$$

Eliminating the intermediate vectors $\bar{w}^{(i)}$ and $\bar{y}^{(i)}$ we obtain the algorithm in the form given by Craig (1955), p. 72, except that Craig used the opposite choice of the sign of the residual vector.

(5.6) Algorithm [Craig (1955)]

$$\bar{x}^{(0)} := 0, \bar{r}^{(0)} := b, \bar{q}^{(1)} := b$$

If $\bar{r}^{(0)} = 0$ terminate.

Do for $i:=1, 2, \dots$, until $\bar{r}^{(i)} = 0$

$$\bar{p}_i := \|\bar{r}^{(i-1)}\|^2 / \|A^T \bar{q}^{(i)}\|^2$$

$$\bar{x}^{(i)} := \bar{x}^{(i-1)} + A^T \bar{q}^{(i)} \bar{p}_i$$

$$\bar{r}^{(i)} := b - A \bar{x}^{(i)} [\equiv \bar{r}^{(i-1)} - AA^T \bar{q}^{(i)} \bar{p}_i]$$

$$\bar{\beta}_i := \|\bar{r}^{(i)}\|^2 / \|\bar{r}^{(i-1)}\|^2$$

$$\bar{q}^{(i+1)} := \bar{r}^{(i)} + \bar{q}^{(i)} \bar{\beta}_i$$

As is noted in Faddeev and Faddeeva (1963, pp. 403-405) this algorithm can be further simplified by introducing the substitution

$$\bar{v}^{(i)} = A^T \bar{q}^{(i)}$$

Note that the vectors $\{\bar{v}^{(1)}, \dots, \bar{v}^{(k)}\}$ are mutually orthogonal since the vectors $\{\bar{q}^{(1)}, \dots, \bar{q}^{(k)}\}$ are mutually (AA^T) -conjugate.

We call the resulting algorithm CGC.

(5.7) Algorithm CGC Conjugate Gradient Algorithm for the Consistent System $Ax=b$

$$\bar{x}^{(0)} := 0, \bar{r}^{(0)} := b, \bar{v}^{(1)} := A^T b$$

If $\bar{r}^{(0)} = 0$ terminate

Do for $i:=1, 2, \dots$, until $\bar{r}^{(i)} = 0$

$$\bar{p}_i = \|\bar{r}^{(i-1)}\|^2 / \|\bar{v}^{(i)}\|^2$$

$$\bar{x}^{(i)} = \bar{x}^{(i-1)} + \bar{v}^{(i)} \bar{p}_i$$

$$\bar{r}^{(i)} = \bar{r}^{(i-1)} - A \bar{v}^{(i)} \bar{p}_i$$

$$\bar{b}_i = \|\bar{r}^{(i)}\|^2 / \|\bar{r}^{(i-1)}\|^2$$

$$\bar{v}^{(i-1)} = A^T \bar{r}^{(i)} + \bar{v}^{(i)} \bar{b}_i$$

PART III
OTHER ALGORITHMS

Chapter 6 Solving a Least Squares Problem, $Ax \cong b$

Let A be an $m \times n$ matrix and let b be an m -vector. We wish to find an n -vector x which minimizes $\|b-Ax\|$. We denote this least squares problem by the notation

$$(6.1) \quad Ax \cong b$$

Generally one would have $m \geq n$ and $\text{Rank}(A) = n$ in such a problem. We will not make these assumptions however as they are not necessary to the mathematical development of the algorithm to be described.

If $\text{Rank}(A) < n$ the solution to Problem (6.1) is not unique. In this case however the problem has a unique solution vector of least euclidean length. It can easily be verified that this unique minimal length solution vector lies in the row space of A and in fact is the only solution vector for Problem (6.1) lying in the row space of A . The algorithm to be described constructs the solution vector in the row space of A and thus finds the minimal length solution vector if $\text{Rank}(A) < n$.

Let \hat{x} be the unique minimal length solution vector for Problem (6.1). Define the residual vector

$$(6.2) \quad \hat{r} = b - A\hat{x}$$

and

$$(6.3) \quad \hat{b} = A\hat{x} = b - \hat{r}$$

Thus b can be written as

$$(6.4) \quad b = \hat{b} + \hat{r}$$

where \hat{b} lies in the column space of A and \hat{r} is orthogonal to the column space of A. The vector \hat{b} is the orthogonal projection of b into the column space of A and will be referred to as the projected right-side vector.

Suppose there exists an integer k ($1 \leq k \leq \min\{m, n\}$) and matrices

$$(6.5) \quad U_{m \times k} = [u^{(1)}, \dots, u^{(k)}]$$

$$(6.6) \quad V_{n \times k} = [v^{(1)}, \dots, v^{(k)}]$$

$$(6.7) \quad R_{k \times k} = \begin{bmatrix} \alpha_1 & \beta_2 & & & 0 \\ & \cdot & \cdot & \cdot & \\ & \alpha_2 & & \cdot & \beta_k \\ & & \cdot & \cdot & \\ 0 & & & & \alpha_k \end{bmatrix} \quad (\text{all } \alpha_i > 0 \text{ and } \beta_i > 0)$$

and a k-vector, \hat{p} , such that

$$(6.8) \quad U^T U = I_k$$

$$(6.9) \quad V^T V = I_k$$

$$(6.10) \quad AV = UR$$

$$(6.11) \quad A^T U = VR^T$$

and

$$(6.12) \quad \hat{x} = V\hat{p}$$

Since R is nonsingular Equations (6.9) and (6.11) imply that the column vectors of V form an orthonormal basis for a subspace, \mathcal{V} , of the row space of A. Similarly Equations (6.8) and (6.10) imply that the column vectors of U form an orthonormal basis for a subspace \mathcal{U} of the column space (range space) of A. Equation (6.12) shows that the subspace \mathcal{V} contains the solution vector \hat{x} . From Equations (6.3), (6.10), and (6.12) it follows that the subspace \mathcal{U} contains the projected right-side vector \hat{b} .

Assuming the availability of the matrices U, V, and R, Problem (6.1) can be approached as follows: Introduce the change of variables

$$(6.13) \quad x = Vp$$

in Problem (6.1) and use Equation (6.10) obtaining the equivalent least squares problem

$$(6.14) \quad URp \cong b$$

Let $\bar{U}_{m \times (m-k)}$ be a matrix which when adjoined to U forms an $m \times m$ orthogonal matrix:

$$(6.15) \quad [U:\bar{U}]^T [U:\bar{U}] = I_m$$

Left multiply Equation (6.14) by $[U:\bar{U}]^T$ obtaining the equivalent least squares problem.

$$(6.16) \quad \begin{bmatrix} R \\ 0 \end{bmatrix} p \cong g \equiv \begin{bmatrix} \tilde{g} \\ \bar{g} \end{bmatrix} \left\{ \begin{array}{l} k \\ m-k \end{array} \right.$$

where

$$(6.17) \quad \mathbf{g} \equiv \begin{bmatrix} \tilde{\mathbf{g}} \\ \mathbf{g} \end{bmatrix} = [\mathbf{U}:\bar{\mathbf{U}}]^T \mathbf{b} \equiv \begin{bmatrix} \mathbf{U}^T \mathbf{b} \\ \bar{\mathbf{U}}^T \mathbf{b} \end{bmatrix}$$

The least squares solution vector $\hat{\mathbf{p}}$ for Problem (6.16) may be obtained as the solution of the upper bidiagonal nonsingular system of equations

$$(6.18) \quad \mathbf{R}\mathbf{p} = \tilde{\mathbf{g}}$$

Obtaining $\hat{\mathbf{p}}$ from Equation (6.18) the solution vector $\hat{\mathbf{x}}$ can then be computed from Equation (6.13).

