

National Library of Canada Bibliothèque nationale du Canada

Acquisitions and Direction des acquisitions et Bibliographic Services Branch des services bibliographiques

395 Wellington Street Ottawa, Ontano K1A 0N4 395, rue Wellington Ottawa (Ontario) K1A 0N4

Your Ne - Votre rélérence Our Ne - Notre rélérence

NOTICE

The quality of this microform is heavily dependent upon the quality of the original thesis submitted for microfilming. Every effort has been made to ensure the highest quality of reproduction possible.

If pages are missing, contact the university which granted the degree.

Some pages may have indistinct print especially if the original pages were typed with a poor typewriter ribbon or if the university sent us an inferior photocopy.

Reproduction in full or in part of this microform is governed by the Canadian Copyright Act, R.S.C. 1970, c. C-30, and subsequent amendments. La qualité de cette microforme dépend grandement de la qualité de la thèse soumise au microfilmage. Nous avons tout fait pour assurer une qualité supérieure de reproduction.

AVIS

S'il manque des pages, veuillez communiquer avec l'université qui a conféré le grade.

La qualité d'impression de certaines pages peut laisser à désirer, surtout si les pages originales ont été dactylographiées à l'aide d'un ruban usé ou si l'université nous a fait parvenir une photocopie de qualité inférieure.

La reproduction, même partielle, de cette microforme est soumise à la Loi canadienne sur le droit d'auteur, SRC 1970, c. C-30, et ses amendements subséquents.





# Matrix Based Derivations and Representations of Krylov Subspace Methods

Penelope L. Anderson

A thesis submitted to the Faculty of Graduate Studies and Research in partial fulfillment of the requirements of the degree of Master of Science

Copyright, Penelope L. Anderson, October 1995

1



National Library of Canada

Acquisitions and Bibliographic Services Branch Bibliothèque nationale du Canada

Direction des acquisitions et des services bibliographiques

395 Wellington Street Ottawa, Ontario K1A 0N4 395, rue Wellington Ottawa (Ontario) K1A 0N4

Your Ne - Voire rélérence

Our Ne Notre référence

The author has granted an irrevocable non-exclusive licence allowing the National Library of Canada to reproduce, loan, distribute or sell copies of his/her thesis by any means and in any form or format, making this thesis available to interested persons.

L'auteur a accordé une licence irrévocable et non exclusive permettant à la Bibliothèque du nationale Canada de reproduire, prêter, distribuer ou vendre des copies de sa thèse de quelque manière et sous quelque forme que ce soit pour mettre des exemplaires de cette thèse la disposition à des personnes intéressées.

The author retains ownership of the copyright in his/her thesis. Neither the thesis nor substantial extracts from it may be printed or otherwise reproduced without his/her permission. L'auteur conserve la propriété du droit d'auteur qui protège sa thèse. Ni la thèse ni des extraits substantiels de celle-ci ne doivent être imprimés ou autrement reproduits sans son autorisation.

ISBN 0-612-12154-2



### Acknowledgements

I would like to thank Professor Chris Paige for his guidance, support, patience and encouragement. My office-mate Xiaowen Chang for his help, friendship and the occasional dim sum. My family for their love and support whether near or far. My friends and room-mates for putting up with me throughout this process.

This thesis is dedicated to my step-father, John L. Hastings.

#### Abstract

This thesis is based on recent work by Paige which gave a formalism for presenting and analyzing the class of algorithms which manipulate an appropriate Krylov subspace in solving large sparse systems of linear equations. This formalism — a way of dividing a method of solution into a Krylov process and an associated subproblem - is described and then applied to several of the more popular algorithms in use today including the methods of Conjugate Gradients and Bi-Conjugate Gradients. The aim is to clarify these algorithms to make them easier to understand, analyze and use. Several of the methods presented in this thesis were developed in exactly this way - notably the Symmetric LQ method and the Generalized Minimum Residual method — and required little or no effort to characterize using the formalism. It was successfully applied to Conjugate Gradients and BiConjugate Gradients, already recognized as being closely related to the symmetric and unsymmetric Lanczos processes respectively. The newer algorithms such as Conjugate Gradients Squared and BiConjugate Gradients Stabilized, with less obvious relation to a specific Krylov process, provided more difficulty in their clarification.

#### Résumé

Ce mémoire est basé sur le travail récent de Paige sur un formalisme pour présenter et analyser une classe d'algorithmes manipulant un sous-espace de Krylov approprié dans la résolution de larges systèmes creux d'équations linéaires. Ce formalisme - une façon de diviser une méthode de résolution en un processus de Krylov et un sous-problème associé — est décrit et ensuite appliqué à plusieurs des algorithmes les plus populaires en ce moment, dont la méthode des Gradients Conjugués et celle des Gradients BiConjugués. Le but de ce mémoire est de clarifier ces algorithmes pour les rendre plus faciles à comprendre, analyser et utiliser. Plusieurs des méthodes ici présentées ont justement été développées de cette façon - la méthode des LQ Symétriques et la méthode des Résidus Minimums Généralisés en particulier - et leur caractérisation à l'aide du formalisme n'a requis que peu d'efforts. Le formalisme a été appliqué aux Gradients Conjugués et aux Gradients BiConjugués, connus pour être liés respectivement aux processus symétriques et asymétriques de Lanczos. Les algorithmes plus récents, tel que les Carrés des Gradients Conjugués et les Gradients BiConjugués Stabilisés, dont le lien avec un processus de Krylov spécifique est moins évidente, ont été plus difficiles à clarifier.

# Contents

1	Introduction		6
	1.1	Formalism	6
	1.2	Notation	8
	1.3	Table of Abbreviations	9
	1.4	A note on the literature review	9
2	Krylov Processes		
	2.1	The general Krylov process	10
	2.2	The Arnoldi algorithm	11
	2.3	The Lanczos process for symmetric A	12
	2.4	The Lanczos process for unsymmetric $A$	15
3	Pre	reliminaries - The QR decomposition 22	
4	Algorithms for Solving Linear Systems		
	4.1	The Lanczos (Conjugate Gradients) Method	27
	4.2	Symmetric LQ Method (SYMMLQ)	37
	4.3	Minimum Residual Method (MINRES)	42
	-1.4	Conjugate Gradients on the Normal Equations (CGNE) -	
		LSQR implementation	46
	4.5	Generalized Minimum Residual Method with restarts (GM-	
		RES(m))	53
	4.6	The Unsymmetric Lanczos (BiConjugate Gradients) Method	56
	4.7	Quasi-Minimal Residuals (QMR)	66
	4.8	Conjugate Gradients Squared (CGS)	70
	4.9	BiConjugate Gradients Stabilized (BiCGSTAB)	83
5	Co	Conclusion 9	

.



## 1 Introduction

### 1.1 Formalism

Based on ideas in [PS75], Paige [Pai94] states a simple formalism for motivating, presenting and studying Krylov subspace methods for problems involving a matrix A. These methods are useful for example in solving the eigenproblem  $Ax = \lambda x$  or solution of equations Ax = b when A is large and sparse.

Given a matrix  $A \in \mathbb{R}^{n \times n}$ , and a starting vector v, a Krylov Process produces vectors  $\{v_1, v_2, \ldots, v_k\}$  which span the k-th Krylov subspace

$$\mathcal{K}^{k}(A;v) \equiv \operatorname{span}\{v, Av, \dots, A^{k-1}v\}.$$

Krylov subspace methods seek elements of these subspaces which are in some sense good approximations to the solutions of the original problems.

The formalism for Krylov subspace methods is as follows:

A Krylov subspace process + Solving a subsidiary problem  $\rightarrow$  A Krylov subspace method.

The terms used above were chosen carefully to distinguish the parts of the algorithm.

*Process* refers to the particular process we are considering, for example, the Lanczos process, i.e. the way in which the vectors are formed to span the Krylov subspace in question. Many processes are presented in § 2.

Method refers to the theoretical method used to solve the problem. For example [BBC<sup>+</sup>94] gives one implementation of the Conjugate Gradients, or CG method for solving Ax = b; another implementation using the Lanczos process and then considering iterates  $x_k \in \mathcal{K}^k(A; b)$  obtained by manipulating the Lanczos vectors is well known, see for example [GL89, § 9.3.1] and [PS75] and § 4.1 in this thesis. Although the implementations are different, the resulting methods are theoretically the same (using exact arithmetic).

For any underlying Krylov process, the manipulation of the vectors obtained in that process is known as the *subproblem* for this process used to obtain this method.

A significant advantage of this approach is that both the *process* and the *subsidiary problem* may be presented as clear matrix formulations, and these make the *method* casy to understand and analyse. The less transparent vector representations (i.e. the implementations) of the method can then be obtained from these matrix formulations. The vector representation is important for implementation, but the matrix representation is a far more powerful tool for teaching, for understanding, for relating and deriving different methods, and for analysing their properties.

Paige and Saunders [PS75] implicitly used this formalism to understand old methods, and to develop new methods. However, several of the more recent methods in this area appear to be significantly more complicated than those in [PS75], and not nearly as well motivated. The result is that many users of these algorithms have little understanding of them, a fact which does not help either their use or implementation.

In this thesis the formalism is applied to the more successful Krylov subspace methods for Ax = b used today – for example GMRES (Generalized Minimum Residual), BiCG (BiConjugate Gradients), QMR (Quasi-Minimal Residual), CGS (Conjugate Gradients Squared) – all of which are summarized, usually in vector form, in [BBC+94].

The aim is to see if the formalism described here, which was so successful in clarifying earlier methods, can also help to simplify the understanding and development of today's more complicated methods.

### 1.2 Notation

\_•.

- Matrices are represented by upper case Roman letters such as A or  $B_k$ .
- Vectors are represented by lower case Roman letters such as  $v_k$ , including the columns of matrices, for example  $w_j$  is the *j*-th column of  $W_k$ . Usually columns of a matrix use the same letter as the matrix itself; one prominent exception is  $e_j$  which denotes the *j*-th column of a unit matrix *I* of appropriate dimension. Occasionally, for clarity, superscripts such as the *k* in  $e_i^{(k)}$  will be used to specify the dimension.
- Lower case Greek letters represent scalars. The scalar entries of matrices or vectors may be represented by lower case Greek letters or by lower case Roman letters the same as the entity itself, with appropriate subscripts. For example, the entries of the vector  $z_k$  are  $\zeta_j$ ; the entries of the matrix  $H_k$  are  $h_{ij}$ .
- $c_k$  and  $s_k$  represent  $cos(\theta_k)$  and  $sin(\theta_k)$  respectively in those algorithms which involve a QR or LQ decomposition.
- $\|\cdot\|$  is used throughout to denote the 2-norm.

### 1.3 Table of Abbreviations

QR	referring to the decomposition of a matrix into
	an orthogonal matrix and an upper triangular matrix
LQ	referring to the decomposition of a matrix into
	a lower triangular matrix and an orthogonal matrix
CG	Conjugate Gradients
SYMMLQ	Symmetric LQ
MINRES	Minimum Residual
CGNE	Conjugate Gradients on the Normal Equations
LSQR	Least Squares with the QR decomposition
GMRES	Generalized Minimum Residual
BiCG	BiConjugate Gradients
QMR	Quasi-Minimal Residual
CGS	Conjugate Gradients Squared
BICGSTAB	<b>BiConjugate Gradients Stabilized</b>
SPD	Symmetric Positive Definite

### 1.4 A note on the literature review

This thesis is required to have a review of the literature. In that the body of this thesis describes popular methods of solution of large sparse systems of linear equations and clearly refers to the original papers and subsequent articles written on these methods, it is our belief that this review is implicit and need not be presented separately.

# 2 Krylov Processes

### 2.1 The general Krylov process

We first state a very general Krylov process, then derive some of the more popular processes from this by manipulating its matrix formulation. We feel this is briefer and clearer than the more usual approaches to deriving these processes.

The general Krylov process with a given matrix  $A \in \mathbb{R}^{n \times n}$  and starting vector  $v \in \mathbb{R}^n$  forms a sequence of vectors  $v_1, v_2, \ldots \in \mathbb{R}^n$  as follows

$$h_{1,0}v_1 = v$$
  
$$h_{j+1,j}v_{j+1} = Av_j - h_{j,j}v_j - \dots - h_{1,j}v_1, \quad j = 1, 2, \dots, \quad (1)$$

where if  $h_{j+1,j}v_{j+1} = 0$  the process is stopped. Later the coefficients  $h_{i,j}$  will be chosen to give different processes. Clearly after k - 1 steps  $\{v_j\}_{1}^{k}$  spans the k-th Krylov subspace

$$\mathcal{K}^{k}(A;v) \equiv \operatorname{span}\{v, Av, \dots, A^{k-1}v\}.$$
 (2)

After k steps (1) corresponds to the columns of

$$AV_{k} = V_{k}H_{k} + h_{k+1,k}v_{k+1}e_{k}^{T} = V_{k+1}H_{k+1,k}$$
(3)  

$$V_{k} \equiv \begin{bmatrix} v_{1} & v_{2} & \cdots & v_{k} \end{bmatrix},$$

$$H_{1,1} \quad h_{1,2} & \cdots & h_{1,k} \\ h_{2,1} \quad h_{2,2} & \cdots & h_{2,k} \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ h_{k,k-1} \quad h_{k,k} \end{bmatrix}, \quad H_{k+1,k} \equiv \begin{bmatrix} H_{k} \\ h_{k+1,k}e_{k}^{T} \end{bmatrix}.$$

Notice that  $H_{k+1,k}$  is just the upper Hessenberg  $H_k$  supplemented by a row of zeros except for the last element. If  $h_{k+1,k}v_{k+1} = 0$  then  $AV_k = V_kH_k$ , and

 $\mathcal{K}^k(A; v)$  is an invariant subspace of A, which we usually want. However the reverse need not hold for the general process, so the specific Krylov processes we consider will (in theory, that is in the absence of rounding errors) ensure termination, i.e.  $h_{k+1,k}v_{k+1} = 0$ , immediately the columns of  $V_k$  span an invariant subspace of A.

Householder [Hou74] gives a detailed treatment of Krylov's method [Kry31] and its relation to other algorithms for the algebraic eigenvalue problem.

### 2.2 The Arnoldi algorithm

One obvious choice to ensure termination is to derive orthogonal vectors  $v_1, v_2, \ldots$  in (1). Here and later we will use simple inductive arguments to obtain expressions for coefficients. Suppose we have nonsingular

$$V_k^T V_k = D_k = \operatorname{diag}(\rho_1, \dots, \rho_k), \tag{4}$$

which is trivially obtained for k = 1, then from (3)

$$V_k^T A V_k = D_k H_k + h_{k+1,k} V_k^T v_{k+1} e_k^T,$$

and to ensure  $V_k^T v_{k+1} = 0$  in theory we take

$$h_{i,k} = v_i^T A v_k / \rho_i, \quad i = 1, \dots, k.$$
<sup>(5)</sup>

If we were to compute the vectors via (5) and then (1), this would correspond to the numerically inadequate classical Gram-Schmidt approach, see for example [GLS9, § 5.2], so a more numerically reliable computation is via a modified Gram-Schmidt approach

$$u_k = Av_k$$
  
for  $i = 1, \dots, k$ 

 $h_{i,k} = v_i^T u_k / \rho_i$  $u_k = u_k - h_{i,k} v_i$ and then  $h_{k+1,k} v_{k+1} = u_k$ .

This gives (3) with

$$V_k^T V_k = D_k, \quad V_k^T v_{k+1} = 0, \quad V_k^T A V_k = D_k H_k.$$
 (6)

If  $h_{k+1,k}v_{k+1}$  is then nonzero we take  $h_{k+1,k}$  to be whatever we want. For orthonormal vectors we choose  $h_{k+1,k}$  to give  $||v_{k+1}|| = 1$ , and we can continue the process. So in this case (5) and (6) become

$$h_{i,k} = v_i^T A v_k, \quad i = 1, \dots k, \tag{7}$$

$$V_k^T V_k = I, \quad V_k^T v_{k+1} = 0, \quad V_k^T A V_k = H_k.$$
 (8)

This was proposed by Arnoldi for solving the eigenproblem: on termination all the eigenvalues of  $H_k$  are eigenvalues of A, see [Arn51].

### 2.3 The Lanczos process for symmetric A

If A is symmetric in the Arnoldi algorithm with arbitrary  $h_{j+1,j}$  leading to (6), we see  $V_k^T A V_k = D_k H_k = H_k^T D_k$  is tridiagonal and symmetric. That means  $h_{i,j} = 0$ , i = 1, 2, ..., j-2 in (1), and from the j-1, j element, see (4),

$$\rho_{j-1}h_{j-1,j} = \rho_j h_{j,j-1}.$$
 (9)

Let us define  $\alpha_j \equiv h_{j,j}$ ,  $\beta_j \equiv h_{j,j-1}$ ,  $\gamma_j \equiv h_{j-1,j}$ . This leads to Lanczos' 3-term process with arbitrary normalization

$$\gamma_0 v_0 \equiv 0$$
  
 $\beta_1 v_1 = v, \quad \beta_1 \neq 0$  arbitrary,  
for  $j = 1, \dots, k$ 

$$\alpha_j = v_j^T A v_j / \rho_j$$
  

$$\gamma_j = \beta_j \rho_j / \rho_{j-1}$$
(10)

$$\beta_{j+1}v_{j+1} = Av_j - \alpha_j v_j - \gamma_j v_{j-1}, \qquad (11)$$

$$eta_{j+1} 
eq 0$$
 arbitrary  
stop if  $eta_{j+1}v_{j+1} = 0.$ 

Notice the definition (10) of  $\gamma_j$  agrees with the original (5) both from symmetry (9) which says

$$\rho_{j-1}\gamma_j = \rho_j\beta_j,\tag{12}$$

or from the use of the symmetry of A and orthogonality of the  $v_i$  in Arnoldi's algorithm

$$\gamma_{j} \equiv h_{j-1,j}$$

$$= v_{j-1}^{T} A v_{j} / \rho_{j-1} = v_{j}^{T} A v_{j-1} / \rho_{j-1}$$

$$= v_{j}^{T} (\beta_{j} v_{j} + \alpha_{j-1} v_{j-1} + \gamma_{j-1} v_{j-2}) / \rho_{j-1}$$

$$= \beta_{j} \rho_{j} / \rho_{j-1}.$$
(13)

Note by using the matrix representation (6), the symmetry of A gave (12) immediately, while the usual vector approach required more effort. It was shown in [Pai72] that (10) is an efficient and more numerically reliable way of computing  $\gamma_j$  than (13).

In matrix form the result is, using tridiagonal  $T_k \equiv H_k$ ,

$$AV_{k} = V_{k}T_{k} + \beta_{k+1}v_{k+1}e_{k}^{T} = V_{k+1}T_{k+1,k}$$
(14)

$$V_k^T V_k = D_k, \quad V_k^T v_{k+1} = 0, \quad V_k^T A V_k = D_k T_k \equiv S_k, \quad \text{say,} \quad (15)$$
$$\begin{bmatrix} \alpha_1 & \gamma_2 \end{bmatrix}$$

$$T_{k} \equiv \begin{bmatrix} \alpha & \beta_{k} & \alpha_{k} \\ \beta_{2} & \alpha_{2} & \cdot \\ & \cdot & \cdot & \gamma_{k} \\ & & \beta_{k} & \alpha_{k} \end{bmatrix}, \quad T_{k+1,k} \equiv \begin{bmatrix} T_{k} \\ \beta_{k+1}e_{k}^{T} \end{bmatrix}, \quad (16)$$

$$S_{k} \equiv \begin{bmatrix} \rho_{1}\alpha_{1} & \rho_{1}\gamma_{2} & & \\ \rho_{2}\beta_{2} & \rho_{2}\alpha_{2} & \cdot & \\ & \cdot & \cdot & \rho_{k-1}\gamma_{k} \\ & & \rho_{k}\beta_{k} & \rho_{k}\alpha_{k} \end{bmatrix}, \qquad (17)$$

where obviously  $S_k \equiv V_k^T A V_k$  is symmetric by (12). Again, if we choose  $\beta_{j+1}$  to enforce orthonormality of the vectors, we notice  $\rho_j = 1$  for j = 1, 2, ... gives  $\gamma_j = \beta_j$ . The normalized symmetric Lanczos process is now

$$\beta_0 v_0 \equiv 0$$
  

$$\beta_1 v_1 = v, \quad v_1^T v_1 = 1,$$
  
for  $j = 1, \dots, k$   

$$\alpha_j = v_j^T A v_j$$
  

$$\beta_{j+1} v_{j+1} = A v_j - \alpha_j v_j - \beta_j v_{j-1},$$
  
stop if  $\beta_{j+1} v_{j+1} = 0$  else make  $v_{j+1}^T v_{j+1} = 1.$ 
(18)

And in matrix form this is

$$AV_{k} = V_{k}T_{k} + \beta_{k+1}v_{k+1}e_{k}^{T} = V_{k+1}T_{k+1,k}$$
(19)

$$V_{k}^{T}V_{k} = I, \quad V_{k}^{T}v_{k+1} = 0, \quad V_{k}^{T}AV_{k} = T_{k}, \quad (20)$$

$$T_{k} \equiv \begin{bmatrix} \alpha_{1} & \beta_{2} & & \\ \beta_{2} & \alpha_{2} & \cdot & \\ & \beta_{2} & \alpha_{2} & \cdot & \\ & & \beta_{k} & \alpha_{k} \end{bmatrix}, \quad T_{k+1,k} \equiv \begin{bmatrix} T_{k} \\ \beta_{k+1}e_{k}^{T} \end{bmatrix}.$$

The Lanczos process (which is closely related to the 3-term recurrence for orthogonal polynomials) was suggested by Lanczos [Lan50] for finding eigenvalues of symmetric A, and also used in Lanczos' [Lan52] variant of the Conjugate Gradients method (see also [HS52]) for solving Ax = b with symmetric positive definite A.

Since the difference between the above two versions is only one of normalization, they will have the same sort of numerical behaviour — which is surprisingly effective. Orthogonality of the  $v_j$  is soon lost, but eigenvalues are still found to almost machine precision by the Lanczos method (see for example [Pai71], and later works by him and many others), and solution of equations with symmetric positive definite A is also very effective using either the [HS52] CG algorithm (see Reid [Rei71] and many later works) or the [Lan52] variant, see [PS75].

#### 2.4 The Lanczos process for unsymmetric A

Clearly (18) is far more efficient than (1) with (7), and we would like a similar simplification with unsymmetric A. Lanczos [Lan50] suggested such a process. This produces two mutually orthogonal sequences of vectors,  $\{v_j\}_{1}^{k}$  and  $\{\tilde{v}_j\}_{1}^{k}$  spanning the subspaces  $\mathcal{K}^{k}(A; v)$  and  $\mathcal{K}^{k}(A^{T}; \tilde{v})$  respectively. These vectors will be seen to produce a tridiagonal matrix from unsymmetric A. This is achieved by considering the process (1) applied to  $A^{T}$  and  $\tilde{v}$ , as well as to A and v.