In analogy with the conjugate gradient algorithm we wish to cast this procedure into a sequential form so that a succession of approximate solution vectors $\mathbf{x}^{(i)}$ and associated residual vectors

$$(6.19) \quad \mathbf{r}^{(i)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(i)}$$

can be computed as the successive vectors $\mathbf{u}^{(i)}$ and $\mathbf{v}^{(i)}$ are computed. This obviates the need to store old $\mathbf{u}^{(i)}$ and $\mathbf{v}^{(i)}$ vectors.

Equation (6.18) is not directly suitable for such a (forward) sequential procedure since the last (lower right) element of \mathbf{R} must be determined before any components of $\hat{\mathbf{p}}$ can be computed. However using Equations (6.13) and (6.18) we can write

$$(6.20) \quad \mathbf{x} = \mathbf{V}\mathbf{p} = \mathbf{V}\mathbf{R}^{-1}\mathbf{R}\mathbf{p} = (\mathbf{V}\mathbf{R}^{-1})\tilde{\mathbf{g}} \equiv \mathbf{W}\tilde{\mathbf{g}}$$

where the $n \times k$ matrix

$$(6.21) \quad W = [w^{(1)}, \dots, w^{(k)}]$$

satisfies the linear matrix equation

$$(6.22) \quad R^T W^T = V^T$$

or

$$(6.23) \quad \begin{bmatrix} \alpha_1 & & & & 0 \\ & \beta_2 & & & \\ & & \alpha_2 & & \\ & & & \ddots & \\ & & & & \beta_k & \alpha_k \\ 0 & & & & & \end{bmatrix} \cdot \begin{bmatrix} w^{(1)T} \\ \vdots \\ w^{(k)T} \end{bmatrix} = \begin{bmatrix} v^{(1)T} \\ \vdots \\ v^{(k)T} \end{bmatrix}$$

From this matrix equation one obtains the following expressions for sequential computation of the vectors $w^{(i)}$.

$$(6.24) \quad w^{(1)} = v^{(1)} / \alpha_1$$

$$(6.25) \quad w^{(i)} = (v^{(i)} - \beta_i w^{(i-1)}) / \alpha_i \quad i=2, \dots, k$$

Let g_i denote the i^{th} component of the m -vector g . Then g_i for $i \leq k$, is also the i^{th} component of the k -vector \tilde{g} of Equation (6.18). Define the m -vectors

$$g^{(i)} = [g_1, \dots, g_i, \underbrace{0, \dots, 0}_{m-i}]^T \quad i=0, \dots, m$$

and the k-vectors

$$\tilde{g}^{(i)} = [g_1, \dots, g_i, \underbrace{0, \dots, 0}_{k-i}]^T \quad i=0, \dots, k$$

Motivated by Equation (6.20) define a sequence of approximate solution vectors $x^{(i)}$ by

$$(6.26) \quad x^{(0)} = 0$$

$$(6.27) \quad \begin{aligned} x^{(i)} &= W \tilde{g}^{(i)} = \sum_{j=1}^i w^{(j)} g_j \\ &= x^{(i-1)} + w^{(i)} g_i \quad i=1, \dots, k \end{aligned}$$

The associated residual vectors, $r^{(i)}$, are defined as

$$(6.28a) \quad r^{(0)} = b$$

and

$$(6.28b) \quad \begin{aligned} r^{(i)} &= b - Ax^{(i)} \\ &= b - AVR^{-1} \tilde{g}^{(i)} \\ &= b - URR^{-1} \tilde{g}^{(i)} \quad [\text{Using Equation (6.10)}] \\ &= b - U \tilde{g}^{(i)} \end{aligned}$$

$$\begin{aligned}
&= b - \sum_{j=1}^i u^{(j)} g_i \\
&= r^{(i-1)} - u^{(i)} g_i \quad i=1, \dots, k
\end{aligned}$$

We next consider formulas for computing \tilde{g} and $\|A^T r^{(i)}\|$, $i=1, \dots, k$, which depend upon a particular choice of $v^{(1)}$. The choice of $v^{(1)}$ is somewhat arbitrary as long as it is chosen to lie in the row space of A . We follow Paige (1972) in defining β_1 and $v^{(1)}$ by

$$(6.29) \quad v^{(1)} \beta_1 = A^T b$$

where $\beta_1 = \|A^T b\|$ so that $\|v^{(1)}\| = 1$.

From Equation (6.17) we have

$$(6.30) \quad \tilde{g} = U^T b$$

Left multiplying this equation by R^T and using Equation (6.10) and (6.29) gives

$$\begin{aligned}
(6.31) \quad R^T \tilde{g} &= R^T U^T b = V^T A^T b \\
&= \beta_1 V^T v^{(1)} = \beta_1 e^{(1)}
\end{aligned}$$

where $e^{(1)}$ denotes the first column of the $k \times k$ identity matrix.

Writing the equation $R^T \tilde{g} = \beta_1 e^{(1)}$ in terms of its components gives

$$(6.32) \quad \begin{bmatrix} \alpha_1 & & & & 0 \\ & \beta_2 & & & \\ & & \alpha_2 & & \\ & & & \ddots & \\ 0 & & & & \beta_k & \\ & & & & & \alpha_k \end{bmatrix} \cdot \begin{bmatrix} g_1 \\ \vdots \\ g_k \end{bmatrix} = \begin{bmatrix} \beta_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

from which the components g_i can be computed as

$$(6.33) \quad g_1 = \beta_1 / \alpha_1$$

$$(6.34) \quad g_i = -(\beta_i / \alpha_i) g_{i-1} \quad i=2, \dots, k$$

Note that it would also be possible to compute the quantities g_i sequentially as $g_i = u^{(i)T} b$ (see Equation (6.30)) but Paige (1972) reports that Equations (6.33)-(6.34) were found to give better numerical accuracy in test cases.

In a least squares problem one does not generally expect the final residual vector, $\hat{r} = b - A\hat{x}$, to be zero. The residual vector at the solution is characterized however by the property that it is orthogonal to all of the column vectors of A . Thus the vector $h = A^T(b - Ax)$ is zero if and only if x is a solution vector for the least squares problem $Ax \approx b$. It is also true that h is the negative gradient vector with respect to x of the function $\frac{1}{2} \|b - Ax\|^2$. Define

$$(6.35) \quad \begin{aligned} h^{(i)} &= A^T r^{(i)} = A^T (b - Ax^{(i)}) \\ &= A^T (b - U \tilde{g}^{(i)}) \\ &= v^{(1)T} \beta_1 - v R^T g^{(i)} \end{aligned}$$

$$\begin{aligned}
&= v(e^{(1)}\beta_1 - R^T g^{(i)}) \\
&= -v e^{(i+1)}\beta_{i+1}g_i \\
&= -v^{(i+1)}\beta_{i+1}g_i \quad i=1, \dots, k-1
\end{aligned}$$

while for $i=0$ and $i=k$ we have $h^{(0)} = A^T b = v^{(1)}\beta_1$ and $h^{(k)} = 0$. It is of interest to note that the vectors $h^{(0)}, \dots, h^{(k-1)}$ are mutually orthogonal.

The quantities $\gamma_i = \|h^{(i)}\|$ which may be useful in monitoring the progress of the algorithm are thus expressible as

$$(6.36) \quad \gamma_0 = \beta_1$$

$$(6.37) \quad \gamma_i = |\beta_{i+1}g_i| \quad i=1, \dots, k-1$$

$$(6.38) \quad \gamma_k = 0$$

In practice k is generally not known in advance and will in fact be defined as the first value of i for which $\beta_{i+1} = 0$.

We turn now to the determination of the quantities $u^{(i)}$ and α_i for $i=1, \dots, k$ and $v^{(i)}$ and β_i for $i=2, \dots, k$. From Equation (6.10) we obtain the equations

$$(6.39) \quad u^{(1)}\alpha_1 = Av^{(1)}$$

$$(6.40) \quad u^{(i)}\alpha_i = Av^{(i)} - u^{(i-1)}\beta_i \quad i=2, \dots, k$$

and from Equation (6.11) the equations

$$(6.41) \quad v^{(i)}\beta_i = A^T u^{(i-1)} - v^{(i-1)}\alpha_{i-1} \quad i=2, \dots, k$$

$$(6.42) \quad 0 = A^T u^{(k)} - v^{(k)}\alpha_k$$

If β_1 and $v^{(1)}$ are defined by Equation (6.29) and the scalar quantities β_i and α_i are determined so that the vectors $v^{(i)}$ and $u^{(i)}$ respectively have unit euclidean length (as is required by Equations (6.8) and (6.9)) then Equations (6.39)-(6.41)

determine all of the remaining vectors $v^{(i)}$, $i=2, \dots, k$, and $u^{(i)}$, $i=1, \dots, k$.

Collecting these various formulas together leads to the following algorithm.

(6.43) Algorithm ITLS for iterative solution of the least squares problem

$Ax \cong b$. [Due to C. C. Paige (1972) pp. 21-22.]

<u>Step No.</u>	<u>Description</u>
1	$x^{(0)} := 0$
2	$g_0 := -1$
3	$i := 1$
4	$\tilde{v}^{(i)} := \begin{cases} A^T b & \text{if } i=1 \\ A^T u^{(i-1)} - v^{(i-1)} \alpha_{i-1} & \text{if } i > 1 \end{cases}$
5	$\beta_i := \ \tilde{v}^{(i)}\ $
6	$\gamma_{i-1} := \beta_i g_{i-1} $
7	Theoretical termination test: If $\beta_i = 0$ go to Step 17. Practical termination test: If either β_i or γ_{i-1} is sufficiently small go to Step 17.
8	$v^{(i)} := \tilde{v}^{(i)} / \beta_i$
9	$\tilde{u}^{(i)} := \begin{cases} Av^{(1)} & \text{if } i=1 \\ Av^{(i)} - u^{(i-1)} \beta_i & \text{if } i > 1 \end{cases}$
10	$\alpha_i := \ \tilde{u}^{(i)}\ $
11	$u^{(i)} := \tilde{u}^{(i)} / \alpha_i$
12	$w^{(i)} := \begin{cases} v^{(1)} / \alpha_1 & \text{if } i=1 \\ (v^{(i)} - w^{(i-1)} \beta_i) / \alpha_i & \text{if } i > 1 \end{cases}$
13	$g_i := -(\beta_i / \alpha_i) g_{i-1}$
14	$x^{(i)} := x^{(i-1)} + w^{(i)} g_i$
15	$i := i+1$
16	Go to Step 4
17	$k := i-1$
18	Stop

It must be verified that the vectors $v^{(i)}$ and $u^{(i)}$ produced by this algorithm have the orthogonality properties specified by Equations (6.8) and (6.9) and that all of the numbers α_i defined by this algorithm are positive.

Assume that $\ell-1$ iterations of Algorithm (6.43) have been executed producing positive numbers $\beta^{(i)}$ and $\alpha^{(i)}$, $i=1, \dots, \ell-1$, a set of mutually orthogonal unit n -vectors $\{v^{(1)}, \dots, v^{(\ell-1)}\}$, and a set of mutually orthogonal unit m -vectors $\{u^{(1)}, \dots, u^{(\ell-1)}\}$. It is obvious from Steps 4 and 8 that all of the vectors $v^{(i)}$ lie in the row space of A and from Steps 9 and 11 that all of the vectors $u^{(i)}$ lie in the column space (range space) of A .

Consider the quantities computed during the ℓ^{th} iteration.

If β_ℓ , computed at Step 5 is zero then the (theoretical) algorithm terminates, setting $k=\ell-1$. In this case we have $\beta_{k+1}=0$ and $\gamma_k=0$ which means (see Equation (6.35)) that the most recently computed approximate solution $x^{(k)}$ was in fact the unique minimum length solution \hat{x} for the least squares problem, $Ax \cong b$.

If β_ℓ , computed at Step 5 is not zero, i. e., $\beta_\ell > 0$, then we must verify that the vector $\tilde{v}^{(\ell)}$ previously computed at Step 4 is orthogonal to $v^{(i)}$, $i=1, \dots, \ell-1$. Using the formula of Step 4 with $i > 1$ the inner products to be investigated are

$$\begin{aligned}
 (6.44) \quad v^{(i)T} \tilde{v}^{(\ell)} &= v^{(i)T} A T_u^{(\ell-1)} - v^{(i)T} v^{(\ell-1)} \alpha_{\ell-1} \\
 &= [A v^{(i)}]^T T_u^{(\ell-1)} - v^{(i)T} v^{(\ell-1)} \alpha_{\ell-1} \\
 &= [u^{(i)} \alpha_i + u^{(i-1)} \beta_i]^T T_u^{(\ell-1)} - v^{(i)T} v^{(\ell-1)} \alpha_{\ell-1} \\
 &= u^{(i)T} T_u^{(\ell-1)} \alpha_i + u^{(i-1)T} T_u^{(\ell-1)} \beta_i - v^{(i)T} v^{(\ell-1)} \alpha_{\ell-1}
 \end{aligned}$$

For $i < \ell-1$ each of the three right side terms in Equation (6.44) is zero because of the assumed mutual orthogonality of $\{u^{(1)}, \dots, u^{(\ell-1)}\}$ and the assumed mutual orthogonality of $\{v^{(1)}, \dots, v^{(\ell-1)}\}$.

For $i=l-1$ Equation (6.44) becomes

$$v^{(l-1)T} \tilde{v}^{(l)} = \alpha_{l-1} + 0 - \alpha_{l-1} = 0$$

which completes the verification that $\tilde{v}^{(l)}$ is orthogonal to $v^{(i)}$, $i=2, \dots, l-1$. The verification that $v^{(1)T} \tilde{v}^{(l)} = 0$ is equally straightforward.

Similarly it can be shown that $\tilde{u}^{(l)}$ computed at Step 9 satisfies $u^{(i)T} \tilde{u}^{(l)} = 0$ for $i=1, \dots, l-1$. We further assert that $\tilde{u}^{(l)}$ is not zero and thus that $\alpha_l > 0$.

Assume the contrary. Then

$$\begin{aligned} (6.45) \quad 0 &= \tilde{u}^{(l)} = A v^{(l)} - u^{(l-1)} \beta_l \\ &= A v^{(l)} - [A v^{(l-1)} - u^{(l-2)} \beta_{l-1}] \beta_l / \alpha_{l-1} \\ &= \dots = \sum_{i=1}^l c_i A v^{(i)} \\ &= A \sum_{i=1}^l c_i v^{(i)} \end{aligned}$$

where the coefficients c_i , $i=1, \dots, l$ are nonzero. Since the vectors $v^{(i)}$, $i=1, \dots, l$ constitute an orthonormal basis for a subspace of the row space of A the vector $z = \sum_{i=1}^l c_i v^{(i)}$ must be a nonzero vector lying in the row space of A . Such a vector must satisfy $Az \neq 0$ contradicting Equation (6.45). We conclude that $\tilde{u}^{(l)} \neq 0$ and $\alpha_l > 0$.

This completes the verification of the theoretical algorithm (6.43). In practice there remains the problem of fully specifying a satisfactory termination test at Step 7.