The resulting equivalents of the Krylov process (3) are, with obvious notation, after k steps

$$AV_{k} = V_{k}H_{k} + h_{k+1,k}v_{k+1}e_{k}^{T}, \qquad (21)$$

$$A^T \tilde{V}_k = \tilde{V}_k \widetilde{H}_k + \tilde{h}_{k+1,k} \tilde{v}_{k+1} e_k^T, \qquad (22)$$

2

but now we seek elements  $h_{i,j}$  and  $\tilde{h}_{i,j}$  which make  $\tilde{V}_{k+1}^T V_{k+1}$  diagonal. Again we assume we stop whenever we reach an invariant subspace. However we will see this process will actually break down at the first case of  $\tilde{v}_i^T v_i = 0$ , so we will assume v and  $\tilde{v}$  are (fortuitously) chosen to avoid this.

An inductive argument will show how this mutual orthogonality may be

obtained. Suppose after step k-1 we have diagonal and nonsingular

$$\tilde{V}_k^T V_k = D_k = \operatorname{diag}(\rho_1, \dots, \rho_k), \qquad (23)$$

(almost all choices of v and  $\tilde{v}$  will ensure this for k = 1), then with (21) and (22)

$$\tilde{V}_k^T A V_k = D_k H_k + h_{k+1,k} \tilde{V}_k^T v_{k+1} e_k^T$$
(24)

$$= \widetilde{H}_k^T D_k + e_k \widetilde{h}_{k+1,k} \widetilde{v}_{k+1}^T V_k.$$
(25)

Since  $H_k$  and  $\widetilde{H}_k$  are upper Hessenberg, equating the right hand sides shows that for j = 1, ..., k - 2, taking  $h_{j,k} = 0$  gives  $\tilde{v}_j^T v_{k+1} = 0$ , and taking  $\tilde{h}_{j,k} = 0$  gives  $v_j^T \tilde{v}_{k+1} = 0$ . Next the (k, k) elements of (24) and (25) show  $\tilde{v}_k^T v_{k+1} = \tilde{v}_{k+1}^T v_k = 0$  if we take

$$h_{k,k} = \tilde{h}_{k,k} = \tilde{v}_k^T A v_k / \rho_k \equiv \alpha_k, \quad \text{say.}$$
(26)

Next the (k, k-1) element of (25) shows  $\tilde{v}_{k+1}^T v_{k-1} = 0$  if we take

$$\tilde{h}_{k-1,k} = \tilde{v}_k^T A v_{k-1} / \rho_{k-1} \equiv \tilde{\gamma}_k, \quad \text{say}, \tag{27}$$

while the (k, k-1) element of (24) shows for the normalizer of  $v_k$  that

$$h_{k,k-1} = \tilde{v}_k^T A v_{k-1} / \rho_k \equiv \beta_k, \quad \text{say.}$$
(28)

Finally the (k-1,k) element of (24) shows  $\tilde{v}_{k-1}^T v_{k+1} = 0$  if we take

$$h_{k-1,k} = \tilde{v}_{k-1}^T A v_k / \rho_{k-1} \equiv \gamma_k, \quad \text{say}, \tag{29}$$

while the (k-1,k) element of (25) shows for the normalizer of  $\tilde{v}_k$  that

$$\tilde{h}_{k,k-1} = \tilde{z}_{k-1}^T A v_k / \rho_k \equiv \tilde{\beta}_k, \quad \text{say.}$$
(30)

These show that  $T_k \equiv H_k$  and  $\tilde{T}_k \equiv \tilde{H}_k$  are tridiagonal, and also show why nonzero  $\tilde{v}_i^T v_i$  is required to avoid breakdowns. If now  $\tilde{v}_{k+1}^T v_{k+1}$  is nonzero we

have  $\tilde{V}_{k+1}^T V_{k+1}$  diagonal and nonsingular, and we can continue the process. Note the process has not been fully defined, as (at each step) the normalizers  $\beta_{k+1} \equiv h_{k+1,k}$  of  $v_{k+1}$ , and  $\tilde{\beta}_{k+1} \equiv \tilde{h}_{k+1,k}$  of  $\tilde{v}_{k+1}$  may be chosen arbitrarily nonzero. Here we consider two possibilities.

By the standard unsymmetric Lanczos process we will mean the Lanczos process where  $\beta_j$  is arbitrarily chosen, for example to give  $v_j^T v_j = 1$ ,  $j = 1, 2, \ldots$ , but whatever the choice of  $\beta_j$ ,  $\tilde{\beta}_j$  must be chosen to give

$$\tilde{v}_j^T v_j = 1, \quad j = 1, 2, \dots$$

so that  $D_k = I$  and from (24) and (2<sup>F</sup>)  $T_k \equiv H_k = \widetilde{H}_k^T$ , giving (see also (27) and (28), and (29) and (30))

$$\tilde{\gamma}_j = \beta_j, \qquad \tilde{\beta}_j = \gamma_j, \qquad j = 1, 2, \dots,$$

and for (21) and (22)

$$V_{k} \equiv \begin{bmatrix} v_{1} & v_{2} & \cdots & v_{k} \end{bmatrix}, \quad \tilde{V}_{k} \equiv \begin{bmatrix} \tilde{v}_{1} & \tilde{v}_{2} & \cdots & \tilde{v}_{k} \end{bmatrix},$$
$$T_{k} \equiv \begin{bmatrix} \alpha_{1} & \gamma_{2} & & & \\ \beta_{2} & \alpha_{2} & \cdot & & \\ & \ddots & & \gamma_{k} & & \\ & & \beta_{k} & \alpha_{k} \end{bmatrix},$$
$$T_{k+1,k} \equiv \begin{bmatrix} T_{k} \\ \beta_{k+1}e_{k}^{T} \end{bmatrix}, \quad T_{k,k+1} \equiv \begin{bmatrix} T_{k} | \gamma_{k+1}e_{k} \end{bmatrix},$$

$$AV_{k} = V_{k}T_{k} + \beta_{k+1}v_{k+1}e_{k}^{T} = V_{k+1}T_{k+1,k}, \qquad (31)$$

$$A^{T}\tilde{V}_{k} = \tilde{V}_{k}T_{k}^{T} + \gamma_{k+1}\tilde{v}_{k+1}e_{k}^{T} = \tilde{V}_{k+1}T_{k,k+1}^{T}, \qquad \tilde{V}_{k}^{T}V_{k} = I, \quad \tilde{V}_{k}^{T}v_{k+1} = V_{k}^{T}\tilde{v}_{k+1} = 0, \qquad \tilde{V}_{k}^{T}AV_{k} = T_{k}, \quad V_{k}^{T}A^{T}\tilde{V}_{k} = T_{k}^{T}.$$

The detailed algorithm is then (where  $\beta_1, \ldots, \beta_{j+1}$  may be chosen arbitrarily)

$$v_{0} \equiv 0, \quad \tilde{v}_{0} \equiv 0$$

$$\beta_{1}v_{1} = v, \quad \gamma_{1}\tilde{v}_{1} = \tilde{v} \quad \text{so that} \quad \tilde{v}_{1}^{T}v_{1} = 1,$$
for  $j = 1, \dots, k$ 

$$\alpha_{j} = \tilde{v}_{j}^{T}Av_{j}$$

$$\beta_{j+1}v_{j+1} = Av_{j} - \alpha_{j}v_{j} - \gamma_{j}v_{j-1}$$

$$\gamma_{j+1}\tilde{v}_{j+1} = A^{T}\tilde{v}_{j} - \alpha_{j}\tilde{v}_{j} - \beta_{j}\tilde{v}_{j-1}$$
stop if  $\beta_{j+1}v_{j+1} = 0 \quad \text{or} \quad \gamma_{j+1}\tilde{v}_{j+1} = 0 \quad \text{or}$ 

$$\beta_{j+1}v_{j+1}^{T}\tilde{v}_{j+1}\gamma_{j+1} = 0,$$
otherwise choose  $\gamma_{j+1}$  to give  $\tilde{v}_{j+1}^{T}v_{j+1} = 1.$  (32)

By the alternative unsymmetric Lanczos process we will mean the (apparently new variant of the) Lanczos process where  $\beta_j$  is arbitrarily chosen, but then we take

$$\tilde{\beta}_j = \beta_j, \quad j = 1, 2, \dots \tag{33}$$

With this choice it follows from (28) and (30) that

$$\tilde{v}_{k-1}^T A v_k = \tilde{v}_k^T A v_{k-1}, \tag{34}$$

so from (27) and (29)

$$\tilde{\gamma}_k = \tilde{v}_{k-1}^T A v_k / \rho_{k-1} = \gamma_k, \tag{35}$$

giving

$$T_{k} \equiv H_{k} = \widetilde{H}_{k} = \begin{bmatrix} \alpha_{1} & \gamma_{2} & & \\ \beta_{2} & \alpha_{2} & \cdot & \\ & \ddots & \ddots & \gamma_{k} \\ & & & \beta_{k} & \alpha_{k} \end{bmatrix}$$
(36)

۰.

۰.

$$T_{k+1,k} = \begin{bmatrix} T_k \\ \beta_{k+1}e_k^T \end{bmatrix} \cdot \cdot \\ AV_k = V_k T_k + \beta_{k+1}v_{k+1}e_k^T = V_{k+1}T_{k+1,k}, \quad (37)$$

$$A^{T}\tilde{V}_{k} = \tilde{V}_{k}T_{k} + \beta_{k+1}\tilde{v}_{k+1}e_{k}^{T} = \tilde{V}_{k+1}T_{k+1,k}, \qquad (38)$$

$$\tilde{V}_{k}^{T}V_{k} = D_{k} = \operatorname{diag}(\rho_{1}, \dots, \rho_{k}).$$

$$\tilde{V}_{k}^{T}AV_{k} = D_{k}T_{k} \equiv S_{k} \text{ is symmetric.}$$
(39)

The surprising outcome is that the two Krylov processes (with A and  $A^{T}$ ) use the coefficients in the same way. At the *j*-th step we could take, see (26), and (34) and (35),

$$\rho_{j} = \tilde{v}_{j}^{T} v_{j}$$

$$\alpha_{i} = \tilde{v}_{i}^{T} A v_{i} / \rho_{i} \qquad (40)$$

$$\gamma_j = \tilde{v}_{j-1}^T A v_j / \rho_{j-1} \tag{41}$$

$$\beta_{j+1}v_{j+1} = Av_j - \alpha_j v_j - \gamma_j v_{j-1}$$

$$\tag{42}$$

$$\beta_{j+1}\tilde{v}_{j+1} = A^T\tilde{v}_j - \alpha_j\tilde{v}_j - \gamma_j\tilde{v}_{j-1}$$
(43)

 $\beta_{j+1} \neq 0$  arbitrary.

Again, as with (5) in our first approach to the Arnoldi algorithm, and our comments on the Lanczos process for symmetric A, we find that the above computation is unnecessarily numerically inaccurate. Thus we use the following alternative to (41).

$$\gamma_j = \beta_j \rho_j / \rho_{j-1}$$
 from (29), (30) and (33), (44)

which is identical to the choice (10) recommended for the symmetric case.

This alternative unsymmetric Lanczos process is actually just a diagonally scaled version of the standard unsymmetric Lanczos process. This is clear since both processes produce mutually orthogonal vectors spanning the same subspaces — the only difference is the first process produces  $\tilde{V}_k^T V_k = I$ , the second one does not. Although the standard process is usually considered the underlying process in the BiCG method, we will use the alternative process to produce a pleasingly simple and clear derivation of BiCG in § 4.6.

Apart from this surprising new variant of the unsymmetric Lanczos proccss, so far the only original work here has been the matrix based derivation of these algorithms. This is much briefer and, we think, easier to follow than the usual approach, see for example [GL89, § 10.2], but the originality is in the approach alone. However we can now obtain a minor insight into the computational behaviour of these two variants of the unsymmetric Lanczos process.

If A is symmetric, the unsymmetric Lanczos process becomes the symmetric Lanczos process (taking  $\tilde{v} = v$  of course, which appears to be the only efficient choice – halving the work). In this case the main step of the unsymmetric Lanczos process with orthonormal vectors (32) becomes

$$\alpha_j = v_j^T A v_j$$
  

$$\beta_{j+1} v_{j+1} = A v_j - \alpha_j v_j - \beta_j v_{j-1}$$
  
since  $\gamma_j = \beta_j$ ,

which is just the normalized symmetric Lanczos process (18).

Similarly the alternative unsymmetric Lanczos process becomes

$$\rho_{j} = v_{j}^{T} v_{j}$$

$$\alpha_{j} = v_{j}^{T} A v_{j} / \rho_{j}$$

$$\gamma_{j} = \beta_{j} \rho_{j} / \rho_{j-1}$$

$$\beta_{j+1} v_{j+1} = A v_{j} - \alpha_{j} - \gamma_{j} v_{j-1}$$

$$\beta_{j+1} \neq 0 \text{ arbitrary,}$$

which is the same as the arbitrarily normalized Lanczos process for sym-

metric A. Thus we can hope that these two unsymmetric processes will maintain some of the well-known good numerical properties of their equivalent symmetric processes.

.

. 4

# 3 Preliminaries – The QR decomposition

Some of the algorithms presented in this paper, for example LSQR, GM-RES(m) and QMR, require the reduction of a sequence of matrices, each of which contains the previous as leading submatrix, to upper triangular form. To this end, one may perform a sequence of QR decompositions. The notation and method presented in this section will be used in those algorithms for this purpose. For the decomposition one may choose from  $2 \times 2$  rotation matrices,  $2 \times 2$  reflection matrices ( $2 \times 2$  Householder transformations), or fast Givens transformations. For simplicity,  $2 \times 2$  rotation matrices were used here. The matrices arising in these algorithms are of a specific shape – in general upper Hessenberg, more specifically tridiagonal or lower bidiagonal, with an additional row containing a nonzero element in the last column only. The upper triangular result may in these special cases be upper tridiagonal or upper bidiagonal. For the general matrix

$$H_{k+1,k} \equiv \begin{bmatrix} h_{11} & h_{12} & \cdots & h_{1k} \\ h_{21} & h_{22} & \cdots & h_{2k} \\ & \ddots & & \ddots \\ & & & h_{k,k-1} & \ddots & h_{kk} \\ & & & & & h_{k+1,k} \end{bmatrix} \in \mathbb{R}^{(k+1) \times k}$$

 $(h_{ij} \text{ may be zero for } j > i + 1 \text{ or } j > i)$  the decomposition is

$$Q_k^T H_{k+1,k} = \begin{bmatrix} R_k \\ 0 \end{bmatrix},$$
  
where  $Q_k \in R^{(k+1)\times(k+1)}$  is orthogonal,  
and  $R_k \in R^{k\times k}$  is upper triangular.

22

This is achieved using rotations embedded at the i, i position of k + 1dimensional identity matrices, as

$$Q_{i,k} \equiv \begin{bmatrix} I_{i-1} & & & \\ & C_i & s_i \\ & -s_i & c_i \end{bmatrix}$$
  
such that  $Q_k^T \equiv Q_{k,k}^T \cdots Q_{1,k}^T$   
 $\equiv Q_{k,k}^T \begin{bmatrix} Q_{k-1}^T & \\ & 1 \end{bmatrix}.$ 

The effect of the appropriately embedded  $Q_{k-1}^T$  on  $H_{k+1,k}$  is

$$\begin{bmatrix} Q_{k-1}^T \\ & 1 \end{bmatrix} \begin{bmatrix} H_{k,k-1} & \vdots \\ & h_{kk} \\ \hline & 0 & h_{k+1,k} \end{bmatrix} = \begin{bmatrix} h_{1k} \\ R_{k-1} & \vdots \\ & r_{k-1,k} \\ \hline & 0 & \bar{r}_{kk} \\ \hline & 0 & h_{k+1,k} \end{bmatrix},$$

while  $Q_{k,k}^T$ , which alters only rows k and k + 1, zeros  $h_{k+1,k}$ 

$$\begin{bmatrix} c_k & -s_k \\ s_k & c_k \end{bmatrix} \begin{bmatrix} \tilde{r}_{kk} \\ h_{k+1,k} \end{bmatrix} = \begin{bmatrix} r_{kk} \\ 0 \end{bmatrix},$$
$$\Rightarrow r_{kk} = \sqrt{\tilde{r}_{kk}^2 + h_{k+1,k}^2}$$
$$c_k = \tilde{r}_{kk}/r_{kk}$$
$$s_k = -h_{k+1,k}/r_{kk}.$$

The orthogonality of the  $Q_{i,k}$  leads to that of  $Q_k$ .

A scalar multiple of the first column of  $Q_k^T$  is given by

$$g_{k} \equiv \beta Q_{k}^{T} e_{1} \equiv \beta Q_{k,k}^{T} \cdots Q_{1,k}^{T} e_{1} \equiv \begin{bmatrix} f_{k} \\ \bar{\phi}_{k+1} \end{bmatrix} \equiv \begin{bmatrix} \phi_{1} \\ \vdots \\ \phi_{k} \\ \bar{\phi}_{k+1} \end{bmatrix}, \quad (45)$$

23

where

$$f_{k} = \begin{bmatrix} f_{k-1} \\ \phi_{k} \end{bmatrix},$$

$$f_{1} = \phi_{1} = \beta c_{1}$$

$$\bar{\phi}_{2} = \beta s_{1}$$

$$\phi_{k} = c_{k} \bar{\phi}_{k} = \beta c_{k} s_{k-1} \cdots s_{1}$$

$$\bar{\phi}_{k+1} = s_{k} \bar{\phi}_{k} = \beta s_{k} \cdots s_{1}.$$
(46)

<u>Note</u>: This notation  $-Q_k, Q_{i,k}, g_k, f_k, \overline{\phi}_k, \phi_k, c_k, s_k$  – is used throughout this document with exactly these meanings, except in algorithms SYMMLQ and MINRES where the  $Q_k, c_k, s_k$  refer to the LQ decomposition of the matrices involved.

Reference: for example [GL89, § 5.1 5.2]

## 4 Algorithms for Solving Linear Systems

We now consider most of the more common Krylov subspace methods for solving linear systems of equations. Each of the following algorithms has been converted into the matrix formalism presented in § 1.1. Each subsection is begun with the <u>problem</u> statement, either large systems of linear equations, or the linear least squares problem, that is

$$Ax = b$$
,  $A \in \mathbb{R}^{n \times n}$  or  $\min_{x \in \mathbb{R}^n} ||b - Ax||$ ,  $A \in \mathbb{R}^{m \times n}$ .

The type of matrix A, for example symmetric or not, positive definite or indefinite, is given.

We examine different methods, where sometimes these are just different variants of one method. We continue to use the word method rather than implementation or algorithm, as there could also be more than one implementation for a given variant. For each such method the process used to form vectors  $\{v_1, \ldots, v_k\}$  spanning some associated subspace is named, as is the subspace.

Next for each method the <u>subproblem</u> arises as a result of considering successive solutions of the form

$$x_k = x_0 + V_k y_k$$

for some initial guess  $x_0$ . For simplicity,  $x_0$  will usually<sup>1</sup> be taken as zero. So the iterates look like

$$x_{k} = V_{k}y_{k} \in \mathcal{R}(V_{k}), \qquad (47)$$
  
where  $V_{k} \equiv \begin{bmatrix} v_{1} & v_{2} & \cdots & v_{k} \end{bmatrix},$ 

and the way in which this subproblem arises and is solved is described.

<sup>&</sup>lt;sup>1</sup>Except in GMRES(m) where with each restart a new nonzero  $x_0$  is used.

For some methods, where it is important for the development here, we will describe the basic <u>implementation</u> of the method, showing how the vectors from the <u>process</u> can be combined via the solution of the <u>subproblem</u> to produce an efficient (although not yet optimized) algorithm.

Then, to facilitate the most likely implementation of a stopping condition for each algorithm, the norm of the *residual* of the current iterate

$$r_k \equiv b - Ax_k = b - AV_k y_k \tag{48}$$

is considered.

Finally (when we can, and when it is appropriate) we relate this matrix derivation to one or more of the standard vector implementations used today, thus closing the circle by showing how what are usually confusing (to the non-expert) vector presentations can be easily derived and wellmotivated using the matrix approach recommended here.

7

### 4.1 The Lanczos (Conjugate Gradients) Method

<u>Problem</u>: Solve Ax = b for x when A is symmetric positive definite (SPD).

First we derive what will later be shown to be the Conjugate Gradients approximations  $x_k$  from the normalized symmetric Lanczos process (18). We refer to this as:

The unit norm Lanczos-CG method.

<u>Process</u>: the normalized Lanczos process (18) for symmetric A to produce orthonormal  $\{v_1, \ldots, v_k\}$  spanning  $\mathcal{K}^k(A; b)$ , starting with  $v_1 = b/\beta_1$ , where  $\beta_1 \equiv ||b||$ .

To derive the subproblem we use (19) and consider iterates of the form  $x_k = V_k y_k$  as suggested in (47), and take  $T_k y_k = \beta_1 e_1$ , so that

$$r_{k} = b - Ax_{k} = b - AV_{k}y_{k}$$
  
=  $b - V_{k}T_{k}y_{k} - \beta_{k+1}v_{k+1}e_{k}^{T}y_{k}$   
=  $V_{k}(\beta_{1}e_{1} - T_{k}y_{k}) - \beta_{k+1}v_{k+1}e_{k}^{T}y_{k}$   
=  $-\beta_{k+1}v_{k+1}e_{k}^{T}y_{k}.$  (49)

In theory  $\beta_{k+1}v_{k+1} = 0$  for some  $k \leq n$ , the dimension of A, and the problem will be solved.

Subproblem: Solve  $T_k y_k = \beta_1 c_1$ , where  $T_k$  is symmetric tridiagonal.

This has shown in simple matrix terms how the Lanczos process (19)  $AV_k = V_kT_k + \beta_{k+1}v_{k+1}e_k^T$  and the subproblem  $T_ky_k = \beta_1e_1$  give approximations  $x_k = V_ky_k$  to the solution of Ax = b. However we now have to combine these two (Process and Subproblem) to produce a useful *implementation*. <u>Implementation</u>: For A SPD,  $T_k = V_k^T AV_k$  is as well, therefore the subproblem always has a solution and it is possible to perform a Cholesky factorisation of  $T_k = L_k L_k^T$ , where  $L_k$  is lower bidiagonal. In order to obtain objects  $\zeta_k$ ,  $w_k$  that can be computed sequentially and discarded after

5

use — an important efficiency for large sparse problems — define

$$z_{k} \equiv \begin{bmatrix} \zeta_{1} \\ \zeta_{2} \\ \vdots \\ \zeta_{k} \end{bmatrix} \equiv L_{k}^{T} y_{k}, \quad W_{k} \equiv \begin{bmatrix} w_{1} & w_{2} & \cdots & w_{k} \end{bmatrix} \equiv V_{k} L_{k}^{-T},$$

so the system is transformed to

$$T_k y_k = L_k L_k^T y_k = L_k z_k \implies L_k z_k = \beta_1 e_1,$$
$$V_k y_k = V_k L_k^{-T} L_k^T y_k \implies x_k = W_k z_k.$$

.

Now,  $L_k$  is formed one row at a time:

$$T_{k} = \begin{bmatrix} L_{k-1} & 0 \\ 0 & 0 \\ 0 & \delta_{k} & \gamma_{k} \end{bmatrix} \begin{bmatrix} L_{k-1}^{T} & 0 \\ \delta_{k} \\ 0 & 0 & \gamma_{k} \end{bmatrix}$$
$$= \begin{bmatrix} L_{k-1}L_{k-1}^{T} & 0 \\ 0 & \delta_{k}\gamma_{k-1} & \delta_{k}^{2} + \gamma_{k}^{2} \end{bmatrix} = \begin{bmatrix} T_{k-1} & 0 \\ 0 & \beta_{k} \\ 0 & \beta_{k} & \alpha_{k} \end{bmatrix}$$
$$\Rightarrow \gamma_{1} = \sqrt{\alpha_{1}},$$
$$\delta_{k} = \beta_{k}/\gamma_{k-1}, \quad \gamma_{k} = \sqrt{\alpha_{k} - \delta_{k}^{2}}, \quad k > 1.$$
(50)

And so  $z_k$  is formed one entry at a time:

$$\begin{bmatrix} L_{k-1} & 0 \\ 0 & 0 \\ 0 & \delta_k & \gamma_k \end{bmatrix} \begin{bmatrix} z_{k-1} \\ \zeta_k \end{bmatrix} = \begin{bmatrix} L_{k-1}z_{k-1} \\ \delta_k\zeta_{k-1} + \gamma_k\zeta_k \end{bmatrix} = \beta_1 e_1$$
$$\Rightarrow \zeta_1 = \beta_1/\gamma_1, \quad \zeta_k = -\delta_k\zeta_{k-1}/\gamma_k, \quad k > 1.$$
(51)

 $W_k$  is also formed one column at a time  $(W_k L_k^T = V_k)$ :

$$\begin{bmatrix} W_{k-1} & w_k \end{bmatrix} \begin{bmatrix} L_{k-1}^T & 0 \\ & \delta_k \\ \hline 0 & 0 & \gamma_k \end{bmatrix} = \begin{bmatrix} W_{k-1}L_{k-1}^T & \delta_k w_{k-1} + \gamma_k w_k \end{bmatrix}$$

$$\Rightarrow w_1 = v_1/\gamma_1, \quad w_k = (v_k - \delta_k w_{k-1})/\gamma_k, \quad k > 1.$$

Finally,

$$x_{k} = W_{k}z_{k} = \begin{bmatrix} W_{k-1} & w_{k} \end{bmatrix} \begin{bmatrix} z_{k-1} \\ \zeta_{k} \end{bmatrix}$$
$$= W_{k-1}z_{k-1} + \zeta_{k}w_{k} = x_{k-1} + \zeta_{k}w_{k}, \quad k \ge 1.$$

Again we see how each step of the implementation makes obvious sense when approached via this matrix formulation.

For large sparse problems we would like to stop in  $k \ll n$  steps, and since a solution with sufficiently small residual is what is often wanted, we look at how the residual norm behaves here.

<u>Residual Norm</u>: The following estimation of the norm of the residual  $r_k$ , as defined in (49) can be used.

$$r_{k} = -\beta_{k+1}v_{k+1}e_{k}^{T}y_{k} = -\beta_{k+1}v_{k+1}y_{kk}$$
$$||r_{k}|| = \beta_{k+1}|y_{kk}|, \text{ where } y_{kk} \text{ is the } k\text{-th entry of } y_{k}.$$

Since  $y_k$  is not computed directly, one must consider

$$L_{k}^{T}y_{k} = z_{k}$$

$$\begin{bmatrix} L_{k-1}^{T} & 0 \\ & \delta_{k} \\ \hline 0 & 0 & \gamma_{k} \end{bmatrix} \begin{bmatrix} y_{k1} \\ \vdots \\ y_{kk} \end{bmatrix} = \begin{bmatrix} \vdots \\ \gamma_{k}y_{kk} \end{bmatrix} = \begin{bmatrix} z_{k-1} \\ \zeta_{k} \end{bmatrix}$$

$$y_{kk} = \zeta_{k}/\gamma_{k}$$

$$\Rightarrow ||r_{k}|| = \beta_{k+1}|\zeta_{k}/\gamma_{k}|$$

$$= \beta_{k+1}|\delta_{k}\zeta_{k-1}/\gamma_{k}^{2}| \text{ by (51)}$$

$$= \beta_{k+1}|\beta_{k}\zeta_{k-1}/\gamma_{k-1}|/\gamma_{k}^{2} \text{ by (50)}$$

$$= \beta_{k+1}||r_{k-1}||/\gamma_{k}^{2}, \quad k \ge 1.$$

This matrix approach also tends to make analysis of such algorithms fairly easy. For example, since  $r_j$  is a scalar multiple of  $v_{j+1}$  and  $V_k^T V_k = I$ , the residuals  $r_j$  are orthogonal. Next

$$W_{k}^{T}AW_{k} = L_{k}^{-1}V_{k}^{T}AV_{k}L_{k}^{-T} = L_{k}^{-1}T_{k}L_{k}^{-T} = I,$$

so the  $w_j$  (the search vectors in moving from  $x_{j-1}$  to  $x_j$ ) are A-orthogonal.

That  $x_k$  is the same as the ordinary CG solution follows for example from [PPdV95, § 3], since  $x_k = V_k y_k \in \mathcal{K}^k(A; b)$  and  $r_k \equiv b - A x_k \perp \mathcal{K}^k(A; b)$ , and so  $x_k$  is the Galerkin solution from  $\mathcal{K}^k(A; b)$  to Ax = b, which is exactly what CG delivers.

This unit norm Lanczos-CG method, based upon the Lanczos process, has the nice property that the vectors  $v_1, v_2, \ldots$  are orthonormal. However, let us derive a method which more closely mimics the standard CG method in that  $v_1, v_2, \ldots$  are the residual vectors themselves. Because it gives a more direct analogy with the standard CG method, we call it:

#### The Lanczos-CG method

<u>Process</u>: the Lanczos process (11) for symmetric A, starting with  $v_1 = b$ . The normalization  $\beta_{j+1}$  is unspecified as yet.

Using (14) and considering iterates of the form  $x_k = V_k y_k$  as suggested in (47) and solving for  $T_k y_k = e_1$ ,

$$\tau_{k} = b - Ax_{k} = b - AV_{k}y_{k}$$
  
=  $V_{k}(e_{1} - T_{k}y_{k}) - \beta_{k+1}v_{k+1}e_{k}^{T}y_{k}$   
=  $-\beta_{k+1}y_{kk}v_{k+1}$ . (52)

Subproblem: Solve  $T_k y_k = e_1$ , where  $T_k$  is tridiagonal, but not necessarily symmetric.

Implementation: We choose to perform an LDU decomposition of  $T_k = D_k^{-1}S_k$ , see (15).  $S_k = V_k^T A V_k$  is SPD, and  $D_k$  is a diagonal matrix with

positive elements, so the LDU decomposition without pivoting appears to be a safe choice. There is considerable freedom in such a decomposition, and the one which suits our purpose is as follows. With the as yet unknown  $\mathcal{D}_k \equiv -\text{diag}(\beta_2, \ldots, \beta_{k+1})$ , write

$$T_{k} = L_{k}\mathcal{D}_{k}U_{k}$$

$$T_{k} = \left[\begin{array}{c|c} L_{k-1} & 0 \\ \hline 0 & -1 & 1 \end{array}\right] \left[\begin{array}{c|c} \mathcal{D}_{k-1} & 0 \\ \hline 0 & 0 & -\beta_{k+1} \end{array}\right] \left[\begin{array}{c|c} U_{k-1} & 0 \\ \hline -\frac{\gamma_{k}}{\beta_{k}} \\ \hline 0 & 0 & 1 \end{array}\right]$$

$$= \left[\begin{array}{c|c} L_{k-1}\mathcal{D}_{k-1}U_{k-1} & 0 \\ \hline L_{k-1}\mathcal{D}_{k-1}U_{k-1} & \gamma_{k} \\ \hline 0 & \beta_{k} & -(\gamma_{k}+\beta_{k+1}) \end{array}\right].$$

Note that a simpler method could be obtained by combining  $\mathcal{D}$  and U, and using for example  $\mathcal{U} \equiv \mathcal{D}U$  instead, but the approach here will produce the equivalents of various elements of CG that we want.

We can now determine our normalizers  $\beta_{j+1}$  by equating the right side of the above equation with  $T_k$  in (16). Clearly

$$\alpha_1 = -\beta_2; \qquad \alpha_k = -(\gamma_k + \beta_{k+1}), \quad k > 1,$$
  
$$\Rightarrow \beta_2 \equiv -\alpha_1; \qquad \beta_{k+1} \equiv -(\alpha_k + \gamma_k), \quad k > 1.$$
(53)

We will show below that this  $\beta_{j+1} \neq 0$ , and so is a satisfactory normalizer for the process, and  $U_{k+1}$  exists, but first we complete this implementation. Define

$$z_{k} \equiv \begin{bmatrix} \zeta_{1} \\ \zeta_{2} \\ \vdots \\ \zeta_{k} \end{bmatrix} \equiv U_{k}y_{k}, \quad W_{k} \equiv \begin{bmatrix} w_{1} & w_{2} & \cdots & w_{k} \end{bmatrix} \equiv V_{k}U_{k}^{-1}, \quad (54)$$

so  $L_k \mathcal{D}_k z_k = c_1$  is our subproblem, and  $x_k = W_k z_k$  is our solution. Solving

$$L_k \mathcal{D}_k z_k = \begin{bmatrix} -\beta_2 \zeta_1 \\ \beta_2 \zeta_1 - \beta_3 \zeta_2 \\ \vdots \\ \beta_k \zeta_{k-1} - \beta_{k+1} \zeta_k \end{bmatrix} = e_1$$
  
$$\Rightarrow \zeta_1 = -\frac{1}{\beta_2}$$
  
$$\zeta_k = \frac{\beta_k}{\beta_{k+1}} \zeta_{k-1} = \frac{\beta_{k-1}}{\beta_{k+1}} \zeta_{k-2} = \dots = -\frac{1}{\beta_{k+1}}, \quad k > 1. \quad (55)$$

Solving  $W_k U_k = V_k$  gives

$$w_1 = v_1, \quad w_k = v_k + \frac{\gamma_k}{\beta_k} w_{k-1}, \quad k > 1.$$

Finally

$$x_k = W_{k-1} z_{k-1} + \zeta_k w_k = x_{k-1} + \zeta_k w_k, \quad k \ge 1.$$

Residual Norm:

$$r_{k} = -\beta_{k+1}y_{kk}v_{k+1} \quad \text{by (52)}$$

$$= -\beta_{k+1}\zeta_{k}v_{k+1} \quad \text{by (54)}$$

$$= -\beta_{k+1} \times -\frac{1}{\beta_{k+1}}v_{k+1} \quad \text{by (55)}$$

$$= v_{k+1}^{l} \qquad (56)$$

$$\Rightarrow ||r_{k}|| = ||v_{k+1}||.$$

So the Lanczos process (11) for symmetric A, with  $v_1 = b$ , can be used to produce the (CG) residual vectors  $\tau_k = b - Ax_k$  themselves, as long as  $\beta_{j+1}$  is chosen as in (53).

It remains for us to prove that no  $\beta_{j+1}$  found this way can be zero. Since A is SPD,  $\alpha_1$  is nonzero and so  $\beta_2$  is nonzero. Suppose  $\beta_2, \ldots, \beta_j$  obtained from (53) are nonzero. Then with  $e^T \equiv [1, \ldots, 1], e^T T_j = [0, \ldots, 0, \alpha_j + \gamma_j]$ , which must be nonzero since  $T_j$  is nonsingular. Thus no  $\beta_{j+1}$  will be zero, and this algorithm will not break down if A is SPD. It is also clear that the  $L_k \mathcal{U}_k$  'factorization' is numerically well behaved, with the only computation being  $\beta_{j+1} = -(\alpha_j + \gamma_j)$ , and there is no large element growth in  $\mathcal{U}_k$ . Equivalence of the Lanczos-CG method with CG:

This Lanczos-CG method is just a scaling of the unit norm Lanczos-CG method, and provides the same solution in a marginally different way. It is therefore also the CG solution. However we would like to compare the resulting algorithms closely, so we give a proof which shows clearly what equivalences occur.

The CG Method — slightly altered from 'Templates' [BBC+94] — is as follows.

$$r_{0} = b$$
  
for  $i = 1, 2, ...$   
 $\rho_{i-1} = r_{i-1}^{T} r_{i-1}$   
if  $(i = 1)$   
 $p_{i} = r_{i-1}$   
else  
 $\beta_{i-1} = \rho_{i-1}/\rho_{i-2}$   
 $p_{i} = r_{i-1} + \beta_{i-1}p_{i-1}$   
 $q_{i} = Ap_{i}$   
 $\alpha_{i} = \rho_{i-1}/p_{i}^{T}q_{i}$   
 $x_{i} = x_{i-1} + \alpha_{i}p_{i}$   
 $r_{i} = r_{i-1} - \alpha_{i}q_{i}$ 

check convergence; continue if necessary

end

To avoid confusion we will use superscript  $^{C}$  to distinguish 'Templates' terms from those here, where it will also make life easier if we re-index two

of the 'Templates' terms as follows:

$$\rho_i^C \equiv \text{`Templates' } \rho_{i-1},$$
  
$$\beta_i^C \equiv \text{`Templates' } \beta_{i-1}.$$

We now write the algorithms side by side for comparison, where in the Lanczos-CG algorithm we replace  $v_{i+1}$  by  $r_i$ , see (56), and introduce  $u_i \equiv Av_i$ . To handle the i = 1 case we define

CGLanczos-CG
$$\rho_0^C = 1$$
 $\rho_0 = 1, \quad \beta_1 = 0$  $x_0^C = 0, \quad p_0^C = 0$  $x_0 = 0, \quad w_0 = 0, \quad r_{-1} = 0$ 

Then for  $i = 1, 2, 3, \ldots$  we have

CG Lanczos-CG  

$$\rho_i^C = r_{i-1}^{CT} r_{i-1}^C \qquad \rho_i = r_{i-1}^T r_{i-1}$$

$$\beta_i^C = \rho_i^C / \rho_{i-1}^C \qquad \gamma_i = \beta_i \rho_i / \rho_{i-1}$$

$$p_i^C = r_{i-1}^C + \beta_i^C p_{i-1}^C \qquad w_i = r_{i-1} + (\gamma_i / \beta_i) w_{i-1}$$

$$q_i^C = A p_i^C \qquad u_i = A r_{i-1}$$

$$\alpha_i^C = \rho_i^C / p_i^{CT} q_i^C \qquad \alpha_i = r_{i-1}^T u_i / \rho_i$$

$$\beta_{i+1} = -(\alpha_i + \gamma_i)$$

$$x_i^C = x_{i-1}^C + \alpha_i^C p_i^C \qquad x_i = x_{i-1} - (1/\beta_{i+1}) w_i$$

$$r_i^C = r_{i-1}^C - \alpha_i^C q_i^C \qquad r_i = (u_i - \alpha_i r_{i-1} - \gamma_i r_{i-2}) / \beta_{i+1}$$

The only significant difference in cost occurs in the last line — n extra multiplications, subtractions and divisions in the Lanczos-CG algorithm. However note the Lanczos-CG here was written for comparison, not efficiency. The equivalences are

$$\rho_i = \rho_i^C, \quad \gamma_i / \beta_i = \beta_i^C, \quad w_i = p_i^C$$
$$-1\beta_{i+1} = \alpha_i^C, \quad x_i = x_i^C, \quad r_i = r_i^C.$$
(57)
A straightforward examination shows this is true for i = 1. Suppose it is true for i = 1, 2, ..., k - 1. Then it is clear that

$$\rho_k = \rho_k^C, \quad \gamma_k / \beta_k = \beta_k^C, \quad w_k = p_k^C.$$

We now need only prove

$$-1/\beta_{k+1} = \alpha_k^C, \tag{58}$$

so that  $x_k = x_k^C$ , and their residuals  $r_k$  and  $r_k^C$  must then also be equal, since  $r_i^C$  is clearly the residual for  $x_i^C$  by examining the last two lines of CG, and we have shown  $r_i$  is the residual of  $x_i$  in the Lanczos-CG. We can prove (58) using our matrix formulation of the Lanczos-CG algorithm. We know  $W_k = V_k U_k^{-1}$ , so

$$W_{k}^{T}AW_{k} = U_{k}^{-T}V_{k}^{T}AV_{k}U_{k}^{-1} = U_{k}^{-T}D_{k}T_{k}U_{k}^{-1}$$
$$= U_{k}^{-T}D_{k}(L_{k}\mathcal{D}_{k}U_{k})U_{k}^{-1} = U_{k}^{-T}D_{k}L_{k}\mathcal{D}_{k},$$

which, by the way, is lower triangular and symmetric, and therefore diagonal, showing A-orthogonality again. The (i, i)-element of this is

$$w_i^T A w_i = e_i^T U_k^{-T} (\rho_i e_i - \rho_{i+1} e_{i+1}) (-\beta_{i+1}).$$

But  $e_i^T U_k^{-T} = (x, \dots, x, 1, 0, \dots, 0)$  is zero after the *i*-th element, so since  $p_i^C = w_i, i = 1, \dots, k$ ,

$$p_i^{CT}Ap_i^C = w_i^TAw_i = -\rho_i\beta_{i+1}, \qquad i = 1, \dots, k.$$

From the CG algorithm we then see

$$\alpha_{i}^{C} = \rho_{i}^{C} / p_{i}^{CT} q_{i}^{C} = \rho_{i} / p_{i}^{CT} A p_{i}^{C} = -1 / \beta_{i+1}, \quad i = 1, \dots, k$$

and so (58) holds, and the induction proof of (57) is complete, showing that in theory these algorithms not only compute the same iterates (as was obvious earlier), but also several of the intermediate quantities are identical. However the algorithms are clearly not identical, as the last line of each clearly highlights, and may have different numerical behaviours. This suggests a matrix development of CG would be useful, but we leave this for future work.



36

Apple
 Apple

### 4.2 Symmetric LQ Method (SYMMLQ)

<u>Problem</u>: Solve Ax = b with A nonsingular symmetric, but indefinite. <u>Process</u>: the normalized Lanczos process (18) for symmetric A on  $\mathcal{K}^k(A; b)$ starting with  $v_1 = b/||b||$ .

In this section iterates of the form  $x_k^C = V_k y_k^C$  are produced as well as iterates  $x_k^L = V_{k+1} y_k^L$ . (The superscript C denotes the solution that would be obtained by the method of Conjugate Gradients, while the superscript L refers to the SYMMLQ method.) Using (19) and solving  $T_k y_k^C = \beta_1 e_1$ ,  $\beta_1 \equiv ||b||$ , one obtains the result (49) as for the Lanczos method.

Subproblem (Conjugate Gradients method):  $T_k y_k^C = \beta_1 e_1$ , where the matrix  $T_k$  is symmetric tridiagonal.  $T_k$  may be singular for some k, in which case this subproblem will not have a (unique) solution, see later, (p. 38, 39). Implementation: Since A is indefinite, there is a possibility that  $T_k$  is as well. Therefore a Cholesky factorisation may break down. Instead, consider an LQ decomposition.

$$T_k \bar{Q}_k^T = \bar{L}_k, \quad \bar{Q}_k^T \bar{Q}_k = I,$$

where one method of factorising has the  $\bar{Q}_k^T = Q_{1,k} \cdots Q_{k-1,k}$  as orthogonal products of (reflection) matrices embedded in the k-dimensional identity starting in the (i, i) entry.

$$Q_{i,k} \equiv \begin{bmatrix} I_{i-1} & & & \\ & C_i & S_i & \\ & S_i & -C_i & \\ \hline & & & I_{k-i-1} \end{bmatrix} \in \mathcal{R}^{k \times k}$$
(59)

These are designed to zero the superdiagonal (i, i+1) entry of  $T_k$ .  $\bar{L}_k$  is the lower tridiagonal result. The notation  $\bar{L}_k$  is used to distinguish it from  $L_k$  — the leading  $k \times k$  part of  $\bar{L}_{k+1}$  — from which it differs only in the (k, k)

element ( $\tilde{\gamma}_k$  becomes  $\gamma_k$ ). The first two steps illustrate the process:

$$\begin{bmatrix} T_{2} \\ 0 & \beta_{3} \end{bmatrix} \bar{Q}_{2}^{T} = \begin{bmatrix} \alpha_{1} & \beta_{2} \\ \beta_{2} & \alpha_{2} \\ 0 & \beta_{3} \end{bmatrix} \begin{bmatrix} c_{1} & s_{1} \\ s_{1} & -c_{1} \end{bmatrix} = \begin{bmatrix} \gamma_{1} \\ \delta_{2} & \bar{\gamma}_{2} \\ \epsilon_{3} & \bar{\delta}_{3} \end{bmatrix} = \begin{bmatrix} \bar{L}_{2} \\ \epsilon_{3} & \bar{\delta}_{3} \end{bmatrix}$$
$$T_{3}\bar{Q}_{3}^{T} = \begin{bmatrix} \gamma_{1} \\ \delta_{2} & \bar{\gamma}_{2} & \beta_{3} \\ \epsilon_{3} & \bar{\delta}_{3} & \alpha_{3} \end{bmatrix} \begin{bmatrix} 1 \\ c_{2} & s_{2} \\ s_{2} & -c_{2} \end{bmatrix} = \begin{bmatrix} \gamma_{1} \\ \delta_{2} & \gamma_{2} \\ \epsilon_{3} & \delta_{3} & \bar{\gamma}_{3} \end{bmatrix}$$
$$= \begin{bmatrix} L_{2} & 0 \\ 0 \\ \epsilon_{3} & \delta_{3} & \bar{\gamma}_{3} \end{bmatrix} = \bar{L}_{3}.$$

Thus after k steps

$$\bar{L}_{k} = \begin{bmatrix} \gamma_{1} & & & \\ \delta_{2} & \gamma_{2} & & \\ \epsilon_{3} & \delta_{3} & \gamma_{3} & & \\ & \ddots & \ddots & & \\ & & \epsilon_{k-1} & \delta_{k-1} & \gamma_{k-1} \\ & & & \epsilon_{k} & \delta_{k} & \bar{\gamma}_{k} \end{bmatrix} = \begin{bmatrix} & & 0 \\ L_{k-1} & \vdots \\ & 0 \\ \hline 0 & \epsilon_{k} & \delta_{k} & \bar{\gamma}_{k} \end{bmatrix} .$$
(60)