Some estimate, say ϵ_i , of the norm of the round-off error vector associated with the computed vector $\tilde{v}^{(i)}$ could be computed. Then the number β_i could be regarded as being sufficiently small for termination if $\beta_i \leq \epsilon_i$.

Since γ_{i-1} represents an evaluation of $A^T(b - Ax^{(i-1)})$ one might define

$$(6.46) \quad \omega_i = \|A\|(\|b\| + \|A\| \cdot \|x^{(i)}\|)$$

and terminate the algorithm when $\gamma_{i-1} \leq \eta \omega_{i-1}$ where η denotes the relative machine precision. (Define η to be the smallest number such that the computed value of $1 + \eta$ is distinguishable from 1.)

More complex algorithmic logic might be needed. Thus if $\beta_i \leq \epsilon_i$ but $\gamma_{i-1} \gg \eta \omega_{i-1}$ then rather than accepting $x^{(i-1)}$ as the solution it might be useful to restart the algorithm, attacking the modified problem

$$(6.47) \quad A dx \cong b - Ax^{(i-1)}$$

to obtain a correction vector dx to be added to $x^{(i-1)}$.

We intend to study the problem of termination tests for this algorithm and treat the subject in a subsequent report.

Chapter 7 The Theoretical Equivalence of Algorithms CGLS and ITLS

We will show that the sequence of approximate solution vectors $\bar{x}^{(i)}$ generated by Algorithm CGLS (4.3) is identical (theoretically) with the sequence of approximate solution vectors $x^{(i)}$ generated by Algorithm ITLS (6.43).

It will be convenient to state the relationships represented by Algorithm CGLS in matrix form. In terms of the quantities defined in Algorithm CGLS define the matrices

$$(7.1) \quad \bar{V} = [\bar{v}^{(1)}, \dots, \bar{v}^{(k)}]$$

$$(7.2) \quad \bar{H} = [\bar{h}^{(0)}, \dots, \bar{h}^{(k-1)}]$$

$$(7.3) \quad D = \text{Diag} \{ \|\bar{A}\bar{v}^{(1)}\|, \dots, \|\bar{A}\bar{v}^{(k)}\| \} \\ \equiv \text{Diag} \{ d_1, \dots, d_k \}$$

and

$$(7.4) \quad F = \text{Diag} \{ \|\bar{h}^{(0)}\|, \dots, \|\bar{h}^{(k-1)}\| \} \\ \equiv \text{Diag} \{ f_0, \dots, f_{k-1} \}$$

Then from the orthogonality of the vectors $\bar{h}^{(i)}$ we have

$$(7.5) \quad \bar{H}^T \bar{H} = F^2$$

and from the $(A^T A)$ -conjugacy of the vectors $\bar{v}^{(i)}$ we have

$$(7.6) \quad \bar{V}^T A^T A \bar{V} = D^2$$

Following Householder ((1964) pp. 2 and 139-141) define

$$(7.7) \quad J = \begin{bmatrix} 0 & & & & 0 \\ & 1 & & & \\ & & 0 & & \\ & & & \ddots & \\ & & & & 1 & \\ 0 & & & & & 0 \end{bmatrix}$$

Then Equations (4.8) and (4.10) can be written respectively as

$$(7.8) \quad \bar{H}(I-J) = A^T A \bar{V} F^2 D^{-2}$$

and

$$(7.9) \quad \bar{V} F^{-2} (I-J^T) F^2 = \bar{H}$$

In terms of these quantities we now define quantities, distinguished by a caret, which will be shown to satisfy the relationships of Algorithm ITLS.

$$(7.10) \quad \hat{V} = \bar{H} F^{-1}$$

$$(7.11) \quad \hat{U} = A \bar{V} D^{-1}$$

$$(7.12) \quad \hat{R} = D F^{-2} (I-J^T) F$$

$$(7.13) \quad \hat{W} = \bar{V} D^{-1}$$

Note that \hat{R} is upper bidiagonal.

Using the definitions (7.10) - (7.12) and the Equations (7.5) - (7.9) it can be directly verified that the equations

$$(7.14) \quad \hat{V}^T \hat{V} = I$$

$$(7.15) \quad \hat{U}^T \hat{U} = I$$

$$(7.16) \quad A \hat{V} = \hat{U} \hat{R}$$

and

$$(7.17) \quad A^T \hat{U} = \hat{V} \hat{R}^T$$

are satisfied.

From Equation (7.10) the first column vector of the matrix \hat{V} is equal to $\bar{h}^{(0)} / \|\bar{h}^{(0)}\|$ or equivalently $A^T b / \|A^T b\|$. This condition along with Equations (7.16) - (7.17) assure that the matrices \hat{U} , \hat{V} , and \hat{R} are identical with the matrices U , V , and R which would be generated by Algorithm ITLS in solving the least squares problem $Ax \cong b$.

Furthermore by use of Equation (7.9) it can be directly verified that the matrix \hat{W} of Equation (7.13) satisfies

$$(7.18) \quad \hat{W} \hat{R} = \hat{V}$$

In accordance with Equation (6.30) define

$$(7.19) \quad \hat{g} = \hat{U}^T b$$

Let $\hat{w}^{(i)}$ denote the i^{th} column vector of \hat{W} and let \hat{g}_i denote the i^{th} component of \hat{g} . We wish to show that the correction $\bar{v}^{(i)}_{p_i}$ which is added to the approximate solution vector at Equation (4.7) of Algorithm CGLS is identical with the correction vector $\hat{w}^{(i)}_{g_i}$ which is added to the approximate solution vector at Step 14 of Algorithm ITLS. Thus we wish to show that

$$(7.20) \quad \bar{v}^{(i)}_{p_i} = \hat{w}^{(i)}_{g_i} \quad i=1, \dots, k$$

From Equation (7.13) we have $\hat{w}^{(i)} = \bar{v}^{(i)}/d_i$ and from Equations (4.6), (7.3), and (7.4) we have $\bar{p}_i = f_{i-1}^2/d_i^2$. Thus Equation (7.20) may be established by proving that

$$(7.30) \quad f_{i-1}^2 = d_i \hat{g}_i \quad i=1, \dots, k$$

Introduce the k -vector $e = [1, \dots, 1]^T$ so that Equation (7.30) can be written as

$$(7.31) \quad \begin{aligned} F^2 e &= D \hat{g} \\ &= D \hat{U}^T b && \text{[Using Equation (7.19)]} \\ &= \bar{V}^T A^T b && \text{[Using Equation (7.11)]} \\ &= \bar{V}^T \bar{h}^{(c)} && \text{[Using Equation (4.4)]} \end{aligned}$$

Left multiply Equation (7.8) by \bar{V}^T and use Equation (7.6) obtaining

$$(7.32) \quad \bar{V}^T \bar{H}(I-J) = F^2$$

Substitute this expression for F^2 into Equation (7.31) obtaining

$$(7.33) \quad \bar{V}^T \bar{H}(\mathbf{I}-\mathbf{J})\mathbf{e} = \bar{V}^T \bar{h}^{(0)}$$

as the equation to be verified. This equation is clearly true since

$$(\mathbf{I}-\mathbf{J})\mathbf{e} = [1, 0, \dots, 0]^T.$$

This completes the verification that the algorithms CGLS and ITLS produce the same sequence of approximate solutions. It follows that the vector $\bar{h}^{(i)} = \mathbf{A}^T [b - \mathbf{A}\bar{x}^{(i)}]$ of Equation (4.8) is identical with the vector $h^{(i)} = \mathbf{A}^T [b - \mathbf{A}x^{(i)}]$ of Equation (6.35).