Note that  $T_k$  in (20) may be singular, in which case  $\overline{L}_k$  is singular. But since  $\beta_2 \cdots \beta_k \neq 0$ ,  $L_{k-1}$  is nonsingular. For ease of calculation define

$$\bar{z}_{k} \equiv \begin{bmatrix} \zeta_{1} \\ \vdots \\ \zeta_{k-1} \\ \bar{\zeta}_{k} \end{bmatrix} \equiv \begin{bmatrix} z_{k-1} \\ \bar{\zeta}_{k} \end{bmatrix} \equiv \bar{Q}_{k} y_{k}^{C}$$
$$\bar{W}_{k} \equiv \begin{bmatrix} w_{1} \cdots w_{k-1} & \bar{w}_{k} \end{bmatrix} \equiv \begin{bmatrix} W_{k-1} & \bar{w}_{k} \end{bmatrix} \equiv V_{k} \bar{Q}_{k}^{T}, \quad (61)$$

so the system

$$T_k y_k^C = \beta_1 e_1, \quad x_k^C = V_k y_k^C$$

is transformed to

$$\tilde{L}_k \tilde{z}_k = \beta_1 e_1, \quad x_k^C = \tilde{W}_k \tilde{z}_k.$$

At each step,  $\gamma_k$ ,  $c_k$  and  $s_k$  are calculated:

$$\gamma_k = \sqrt{\bar{\gamma}_k^2 + \beta_{k+1}^2}, \quad c_k = \bar{\gamma}_k / \gamma_k, \quad s_k = \beta_{k+1} / \gamma_k, \tag{62}$$

while the last row of  $\vec{L}_{k+1}$  is developed at the same time:

$$\bar{\gamma}_{k+1} = -(\beta_{k+1}c_{k-1}s_k + \alpha_{k+1}c_k), \quad \delta_{k+1} = \alpha_{k+1}s_k - \beta_{k+1}c_{k-1}c_k,$$
  

$$\epsilon_{k+1} = \beta_{k+1}s_{k-1}.$$
(63)

By comparing the last entry of  $\bar{L}_k \bar{z}_k$  with the k-th entry of  $\bar{L}_{k+1} \bar{z}_{k+1}$  (both should be zero):

$$\bar{\gamma}_k \bar{\zeta}_k = \gamma_k \zeta_k \quad \Rightarrow \quad \zeta_k = \bar{\gamma}_k \bar{\zeta}_k / \gamma_k = c_k \bar{\zeta}_k.$$
(64)

Using (59) and (61):

$$\bar{w}_{1} = v_{1}$$

$$\begin{bmatrix} \bar{w}_{k} & v_{k+1} \end{bmatrix} \begin{bmatrix} c_{k} & s_{k} \\ s_{k} & -c_{k} \end{bmatrix} = \begin{bmatrix} w_{k} & \bar{w}_{k+1} \end{bmatrix},$$
(65)

so passing from  $\{\bar{w}_k, v_{k+1}\}$  to  $\{w_k, \bar{w}_{k+1}\}$  is painless. The above algorithm is not implemented, since if  $\bar{L}_k$  is singular, then  $\bar{z}_k$  is undefined. This led Paige and Saunders [PS75] to consider the following alternative subproblem. Subproblem (SYMMLQ): The minimum 2-norm solution to the problem

 $T_{k+1,k}^T y_k^L = \beta_1 e_1$ , where  $T_{k+1,k}$  is as in (19). The matrix  $T_{k+1,k}$  has full column rank, and the subproblem always has a unique solution.

This is solved by performing an LQ decomposition of  $T_{k+1,k}^T$ ,

$$T_{k+1,k}^{T}\bar{Q}_{k+1}^{T} = \begin{bmatrix} T_{k} \mid \beta_{k+1}e_{k} \end{bmatrix} \begin{bmatrix} Q_{k}^{T} \mid \bar{q}_{k+1}^{T} \end{bmatrix} = \begin{bmatrix} L_{k} \mid 0 \end{bmatrix}$$
  
solving  $L_{k}z_{k} = \beta_{1}e_{1}$   
and setting  $y_{k}^{L} = Q_{k}^{T}z_{k}$ .  
Finally  $x_{k}^{L} = V_{k+1}y_{k}^{L} = V_{k+1}Q_{k}^{T}z_{k} = W_{k}z_{k} = x_{k-1}^{L} + \zeta_{k}w_{k}$ .

Now  $L_k$  is nonsingular if  $\beta_{k+1} \neq 0$ , so  $z_k$  does exist and we can safely compute the SYMMLQ solution. Notice  $x_k^c$  can easily be computed using

$$x_k^C = \bar{W}_k \bar{z}_k = W_{k-1} z_{k-1} + \bar{\zeta}_k \bar{w}_k = x_{k-1}^L + \bar{\zeta}_k \bar{w}_k.$$
(66)

The SYMMLQ solution is more easily updated than  $x_k^C$  (when  $x_k^C$  exists), but the latter may be more desirable at the finish.

Residual Norms: For the Conjugate Gradients solution

$$r_{k}^{C} = b - A\bar{W}_{k}\bar{z}_{k} = b - AV_{k}\bar{Q}_{k}^{T}\bar{z}_{k}$$

$$= b - V_{k+1}T_{k+1,k}\bar{Q}_{k}^{T}\bar{z}_{k}$$

$$= b - V_{k+1}\left[\frac{\bar{L}_{k}}{0 \ 0 \ \beta_{k+1}s_{k-1} \ -\beta_{k+1}c_{k-1}}\right]\bar{z}_{k} \quad (67)$$

$$= b - V_{k+1}\{\beta_{1}e_{1} + \beta_{k+1}(s_{k-1}\zeta_{k-1} - c_{k-1}\bar{\zeta}_{k})e_{k+1}\}$$

$$= -\beta_{k+1}(s_{k-1}\zeta_{k-1} - c_{k-1}\bar{\zeta}_{k})v_{k+1} \quad (68)$$

$$\||r_{k}^{C}\| = \beta_{k+1}|s_{k-1}\zeta_{k-1} - c_{k-1}\bar{\zeta}_{k}|.$$

Now, for the SYMMLQ residual,

⇒

$$\begin{aligned} r_k^L &= b - Ax_k^L = b - A(x_{k+1}^C - \bar{\zeta}_{k+1}\bar{w}_{k+1}) & \text{by (66)} \\ &= r_{k+1}^C + \bar{\zeta}_{k+1}A\bar{w}_{k+1} \\ &= r_{k+1}^C + \bar{\zeta}_{k+1}AV_{k+1}Q_{k+1}^Te_{k+1} & \text{by (61)} \\ &= r_{k+1}^C + \bar{\zeta}_{k+1}V_{k+2}T_{k+2,k+1}Q_{k+1}^Te_{k+1} \\ &= r_{k+1}^C + \bar{\zeta}_{k+1}V_{k+2}(\bar{\gamma}_{k+1}e_{k+1} - \beta_{k+2}c_ke_{k+2}) & \text{by (60) and (67)} \\ &= r_{k+1}^C + \bar{\zeta}_{k+1}(\bar{\gamma}_{k+1}v_{k+1} - \beta_{k+2}c_kv_{k+2}) \\ &= \bar{\zeta}_{k+1}\bar{\gamma}_{k+1}v_{k+1} - \\ &\quad \{\beta_{k+2}(s_k\zeta_k - c_k\bar{\zeta}_{k+1}) + \bar{\zeta}_{k+1}\beta_{k+2}c_k\}v_{k+2} & \text{by (68)} \\ &= \zeta_{k+1}\gamma_{k+1}v_{k+1} - \epsilon_{k+2}\zeta_kv_{k+2} & \text{by (63) and (64)} \\ \Rightarrow \|r_k^L\|^2 &= \zeta_{k+1}^2\gamma_{k+1}^2 + \epsilon_{k+2}^2\zeta_k^2. \end{aligned}$$

Because the implementation of SYMMLQ was originally obtained from this sort of approach in [PS75], there is no need to show any equivalence of this with a standard "vector" form here.

### 4.3 Minimum Residual Method (MINRES)

<u>Problem</u>:  $\min_{x \in R^n} ||b - Ax||$  for A symmetric indefinite, where A can be singular, and where Ax = b need not be a consistent set of equations.

<u>Process</u>: the normalized Lanczos process for symmetric A on  $\mathcal{K}^k(A; b)$  starting with  $v_1 = b/\beta_1$ ,  $\beta_1 \equiv ||b||$ .

Iterates of the form  $x_k = V_k y_k$  are considered. Using (18),(19),(20) and (48),

$$r_{k} = b - AV_{k}y_{k} = V_{k+1}(\beta_{1}e_{1} - T_{k+1,k}y_{k})$$
$$||r_{k}|| = ||\beta_{1}e_{1} - T_{k+1,k}y_{k}||.$$

<u>Subproblem</u>:  $\min_{y \in \mathbb{R}^k} ||\beta_1 e_1 - T_{k+1,k} y_k||$ , that is in theory solve the normal equations  $(T_k^2 + \beta_{k+1}^2 e_k e_k^T) y_k^M = \beta_1 T_k e_1$ , for  $T_k$  symmetric tridiagonal. (The superscript M indicates the iterates  $x_k^M = V_k y_k^M$  are of the MINRES method.) The matrix  $T_{k+1,k}$  has full column rank, so the subproblem always has a unique solution.

Implementation: As for SYMMLQ, consider an LQ decomposition of  $T_k$ .

$$T_{k}\bar{Q}_{k}^{T} = \bar{L}_{k}, \quad \bar{Q}_{k}^{T}\bar{Q}_{k} = I$$

$$T_{k+1,k}^{T}T_{k+1,k} = T_{k}^{2} + \beta_{k+1}^{2}e_{k}e_{k}^{T} = T_{k}T_{k}^{T} + \beta_{k+1}^{2}e_{k}e_{k}^{T}$$

$$= \bar{L}_{k}\bar{L}_{k}^{T} + \beta_{k+1}^{2}e_{k}e_{k}^{T} = L_{k}L_{k}^{T}.$$

The last equality follows by considering the (k, k)-element of the matrices and (62). So the subproblem becomes:

$$L_k L_k^T y_k^M = \beta_1 \bar{L}_k \bar{Q}_k e_1$$
  
or  $L_k^T y_k^M = \beta_1 D_k \bar{Q}_k e_1$ ,

since  $L_k$  is nonsingular if  $\beta_{k+1} \neq 0$ , and from (62)

$$\bar{L}_k = L_k D_k$$
 where  $D_k \equiv \text{diag}(1, \dots, 1, c_k)$ .

For case of computation define

$$\iota_{k} \equiv \begin{bmatrix} \tau_{1} \\ \tau_{2} \\ \vdots \\ \tau_{k} \end{bmatrix} \equiv \beta_{1} D_{k} \bar{Q}_{k} e_{1}, \quad M_{k} \equiv \begin{bmatrix} m_{1} & m_{2} & \dots & m_{k} \end{bmatrix} \equiv V_{k} L_{k}^{-T}, \quad (69)$$

which transforms the system to

$$L_{k}^{T}y_{k}^{M} = t_{k}, \quad x_{k}^{M} = V_{k}y_{k}^{M} = V_{k}L_{k}^{-T}L_{k}^{T}y_{k}^{M} = M_{k}t_{k}.$$

By continuing the sequence

$$\bar{Q}_2 e_1 = Q_{1,2} e_1 = \begin{bmatrix} c_1 \\ s_1 \end{bmatrix}, \quad \bar{Q}_3 e_1 = Q_{2,3} Q_{1,3} e_1 = \begin{bmatrix} c_1 \\ c_2 s_1 \\ s_2 s_1 \end{bmatrix},$$

we see the entries of  $t_k$  are given by:

$$\tau_1 \equiv \beta_1 c_1, \quad \tau_i \equiv \beta_1 s_1 \cdots s_{i-1} c_i, \quad i = 2, \dots k.$$
(70)

The columns of  $M_k$  satisfy a three term recurrence  $(M_k L_k^T = V_k)$  (see (60)):

$$\begin{bmatrix} M_{k-1} & m_k \end{bmatrix} \begin{bmatrix} L_{k-1}^T & 0 \\ \epsilon_k \\ 0 & \cdots & 0 \end{bmatrix} = \begin{bmatrix} \epsilon_k m_{k-2} + 1 \\ M_{k-1} L_{k-1}^T & \delta_k m_{k-1} + 1 \\ \gamma_k m_k \end{bmatrix}$$

$$\Rightarrow \quad m_1 = v_1/\gamma_1, \quad m_2 = (v_2 - \delta_2 m_1)/\gamma_2$$

$$m_k = (v_k - \delta_k m_{k-1} - \epsilon_k m_{k-2})/\gamma_k, \quad k > 2.$$

Finally:

$$x_k^M = M_k t_k = x_{k-1}^M + \tau_k m_k.$$

One may also consider the Conjugate Gradient solution,  $x_k^c$ .

 $x_{k}^{c} = V_{k}y_{k}^{c}$ where  $T_{k}y_{k}^{c} = \beta_{1}e_{1}$   $\bar{Q}_{k}T_{k}y_{k}^{c} = \bar{L}_{k}^{T}y_{k}^{c} = \bar{Q}_{k}\beta_{1}e_{1}$   $x_{k}^{c} = V_{k}L_{k}^{-T}L_{k}^{T}y_{k}^{c} = M_{k}L_{k}^{T}y_{k}^{c}$   $= M_{k}(\bar{L}_{k}D_{k}^{-1})^{T}y_{k}^{c} = M_{k}D_{k}^{-1}\bar{L}_{k}^{T}y_{k}^{c}$   $= M_{k}D_{k}^{-1}\bar{Q}_{k}\beta_{1}e_{1} = M_{k}D_{k}^{-2}(D_{k}\bar{Q}_{k}\beta_{1}e_{1})$   $= M_{k}D_{k}^{-2}t_{k}, \quad \text{by (69)}$   $= x_{k}^{M} + \tau_{k}(c_{k}^{-2} - 1)m_{k}$   $= x_{k}^{M} + \tau_{k}(s_{k}/c_{k})^{2}m_{k}.$ 

So  $x_k^c$  is available from  $x_k^M$ , although the reverse does not appear to be true, that is, we have not derived a simple way of computing  $x_k^M$  from the Conjugate Gradient solutions method.

<u>Residual Norm</u>: To calculate  $||r_k^c||$  or  $||r_k^M||$ , it is helpful to consider the relationship between the  $m_k$ , and the  $w_k$  as defined in (61) in SYMMLQ. Since we are considering the theoretical Lanczos process, it must stop in at most n steps. If it stops at  $\beta_{m+1} = 0$ , then

$$AV_m = V_m T_m, \quad T_m = L_m Q_m, \quad M_m = V_m L_m^{-T}$$
$$AM_m = V_m T_m L_m^{-T} = V_m Q_m^T = W_m,$$
$$Am_k = w_k, \quad k = 1, \dots, m.$$

Now for the residuals

$$r_{k}^{M} - r_{k}^{c} \equiv A(x_{k}^{c} - x_{k}^{M}) = \tau_{k}(s_{k}/c_{k})^{2}w_{k}$$
$$r_{k}^{M} \equiv b - Ax_{k}^{M} = r_{k}^{c} + \tau_{k}(s_{k}/c_{k})^{2}w_{k}$$
$$= \beta_{1}s_{1}\cdots s_{k}(s_{k}w_{k} - v_{k+1})/c_{k} \quad \text{by (70)}$$

$$= \beta_1 s_1 \cdots s_k \overline{w}_{k+1} \quad \text{by (65)}$$
$$\|r_k^M\| = \beta_1 |s_1 \cdots s_k|$$
$$= |c_k| \|r_k^c\|.$$

As with SYMMLQ, the implementation of MINRES was originally obtained from this sort of approach in [PS75]. Thus there is no need to show any equivalence of this with a standard "vector" form here.

# 4.4 Conjugate Gradients on the Normal Equations (CGNE) – LSQR implementation

<u>Problem</u>: Solve Ax = b if  $A \in \mathbb{R}^{m \times n}$  is consistent, otherwise  $\min_{x \in \mathbb{R}^n} ||b - Ax||$  for general A, by (theoretically) performing the Conjugate Gradient algorithm on the Normal Equations<sup>2</sup>,

$$A^T A x = A^T b, (71)$$

which is a SPD system when A has rank n.

<u>Process</u>: the normalized Lanczos process for symmetric matrices applied to give a basis for  $\mathcal{K}^k(A^T A; A^T b)$ .

The practical implementation does not use  $A^T A$  to compute the Lanczos vectors. Instead it considers two bidiagonalization procedures for the unsymmetric A called Bidiag1 and Bidiag2 respectively.

Bidiag1: This reduces A to lower bidiagonal form using

$$\beta_1 u_1 = b$$

$$\alpha_1 v_1 = A^T u_1$$

$$\beta_{k+1} u_{k+1} = A v_k - \alpha_k u_k$$
(72)

$$\alpha_{k+1}v_{k+1} = A^{T}u_{k+1} - \beta_{k+1}v_{k}$$
(73)

where the  $\alpha_k$  and  $\beta_k$  are chosen to normalize  $v_k$  and  $u_k$  respectively, so  $||u_k|| = ||v_k|| = 1$ . The process stops when either (72) or (73) is zero. With the definitions

$$U_k \equiv \begin{bmatrix} u_1 & u_2 & \dots & u_k \end{bmatrix}, \quad V_k \equiv \begin{bmatrix} v_1 & v_2 & \dots & v_k \end{bmatrix}$$

<sup>2</sup>See Claim at the end of this subsection, p. 51.

$$B_{k} \equiv \begin{bmatrix} \alpha_{1} & & \\ \beta_{2} & \alpha_{2} & \\ & \beta_{3} & \cdot \\ & & \cdot & \alpha_{k} \end{bmatrix}, \quad B_{k+1,k} \equiv \begin{bmatrix} B_{k} \\ & \beta_{k+1}e_{k}^{T} \end{bmatrix},$$

the recurrence relations may be written as

$$U_{k+1}(\beta_1 e_1) = b$$
  

$$AV_k = U_{k+1}B_{k+1,k}$$
(74)

$$A^T U_{k+1} = V_{k+1} B_{k+1}^T. (75)$$

The beauty of this is that the choice of coefficients ensures in theory (prior to stopping)

 $U_{k+1}^T U_{k+1} = I = V_{k+1}^T V_{k+1}, \quad U_{k+1}^T A V_{k+1} = B_{k+1}.$ 

Proof by induction:

$$B_{1} = \alpha_{1} = v_{1}^{T} A^{T} u_{1} = U_{1}^{T} AV_{1}$$

$$U_{1}^{T} U_{1} = I = V_{1}^{T} V_{1}.$$
Suppose  $U_{k}^{T} U_{k} = I = V_{k}^{T} V_{k}, \quad U_{k}^{T} AV_{k} = B_{k}.$ 
From (72)  $\beta_{k+1} U_{k}^{T} u_{k+1} = U_{k}^{T} Av_{k} - \alpha_{k} e_{k} = B_{k} e_{k} - \alpha_{k} e_{k} = 0,$ 
so  $U_{k+1}^{T} U_{k+1} = I$ 
and from (74)  $U_{k+1}^{T} AV_{k} = B_{k+1,k}.$ 
From (73)  $\alpha_{k+1} V_{k}^{T} v_{k+1} = V_{k}^{T} A^{T} u_{k+1} - \beta_{k+1} e_{k}$ 
 $= B_{k+1,k}^{T} e_{k} - B_{k+1,k}^{T} e_{k} = 0,$ 
so  $V_{k+1}^{T} V_{k+1} = I$ 
and from (75)  $U_{k+1}^{T} AV_{k+1} = B_{k+1}.$ 

Bidiag1 is related to two Lanczos processes, as can be seen by combining (74) and (75).

$$AA^T U_k = AV_k B_k^T = U_{k+1} B_{k+1,k} B_k^T,$$

$$A^{T}AV_{k} = A^{T}U_{k+1}B_{k+1,k} = V_{k+1}B_{k+1}^{T}B_{k+1,k}.$$
 (76)

Clearly  $\{u_i\}_{i=1}^{k}$  spans the krylov subspace  $\mathcal{K}^k(AA^T; b)$ , while  $\{v_i\}_{i=1}^{k}$  spans the krylov subspace  $\mathcal{K}^k(A^TA; A^Tb)$ .

Bidiag2: This reduces A to upper bidiagonal form using

$$\theta_1 v_1 = A^T b$$

$$\rho_1 p_1 = A v_1$$

$$\theta_{k+1} v_{k+1} = A^T p_k - \rho_k v_k$$

$$\rho_{k+1} p_{k+1} = A v_{k+1} - \theta_{k+1} p_k$$
(77)

where the  $\theta_k$  and  $\rho_k$  are chosen to give  $v_k$  and  $p_k$  unit norm respectively. With the definitions

$$P_{k} \equiv \begin{bmatrix} p_{1} & p_{2} & \dots & p_{k} \end{bmatrix}, \quad V_{k} \equiv \begin{bmatrix} v_{1} & v_{2} & \dots & v_{k} \end{bmatrix}$$
$$R_{k} \equiv \begin{bmatrix} \rho_{1} & \theta_{2} & & & \\ & \rho_{2} & \theta_{3} & & \\ & & \ddots & & \\ & & & \rho_{k-1} & \theta_{k} \\ & & & & \rho_{k} \end{bmatrix}, \quad R_{k,k+1} \equiv \begin{bmatrix} R_{k} & \theta_{k+1}e_{k} \end{bmatrix},$$

the recurrence relations may be written as

$$V_{k}(\theta_{1}e_{1}) = A^{T}b$$

$$A^{T}P_{k} = V_{k}R_{k}^{T} + \theta_{k+1}v_{k+1}e_{k}^{T} = V_{k+1}R_{k,k+1}^{T}$$

$$AV_{k+1} = P_{k+1}R_{k+1}$$

$$P_{k+1}^{T}P_{k+1} = V_{k+1}^{T}V_{k+1} = I$$

$$P_{k+1}^{T}AV_{k+1} = R_{k+1}.$$

The last two lines follow since this recurrence is essentially Bidiag1 with b replaced by  $A^T b$ , and  $A^T$  replaced by A. Again we have two Lanczos

processes

$$AA^{T}P_{k} = P_{k+1}R_{k+1}R_{k,k+1}^{T}, \quad A^{T}AV_{k} = V_{k+1}R_{k,k+1}^{T}R_{k},$$
(78)

so  $\{p_i\}_{i}^{k}$  spans  $\mathcal{K}^{k}(AA^{T}; AA^{T}b)$ , and  $\{v_i\}_{i}^{k}$  spans  $\mathcal{K}^{k}(A^{T}A; A^{T}b)$ .

<u>Relationship between Bidiag1 and Bidiag2</u>: The  $V_k$  are the same in Bidiag1 and Bidiag2, that is they are the result of applying the Lanczos process to develop  $\mathcal{K}^k(A^T A; A^T b)$ . This follows since  $\{v_i\}_1^k$  spans this same Krylov subspace in Bidiag1 and Bidiag2. The Lanczos processes with  $A^T A$  in (76) and (78) are therefore identical, and

$$B_{k+1}^T B_{k+1,k} = R_{k,k+1}^T R_k.$$

Dropping the last row of this gives

$$B_{k+1,k}^T B_{k+1,k} = R_k^T R_k.$$

This equality shows that  $R_k$  is the upper bidiagonal factor of  $B_{k+1,k}$  in its QR factorisation.

$$Q_k^T B_{k+1,k} = \begin{bmatrix} R_k \\ 0 \end{bmatrix} \quad \text{for } Q_k \text{ as in } \S 3,$$
$$g_k \equiv \beta_1 Q_k^T e_1 \quad \text{as in } (45).$$

The problem could be solved using either Bidiag1 or Bidiag2. Because of its greater numerical reliability in certain circumstances, LSQR was based on Bidiag1, and it is given here. It carries out the above QR factorization of  $B_{k+1,k}$  as the process progresses.

To solve  $\min_{x \in \mathbb{R}^n} ||b - Ax||$  using quantities from Bidiag1, let

$$x_k = V_k y_k$$
  
$$t_{k+1} \equiv \beta_1 e_1 - B_{k+1,k} y_k$$
  
$$r_k \equiv b - A x_k$$

$$= U_{k+1}\beta_1 e_1 - U_{k+1}B_{k+1,k}y_k = U_{k+1}t_{k+1} \quad \text{by (74)}$$
$$\|r_k\| = \|U_{k+1}t_{k+1}\| = \|t_{k+1}\|$$
$$\Rightarrow \min_{x \in \mathcal{K}^k} \|r_k\| = \min_{y_k \in R^k} \|t_{k+1}\|$$

<u>Subproblem</u>:  $\min_{y_k \in \mathbb{R}^k} \|\beta_1 e_1 - B_{k+1,k} y_k\|$ , for  $B_{k+1,k} \in \mathbb{R}^{(k+1) \times k}$  lower bidiagonal supplemented by an extra row. This always has a solution. <u>Implementation</u>: This is solved using the QR factorisation of  $B_{k+1,k}$  which with (45) transforms the subproblem to

$$\min_{\substack{y_k \in R^k \\ y_k \in R^k}} \|Q_k^T \beta_1 e_1 - Q_k^T B_{k+1,k} y_k\|$$

$$= \min_{\substack{y_k \in R^k \\ y_k \in R^k}} \|g_k - \begin{bmatrix} R_k \\ 0 \end{bmatrix} y_k\|$$

$$= \min_{\substack{y_k \in R^k \\ \phi_{k+1}}} \| \begin{bmatrix} f_k \\ \bar{\phi}_{k+1} \end{bmatrix} - \begin{bmatrix} R_k \\ 0 \end{bmatrix} y_k\|.$$

This is equivalent to

$$R_{k}y_{k} = f_{k}$$

$$l_{k+1} = \beta_{1}e_{1} - Q_{k}\begin{bmatrix} R_{k} \\ 0 \end{bmatrix} y_{k} = Q_{k}g_{k} - Q_{k}\begin{bmatrix} f_{k} \\ 0 \end{bmatrix}$$

$$= Q_{k}\begin{bmatrix} 0 \\ \bar{\phi}_{k+1} \end{bmatrix}.$$

For case of calculation define

$$D_k \equiv \begin{bmatrix} d_1 & d_2 & \dots & d_k \end{bmatrix} = V_k R_k^{-1}$$
$$W_k \equiv \begin{bmatrix} w_1 & w_2 & \dots & w_k \end{bmatrix} = D_k \operatorname{diag}(\rho_1, \dots, \rho_k)$$
$$x_k = V_k y_k = V_k R_k^{-1} f_k = D_k f_k.$$

Now the  $d_k$  (and their scalar multiples, the  $w_k$ ) are easily calculated from

 $V_k^T = R_k^T D_k^T:$ 

$$\begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_k^T \end{bmatrix} = \begin{bmatrix} \rho_1 \\ \theta_2 \\ \cdot \\ \rho_{k-1} \\ \theta_k \\ \rho_k \end{bmatrix} \begin{bmatrix} d_1^T \\ d_2^T \\ \vdots \\ d_k^T \end{bmatrix}$$
$$(d_0 \equiv 0, \quad w_0 \equiv 0)$$
$$d_k = (v_k - \theta_k d_{k-1})/\rho_k$$
$$w_k = v_k - \theta_k w_{k-1}/\rho_{k-1}$$

and hence the  $x_k$ :

$$x_k = D_k f_k = D_{k-1} f_{k-1} + \phi_k d_k = x_{k-1} + \phi_k d_k.$$

Residual Norm:

$$r_{k} = U_{k+1}t_{k+1} = U_{k+1}Q_{k}\begin{bmatrix} 0\\ \bar{\phi}_{k+1}\end{bmatrix} = \bar{\phi}_{k+1}U_{k+1}Q_{k}e_{k+1}$$
$$\|r_{k}\| = |\bar{\phi}_{k+1}| = \beta_{1}|s_{1}\cdots s_{k}| \quad \text{by (46)}$$
$$\Rightarrow \|r_{k}\| = |s_{k}|\|r_{k-1}\|.$$

<u>Claim</u>: LSQR is equivalent to CG on the Normal Equations (NE). <u>Proof</u>: By <sup>•</sup> denote the quantities from CG on the NE, that is

$$\hat{A}x = \hat{b},$$

where  $\hat{A} \equiv A^T A$  is SPD and  $\hat{b} \equiv A^T b$ . Now it has already been shown that  $V_k$  produced in Bidiagl is the same as  $V_k$  produced in the Lanczos process on  $\mathcal{K}^k(\hat{A}, \hat{b})$ . To show the theoretical equivalence of the two methods, we must prove

$$\hat{x}_k = V_k \hat{y}_k$$
 is equivalent to  $x_k = V_k y_k$   
or  $\hat{y}_k = y_k$ .

Now the  $y_k$  satisfy  $R_k y_k = f_k$ , while the  $\hat{y}_k$  are obtained from

$$\hat{T}_k \hat{y}_k = \hat{\beta}_1 e_1^{(k)}$$

$$R_k^T R_k \hat{y}_k = \theta_1 e_1^{(k)} \text{ by (77) and (78)}$$

$$R_k \hat{y}_k = R_k^{-T} \theta_1 e_1^{(k)}.$$

So it remains to prove the equivalence of  $\hat{y}_k$  and  $y_k$  using

$$R_{k}\hat{y}_{k} = R_{k}^{T}\theta_{1}e_{1}^{(k)} \text{ and } R_{k}y_{k} = f_{k}.$$
Now  $R_{k}^{T}f_{k} = \text{ first } k \text{ entries of } \left[\begin{array}{c|c} R_{k}^{T} \\ \hline \end{array}\right]g_{k}.$ 

$$\left[\begin{array}{c|c} R_{k}^{T} \\ \hline \end{array}\right]g_{k} = \left[\begin{array}{c|c} R_{k}^{T} \\ \hline \end{array}\right]\beta_{1}Q_{k}^{T}e_{1}^{(k+1)}$$

$$= \left(Q_{k}\left[\begin{array}{c|c} R_{k} \\ \hline \end{array}\right]\right)^{T}\beta_{1}e_{1}^{(k+1)}$$

$$= \left[\begin{array}{c|c} B_{k+1,k}^{T} \\ 0 \end{array}\right]\beta_{1}e_{1}^{(k+1)}$$

$$= \alpha_{1}\beta_{1}e_{1}^{(k+1)} = \theta_{1}e_{1}^{(k+1)}$$

$$\Rightarrow f_{k} = R_{k}^{T}\theta_{1}e_{1}^{(k)}.$$

The two methods produce the same iterates and therefore are the same. However naively applying the CG method to the Normal Equations has unsatisfactory behaviour on ill-conditioned systems. Thus it is more favourable to use LSQR.

The LSQR implementation of CGNE was obtained from this sort of approach in [PS82]. The mathematical equivalence of LSQR with CGNE is shown above.

References: [GK65, PS82]

## 4.5 Generalized Minimum Residual Method with restarts (GMRES(m))

<u>Problem</u>: Solve Ax = b for x when A is general nonsingular.

<u>Process</u>: Arnoldi algorithm applied to  $\mathcal{K}^k(A; r_0)$ , where  $r_0 \equiv b - Ax_0$  for some initial guess  $x_0$ . Define  $h_{1,0} \equiv ||r_0||$ ,  $v_1 = r_0/h_{1,0}$ .

This algorithm considers iterates of the form  $x_k = x_0 + V_k y_k$  so that  $||r_k||$ is minimized over  $(x_k - x_0) \in \mathcal{K}^k$ , that is

$$\begin{array}{ll}
\min_{x_k - x_0 \in \mathcal{K}^k} \|b - Ax_k\| &= \min_{y_k \in \mathcal{R}^k} \|b - A(x_0 + V_k y_k)\| = \\
\min_{y_k \in \mathcal{R}^k} \|r_0 - AV_k y_k\| &= \min_{y_k \in \mathcal{R}^k} \|V_{k+1}(h_{1,0}e_1 - H_{k+1,k} y_k)\| \quad \text{by (3).}
\end{array}$$

Subproblem:  $\min_{y_k \in \mathbb{R}^k} ||h_{1,0}e_1 - H_{k+1,k}y_k||$ , for  $H_{k+1} \in \mathbb{R}^{(k+1)\times k}$  upper Hessenberg supplemented by an extra row. This clearly always has a solution. Implementation: This may be solved by performing a QR factorisation on  $H_{k+1,k}$ 

$$Q_k^T I_{k+1,k} = \begin{bmatrix} R_k \\ 0 \end{bmatrix}, \quad g_k \equiv h_{1,0} Q_k^T e_1 = \begin{bmatrix} f_k \\ \bar{\phi}_{k+1} \end{bmatrix} \quad \text{as in (45).}$$

So by solving for  $y_k = R_k^{-1} f_k$  by backward substitution, the subproblem becomes

$$\min_{y_k \in R^k} \|g_k - \begin{bmatrix} R_k y_k \\ 0 \end{bmatrix} \| = |\bar{\phi}_{k+1}|.$$

Residual Norm:

$$||r_k|| = |\bar{\phi}_{k+1}| = h_{1,0}|s_1 \cdots s_k| \quad \text{by (46)}$$
  
$$\Rightarrow ||r_k|| = |s_k|||r_{k-1}||.$$

As a practical note, since the  $v_1, v_2, \ldots$  are kept, it is not necessary to compute  $x_k$  until iteration m, at which time

$$x_m = V_m y_m$$
$$r_m = b - A x_n$$

are calculated and the algorithm is restarted with  $x_0 \equiv x_m$  and  $r_0 \equiv r_m$ . At any iteration, based on the magnitude of  $||r_k||$ , the algorithm may be stopped and  $x_k$  calculated.

The drawback of GMRES(n) requiring an increasing amount of work and storage as the iteration count rises is overcome by restarting every m iterations, for some  $m \ll n$ . Unfortunately the optimal value of m is difficult to determine and in certain cases the method will not converge for m < n. Too large a value of m results in unneccessary work and storage. On the other hand, too small a value may give slow convergence or fail to converge at all.

Because the implementation of GMRES was originally obtained from this sort of approach in [SS86], there is no need to show any equivalence of this with a standard "vector" form here. The GMRES(m) method slightly altered from 'Templates' — is as follows.

 $x_{0} \text{ initial guess}$ for j = 1, 2, ... $h_{1,0} = ||b - Ax_{0}||$  $v_{1} = (b - Ax_{0})/h_{1,0}$  $Q_{0}^{T} = [1]$  $\tilde{\phi}_{1} = h_{1,0}$ for i = 1, 2, ..., m $u_{i} = Av_{i}$ for k = 1, 2, ..., i $h_{k,i} = v_{k}^{T}u_{i}$  (or, k-th entry  $h_{i}$  gets  $h_{k,i}$ )  $u_{i} = u_{i} - h_{k,i}v_{k}$  $h_{i+1,i} = ||u_{i}||$  $v_{i+1} = u_{i}/h_{i+1,i}$ 

$$\begin{split} R_i &= \begin{bmatrix} R_{i-1} \\ 0 \end{bmatrix} Q_{i-1}^T h_i \\ \hat{r}_{i,i} &= r_{i,i} \\ r_{i,i} &= \sqrt{\hat{r}_{i,i}^2 + h_{i+1,i}^2} \\ c_i &= \hat{r}_{i,i}/r_{i,i}, \quad s_i &= -h_{i+1,i}/r_{i,i} \\ Q_i^T &= Q_{i,i}^T \begin{bmatrix} Q_{i-1}^T \\ 0 \end{bmatrix} \\ \phi_i &= c_i \bar{\phi}_i \\ f_i &= \begin{bmatrix} f_{i-1} \\ \phi_i \end{bmatrix} \\ \bar{\phi}_{i+1} &= s_i \bar{\phi}_i \\ \text{if } |\bar{\phi}_{i+1}| < \text{tol} \\ \text{update}(x, i) \text{ and quit} \\ \text{update}(x, m) \end{split}$$

.

end

update(x, j): solve  $R_j y_j = f_j$  by back substitution form  $x_j = V_j y_j$ 

end

## 4.6 The Unsymmetric Lanczos (BiConjugate Gradients) Method

<u>Problem</u>: Solve Ax = b for x when A is general nonsingular.

First we derive what will later be shown to be the BiConjugate Gradients approximations  $x_k$  from the standard unsymmetric Lanczos process. We refer to this as:

The standard unsymmetric Lanczos-BiCG method.

<u>Process</u>: the standard unsymmetric Lanczos process (32) to produce mutually orthogonal  $\{v_1, \ldots, v_k\}$  and  $\{\tilde{v}_1, \ldots, \tilde{v}_k\}$  spanning the krylov subspaces  $\mathcal{K}^k(A; b)$  and  $\mathcal{K}^k(A^T; b)$  respectively, starting with  $v_1 = b/\beta_1$  and  $\tilde{v}_1 = b/\gamma_1$ , where  $\beta_1 = \gamma_1 = ||b||$ , where  $\beta_{j+1}$  is chosen to normalize  $v_{j+1}$ .

To derive the subproblem we use (31) and consider iterates of the form  $x_k = V_k y_k$  as suggested in (47), and take  $T_k y_k = \beta_1 e_1$ , so that

$$r_{k} = b - Ax_{k} = b - AV_{k}y_{k}$$
  
=  $b - V_{k}T_{k}y_{k} - \beta_{k+1}v_{k+1}e_{k}^{T}y_{k}$   
=  $V_{k}(\beta_{1}e_{1} - T_{k}y_{k}) - \beta_{k+1}v_{k+1}e_{k}^{T}y_{k}$   
=  $-\beta_{k+1}v_{k+1}e_{k}^{T}y_{k},$  (79)

much the same as (49) in the unit norm Lanczos method. In theory we have  $\beta_{k+1}v_{k+1} = 0$  for some  $k \leq n$ , the dimension of A, and the problem will be solved.

Subproblem: Solve  $T_k y_k = \beta_1 e_1$ , where  $T_k$  is tridiagonal.

When  $T_k$  is singular this subproblem does not have a (unique) solution, but as we are only intent on giving a new approach to the general BiCG idea, we do not go into such aspects here. One very satisfactory approach to avoiding such difficulties is to use the Quasi-Minimal Residuals method discussed later in § 4.7. This work on BiCG will motivate that. <u>Implementation</u>: This may be solved using an LU factorisation (when it exists) of  $T_k = L_k U_k$ , where  $L_k$  is lower bidiagonal and  $U_k$  is chosen to be unit upper bidiagonal. There is nothing special about  $T_k$ , so the use of this decomposition without pivoting may not be stable. For further insights on this, see the following alternative unsymmetric Lanczos-BiCG method.

For case of calculation and savings on storage, define

$$z_{k} \equiv \begin{bmatrix} \zeta_{1} \\ \zeta_{2} \\ \vdots \\ \zeta_{k} \end{bmatrix} \equiv U_{k}y_{k}, \quad W_{k} \equiv \begin{bmatrix} w_{1} & w_{2} & \cdots & w_{k} \end{bmatrix} \equiv V_{k}U_{k}^{-1},$$

so the system is transformed to

$$T_k y_k = L_k U_k y_k = L_k z_k \implies L_k z_k = \beta_1 e_1,$$
$$V_k y_k = V_k U_k^{-1} U_k y_k \implies x_k = W_k z_k.$$

Now,  $L_k$  is formed one row at a time;  $U_k$  one column at a time:

$$T_{k} = \begin{bmatrix} L_{k-1} & 0 \\ 0 & \epsilon_{k} & \delta_{k} \end{bmatrix} \begin{bmatrix} U_{k-1} & 0 \\ \eta_{k} \\ 0 & 0 & 1 \end{bmatrix}$$
$$= \begin{bmatrix} L_{k-1}U_{k-1} & 0 \\ 0 & \epsilon_{k} & \delta_{k-1}\eta_{k} \\ 0 & \epsilon_{k} & \epsilon_{k}\eta_{k} + \delta_{k} \end{bmatrix} = \begin{bmatrix} T_{k-1} & 0 \\ \eta_{k} & \eta_{k} \\ \eta_{k} & \eta_{k} \end{bmatrix}$$
$$\Rightarrow \delta_{1} = \alpha_{1}$$
$$\epsilon_{k} = \beta_{k}, \quad \eta_{k} = \gamma_{k}/\delta_{k-1}, \quad \delta_{k} = \alpha_{k} - \epsilon_{k}\eta_{k}, \quad k > 1.$$
(80)

And so  $z_k$  is formed one entry at a time:

$$\begin{bmatrix} L_{k-1} & 0 \\ 0 \\ \hline 0 & \epsilon_k & \delta_k \end{bmatrix} \begin{bmatrix} z_{k-1} \\ \zeta_k \end{bmatrix} = \begin{bmatrix} L_{k-1} z_{k-1} \\ \epsilon_k \zeta_{k-1} + \delta_k \zeta_k \end{bmatrix} = \beta_1 e_1$$
$$\Rightarrow \zeta_1 = \beta_1 / \delta_1, \quad \zeta_k = -\epsilon_k \zeta_{k-1} / \delta_k, \quad k > 1.$$
(81)

 $W_k$  is also formed one column at a time  $(W_k U_k = V_k)$ :

$$\begin{bmatrix} W_{k-1} & w_k \end{bmatrix} \begin{bmatrix} U_{k-1} & 0 \\ & \eta_k \\ \hline 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} W_{k-1}U_{k-1} & \eta_k w_{k-1} + w_k \end{bmatrix}$$
$$\Rightarrow w_1 = v_1, \quad w_k = v_k - \eta_k w_{k-1}, \quad k > 1.$$

Finally,

$$x_k = W_k z_k = W_{k-1} z_{k-1} + \zeta_k w_k = x_{k-1} + \zeta_k w_k, \quad k \ge 1.$$

<u>Residual Norm</u>: The following estimation of the norm of the residual  $r_k$ , as defined in (79) can be used.

$$r_k = -\beta_{k+1}v_{k+1}e_k^T y_k = -\beta_{k+1}v_{k+1}y_{kk}$$
$$||r_k|| = \beta_{k+1}|y_{kk}|, \text{ where } y_{kk} \text{ is the } k\text{-th entry of } y_k$$

Since  $y_k$  is not computed directly, one must consider

$$U_{k}y_{k} = z_{k}$$

$$\begin{bmatrix} U_{k-1} & 0 \\ \eta_{k} \\ \hline 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_{k1} \\ \vdots \\ y_{kk} \end{bmatrix} = \begin{bmatrix} \vdots \\ y_{kk} \end{bmatrix} = \begin{bmatrix} z_{k-1} \\ \zeta_{k} \end{bmatrix}$$

$$y_{kk} = \zeta_{k}$$

$$\Rightarrow ||r_{k}|| = \beta_{k+1}|\zeta_{k}|$$

$$= \beta_{k+1}|\epsilon_{k}\zeta_{k-1}/\delta_{k}| \text{ by (81)}$$

$$= \beta_{k+1}\beta_{k}|\zeta_{k-1}/\delta_{k}| \text{ by (80)}$$

$$||r_{k}|| = \beta_{k+1}||r_{k-1}||/|\delta_{k}|, \quad k \ge 1.$$

Although this solves the problem in much the same way as BiCG in 'Templates'. [BBC+94], it gives different scalings for equivalent vectors. Let us now derive a method which, as with the BiCG method of [BBC+94], produces the residual vectors themselves. The alternative unsymmetric Lanczos-BiCG method

<u>Process</u>: the alternative unsymmetric Lanczos process, see (33) to (44), starting with  $v_1 = \tilde{v}_1 = b$  ( $\beta_1 \equiv 1$ ). The normalization  $\beta_{j+1}$  is unspecified as yet.

Using (37) and and considering iterates of the form  $x_k = V_k y_k$  as suggested in (47) and solving for  $T_k y_k = e_1$ ,

$$r_{k} = b - Ax_{k} = b - AV_{k}y_{k}$$
  
=  $b - V_{k}T_{k}y_{k} - \beta_{k+1}v_{k+1}e_{k}^{T}y_{k}$   
=  $V_{k}(e_{1} - T_{k}y_{k}) - \beta_{k+1}v_{k+1}e_{k}^{T}y_{k}$   
=  $-\beta_{k+1}v_{k+1}e_{k}^{T}y_{k}.$  (82)

Subproblem: Solve  $T_k y_k = e_1$ , where  $T_k$  is tridiagonal. Again  $T_k$  may be singular, in which case this approach, and ordinary BiCG, fails.

<u>Implementation</u>: This may be solved using an LDU decomposition of  $T_k$ . The one shown below suits our purposes. With the as yet unkown  $\mathcal{D}_k \equiv -\text{diag}(\beta_2, \ldots, \beta_{k+1})$ , write

$$T_{k} = L_{k} \mathcal{D}_{k} U_{k}$$

$$= \left[ \begin{array}{c|c} L_{k-1} & 0 \\ \hline 0 & -1 & 1 \end{array} \right] \left[ \begin{array}{c|c} \mathcal{D}_{k-1} & 0 \\ \hline 0 & 0 & -\beta_{k+1} \end{array} \right] \left[ \begin{array}{c|c} U_{k-1} & 0 \\ \hline -\frac{\gamma_{k}}{\beta_{k}} \\ \hline 0 & 0 & 1 \end{array} \right]$$
(83)
$$= \left[ \begin{array}{c|c} L_{k-1} \mathcal{D}_{k-1} U_{k-1} \\ \hline 0 & \beta_{k} \end{array} \right] \left[ \begin{array}{c|c} -\gamma_{k} \\ \hline \gamma_{k} \\ \hline 0 & \beta_{k} \end{array} \right]$$

A simpler method could be obtained by combining  $\mathcal{D}$  and U, and using for example  $\mathcal{U} \equiv \mathcal{D}U$  instead, but the approach here will produce the equivalents of various elements of BiCG that we want.