Chapter 8 Solving a Consistent System $Ax=b$

Let A be an $m \times n$ matrix and let b be an m -vector contained in the column space (range space) of A . Consider the problem of finding an n -vector x satisfying

$$(8.1) \quad Ax = b$$

The most common case of a consistent system of linear equations would be the case in which the matrix A is square and nonsingular. Also of practical interest is the case of a full rank underdetermined problem, i. e., $m < n$ and $\text{Rank}(A) = m$. The algorithm to be described permits either $m \leq n$ or $m > n$ and does not require any restriction on the rank of A .

If $\text{Rank}(A) < n$ the solution to Problem (8.1) is nonunique. In this case there is a unique solution vector \hat{x} of minimum euclidean length. This minimum length solution vector is characterized by being the only solution vector for Problem (8.1) lying in the row space of A . The algorithm to be described constructs this solution vector, \hat{x} .

Assume there exists an integer k ($1 \leq k \leq \min\{m, n\}$) and matrices

$$(8.2) \quad U_{m \times k} = [u^{(1)}, \dots, u^{(k)}]$$

$$(8.3) \quad V_{n \times k} = [v^{(1)}, \dots, v^{(k)}]$$

$$(8.4) \quad L_{k \times k} = \begin{bmatrix} \alpha_1 & & & & 0 \\ \beta_2 & \alpha_2 & & & \\ & \cdot & \cdot & & \\ & & \cdot & \cdot & \\ 0 & & & \beta_k & \alpha_k \end{bmatrix} \quad (\text{all } \alpha_i > 0 \text{ and } \beta_i > 0)$$

and a k -vector, \hat{p} , such that

$$(8.5) \quad U^T U = I_k$$

$$(8.6) \quad V^T V = I_k$$

$$(8.7) \quad AV = UL$$

$$(8.8) \quad A^T U = VL^T$$

and

$$(8.9) \quad \hat{x} = V\hat{p}$$

Since L is nonsingular Equations (8.6) and (8.8) imply that the column vectors of V form an orthonormal basis for a subspace \mathcal{Y} of the row-space of A . Similarly Equations (8.5) and (8.7) imply that the column vectors of U form an orthonormal basis for a subspace \mathcal{U} of the column space (range space) of A . Equation (8.9) shows that the subspace \mathcal{Y} contains the solution vector \hat{x} . From Equations (8.1), (8.7), and (8.9) it follows that the subspace \mathcal{U} contains the right-side vector b .

Assuming the availability of the matrices U , V , and L , Problem (8.1) can be approached as follows:

Introduce the change of variables

$$(8.10) \quad x = Vp$$

in Problem (8.1) and use Equation (8.7) obtaining the equivalent problem

$$(8.11) \quad ULp = b$$

Left multiplying this equation by U^T and using Equation (8.5) gives the $k \times k$ nonsingular lower bidiagonal system

$$(8.12) \quad Lp = g$$

where g is the k -vector defined by

$$(8.13) \quad g = U^T b$$

A computational algorithm can thus be based on the computation of g using Equation (8.13), solving for \hat{p} in Equation (8.12), and finally computing \hat{x} using Equation (8.9). These steps are all directly amendable to being organized in a sequential form which uses the vectors $u^{(i)}$ and $v^{(i)}$ as they are produced.

Assuming the nontrivial case of $b \neq 0$ we follow Paige (1972) in defining $u^{(1)}$ by the equations

$$(8.14) \quad \beta_1 = \|b\|$$

$$(8.15) \quad u^{(1)} = b/\beta_1$$

Thus $u^{(1)}$ is in the column space (range space) of A since it was assumed that this was true of b .

The other vectors $u^{(i)}$ and $v^{(i)}$ and the elements α_i and β_i of the matrix L are sequentially determined by the unit length requirements of Equations (8.5) and (8.6) and the following equations which follow directly from Equations (8.7) and (8.8).

$$(8.16) \quad u^{(i)} \beta_i = A v^{(i-1)} - u^{(i-1)} \alpha_{i-1} \quad i=2, \dots, k$$

$$(8.17) \quad 0 = A v^{(k)} - u^{(k)} \alpha_k$$

$$(8.18) \quad v^{(1)} \alpha_1 = A^T u^{(1)}$$

$$(8.19) \quad v^{(i)} \alpha_i = A^T u^{(i)} - v^{(i-1)} \beta_i \quad i=2, \dots, k$$

Comparing Equations (8.16) and (8.17) we note that the integer k is generally not known a priori but is determined as the first value of i for which $A v^{(i)} - u^{(i)} \alpha_i = 0$.

With $u^{(1)}$ chosen in the column space of A it follows from Equations (8.16), (8.18), and (8.19) that all $u^{(i)}$ will be in the column space of A and all $v^{(i)}$ will be in the row space of A . It will subsequently be verified that vectors $u^{(i)}$ and $v^{(i)}$ produced in this way necessarily satisfy the orthogonality conditions of Equations (8.5) and (8.6).

With $u^{(1)}$ defined by Equation (8.15) the vector g defined by Equation (8.13) is representable as

$$(8.20) \quad g = \beta_1 e^{(1)}$$

where $e^{(1)}$ denotes the first column vector of the $k \times k$ identity matrix. Using this expression for g Equation (8.12) becomes

$$(8.21) \quad Lp = \beta_1 e^{(1)}$$

which permits the components, p_i , of the solution vector \hat{p} to be expressed as

$$(8.21) \quad p_1 = \beta_1 / \alpha_1$$

$$(8.22) \quad p_i = -(\beta_i / \alpha_i) p_{i-1} \quad i=2, \dots, k$$

Define the sequence of k-vectors

$$(8.23) \quad p^{(i)} = [p_1, \dots, p_i, \underbrace{0, \dots, 0}_{k-i}]^T \quad i=0, \dots, k$$

Define the sequence of approximate solution vectors

$$(8.24) \quad x^{(0)} = 0$$

$$(8.25) \quad x^{(i)} = V p^{(i)} = \sum_{j=1}^i v^{(j)} p_j = x^{(i-1)} + v^{(i)} p_i \quad i=1, \dots, k$$

Associated residual vectors are defined by

$$(8.26) \quad r^{(i)} = b - A x^{(i)} \quad i=0, \dots, k$$

and may be expressed as

$$(8.27) \quad r^{(0)} = b = u^{(1)} \beta_1$$

$$(8.28) \quad r^{(i)} = b - A V p^{(i)} = b - U L p^{(i)}$$

$$= b - U \begin{bmatrix} \beta_1 \\ 0 \\ \vdots \\ 0 \\ \beta_{i+1} p_i \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \leftarrow \text{row } i+1$$

$$= b - u^{(1)} \beta_1 - u^{(i+1)} \beta_{i+1} p_i$$

$$= -u^{(i+1)} \beta_{i+1} p_i \quad i=1, \dots, k-1$$

$$(8.29) \quad r^{(k)} = 0$$

If we introduce the additional definitions, $p_0 = -1$, and $\beta_{k+1} = 0$ the norm of the residual vector can be expressed as

$$(8.30) \quad \rho_i = \|r^{(i)}\| = |\beta_{i+1} p_i| \quad i=0, \dots, k$$

These considerations may be organized into a computational algorithm as follows:

(8.31) Algorithm ITC For Iterative Solution of Consistent Linear Equations

[Given by Paige (1972) pp. 21-22]

<u>Step Number</u>	<u>Description</u>
1	$x^{(0)} := 0$
2	$p_0 := -1$
3	$i := 1$
4	$\tilde{u}^{(i)} := \begin{cases} b & \text{if } i=1 \\ A v^{(i-1)} - u^{(i-1)} \alpha_{i-1} & \text{if } i > 1 \end{cases}$
5	$\beta_i := \ \tilde{u}^{(i)}\ $
6	$\rho_{i-1} := \beta_i p_{i-1} $
7	Theoretical termination test: If $\beta_i = 0$ go to Step 14. Practical termination test: If either β_i or ρ_{i-1} is sufficiently small go to Step 16.
8	$u^{(i)} := \tilde{u}^{(i)} / \beta_i$
9	$\tilde{v}^{(i)} := \begin{cases} A^T u^{(1)} & \text{if } i=1 \\ A^T u^{(i)} - v^{(i-1)} \beta_i & \text{if } i > 1 \end{cases}$
10	$\alpha_i := \ \tilde{v}^{(i)}\ $
11	$v^{(i)} := \tilde{v}^{(i)} / \alpha_i$
12	$p_i := -(\beta_i / \alpha_i) p_{i-1}$
13	$x^{(i)} := x^{(i-1)} + v^{(i)} p_i$
14	$i := i + 1$
15	Go to Step 4
16	$k := i - 1$
17	Stop

It must be verified that all of the vectors $u^{(i)}$ and $v^{(i)}$ produced by this algorithm have the orthogonality properties specified by Equations (8.5) and (8.6) and that all of the numbers α_i defined by this algorithm are positive.