Our normalizers  $\beta_{j+1}$  may now be determined by equating the right side

of the above equation with the tridiagonal  $T_k$  in (36). So

$$\alpha_1 = -\beta_2; \quad \alpha_k = -(\gamma_k + \beta_{k+1}), \quad k > 1,$$
  
$$\Rightarrow \beta_2 \equiv -\alpha_1; \quad \beta_{k+1} \equiv -(\alpha_k + \gamma_k), \quad k > 1.$$
(84)

We will show later that if  $\beta_{j+1} = 0$  the ordinary BiCG breaks down too, otherwise this is a satisfactory normalizer for the process, and  $U_{k+1}$  exists. In the meantime define

$$z_{k} \equiv \begin{bmatrix} \zeta_{1} \\ \zeta_{2} \\ \vdots \\ \zeta_{k} \end{bmatrix} \equiv U_{k} y_{k}, \quad W_{k} \equiv \begin{bmatrix} w_{1} & w_{2} & \cdots & w_{k} \end{bmatrix} \equiv V_{k} U_{k}^{-1}.$$
(85)

For our later analysis, we also define

$$\widetilde{W}_k \equiv \left[ \begin{array}{ccc} \tilde{w}_1 & \tilde{w}_2 & \cdots & \tilde{w}_k \end{array} \right] \equiv \widetilde{V}_k U_k^{-1},$$

although this is not used in the algorithm itself.

The system is transformed to

$$T_k y_k = L_k \mathcal{D}_k U_k y_k = L_k \mathcal{D}_k z_k \implies L_k \mathcal{D}_k z_k = e_1$$
$$V_k y_k = V_k U_k^{-1} U_k y_k \implies x_k = W_k z_k.$$

Now  $z_k$  is formed one entry at a time:

$$L_{k}\mathcal{D}_{k}z_{k} = \begin{bmatrix} -\beta_{2}\zeta_{1} \\ \beta_{2}\zeta_{1} - \beta_{3}\zeta_{2} \\ \vdots \\ \beta_{k}\zeta_{k-1} - \beta_{k+1}\zeta_{k} \end{bmatrix} = e_{1}$$
  

$$\Rightarrow \zeta_{1} = -\frac{1}{\beta_{2}}$$
  

$$\zeta_{k} = \frac{\beta_{k}}{\beta_{k+1}}\zeta_{k-1} = \frac{\beta_{k-1}}{\beta_{k+1}}\zeta_{k-2} = \dots = -\frac{1}{\beta_{k+1}}, \quad k > 1. \quad (86)$$

Solving  $W_k U_k = V_k$  and  $\widetilde{W}_k U_k = \widetilde{V}_k$ ,

$$w_1 = v_1, \quad w_k = v_k + \frac{\gamma_k}{\beta_k} w_{k-1}, \quad k > 1$$
  
$$\tilde{w}_1 = \tilde{v}_1, \quad \tilde{w}_k = \tilde{v}_k + \frac{\gamma_k}{\beta_k} \tilde{w}_{k-1}, \quad k > 1.$$

Finally,

$$x_k = W_{k-1} z_{k-1} + \zeta_k w_k = x_{k-1} + \zeta_k w_k.$$

<u>Residual Norm</u>: The following estimation of the norm of the residual  $r_k$ , as defined in (82) may be used.

$$r_{k} = -\beta_{k+1}y_{kk}v_{k+1}$$
  
=  $-\beta_{k+1}\zeta_{k}v_{k+1}$  by (85)  
=  $v_{k+1}$  by (86) (87)  
 $\Rightarrow ||r_{k}|| = ||v_{k+1}||.$ 

So the alternative unsymmetric Lanczos process, with  $v_1 = b$ , can be used to produce the (BiCG) residual vectors  $r_k = b - Ax_k$  themselves, as long as  $\beta_{j+1}$  is chosen as in (84).

It remains to examine under what conditions  $\beta_{j+1}$  can be zero. Suppose  $\beta_2, \ldots, \beta_j$  obtained from (84) are nonzero. Then with  $e^T \equiv [1, \ldots, 1]$ , we have  $e^T T_j = [0, \ldots, 0, \alpha_j + \gamma_j]$ , which will be nonzero if  $T_j$  is nonsingular. However if  $\alpha_j + \gamma_j = 0$  then  $T_j$  is singular, and both this and the standard BiCG method fails. Otherwise, it is clear that the  $L_k \mathcal{U}_k$  'factorization' is possible. However small  $\beta_k$  could cause large element growth in  $\mathcal{U}_k$ .

Equivalence of unsymmetric Lanczos-BiCG with BiCG: The BiCG method — slightly altered from 'Templates' — is as follows.

 $r_0 = b$ 

 $\tilde{r}_0 = r_0$ 

for i = 1, 2, ...  $\rho_{i-1} = r_{i-1}^T \tilde{r}_{i-1}$ if  $\rho_{i-1} = 0$ , the method fails, stop if (i = 1)  $p_i = r_{i-1}$   $\tilde{p}_i = \tilde{r}_{i-1}$ else  $\beta_{i-1} = \rho_{i-1}/\rho_{i-2}$ 

 $p_{i-1} = p_{i-1}/p_{i-2}$   $p_i = r_{i-1} + \beta_{i-1}p_{i-1}$   $\tilde{p}_i = \tilde{r}_{i-1} + \beta_{i-1}\tilde{p}_{i-1}$   $q_i = Ap_i$   $\tilde{q}_i = A^T \tilde{p}_i$   $\alpha_i = \rho_{i-1}/\tilde{p}_i^T q_i$   $x_i = x_{i-1} + \alpha_i p_i$   $r_i = r_{i-1} - \alpha_i q_i$   $\tilde{r}_i = \tilde{r}_{i-1} - \alpha_i \tilde{q}_i$ 

check convergence; continue if necessary

end

To avoid confusion we will use the superscript  $^{B}$  to distinguish 'Templates' terms from those here, where it will also make life easier if we reindex two of the 'Templates' terms as follows:

$$\rho_i^B \equiv \text{`Templates'}\rho_{i-1},$$
  
 $\beta_i^B \equiv \text{`Templates'}\beta_{i-1}.$ 

We now write the algorithms side by side for comparison, where in the alternative unsymmetric Lanczos-BiCG algorithm we replace  $v_{i+1}$  by  $r_i$ , see (87). introduce  $u_i \equiv Av_i$  and  $\tilde{u}_i \equiv A^T \tilde{v}_i$  and redefine  $\tilde{r}_i \equiv \tilde{v}_{i+1}$ . For ease of comparison in handling the i = 1 case we assume  $1/\infty = 0$  and define

BiCG  

$$\rho_0^B = 1$$
  
 $x_0^B = 0, \quad p_0^B = 0, \quad \tilde{p}_0^B = 0$ 
  
 $\tilde{r}_0^B = r_0^B = b$ 

unsymmetric Lanczos-BiCG

$$\rho_0 = \infty, \quad \beta_1 = 1$$
  
 $x_0 = 0$ 
  
 $r_{-1} = 0 \quad \tilde{r}_{-1} = 0, \quad \tilde{r}_0 = r_0 = b$ 

Then for i = 1, 2, 3, ... we have

----

BiCGalt. unsym. Lanczos-BiCG
$$\rho_i^B = r_{i-1}^{BT} \tilde{r}_{i-1}^B$$
 $\rho_i = \tilde{r}_{i-1}^T r_{i-1}$  $\beta_i^B = \rho_i^B / \rho_{i-1}^B$  $\gamma_i = \beta_i \rho_i / \rho_{i-1}$  from (44) $p_i^B = r_{i-1}^B + \beta_i^B p_{i-1}^B$  $w_i = r_{i-1} + (\gamma_i / \beta_i) w_{i-1}$  $\tilde{p}_i^B = \tilde{r}_{i-1}^B + \beta_i^B \tilde{p}_{i-1}^B$  $\tilde{w}_i = \tilde{r}_{i-1} + (\gamma_i / \beta_i) \tilde{w}_{i-1}$  $q_i^B = A p_i^B$  $u_i = A r_{i-1}$  $\tilde{q}_i^B = \Lambda^T \tilde{p}_i^B$  $u_i = A^T \tilde{r}_{i-1}$  $\alpha_i^B = \rho_i^B / \tilde{p}_i^{BT} q_i^B$  $\alpha_i = \tilde{r}_{i-1}^T u_i / \rho_i$  from (40) $\beta_{i+1} = -(\alpha_i + \gamma_i)$  $x_i = x_{i-1} - (1/\beta_{i+1}) w_i$  $r_i^B = r_{i-1}^B - \alpha_i^B q_i^B$  $r_i = (u_i - \alpha_i r_{i-1} - \gamma_i r_{i-2})/\beta_{i+1}$  $\tilde{r}_i^B = \tilde{r}_{i-1}^B - \alpha_i^B \tilde{q}_i^B$  $\tilde{r}_i = (\tilde{u}_i - \alpha_i \tilde{r}_{i-1} - \gamma_i \tilde{r}_{i-2})/\beta_{i+1}$ 

The main difference in cost between the two is in the calculation of the residuals and shadow residuals (the  $\tilde{r}_i$ ) — 2n extra multiplications, subtractions and divisions in the alternative unsymmetric Lanczos-BiCG algorithm — which are partially offset by the *n* extra multiplications and additions in the BiCG calculation of  $\tilde{p}_i^B$  ( $\tilde{w}_i$  is not actually calculated). The equivalences are

$$\rho_i = \rho_i^B, \quad \gamma_i / \beta_i = \beta_i^B, \quad w_i = p_i^B, \quad \tilde{w}_i = \tilde{p}_i^B$$
$$-1/\beta_{i+1} = \alpha_i^B, \quad x_i = x_i^B, \quad r_i = r_i^B, \quad \tilde{r}_i = \tilde{r}_i^B.$$
(88)

A straightforward examination shows this is true for i = 1. Suppose it is

true for i = 1, 2, ..., k - 1. Then it is clear that

$$\rho_k = \rho_k^B, \quad \gamma_k / \beta_k = \beta_k^B, \quad w_k = p_k^B, \quad \tilde{w}_k = \tilde{p}_k^B.$$

We now need only prove

$$-1/\beta_{k+1} = \alpha_k^B, \tag{89}$$

for then  $x_k = x_k^B$ , and their residuals  $r_k$  and  $r_k^B$  must then also be equal, since  $r_i^B$  is clearly the residual for  $x_i^B$  by examining the last two lines of BiCG, and we have shown  $r_i$  is the residual of  $x_i$  in the alternative unsymmetric Lanczos-BiCG. In this case  $\tilde{r}_k = \tilde{r}_k^B$  follows from the equality of the  $r_k$  and  $r_k^B$  residuals, since we know that  $r_k = \varphi_k(A)b$  and  $r_k^B = \psi_k(A)b$  for two k degree polynomials  $\varphi_k$  and  $\psi_k$ , and  $\varphi_k = \psi_k$ . But we can also show  $\tilde{r}_k = \varphi_k(A^T)b$  and  $\tilde{r}_k^B = \psi_k(A^T)b$ , so these must also be equal. We can prove (S9) using our matrix formulation of the unsymmetric Lanczos-BiCG algorithm. We know  $W_k = V_k U_k^{-1}$  and  $\widetilde{W}_k = \widetilde{V}_k U_k^{-1}$ , so

$$\widetilde{W}_k^T A W_k = U_k^{-T} \widetilde{V}_k^T A V_k U_k^{-1} = U_k^{-T} D_k T_k U_k^{-1}$$
$$= U_k^{-T} D_k (L_k \mathcal{D}_k U_k) U_k^{-1} = U_k^{-T} D_k L_k \mathcal{D}_k,$$

which, by the way, is lower triangular and symmetric (see (39)), and therefore diagonal, showing "A-biorthogonality" of the  $w_i$  and  $\tilde{w}_j$ , where i, j < k. The (i, i)-element of this is (see (83))

$$\hat{w}_i^T A w_i = e_i^T U_k^{-T} (\rho_i e_i - \rho_{i+1} e_{i+1}) (-\beta_{i+1}).$$

But  $c_i^T U_k^{-T} = (x, \dots, x, 1, 0, \dots, 0)$  is zero after the *i*-th element, so since  $p_i^B = w_i \ \hat{p}_i^B = \hat{w}_i, \ i = 1, \dots, k,$ 

$$\tilde{p}_i^{BT}Ap_i^B = \tilde{w}_i^TAw_i = -\rho_i\beta_{i+1}, \qquad i = 1, \dots, k.$$

From the BiCG algorithm we then see

$$\alpha_i^B = \rho_i^B / \tilde{p}_i^{BT} q_i^B = \rho_i / \tilde{p}_i^{BT} A p_i^B = -1 / \beta_{i+1}, \quad i = 1, \dots, k,$$

and so (89) holds, and the induction proof of (88) is complete, showing that in theory these algorithms not only compute the same iterates, but also several of the intermediate quantities are identical. However the algorithms are clearly not identical, as the last two lines of each clearly highlights, and may have different numerical behaviours. This suggests a matrix development of BiCG would be useful, but we leave this for future work.

#### 4.7 Quasi-Minimal Residuals (QMR)

<u>Problem</u>: Solve Ax = b for x when A is general nonsingular.

<u>Process</u>: an unsymmetric Lanczos process which produces mutually orthonormal  $\{v_1, \ldots, v_k\}$  and  $\{\tilde{v}_1, \ldots, \tilde{v}_k\}$  spanning the krylov spaces  $\mathcal{K}^k(A; b)$ and  $\mathcal{K}^k(A^T; b)$  respectively, starting with  $v_1 = b/\beta_1$  and  $\tilde{v}_1 = b/\tilde{\beta}_1$ , where  $\beta_1 = \tilde{\beta}_1 = ||b||$ . We choose  $\beta_{j+1}$  to make  $v_{j+1}$  unit length and  $\tilde{\beta}_{j+1}$  to make  $\tilde{v}_{j+1}$  unit length. We define  $\rho_{j+1} \equiv \tilde{v}_{j+1}^T v_{j+1}$ .

To derive the subproblem we use (21) (where  $T_k \equiv H_k$  is tridiagonal) and consider iterates of the form  $x_k = V_k y_k$  as suggested in (47).

$$r_{k} = b - Ax_{k} = \beta_{1}v_{1} - AV_{k}y_{k} = V_{k+1}(\beta_{1}e_{1} - T_{k+1,k}y_{k}).$$
(90)

Since  $V_{k+1}$  does not have orthonormal columns, it is difficult to minimize  $||r_k||$ . The reduced system below is considered, and this gives rise to the name 'Quasi', since the residual is not actually minimised.

From (23)  $\tilde{V}_{k+1}^T V_{k+1} = D_{k+1}$ , so  $\tilde{V}_{k+1}^T r_k = D_{k+1}(\beta_1 e_1 - T_{k+1,k} y_k)$ . QMR minimizes  $\|\beta_1 e_1 - T_{k+1,k} y_k\|$ , so in reality (in theory) QMR is minimizing  $\|D_{k+1}^{-1} \tilde{V}_{k+1}^T r_k\|$ , not  $\|r_k\|$ .

Subproblem:  $\min_{y_k \in \mathbb{R}^k} ||\beta_1 e_1 - T_{k+1,k} y_k||$ , for  $T_{k+1,k} \in \mathbb{R}^{(k+1) \times k}$  tridiagonal supplemented by an extra row. This always has a solution.

Implementation: This may be solved using a QR factorisation of  $T_{k+1,k}$  as was done on  $H_{k+1,k}$  in GMRES in § 4.5.

$$Q_k^T T_{k+1,k} = \begin{bmatrix} R_k \\ 0 \end{bmatrix}, \quad g_k \equiv \beta_1 Q_k^T e_1 = \begin{bmatrix} f_k \\ \bar{\phi}_{k+1} \end{bmatrix}.$$

However  $R_k$  is upper tridiagonal, not just upper triangular. Solving for  $y_k = R_k^{-1} f_k$  the subproblem becomes

$$\min_{\mathbf{y}_k \in R^k} \|g_k - \begin{bmatrix} R_k \\ 0 \end{bmatrix} \mathbf{y}_k\| = |\bar{\phi}_{k+1}|.$$
(91)

For case of calculation define

$$W_k \equiv \left[ \begin{array}{ccc} w_1 & w_2 & \cdots & w_k \end{array} 
ight] \equiv V_k R_k^{-1}.$$

The  $w_k$  satisfy a three term recurrence  $(W_k R_k = V_k)$ :

$$\begin{bmatrix} W_{k-1} & w_k \end{bmatrix} \begin{bmatrix} 0 \\ R_{k-1} & r_{k-2,k} \\ & r_{k-1,k} \\ \hline 0 & 0 & r_{kk} \end{bmatrix} = \begin{bmatrix} W_{k-1}R_{k-1} & v_k \end{bmatrix}$$

$$\Rightarrow \quad w_1 = v_1/r_{1,1}, \quad w_2 = (v_2 - r_{1,2}w_1)/r_{2,2}$$

$$w_k = (v_k - r_{k-2,k}w_{k-2} - r_{k-1,k}w_{k-1})/r_{k,k}, \quad k > 2.$$

Finally,

-

$$x_{k} = V_{k}y_{k} = V_{k}R_{k}^{-1}f_{k} = W_{k}f_{k}$$
$$= W_{k-1}f_{k-1} + \phi_{k}w_{k} = x_{k-1} + \phi_{k}w_{k}.$$

<u>Residual and Residual Norm</u>: The following estimation of the norm of the residual  $r_k$ , as defined in (90) is used.

$$||r_{k}|| = ||V_{k+1}(\beta_{1}c_{1} - T_{k+1,k}y_{k})||$$
  

$$\leq \sqrt{k+1}||\beta_{1}c_{1} - T_{k+1,k}y_{k}|| \quad \text{(columns of } V_{k+1} \text{ are unit length)}$$
  

$$= \sqrt{k+1}|\bar{\phi}_{k+1}| \quad \text{by (91)}$$
  

$$= \beta_{1}\sqrt{k+1}|s_{1}\cdots s_{k}| \quad \text{by (46).}$$

Once this bound on the residual norm reaches a sufficiently low value, the following calculation of the actual residual is used for the few remaining iterations.

$$r_{k} = V_{k+1}(\beta_{1}e_{1} - T_{k+1,k}R_{k}^{-1}f_{k})$$
  
=  $\tilde{\phi}_{k+1}V_{k+1}Q_{k}e_{k+1}$  (92)

$$= \bar{\phi}_{k+1} V_{k+1} \left[ \begin{array}{c|c} Q_{k-1} \\ \hline \\ 1 \end{array} \right] (s_k e_k + c_k e_{k+1})$$

$$= \bar{\phi}_{k+1} (s_k V_k Q_{k-1} e_k + c_k v_{k+1})$$

$$= s_k^2 \bar{\phi}_k V_k Q_{k-1} e_k + c_k \bar{\phi}_{k+1} v_{k+1} \quad \text{by (46)}$$

$$= s_k^2 r_{k-1} + c_k \bar{\phi}_{k+1} v_{k+1} \quad \text{by (92).}$$

Finally, here is the unpreconditioned algorithm with no look-ahead or scaling from the original paper [FN91].

$$\begin{aligned} x_{0} &= 0 \\ \beta_{1} &= ||b||, \quad \tilde{\beta}_{1} &= \beta_{1} \\ v_{1} &= b/\beta_{1}, \quad \tilde{v}_{1} &= v_{1} \\ \rho_{1} &= \tilde{v}_{1}^{T} v_{1} \\ u_{1} &= A v_{1}, \quad \tilde{u}_{1} &= A^{T} \tilde{v}_{1} \\ c_{0} &= 1 \\ \bar{\phi}_{1} &= \beta_{1} \\ r_{-1,1} w_{-1} &= 0, \quad r_{0,1} w_{0} &= 0 \\ \text{for } j &= 1, 2, \dots \\ \alpha_{j} &= \tilde{v}^{T} u_{j} / \rho_{j} \\ v_{j+1} &= u_{j} - \alpha_{j} v_{j} \\ \beta_{j+1} &= ||v_{j+1}|| \\ v_{j+1} &= \tilde{v}_{j+1} / \beta_{j+1} \\ \tilde{v}_{j+1} &= \tilde{v}_{j+1} / \beta_{j+1} \\ \tilde{v}_{j+1} &= \tilde{v}_{j+1} / \tilde{\beta}_{j+1} \\ \rho_{j+1} &= \tilde{v}_{j+1} / \tilde{\beta}_{j+1} \\ \rho_{j+1} &= \tilde{v}_{j+1} / \rho_{j} \\ u_{j+1} &= A v_{j+1} - \gamma_{j+1} v_{j} \\ \eta_{j+1} &= \beta_{j+1} \rho_{j+1} / \rho_{j} \end{aligned}$$

$$\begin{split} \hat{u}_{j+1} &= A^T \tilde{v}_{j+1} - \tilde{\gamma}_{j+1} \tilde{v}_j \\ \text{if } j > 2 \\ r_{j-2,j} &= -\gamma_j s_{j-2} \\ \text{if } j > 1 \\ r_{j-1,j} &= \gamma_j c_{j-1} c_{j-2} - \alpha_j s_{j-1} \\ \tilde{r}_{j,j} &= \gamma_j s_{j-1} c_{j-2} + \alpha_j c_{j-1} \\ r_{j,j} &= \sqrt{\tilde{r}_{j,j}^2 + \beta_{j+1}^2} \\ c_j &= \tilde{r}_{j,j}/r_{j,j} \quad s_j = -\beta_{j+1}/r_{j,j} \\ \phi_j &= c_j \bar{\phi}_j, \quad \bar{\phi}_{j+1} = s_j \bar{\phi}_j \\ w_j &= (v_j - r_{j-2,j} w_{j-2} - r_{j-1,j} w_{j-1})/r_{j,j} \\ x_j &= x_{j-1} + \phi_j w_j \end{split}$$

check convergence; continue if necessary

end

.

### 4.8 Conjugate Gradients Squared (CGS)

#### Motivation

The Conjugate Gradients Squared (CGS) method was so named because in a sense (that will become clear) it gives the "square" of some quantities in the BiCG method, and the best introduction to it appears to be via a new CGS algorithm derived from the alternative unsymmetric Lanczos-BiCG algorithm of § 4.6, which we reproduce here for ease of reference. First the initial conditions are (where again we use  $1/\infty = 0$ ):

$$\rho_0 = \infty$$
,  $\beta_1 = 1$ ,  $x_0 = 0$ ,  $r_{-1} = 0 = \tilde{r}_{-1} = 0$ ,  $r_0 = \tilde{r}_0 = b$ .

Then for  $i = 1, 2, 3, \ldots$  we have

alternative unsymmetric Lanczos-BiCG

$$\rho_{i} = \tilde{r}_{i-1}^{T} r_{i-1}$$

$$\gamma_{i} = \beta_{i} \rho_{i} / \rho_{i-1} \text{ from } (44)$$

$$w_{i} = r_{i-1} + (\gamma_{i} / \beta_{i}) w_{i-1}$$

$$\tilde{w}_{i} = \tilde{r}_{i-1} + (\gamma_{i} / \beta_{i}) \tilde{w}_{i-1}$$

$$u_{i} = A r_{i-1}$$

$$\tilde{u}_{i} = A^{T} \tilde{r}_{i-1}$$

$$\alpha_{i} = \tilde{r}_{i-1}^{T} u_{i} / \rho_{i} \text{ from } (40)$$

$$\beta_{i+1} = -(\alpha_{i} + \gamma_{i})$$

$$x_{i} = x_{i-1} - (1 / \beta_{i+1}) w_{i}$$

$$r_{i} = (u_{i} - \alpha_{i} r_{i-1} - \gamma_{i} r_{i-2}) / \beta_{i+1}$$

$$\tilde{r}_{i} = (\tilde{u}_{i} - \alpha_{i} \tilde{r}_{i-1} - \gamma_{i} \tilde{r}_{i-2}) / \beta_{i+1}$$

We see from the initial conditions and last two lines that

------

$$\begin{aligned} r_0 &\equiv b, \\ r_k &= \varphi_k(A)r_0, \quad \tilde{r}_k = \varphi_k(A^T)r_0, \quad k \geq 0, \end{aligned}$$
where the polynomial in A,  $\varphi_k(A)$ , expressed as a scalar polynomial  $\varphi_k \equiv \varphi_k(\theta)$  is given by

$$\varphi_0 \equiv 1, \quad \varphi_1 = \frac{1}{\beta_2} (\theta - \alpha_1) \varphi_0,$$
  
$$\varphi_k = \frac{1}{\beta_{k+1}} [(\theta - \alpha_k) \varphi_{k-1} - \gamma_k \varphi_{k-2}], \quad k > 1.$$
(93)

Now since the BiCG residuals are constructed to be mutually orthogonal,

$$\tilde{r}_i^T r_j = r_0^T \varphi_i(A) \varphi_j(A) r_0 = 0 \quad \text{for} \quad i \neq j,$$

$$= r_0^T \varphi_i^2(A) r_0 \neq 0 \quad \text{for} \quad i = j.$$

It is the squared polynomial in this last line which is important in CGS. Notice in BiCG that multiplication by A and  $A^T$  was required in going from  $r_{i-1}$  to  $r_i$ ; that is from  $\varphi_{i-1}(A)r_0$  to  $\varphi_i(A)r_0$ . Using CGS we will see that it is possible to go from  $\varphi_{i-1}^2(A)r_0$  to  $\varphi_i^2(A)r_0$  using only two multiplications by A. The advantages of this are that multiplications by  $A^T$  are not required (and sometimes it is cumbersome to produce code to do such multiplications), and 2i matrix multiplications are required to produce a polynomial of degree 2i (whereas it is only degree i in BiCG). A disadvantage is that although  $\varphi_i(A)r_0$  has useful properties in BiCG, it is not clear what the useful properties of  $\varphi_i^2(A)r_0$  are. Nevertheless we now derive a new variant of CGS. Define

$$\tilde{r}^{CGS} \equiv r_0,$$

$$r_k^{CGS} \equiv \varphi_k^2(A)r_0, \quad k \ge 0,$$
and
$$s_k^{CGS} \equiv \varphi_k(A)\varphi_{k-1}(A)r_0, \quad k \ge 1.$$
Then
$$(\tilde{r}^{CGS})^T s_k^{CGS} = r_0^T \varphi_k(A)\varphi_{k-1}(A)r_0$$

$$= \tilde{r}_k^T r_{k-1} = 0, \quad k \ge 1.$$
(94)

÷.,

Now we show how these new polynomials  $\varphi_k^2$  and  $\varphi_k \varphi_{k-1}$  may be constructed

recursively, using (93). From the resulting expressions we give the corresponding recurrences for  $r_k^{CGS}$  and  $s_k^{CGS}$ .

$$\varphi_{1}\varphi_{0} = \varphi_{1} = \frac{1}{\beta_{2}}(\theta - \alpha_{1}), \quad \text{or} \quad s_{1}^{CGS} = \varphi_{1} = \frac{1}{\beta_{2}}(A - \alpha_{1}I)r_{0},$$
  

$$\varphi_{k}\varphi_{k-1} = \frac{1}{\beta_{k+1}}[(\theta - \alpha_{k})\varphi_{k-1}^{2} - \gamma_{k}\varphi_{k-1}\varphi_{k-2}], \quad k > 1,$$
  

$$s_{k}^{CGS} = \frac{1}{\beta_{k+1}}[(A - \alpha_{k}I)r_{k-1}^{CGS} - \gamma_{k}s_{k-1}^{CGS}], \quad k > 1. \quad (95)$$
  

$$\varphi_{1}^{2} = \frac{1}{\beta_{2}}(\theta - \alpha_{1})\varphi_{1}\varphi_{0}, \quad \text{or} \quad r_{1}^{CGS} = \frac{1}{\beta_{2}}(A - \alpha_{1}I)s_{1}^{CGS},$$
  

$$(96)$$

and for k > 1

$$\varphi_{k}^{2} = \frac{1}{\beta_{k+1}^{2}} [(\theta - \alpha_{k})^{2} \varphi_{k-1}^{2} - 2\gamma_{k}(\theta - \alpha_{k})\varphi_{k-1}\varphi_{k-2} + \gamma_{k}^{2} \varphi_{k-2}^{2}],$$

$$= \frac{1}{\beta_{k+1}^{2}} [\beta_{k+1}(\theta - \alpha_{k})\varphi_{k}\varphi_{k-1} - \gamma_{k}(\theta - \alpha_{k})\varphi_{k-1}\varphi_{k-2} + \gamma_{k}^{2} \varphi_{k-2}^{2}],$$

$$r_{k}^{CGS} = \frac{1}{\beta_{k+1}^{2}} [\beta_{k+1}(A - \alpha_{k}I)s_{k}^{CGS} - \gamma_{k}(A - \alpha_{k}I)s_{k-1}^{CGS} + \gamma_{k}^{2}r_{k-2}^{CGS}]. \quad (97)$$

The key to this new CGS algorithm is that the  $\alpha_k$ ,  $\gamma_k$  and  $\beta_{k+1}$  are exactly those in the Lanczos-BiCG algorithm above, and may be computed from the CGS vectors as follows.

$$\rho_{i} = \tilde{r}_{i-1}^{T} r_{i-1} = \tilde{r}_{0}^{T} \varphi_{i-1}^{2} (A) r_{0} = (\tilde{r}^{CGS})^{T} r_{i-1}^{CGS}, \quad i \ge 1,$$
(98)  

$$\gamma_{1} = 0, \quad \gamma_{i} = \beta_{i} \rho_{i} / \rho_{i-1}, \quad i > 1,$$
  

$$\alpha_{i} = \tilde{r}_{i-1}^{T} A r_{i-1} / \rho_{i} = \tilde{r}_{0}^{T} \varphi_{i-1} (A) A \varphi_{i-1} (A) r_{0} / \rho_{i}$$
  

$$= (\tilde{r}^{CGS})^{T} A \varphi_{i-1}^{2} (A) r_{0} / \rho_{i} = (\tilde{r}^{CGS})^{T} A r_{i-1}^{CGS} / \rho_{i}, \quad i \ge 1$$
  

$$= (\tilde{r}^{CGS})^{T} (A r_{i-1}^{CGS} - \gamma_{i} s_{i-1}^{CGS}) / \rho_{i}, \quad i \ge 2,$$
(99)

since  $(\tilde{r}^{CGS})^T s_{i-1}^{CGS} = 0$  for  $i \ge 2$ . Finally  $\beta_2 = -\alpha_1$  and  $\beta_{i+1} = -(\alpha_i + \gamma_i)$ , for i > 1, are immediately available.

We can now combine the above expressions to give the full computation for the  $r_k^{CGS}$  and  $s_k^{CGS}$  vectors — we call this the "CGS process" – without yet showing the computation for the solution, which is given later.

#### New CGS method

<u>Problem</u>: Solve Ax = b for x when A is general nonsingular.

<u>CGS Process</u>: We drop the superscript CGS from  $s_j$ , and write  $\tilde{v} \equiv \tilde{r}^{CGS} = r_0$ ;  $v_{j+1} \equiv r_j^{CGS}$ ,  $j \ge 0$ , and for efficient computation define  $\hat{u}_1 \equiv Av_1$ ;  $\hat{u}_j \equiv Av_j - \gamma_j s_{j-1}$  for j > 1, and  $t_j \equiv As_j$  for  $j \ge 1$ . Again,  $1/\infty = 0$ .

$$\rho_{0} \equiv \infty, \quad \beta_{1} = 1, \quad s_{0} \equiv 0, \quad v_{0} \equiv 0$$

$$v_{1} = \tilde{v} = b$$
for  $j = 1, 2, ...$ 

$$\rho_{j} = \tilde{v}^{T} v_{j} \quad \text{from (98)}$$

$$\gamma_{j} = \frac{\rho_{j}}{\rho_{j-1}} \beta_{j}$$

$$\tilde{u}_{j} = A v_{j} - \gamma_{j} s_{j-1}$$

$$\alpha_{j} = \tilde{v}^{T} \hat{u}_{j} / \rho_{j} \quad \text{from (99)}$$

$$\beta_{j+1} = -(\alpha_{j} + \gamma_{j}) \quad (100)$$

$$\beta_{j+1} s_{j} = \hat{u}_{j} - \alpha_{j} v_{j} \quad \text{from (95)} \quad (101)$$

$$p_{j+1}s_j = u_j - u_j v_j \quad \text{from (95)}$$

$$t_j = As_j \quad (101)$$

$$\beta_{j+1}v_{j+1} = t_j - \alpha_j s_j - \frac{\gamma_j}{\beta_{j+1}}(t_{j-1} - \alpha_j s_{j-1}) + \frac{\gamma_j^2}{\beta_{j+1}}v_{j-1}$$
(102)

this last coming from (97). A healthy aspect of this is that, as required in (94),

$$\beta_{j+1}\tilde{v}^T s_j = \tilde{v}^T \hat{u}_j - \alpha_j \tilde{v}^T v_j = \alpha_j \rho_j - \alpha_j \rho_j = 0, \quad j \ge 1,$$
(103)

follows immediately from the *j*th step without any recourse to additional theory or induction.

Despite the elegance of the above motivation for the new CGS method and process, the reader should be warned that the following is essentially "work in progress", and the matrix representation is somewhat complicated, and probably not optimal.

#### Matrix representation

If we use (101) for the odd columns and (102) for the even columns the CGS process becomes, in matrix form

$$\vec{V}_{2k} \equiv \begin{bmatrix} \vec{v}_{1} & \vec{v}_{2} & \cdots & \vec{v}_{2k-1} & \vec{v}_{2k} \end{bmatrix}$$

$$= \begin{bmatrix} v_{1} & s_{1} & \cdots & v_{k} & s_{k} \end{bmatrix}$$

$$A \vec{V}_{2k} G_{2k} = \vec{V}_{2k} F_{2k} + \beta_{k+1} \vec{v}_{2k+1} e_{2k}^{T} = \vec{V}_{2k+1} F_{2k+1,2k}$$

$$(105)$$
where
$$G_{2} \equiv \begin{bmatrix} 1 \\ & 1 \end{bmatrix}, \quad G_{2k} \equiv \begin{bmatrix} G_{2k-2} & & -\frac{\gamma_{k}}{\beta_{k+1}} \\ & 1 & 1 \end{bmatrix}, \quad k > 1$$

$$F_{2} \equiv \begin{bmatrix} \alpha_{1} \\ \beta_{2} & \alpha_{1} \end{bmatrix}, \quad k > 1$$

$$F_{2k} \equiv \begin{bmatrix} r_{2k-2} & & -\frac{\gamma_{k}}{\beta_{k+1}} \\ & \gamma_{k} & -\frac{\alpha_{k}\gamma_{k}}{\beta_{k+1}} \\ & \beta_{k} & \alpha_{k} \\ & \beta_{k+1} & \alpha_{k} \end{bmatrix}, \quad k > 1$$

$$(106)$$

$$F_{2k+1,2k} \equiv \begin{bmatrix} F_{2k} \\ \beta_{k+1} e_{2k}^{T} \end{bmatrix}.$$

Using (105) and solving for  $F_{2k}\bar{y}_{2k} = e_1$ , we consider iterates  $x_k$  not quite the same as in (47),

$$x_{k} = \bar{x}_{2k} = \bar{V}_{2k}G_{2k}\bar{y}_{2k}$$
(107)  

$$\Rightarrow r_{k} = b - Ax_{k} = b - A\bar{x}_{2k} = b - A\bar{V}_{2k}G_{2k}\bar{y}_{2k}$$
  

$$= b - (\bar{V}_{2k}F_{2k} + \beta_{k+1}\bar{v}_{2k+1}e_{2k}^{T})\bar{y}_{2k}$$

$$= \bar{V}_{2k}(e_1 - F_{2k}\bar{y}_{2k}) - \beta_{k+1}v_{k+1}e_{2k}^T\bar{y}_{2k}$$
$$= -\beta_{k+1}\bar{y}_{2k,2k}v_{k+1}.$$
 (108)

<u>Subproblem</u>: Solve  $F_{2k}\bar{y}_{2k} = e_1$ , where  $F_{2k}$  is upper Hessenberg as in (106). This does not have a solution if  $F_{2k}$  is singular, and at that time this implementation of the CGS method would fail. However, we continue assuming that this is not the case.

<u>Implementation</u>: This may be solved using an LDU decomposition of  $F_{2k}$ . The one which suits our purposes is shown below.

$$\begin{split} F_{2} &= L_{2}\bar{D}_{2}\bar{U}_{2} \\ &= \begin{bmatrix} 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} -\beta_{2} \\ -\beta_{2} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -\beta_{2} \\ \beta_{2} & -\beta_{2} \end{bmatrix} \\ F_{2k} &= L_{2k}\bar{D}_{2k}\bar{U}_{2k} \\ L_{2k} &= \begin{bmatrix} L_{2k-2} \\ -1 & 1 \\ -1 & 1 \end{bmatrix}, \quad \bar{D}_{2k} = \begin{bmatrix} \bar{D}_{2k-2} \\ -\beta_{k+1} \\ -\beta_{k+1} \end{bmatrix} \\ \bar{U}_{2k} &= \begin{bmatrix} \bar{U}_{2k-2} \\ -\beta_{k+1} \\ -\beta_{k+1} \\ -\beta_{k+1} \end{bmatrix} \\ = F_{2k} &= \begin{bmatrix} L_{2k-2}\bar{D}_{2k-2} \bar{U}_{2k-2} \\ -\frac{\gamma_{k}}{\beta_{k}} \\ -\frac{\gamma_{k}}{\beta_{k+1}} \\ -\gamma_{k} \\ -\gamma_{k}$$

Our normalizers  $\beta_{j+1}$  are now verified by equating the right sides of the

above equations with  $F_2$  and  $F_{2k}$ . So

$$\alpha_1 = -\beta_2; \quad \alpha_k = -(\gamma_k + \beta_{k+1}), \quad k > 1$$
  
$$\Rightarrow \beta_2 \equiv -\alpha_1; \quad \beta_{k+1} \equiv -(\alpha_k + \gamma_k), \quad k > 1.$$
(109)

This was the choice of  $\beta_{j+1}$  that we made in § 4.6. As long as it is nonzero, it is a satisfactory normalizer for the process, and  $\tilde{U}_{2k+2}$  exists. For ease of calculation, we make the following definitions.

$$\begin{split} \bar{W}_{2k} &\equiv \begin{bmatrix} \bar{w}_1 & \bar{w}_2 & \cdots & \bar{w}_{2k-1} & \bar{w}_{2k} \end{bmatrix} \\ &\equiv \begin{bmatrix} u_1 & q_1 & \cdots & u_k & q_k \end{bmatrix} \equiv \bar{V}_{2k} G_{2k} \bar{U}_{2k}^{-1} \\ \\ \text{and} \quad \bar{z}_{2k} &\equiv \begin{bmatrix} \bar{\zeta}_1 \\ \bar{\zeta}_2 \\ \vdots \\ \bar{\zeta}_{2k} \end{bmatrix} \equiv \bar{U}_{2k} \bar{y}_{2k} \\ &\Rightarrow \bar{\zeta}_{2k} &= \bar{y}_{2k,2k}. \end{split}$$

So our system is transformed to

$$L_{2k}\bar{D}_{2k}\bar{z}_{2k}=e_1, \quad \bar{x}_{2k}=\bar{W}_{2k}\bar{z}_{2k}.$$

Solving

$$L_{2k}\bar{D}_{2k}\bar{z}_{2k} = \begin{bmatrix} -\beta_2\bar{\zeta}_1 \\ \beta_2\bar{\zeta}_1 - \beta_2\bar{\zeta}_2 \\ \beta_2\bar{\zeta}_2 - \beta_3\bar{\zeta}_3 \\ \vdots \\ \beta_k\bar{\zeta}_{2k-2} - \beta_{k+1}\bar{\zeta}_{2k-1} \\ \beta_{k+1}\bar{\zeta}_{2k-1} - \beta_{k+1}\bar{\zeta}_{2k} \end{bmatrix} = e_1$$

$$\Rightarrow \quad \bar{\zeta}_1 = -\frac{1}{\beta_2}$$

$$\zeta_1 = \bar{\zeta}_2 = \bar{\zeta}_1 = -\frac{1}{\beta_2}$$

$$\bar{\zeta}_{2k-1} = \frac{\beta_k}{\beta_{k+1}} \bar{\zeta}_{2k-2}, \quad k > 1$$

$$\zeta_k \equiv \bar{\zeta}_{2k} = \bar{\zeta}_{2k-1} = \frac{\beta_k}{\beta_{k+1}} \bar{\zeta}_{2k-2} \quad (110)$$

$$= \frac{\beta_{k-1}}{\beta_{k+1}} \bar{\zeta}_{2k-4} = \dots = -\frac{1}{\beta_{k+1}}, \quad k > 1.$$

Consider the odd columns of  $\bar{W}_{2k}\bar{U}_{2k}=\bar{V}_{2k}G_{2k}$ .

$$\Rightarrow u_{1} = \bar{w}_{1} = \bar{v}_{1} = v_{1}$$
$$u_{k} = \bar{w}_{2k-1} = \bar{v}_{2k-1} + \frac{\gamma_{k}}{\beta_{k}} \bar{w}_{2k-2} = v_{k} + \frac{\gamma_{k}}{\beta_{k}} q_{k-1}, \quad k > 1.$$
(111)

The even columns give a less pleasing result.

$$q_1 = \bar{w}_2 = s_1$$

$$q_k = \bar{w}_{2k} = s_k - \frac{\gamma_k}{\beta_{k+1}}(s_{k-1} + u_k) + \frac{\gamma_k}{\beta_k}(q_{k-1} - \frac{\gamma_k}{\beta_{k+1}}u_{k-1}), \quad k > 1.$$

Now 
$$r_k = -\beta_{k+1} \bar{y}_{2k,2k} v_{k+1} = -\beta_{k+1} \bar{\zeta}_{2k} v_{k+1} = v_{k+1}$$
 (112)

and 
$$x_k = \bar{x}_{2k-2} + \bar{\zeta}_{2k-1}u_k + \bar{\zeta}_{2k}q_k = x_{k-1} - \frac{1}{\beta_{k+1}}(u_k + q_k)$$
 (113)

so this CGS process produces vectors  $v_{j+1}$  which are the residual vectors  $r_j$  of the method.

It remains to examine under what conditions  $\beta_{j+1}$  can be zero. Supose  $\beta_2, \ldots, \beta_j$  obtained from (109) are nonzero. The with  $e^T \equiv [1, \ldots, 1]$ , we have

$$e^{T}F_{2j-1} = [0, \dots, -\frac{\gamma_{j-1}^{2}}{\beta_{j}} - \frac{\alpha_{j-1}\gamma_{j-1}}{\beta_{j}} + \alpha_{j-1} + \beta_{j}, \alpha_{j} + \gamma_{j}]$$

$$e^{T}F_{2j-1}e_{2j-2} = \frac{-\gamma_{j-1}^{2} - \alpha_{j-1}\gamma_{j-1} - \alpha_{j-1}(\alpha_{j-1} + \gamma_{j-1}) + (\alpha_{j-1} + \gamma_{j-1})^{2}}{\beta_{j}}$$

$$= 0$$

$$\Rightarrow e^{T}F_{2j-1} = [0, \dots, 0, \alpha_{j} + \gamma_{j}],$$

which will be nonzero if  $F_{2j-1}$  is nonsingular. However if  $\alpha_j + \gamma_j = 0$  then  $F_{2j-1}$  is singular, and both this and the standard CGS method fails.

Otherwise it is clear that the  $L_{2k}\bar{D}_{2k}\bar{U}_{2k}$  'factorization' is possible. However small  $\beta_k$  could cause large element growth in  $\bar{U}_{2k}$ .