Assume these conditions are satisfied for $i=1, \dots, \ell-1$. Consider the quantities computed during the ℓ^{th} iteration.

If β_ℓ computed at Step 5 is zero the algorithm terminates setting $k = \ell-1$. Thus $\beta_{k+1}=0$ and $\rho_k=0$ which (see Equation (8.30)) implies that the current approximate solution vector $x^{(k)}$ is actually the unique minimal length solution vector \hat{x} .

If $\beta_\ell \neq 0$ then $\beta_\ell > 0$ and the orthogonality of $\tilde{u}^{(\ell)}$ relative to $u^{(1)}, \dots, u^{(\ell-1)}$ must be verified

$$\begin{aligned}
 (8.32) \quad u^{(i)T} \tilde{u}^{(\ell)} &= u^{(i)T} T_{Av}^{(\ell-1)} - u^{(i)T} T_u^{(\ell-1)} \alpha_{\ell-1} \\
 &= [v^{(i)} \alpha_i + v^{(i-1)} \beta_i] T_v^{(\ell-1)} - u^{(i)T} T_u^{(\ell-1)} \alpha_{\ell-1} \\
 &= v^{(i)T} T_v^{(\ell-1)} \alpha_i - u^{(i)T} T_u^{(\ell-1)} \alpha_{\ell-1} \quad i=2, \dots, \ell-1
 \end{aligned}$$

For $i < \ell-1$ each term of this final expression vanishes while for $i = \ell-1$ the final expression reduces to $\alpha_{\ell-1} - \alpha_{\ell-1} = 0$. The verification that $u^{(1)T} \tilde{u}^{(\ell)} = 0$ is similarly straightforward.

Similarly it can be verified that the vector $\tilde{v}^{(\ell)}$ defined at Step 9 satisfies $v^{(i)T} \tilde{v}^{(\ell)} = 0$ for $i=1, \dots, \ell-1$. It remains to be shown that α_ℓ defined at Step 10 is positive.

Suppose $\alpha_\ell = 0$. Then

$$\begin{aligned}
(8.33) \quad 0 &= \tilde{v}^{(\ell)} = A^T u^{(\ell)} - v^{(\ell-1)} \beta_{\ell} \\
&= A^T u^{(\ell)} - [A^T u^{(\ell-1)} - v^{(\ell-2)} \beta_{\ell-1}] \beta_{\ell} / \alpha_{\ell-1} \\
&= \dots = \sum_{i=1}^{\ell} c_i A^T u^{(i)} = A^T \sum_{i=1}^{\ell} c_i u^{(i)}
\end{aligned}$$

Here the coefficients c_i are all nonzero and the vectors $u^{(1)}, \dots, u^{(\ell)}$ constitute an orthonormal basis for a subspace of the column space (range space) of A .

Thus the vector $z = \sum_{i=1}^{\ell} c_i u^{(i)}$ is a nonzero vector in the column space of A . It follows that $A^T z \neq 0$ which contradicts Equation (8.33). We conclude that $\alpha_{\ell} > 0$.

The practical termination test at Step 7 might be implemented as a comparison of β_i with some computed estimation of the norm of the round-off error vector associated with the computed vector $\tilde{u}^{(i)}$ or a comparison of ρ_{i-1} with $\eta(\|b\| + \|A\| \cdot \|x^{(i)}\|)$ where η denotes the relative machine precision.

Chapter 9 The Theoretical Equivalence of Algorithms CGC and ITC

Using the notation of Algorithm CGC (5.7) define

$$(9.1) \quad v^{(i)} = (-1)^{i-1} \bar{v}^{(i)} / \|\bar{v}^{(i)}\| \quad i=1, \dots, k$$

$$(9.2) \quad u^{(i)} = (-1)^{i-1} \bar{r}^{(i-1)} / \|\bar{r}^{(i-1)}\| \quad i=1, \dots, k$$

$$(9.3) \quad \alpha_i = \|\bar{v}^{(i)}\| / \|\bar{r}^{(i-1)}\| = (\bar{p}_i)^{-1/2} \quad i=1, \dots, k$$

$$(9.4) \quad \beta_1 = \|b\| = \|\bar{r}^{(0)}\|$$

$$(9.5) \quad \beta_i = \|\bar{v}^{(i-1)}\| \cdot \|\bar{r}^{(i-1)}\| / \|\bar{r}^{(i-2)}\|^2 \\ = (\bar{\beta}_{i-1} / \bar{p}_{i-1})^{1/2} \quad i=2, \dots, k$$

$$(9.6) \quad p_0 = -1$$

$$(9.7) \quad p_i = (-1)^{i-1} \|\bar{r}^{(i-1)}\|^2 / \|\bar{v}^{(i)}\| \\ = (-1)^{i-1} \bar{p}_i \|\bar{v}^{(i)}\| \quad i=1, \dots, k$$

$$(9.8) \quad x^{(i)} = \bar{x}^{(i)} \quad i=0, \dots, k$$

Since the sets of vectors $\{\bar{v}^{(1)}, \dots, \bar{v}^{(k)}\}$ and $\{\bar{r}^{(0)}, \dots, \bar{r}^{(k-1)}\}$ are each mutually orthogonal it follows that the sets $\{v^{(1)}, \dots, v^{(k)}\}$ and $\{u^{(1)}, \dots, u^{(k)}\}$ defined by Equations (9.1) and (9.2) are each mutually orthonormal. Using

the equations of Algorithm CGC (5.7) it can be directly verified that the quantities defined by Equations (9.1) - (9.8) satisfy the relations of Algorithm ITC (8.31). Thus the algorithms CGC and ITC theoretically produce the same sequence of approximate solution vectors.

Since the difference between the two algorithms only involves different scale factors it is to be expected that, apart from questions of exponent overflow or underflow, the two algorithms will exhibit essentially the same numerical behaviour also.

The theoretical equivalence of these two algorithms was pointed out by Paige (1972), p. 13.

Chapter 10 Solving a Consistent System $Ax=b$ where A is Symmetric

Let A be an $n \times n$ symmetric matrix and let b be an n -vector contained in the column space (range space) of A . We wish to find a vector x satisfying

$$(10.1) \quad Ax = b$$

In practical problems of this type the matrix A would usually be of rank n . The algorithm to be described does not require that $\text{Rank}(A) = n$. If $\text{Rank}(A) < n$ the solution of Problem (10.1) is nonunique. In this case the algorithm finds the (unique) solution vector \hat{x} lying in the row space of A . This is the minimal length solution vector for Problem (10.1).

Assume the existence of matrices $V_{n \times k}$ and $C_{k \times k}$ and a k -vector \hat{p} such that

$$(10.2) \quad V = [v^{(1)}, \dots, v^{(k)}]$$

$$C = \begin{bmatrix} \alpha_1 & \beta_2 & \dots & 0 \\ & \alpha_2 & \dots & \beta_k \\ & & \dots & \\ 0 & \beta_k & \dots & \alpha_k \end{bmatrix}$$

$$(10.4) \quad V^T V = I_k$$

$$(10.5) \quad AV = VC$$

and

$$(10.6) \quad \hat{x} = V\hat{p}$$

If the matrices V and C are available Problem (10.1) can be attacked as follows. Make the change of variables

$$(10.7) \quad x = Vp$$

in Equation (10.1), then use Equation (10.5) obtaining

$$(10.8) \quad VCp = b$$

Left multiply Equation (10.8) by V^T using Equation (10.4).

$$(10.9) \quad Cp = V^T b \equiv g$$

Thus Problem (10.1) could be solved by first computing $g = V^T b$, next solving $Cp = g$ for \hat{p} , and finally computing $\hat{x} = V\hat{p}$.

Since C is a symmetric tridiagonal matrix the entire matrix C must be determined before any components of \hat{p} can be computed. Thus the equation $Cp = g$ is not directly suitable for use in an algorithm which discards old vectors $v^{(i)}$ as new ones are computed.

Let Q be a $k \times k$ orthogonal matrix which on post-multiplication times C produces a $k \times k$ lower tridiagonal matrix L .

$$(10.10) \quad CQ = L = \begin{bmatrix} l_{11} & & & & \\ l_{21} & l_{22} & & & 0 \\ l_{31} & l_{32} & l_{33} & & \\ \cdot & \cdot & \cdot & \cdot & \\ 0 & l_{k, k-2} & l_{k, k-1} & & l_{kk} \end{bmatrix}$$

Define

$$(10.11) \quad W = VQ$$

and

$$(10.12) \quad z = Q^T p$$

Then

$$(10.13) \quad Lz = CQQ^T p = Cp = g$$

and

$$(10.14) \quad x = Vp = VQQ^T p = Wz$$

Thus Problem (10.1) can be solved by the following sequence of operations assuming that b , V , and C are known a priori

$$(10.15) \quad g = V^T b$$

$$(10.16) \quad L = CQ$$

$$(10.17) \quad \text{Solve } Lz = g \text{ for } \hat{z}$$

$$(10.18) \quad W = VQ$$

$$(10.19) \quad \hat{x} = W\hat{z}$$

Each individual step of this sequence is amenable to being implemented in a sequential manner so that in fact V and C need not be known a priori but rather the column vectors of V and the elements α_i and β_i of C can be computed, used, and discarded sequentially.

This method of solving Problem (10.1) is due to C. C. Paige and M. A. Saunders (Personal correspondence, 1972).

We follow M. Saunders in defining β_1 and $v^{(1)}$ by

$$(10.20) \quad \beta_1 = \|b\|$$

and

$$(10.21) \quad v^{(1)} = b/\beta_1 \quad [\text{assuming } \beta_1 \neq 0]$$

From Equation (10.5) we may write the equations

$$(10.22) \quad v^{(2)}e_2 = Av^{(1)} - v^{(1)}\alpha_1$$

$$(10.23) \quad v^{(i+1)}e_{i+1} = Av^{(i)} - v^{(i)}\alpha_i - v^{(i-1)}e_i \quad i=2, \dots, k-1$$

and

$$(10.24) \quad 0 = Av^{(k)} - v^{(k)}\alpha_k - v^{(k-1)}e_k$$

The numbers β_{i+1} may be computed as normalization factors to assure that the vectors $v^{(i+1)}$ have unit length as required by Equation (10.4). The numbers α_i can be computed as

$$(10.25) \quad \alpha_i = v^{(i)T}Av^{(i)}$$

which is the condition (see Equation (2.1)) which assures that $v^{(i+1)}$ computed by Equation (10.22) or (10.23) will be orthogonal to $v^{(i)}$.

The integer k will generally not be known in advance and may (theoretically) be determined as the first value of i for which the right side of Equation (10.23) is zero. With k so determined it can be verified that the vectors $v^{(1)}, \dots, v^{(k)}$ produced by use of Equations (10.21) - (10.23) form an orthonormal basis for a subspace of the column space of A . This verification makes essential use of the symmetry of A .

It is pointed out by Paige and Saunders that the computation of the orthogonal set of vectors $v^{(i)}, i=1, \dots, k$, using Equations (10.22)-(10.23) is a method due to Lanczos (1950 and 1952).

The orthogonal matrix Q will be (implicitly) constructed as the product

$$(10.26) \quad Q = \tilde{G}_1 \tilde{G}_2 \cdots \tilde{G}_{k-1}$$

where

$$(10.27) \quad \tilde{G}_i = \left. \begin{array}{ccc} \left[\begin{array}{ccc} I & 0 & 0 \\ 0 & G_i & 0 \\ 0 & 0 & I \end{array} \right] & \left. \begin{array}{l} i-1 \\ 2 \\ k-i-1 \end{array} \right\} & i=1, \dots, k-1 \end{array} \right.$$

$$(10.28) \quad G_i = \begin{bmatrix} c_i & s_i \\ s_i & -c_i \end{bmatrix} \quad i=1, \dots, k-1$$

and

$$(10.29) \quad c_i^2 + s_i^2 = 1$$

Each matrix G_i effects a nontrivial transformation on a particular 3×2 or 2×2 submatrix of the appropriate intermediate matrix which arises in the process of transforming the symmetric tridiagonal matrix C to the lower tridiagonal matrix L . This action may be expressed as follows:

$$(10.30) \quad \bar{l}_{11} = \alpha_1$$

$$(10.31) \quad \bar{l}_{21} = \theta_2$$

$$(10.32) \quad \begin{bmatrix} \bar{l}_{i,i} & \theta_{i+1} \\ \bar{l}_{i+1,i} & \alpha_{i+1} \\ 0 & \theta_{i+2} \end{bmatrix} \cdot G_i = \begin{bmatrix} l_{i,i} & 0 \\ l_{i+1,i} & \bar{l}_{i+1,i+1} \\ l_{i+2,i} & \bar{l}_{i+2,i+1} \end{bmatrix} \quad i=1, \dots, k-2$$

$$(10.33) \quad \begin{bmatrix} \bar{l}_{k-1,k-1} & \theta_k \\ \bar{l}_{k,k-1} & \alpha_k \end{bmatrix} \cdot G_{k-1} = \begin{bmatrix} l_{k-1,k-1} & 0 \\ l_{k,k-1} & l_{k,k} \end{bmatrix}$$

As each additional row of L is determined by these equations an additional component z_i of the vector \hat{z} satisfying Equation (10.17) can be determined. Note that

$$(10.34) \quad g = \begin{bmatrix} \theta_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

due to Equations (10.4), (10.15), and (10.21). Thus the equations for the components z_i are

$$(10.35) \quad z_1 = \theta_1 / l_{11}$$

$$(10.36) \quad z_2 = -l_{21}z_1 / l_{22}$$

and

$$(10.37) \quad z_i = -(l_{i,i-2}z_{i-2} + l_{i,i-1}z_{i-1}) / l_{ii} \quad i=3, \dots, k$$

Define the vectors

$$(10.38) \quad z^{(i)} = [z_1, \dots, z_i, \underbrace{0, \dots, 0}_{k-i}] \quad i=0, 1, \dots, k$$

and, motivated by Equation (10.19), define a sequence of approximate solution vectors

$$x^{(i)} = Wz^{(i)} \quad i=0, 1, \dots, k$$

The associated residual vectors are expressible as

$$(10.39) \quad \begin{aligned} r^{(i)} &= b - Ax^{(i)} = b - AWz^{(i)} \\ &= b - AVQz^{(i)} = b - VCQz^{(i)} \\ &= b - VLz^{(i)} = v^{(1)}\beta_1 - VLz^{(i)} \\ &= V[e^{(1)}\beta_1 - Lz^{(i)}] \quad i=0, 1, \dots, k \end{aligned}$$

from which we may write

$$(10.40) \quad r^{(0)} = b = v^{(1)} \theta_1$$

$$(10.41) \quad r^{(i)} = -[v^{(i+1)}; v^{(i+2)}] \cdot \begin{bmatrix} \ell_{i+1, i-1} & \ell_{i+1, i} \\ 0 & \ell_{i+2, i} \end{bmatrix} \cdot \begin{bmatrix} z_{i-1} \\ z_i \end{bmatrix} \quad i=1, \dots, k-1$$

$$(10.42) \quad r^{(k)} = 0$$

Interpretation of Equation (10.41) for $i=1$ requires the definitions $\ell_{2, 0}=0$ and $z_0=0$ while for $i=k-1$ one must define $v^{(k+1)}=0$ and $\ell_{k+1, k-1}=0$.

The norms of the residual vectors are expressible as

$$(10.43) \quad \rho_0 \equiv \|r^{(0)}\| = \theta_1$$

$$(10.44) \quad \rho_i \equiv \|r^{(i)}\| = [(\ell_{i+1, i-1} z_{i-1} + \ell_{i+1, i} z_i)^2 + (\ell_{i+2, i} z_i)^2]^{1/2}$$

for $i=1, \dots, k-1$

$$(10.45) \quad \rho_k \equiv \|r^{(k)}\| = 0$$

where again we define $\ell_{2, 0}=0$, $z_0=0$, and $\ell_{k+1, k-1}=0$.

Combining these equations appropriately one can obtain the following algorithm

(10.46) Algorithm ICSE Iterative Solution of a Consistent Symmetric System of Linear Equations [Originated by C. C. Paige and M. A. Saunders, personal communication, 1972]

<u>Step</u>	<u>Description</u>
1	$x^{(0)} := 0$
2	$\beta_1 := \ b\ $
3	If $\beta_1 = 0$ set $k := 0$ and go to Step 24
4	$v^{(1)} := b / \beta_1$
5	$u^{(1)} := v^{(1)}$
6	<u>QUIT := FALSE</u>
7	$i := 1$
8	$y^{(i)} := Av^{(i)}$
9	$\alpha_i := v^{(i)T} y^{(i)}$
10	$\tilde{v}^{(i+1)} := \begin{cases} y^{(1)} - \alpha_1 v^{(1)} & \text{if } i=1 \\ y^{(i)} - \alpha_i v^{(i)} - \beta_i v^{(i-1)} & \text{if } i > 1 \end{cases}$
11	$\beta_{i+1} := \ \tilde{v}^{(i+1)}\ $
12	Theoretical termination test: If $\beta_{i+1} = 0$ set <u>QUIT := TRUE</u> Practical termination test: If β_{i+1} is sufficiently small set <u>QUIT := TRUE</u>
13	$\left\{ \begin{array}{l} \begin{bmatrix} \bar{l}_{11} \\ \bar{l}_{21} \end{bmatrix} := \begin{bmatrix} \alpha_1 \\ \beta_2 \end{bmatrix} \quad \text{if } i=1 \\ \begin{bmatrix} \bar{l}_{i,i-1} & \bar{l}_{i,i} \\ \bar{l}_{i+1,i-1} & \bar{l}_{i+1,i} \end{bmatrix} := \begin{bmatrix} \bar{l}_{i,i-1} & \alpha_i \\ 0 & \beta_{i+1} \end{bmatrix} G_{i-1} \quad \text{if } i > 1 \end{array} \right.$

<u>Step</u>	<u>Description</u>
14	$\begin{cases} \text{If } i=1 \text{ go to Step 16} \\ \text{If } i=2 \text{ set } \rho_1 := [(\ell_{21} z_1)^2 + (\ell_{31} z_1)^2]^{1/2} \\ \text{If } i > 2 \text{ set } \rho_{i-1} := [(\ell_{i, i-2} z_{i-2} + \ell_{i, i-1} z_{i-1})^2 + (\ell_{i+1, i-1} z_{i-1})^2]^{1/2} \end{cases}$
15	Practical termination test: If ρ_{i-1} is sufficiently small set $k:=i-1$ and go to Step 24.
16	$\ell_{ii} := (\bar{\ell}_{ii}^2 + \theta_{i+1}^2)^{1/2}$
17	$z_i := \begin{cases} \theta_1 / \ell_{11} & \text{if } i=1 \\ -\ell_{21} z_1 / \ell_{22} & \text{if } i=2 \\ -(\ell_{i, i-2} z_{i-2} + \ell_{i, i-1} z_{i-1}) / \ell_{ii} & \text{if } i > 2 \end{cases}$
18	<p>If <u>QUIT=FALSE</u> set $v^{(i+1)} := \tilde{v}^{(i+1)} / \theta_{i+1}$, $c_i := \bar{\ell}_{ii} / \ell_{ii}$,</p> <p>$s_i := \theta_{i+1} / \ell_{ii}$,</p> <p>$G_i := \begin{bmatrix} c_i & s_i \\ s_i & -c_i \end{bmatrix}$, and $[w^{(i)}, u^{(i+1)}] := [u^{(i)}, v^{(i+1)}] G_i$</p>
19	If <u>QUIT=TRUE</u> set $w^{(i)} := u^{(i)}$
20	$x^{(i)} := x^{(i-1)} + w^{(i)} z_i$
21	If <u>QUIT=TRUE</u> set $k:=i$ and go to Step 24
22	$i:=i+1$
23	Go to Step 8
24	Here the algorithm is finished with the solution vector $x^{(k)}$.

References

- F. S. Beckman (1960), The Solution of Linear Equations by the Conjugate Gradient Method, pp. 62-72 of Mathematical Methods for Digital Computers, ed. by Ralston and Wilf, Wiley.
- E. J. Craig (1955), The N-Step Iteration Procedures, Journal of Math. and Physics, 34, (1955), 64-73.
- D. K. Faddeev and V. N. Faddeeva (1963), Computational Methods of Linear Algebra, Freeman.
- L. Fox (1965), An Introduction to Numerical Linear Algebra, Oxford Univ. Press.
- G. Golub and W. Kahan (1965), Calculating the Singular Values and Pseudoinverse of a Matrix, J. SIAM Numer. Anal., Ser B, 2, 205-224.
- M. R. Hestenes and E. Stiefel (1952), Methods of Conjugate Gradients for Solving Linear Systems, J. Res. Nat. Bur. Standards, 49, 409-436.
- A. S. Householder (1964), The Theory of Matrices in Numerical Analysis, Blaisdell
- C. Lanczos (1950), An Iteration Method for the Solution of the Eigenvalue Problem of Linear Differential and Integral Operators, J. Res. Nat. Bur. Standards, 45, 255-282.
- C. Lanczos (1952), Solution of Systems of Linear Equations by Minimized Iterations, J. Res. Nat. Bur. Standards, 49, 33-53.
- C. C. Paige (1972), Bidiagonalization of Matrices and Solution of Linear Equations, Stanford Univ. Report No. CS-295, 28 pp.
- J. K. Reid (1971), On the Method of Conjugate Gradients for the Solution of Large Sparse Systems of Linear Equations, pp. 231-254 in Large Sparse Sets of Linear Equations, ed. by J. K. Reid, Academic Press.