Let us introduce vectors  $P_k \equiv [p_1, p_2, \dots, p_k]$ , not used in the algorithm, but useful for this analysis.

$$P_{k} \equiv \bar{W}_{2k}M_{2k,k}E_{k}^{-1}$$
where  $M_{2k,k} \equiv \begin{bmatrix} M_{2k-2,k-1} & & \\ & & & \\ & & & \\ \hline \end{array} \end{array} \\ \hline \end{array} \\ \hline & & & \\ \hline \hline \\ \hline & & & \\ \hline \end{array} \end{array}$ 

Equivalence of this method with the vector form: The CGS method — slightly altered from 'Templates' -- is as follows.

$$r_{0} = b$$
  

$$\dot{r} = r_{0}$$
  
for  $i = 1, 2, ...$   

$$\rho_{i-1} = \tilde{r}^{T} r_{i-1}$$
  
if  $\rho_{i-1} = 0$ , method fails  
if  $i = 1$   

$$u_{i} = r_{i-1}$$
  

$$p_{i} = u_{i}$$
  
clse  

$$\beta_{i-1} = \rho_{i-1}/\rho_{i-2}$$

$$u_{i} = r_{i-1} + \beta_{i-1}q_{i-1}$$
$$p_{i} = u_{i} + \beta_{i-1}(q_{i-1} + \beta_{i-1}p_{i-1})$$

$$\hat{v} = Ap_i$$

$$\alpha_i = \rho_{i-1} / \hat{r}^T \hat{v}$$

$$q_i = u_i - \alpha_i \hat{v}$$

$$\hat{u} = u_i + q_i$$

$$x_i = x_{i-1} + \alpha_i \hat{u}$$

$$\hat{q} = A\hat{u}$$

$$r_i = r_{i-1} - \alpha_i \hat{q}$$
check convergence; continue if necessary

end

To avoid confusion we will use the superscript  $^{C}$  to distinguish 'Templates' terms from those here, where it will also make life easier if we reindex two of the 'Templates' terms as follows:

$$\rho_i^C \equiv \text{`Templates'}\rho_{i-1},$$
  
$$\beta_i^C \equiv \text{`Templates'}\beta_{i-1}.$$

We now write the algorithms side by side for comparison, where in this new CGS algorithm we replace  $v_{i+1}$  by  $r_i$ , see (112). To handle the i = 1 case we define

CGS	new CGS
$\rho_0^C = 1$	$ \rho_0 = \infty, \beta_1 = 1 $
$p_0^C = 0, q_0^C = 0$	$u_0 = 0, p_0 = 0, s_0 = 0, q_0 = 0$

Then for  $i = 1, 2, 3, \ldots$  we have

CGSnew CGS
$$\rho_i^C = \tilde{r}^{CT} r_{i-1}^C$$
 $\rho_i = \tilde{r}^T r_{i-1}$  $\beta_i^C = \rho_i^C / \rho_{i-1}^C$  $\gamma_i = \beta_i \rho_i / \rho_{i-1}$  $u_i^C = r_{i-1}^C + \beta_i^C q_{i-1}^C$  $u_i = r_{i-1} + \frac{\gamma_i}{\beta_i} q_{i-1}$ 

$$\begin{split} p_i^C &= u_i^C + \beta_i^C (q_{i-1}^C + \beta_i^C p_{i-1}^C) & p_i = u_i + \frac{\gamma_i}{\beta_i} (q_{i-1} + \frac{\gamma_i}{\beta_i} p_{i-1}) \\ \alpha_i^C &= \rho_i^C / \hat{r}^{CT} A p_i^C & \alpha_i = \tilde{r}^T (A r_{i-1} - \gamma_i s_{i-1}) / \rho_i \\ \beta_{i+1} &= -(\alpha_i + \gamma_i) \\ s_i &= (A r_{i-1} - \alpha_i r_{i-1} - \gamma_i s_{i-1}) / \beta_{i+1} \\ q_i^C &= u_i^C - \alpha_i^C A p_i^C & q_i = s_i - \frac{\gamma_i}{\beta_{i+1}} (s_{i-1} + u_i) \\ &+ \frac{\gamma_i}{\beta_i} (q_{i-1} - \frac{\gamma_i}{\beta_{i+1}} u_{i-1}) \\ x_i^C &= x_{i-1}^C + \alpha_i^C (u_i^C + q_i^C) & x_i = x_{i-1} - 1 / \beta_{i+1} (u_i + q_i) \\ r_i^C &= r_{i-1}^C - \alpha_i^C A (u_i^C + q_i^C) & r_i = \frac{1}{\beta_{i+1}} (A s_i - \alpha_i s_i) \\ &- \frac{\gamma_i}{\beta_{i+1}^2} (A s_{i-1} - \alpha_i s_{i-1}) + \frac{\gamma_i^2}{\beta_{i+1}^2} r_{i-2} \end{split}$$

There are several differences between the two algorithms — they don't even use all the same vectors — the new CGS uses the  $s_k$  which are orthogonal to  $\hat{r}$ , while CGS produces the  $p_k^C$ . The other vectors are the same; the exact relationships are

$$\gamma_{i}/\beta_{i} = \beta_{i}^{C}, \quad i > 1$$

$$\rho_{i} = \rho_{i}^{C}, \quad -1/\beta_{i+1} = \alpha_{i}^{C}$$

$$u_{i} = u_{i}^{C}, \quad p_{i} = p_{i}^{C}$$

$$s_{i} = r_{i-1}^{C} - \alpha_{i}^{C}Au_{i}^{C}, \quad q_{i} = q_{i}^{C}$$

$$x_{i} = x_{i}^{C}, \quad r_{i} = r_{i}^{C}.$$
(115)

A straightforward examination shows this is true for i = 1. Suppose it is true for i = 1, 2, ..., k - 1. Then it is clear that

$$\rho_k = \rho_k^C, \quad \gamma_k / \beta_k = \beta_k^C, \quad u_k = u_k^C, \quad p_k = p_k^C.$$

We now need only prove

$$-1/\beta_{k+1} = \alpha_k^C \tag{116}$$

$$s_k = r_{k-1}^C - \alpha_k^C A u_k^C \tag{117}$$

 $q_k = q_k^C$ 

so that  $x_k = x_k^C$ . The equality of  $r_k$  and  $r_k^C$  follows since we have shown that  $r_k$  is the residual of  $x_k$  and  $r_k^C$  is obviously the residual of  $x_k^C$  by examining the last two lines of CGS. We can prove (116) using our matrix formulation of the new CGS algorithm. Consider

$$\begin{split} \tilde{v}^{T}AP_{k} &= \tilde{v}^{T}A\bar{W}_{2k}M_{2k,k}E_{k}^{-1} \\ &= \tilde{v}^{T}A\bar{V}_{2k}G_{2k}\bar{U}_{2k}^{-1}M_{2k,k}E_{k}^{-1} \\ &= \tilde{v}^{T}\bar{V}_{2k+1}F_{2k+1,2k}\bar{U}_{2k}^{-1}M_{2k,k}E_{k}^{-1} \\ &= \tilde{v}^{T}\bar{V}_{2k+1}L_{2k+1,2k}\bar{D}_{2k}M_{2k,k}E_{k}^{-1} \\ &= \left[ \rho_{1} \quad 0 \quad \rho_{2} \quad \cdots \quad 0 \quad \rho_{k+1} \right] L_{2k+1,2k}\bar{D}_{2k}M_{2k,k}E_{k}^{-1}. \end{split}$$

The (i, i)- element of this is

$$\tilde{v}^{T}Ap_{i} = \begin{bmatrix} -\rho_{1}\beta_{2} & \rho_{2}(\gamma_{2} - \beta_{3}) & \cdots & \rho_{k}(\gamma_{k} - \beta_{k+1}) \end{bmatrix} \begin{bmatrix} (\frac{\gamma_{2}\cdots\gamma_{k}}{\beta_{3}\cdots\beta_{k}})^{2} \\ (\frac{\gamma_{1}}{\beta_{3}})^{2} \\ \vdots \\ (\frac{\gamma_{k}}{\beta_{k}})^{2} \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

 $= -\rho_i \beta_{i+1}$ 

after much cancellation. Now since  $\tilde{r}^C = \tilde{v}$  and  $p_i^C = p_i$ , i = 1, ..., k,

$$\tilde{r}^{CT}Ap_i^C = \tilde{v}^T Ap_i = -\rho_i \beta_{i+1}, \quad i = 1, \dots, k.$$

From the CGS algorithm we then see

$$\alpha_i^C = \rho_i^C / \bar{r}^{CT} A p_i^C = \rho_i / \bar{r}^{CT} A p_i^C = -1 / \beta_{i+1}, \quad i = 1, \dots, k,$$

and so (116) holds. Now, to show  $s_k = r_{k-1}^C - \alpha_k^C A u_k^C = r_{k-1} + 1/\beta_{k+1} A u_k$ , look at

$$Au_i = A\bar{W}_{2k}e_{2i-1}$$

$$= \bar{V}_{2k+1} L_{2k+1,2k} \bar{D}_{2k} c_{2i-1}$$
$$= -\beta_{i+1} (v_i - s_i).$$

So  $s_i$  is alternatively (and more simply) given by the above expression and (117) falls out. Finally,  $q_k = q_k^C$ . For lack of a simpler way to show this, consider the matrix equation

$$AP_{k} = A\bar{W}_{2k}M_{2k,k}E_{k}^{-1}$$
$$= \bar{W}_{2k+1}\bar{U}_{2k+1}G_{2k+1}^{-1}L_{2k+1,2k}\bar{D}_{2k}M_{2k,k}E_{k}^{-1}.$$

It is not obvious, but this has *i*-th column

$$Ap_i = \tilde{W}_{2k+1}\beta_{i+1}(e_{2i} - e_{2i-1})$$
$$= \beta_{i+1}(q_i - u_i), \quad i = 1, \dots, k$$
$$\Rightarrow q_k = u_k + 1/\beta_{k+1}Ap_k$$
$$= u_k^C - \alpha_k^C Ap_k^C.$$

The induction proof (115) is complete, showing that in theory these algorithms not only compute the same iterates, but also several of the intermediate quantities are identical. However the algorithms are clearly not identical, as the use of different vectors shows, and may have different numerical behaviours. This suggests a matrix development of CGS would be useful, but we leave this for future work.

Reference: [Son89]

### 4.9 BiConjugate Gradients Stabilized (BiCGSTAB)

#### Motivation

The BiConjugate Gradients Stabilized (BiCGSTAB) method was so named because it increases the convergence of BiCG but in a more stable way than blindly "squaring" the residual vectors  $r_i$  as in CGS. Instead we form  $r_i = \varphi_i(A)\theta_i(A)r_0$ , where  $\theta_i(\theta) \equiv (1-\omega_i\theta)\theta_{i-1}$  is an *i*th degree polynomial describing a steepest descent update —  $\omega_i$  is chosen to minimize the norm of  $r_i$  — in the hope that these new iterates will converge more regularly than the  $\varphi_i^2(A)r_0$  of CGS. Here is our new variant of BiCGSTAB. Using  $\varphi_i \equiv \varphi_i(\theta)$  from § 4.8, define

$$\theta_{0} \equiv 1, \quad \theta_{k} = (1 - \omega_{k}\theta)\theta_{k-1}, \quad k \geq 1, \quad (113)$$

$$\tilde{r}^{BCS} \equiv r_{0},$$

$$r_{k}^{BCS} \equiv \varphi_{k}(A)\theta_{k}(A)r_{0}, \quad k \geq 0,$$

$$s_{k}^{BCS} \equiv \varphi_{k}(A)\theta_{k-1}(A)r_{0}, \quad k \geq 1,$$
and 
$$u_{k}^{BCS} \equiv \varphi_{k-1}(A)\theta_{k}(A)r_{0}, \quad k \geq 1.$$
Then  $(\tilde{r}^{BCS})^{T}s_{k}^{BCS} = r_{0}^{T}\varphi_{k}(A)\theta_{k-1}(A)r_{0}$ 

$$= \tilde{r}_{k}^{T}(\theta_{k-1}(A)r_{0}) = 0 \quad (119)$$
since  $\tilde{r}_{k} \perp \theta_{k-1}(A)r_{0} \in \mathcal{K}^{k}(A;r_{0}).$ 

These new polynomials,  $\varphi_{k-1}\theta_k, \varphi_k\theta_{k-1}$  and  $\varphi_k\theta_k$ , may be constructed recursively using (93) and (118). From these expressions we give the corresponding recurrences for  $r_k^{BCS}, s_k^{BCS}$  and  $u_k^{BCS}$ .

$$\varphi_{k-1}\theta_k = (1 - \omega_k \theta)\varphi_{k-1}\theta_{k-1}, \quad k \ge 1$$
  
or  $u_k^{BCS} = (I - \omega_k A)r_{k-1}^{BCS}, \quad k \ge 1.$   
 $\varphi_1\theta_0 = \frac{1}{\beta_2}(\theta - \alpha_1)\varphi_0\theta_0, \quad \text{or} \quad s_1^{BCS} = \frac{1}{\beta_2}\dots \alpha_1 I)r_0,$   
 $\varphi_k\theta_{k-1} = \frac{1}{\beta_{k+1}}[(\theta - \alpha_k)\varphi_{k-1}\theta_{k-1} - \gamma_k\varphi_{k-2}\theta_{k-1}], \quad k > 1$ 

or 
$$s_k^{BCS} = \frac{1}{\beta_{k+1}} [(A - \alpha_k I) r_{k-1}^{BCS} - \gamma_k u_{k-1}^{BCS}], \quad k > 1.$$
  
 $\varphi_0 \theta_0 = 1, \quad \text{or} \quad r_0^{BCS} = r_0,$   
 $\varphi_k \theta_k = (1 - \omega_k \theta) \varphi_k \theta_{k-1}, \quad \text{or} \quad r_k^{BCS} = (I - \omega_k A) s_k^{BCS}, \quad k \ge 1.$ 

As in CGS,  $\alpha_k$ ,  $\beta_k$  and  $\gamma_k$  are scalars from the unsymmetric Lanczos-BiCG algorithm and may be computed from the BiCGSTAB vectors by expanding the polynomials involved and using the simple fact

$$\tilde{r}_i^T \psi_j(A) r_0 = 0$$

for any polynomial  $\psi_j$  of degree j < i.

Define 
$$\tilde{\beta}_{i} \equiv \prod_{j=1}^{i} \beta_{j+1}, \quad \tilde{\omega}_{i} \equiv \prod_{j=1}^{i} \omega_{j}.$$
  
 $\rho_{i} = \tilde{r}_{i-1}^{T} r_{i-1} = \tilde{r}_{i-1}^{T} (\frac{1}{\beta_{i-1}} A^{i-1} r_{0} + \cdots)$   
 $= \frac{1}{\beta_{i-1}} \tilde{r}_{i-1}^{T} A^{i-1} r_{0}, \quad i \ge 1.$   
 $\rho_{i}^{BCS} = (\tilde{r}^{BCS})^{T} r_{i-1}^{BCS}$   
 $= (\tilde{r}^{BCS})^{T} \varphi_{i-1}(A) \theta_{i-1}(A) r_{0}$   
 $= \tilde{r}_{i-1}^{T} (\tilde{\omega}_{i-1} A^{i-1} r_{0} + \cdots)$   
 $= \tilde{\omega}_{i-1} \tilde{r}_{i-1}^{T} A^{i-1} r_{0}.$   
 $\Rightarrow \rho_{i} = \rho_{i}^{BCS} / (\tilde{\omega}_{i-1} \tilde{\beta}_{i-1}).$   
 $\gamma_{1} = 0$   
 $\gamma_{i} = \beta_{i} \rho_{i} / \rho_{i-1} = \rho_{i}^{BCS} / (\rho_{i-1}^{BCS} \omega_{i-1}), \quad i > 1.$   
Now  $(\tilde{r}^{BCS})^{T} (Ar_{i-1}^{BCS} - \gamma_{i} u_{i-1}^{BCS})$   
 $= (\tilde{r}^{BCS})^{T} (\beta_{i+1} \varphi_{i}(A) + \alpha_{i}(A) \varphi_{i-1}(A)) \theta_{i-1}(A) r_{0}$   
 $= \beta_{i+1} \tilde{r}_{i}^{T} \theta_{i-1}(A) r_{0} + \alpha_{i} \tilde{r}_{i-1}^{T} \theta_{i-1}(A) r_{0}$ 

$$= \alpha_{i} \bar{\omega}_{i-1} \hat{r}_{i-1}^{T} A^{i-1} r_{0}$$

$$= \alpha_{i} \bar{\omega}_{i-1} \bar{\beta}_{i-1} \rho_{i}$$

$$= \alpha_{i} \rho_{i}^{BCS}$$

$$\Rightarrow \alpha_{i} = (\hat{r}^{BCS})^{T} (A r_{i-1}^{BCS} - \gamma_{i} u_{i-1}^{BCS}) / \rho_{i}^{BCS}.$$
(120)

Finally  $\beta_2 = -\alpha_1$  and  $\beta_{i+1} = -(\alpha_i + \gamma_i)$ , for i > 1, are immediately available.

We can now combine the above expressions to give the full computation for the  $r_k^{BCS}$ ,  $s_k^{BCS}$  and  $u_k^{BCS}$  vectors — we call this the "BiCGSTAB process" — without yet showing the computation for the solution, which is given later.

New BiCGSTAB method

<u>Problem</u>: Solve Ax = b for x when A is general nonsingular.

<u>BiCGSTAB Process</u>: We drop the superscript BCS from  $s_j$  and  $u_j$ , and write  $\tilde{v} \equiv \tilde{r}^{BCS} = r_0; v_{j+1} \equiv r_j^{BCS}, \quad j \ge 0$ , and for more efficient computation define  $q_1 \equiv Av_1; \quad q_j \equiv Av_j - \gamma_j u_{j-1}$  for j > 1 and  $t_j \equiv As_j$  for  $j \ge 1$ . Again,  $1/\infty = 0$ .

$$\rho_{0} \equiv \infty, \quad \omega_{0} \equiv \infty, \quad \gamma_{1} \equiv 0, \quad u_{0} \equiv 0$$

$$v_{1} = \tilde{v} = b \qquad (121)$$
for  $j = 1, 2, ...$ 

$$\rho_{j} = \tilde{v}^{T} v_{j}$$

$$\gamma_{j} = -\rho_{j}/(\rho_{j-1}\omega_{j-1})$$

$$q_{j} = Av_{j} - \gamma_{j}u_{j-1} \qquad (122)$$

$$\alpha_{j} = \tilde{v}^{T} q_{j}/\rho_{j},$$

$$\beta_{j+1} = -(\alpha_{j} + \gamma_{j}) \qquad (123)$$

$$\beta_{j+1}s_{j} = q_{j} - \alpha_{j}v_{j} \qquad (124)$$

$$t_{j} = As_{j}$$

$$\omega_{j} = t_{j}^{T}s_{j}/t_{j}^{T}t_{j}$$

$$u_j = v_j - \omega_j A v_j \tag{125}$$

$$v_{j+1} = s_j - \omega_j t_j.$$
 (126)

Notice that, as required in (119),

$$\beta_{j+1}\tilde{v}^T s_j = \tilde{v}^T q_j - \alpha_j \tilde{v}^T v_j = \alpha_j \rho_j - \alpha_j \rho_j = 0, \quad j \ge 1,$$

follows immediately from the jth step.

We are again obliged to warn the reader that what follows is "work in progress", and the matrix representation is somewhat complicated, and probably not optimal.

#### Matrix Representation

If we use (122) and (125) substituted into (124) for the odd columns and (126) for the even columns the BiCGSTAB process becomes, in matrix form

$$\bar{V}_{2k} \equiv \begin{bmatrix} v_1 & s_1 & \cdots & v_k & s_k \end{bmatrix}$$

$$A\bar{V}_{2k}G_{2k} = \bar{V}_{2k}F_{2k} - \frac{1}{\omega_k}v_{k+1}e_{2k}^T = \bar{V}_{2k+1}F_{2k+1,2k} \quad (127)$$
where
$$G_2 \equiv \begin{bmatrix} 1 \\ & 1 \end{bmatrix}, \quad G_{2k} \equiv \begin{bmatrix} G_{2k-2} & \omega_{k-1}\gamma_k \\ & 0 \\ \hline & 1 \\ \hline & 1 \end{bmatrix}, \quad k > 1$$

$$F_2 \equiv \begin{bmatrix} \alpha_1 \\ \beta_2 & \frac{1}{\omega_1} \end{bmatrix},$$

$$F_{2k} \equiv \begin{bmatrix} \sigma_{1} \\ \beta_2 & \frac{1}{\omega_1} \end{bmatrix},$$

$$F_{2k-2} = \begin{bmatrix} \gamma_k \\ 0 \\ \hline & 0 \end{bmatrix}, \quad k > 1 \quad (128)$$

ak

 $\beta_{k+1}$ 

 $\frac{1}{\omega_k}$ 

 $-\frac{1}{\omega_{k-1}}$ 

$$F_{2k+1,2k} \cong \begin{bmatrix} F_{2k} \\ -\frac{1}{\omega_k} c_{2k}^T \end{bmatrix}.$$

Using (127) and solving for  $F_{2k}\bar{y}_{2k} = e_1$ , we consider iterates  $x_k$  as in (107),

$$x_{k} = \bar{x}_{2k} = \bar{V}_{2k}G_{2k}\bar{y}_{2k}$$

$$\Rightarrow r_{k} = b - Ax_{k} = b - A\bar{x}_{2k} = b - A\bar{V}_{2k}G_{2k}y_{2k}$$

$$= b - (V_{2k}F_{2k} - \frac{1}{\omega_{k}}v_{k+1}c_{2k}^{T})\bar{y}_{2k}$$

$$= \bar{V}_{2k}(e_{1} - F_{2k}\bar{y}_{2k}) + \frac{1}{\omega_{k}}v_{k+1}e_{2k}^{T}\bar{y}_{2k}$$

$$= \frac{\bar{y}_{2k,2k}}{\omega_{k}}v_{k+1}.$$
(129)

<u>Subproblem</u>: Solve  $F_{2k}\bar{y}_{2k} = e_1$ , where  $F_{2k}$  is upper Hessenberg as in (128). This does not have a solution if  $F_{2k}$  is singular, and at that time this implementation of the BiCGSTAB method would fail. However, we continue assuming that this is not the case.

<u>Implementation</u>: This may be solved using an LDU decomposition  $F_{2k}$ . The one which suits our purposes is shown below.

$$F_{2} = L_{2}\bar{D}_{2}\bar{U}_{2}$$

$$= \begin{bmatrix} 1\\ -1 & 1 \end{bmatrix} \begin{bmatrix} -\beta_{2} \\ \frac{1}{\omega_{1}} \end{bmatrix} \begin{bmatrix} 1\\ 1 \end{bmatrix} = \begin{bmatrix} -\beta_{2} \\ \beta_{2} & \frac{1}{\omega_{1}} \end{bmatrix}.$$
 (130)
$$F_{2k} = L_{2k}\bar{D}_{2k}\bar{U}_{2k}$$

$$L_{2k} = \begin{bmatrix} L_{2k-2} \\ -1 & 1\\ -1 & 1 \end{bmatrix}, \quad \bar{D}_{2k} = \begin{bmatrix} \bar{D}_{2k-2} \\ -\beta_{k+1} \\ \frac{1}{\omega_{k}} \end{bmatrix},$$

$$\tilde{U}_{2k} = \begin{bmatrix} \tilde{U}_{2k-2} & -\frac{\gamma_k}{\beta_k} \\ & \omega_{k-1}\gamma_k \\ \hline & 1 \\ & 1 \end{bmatrix}$$

$$\Rightarrow F_{2k} = \begin{bmatrix} L_{2k-2}\bar{D}_{2k-2}\bar{U}_{2k-2} & \gamma_k \\ & 0 \\ \hline & -\frac{1}{\omega_{k-1}} & -(\gamma_k + \beta_{k+1}) \\ & \beta_{k+1} & \frac{1}{\omega_k} \end{bmatrix}.$$

Our normalizers  $\beta_{j+1}$  are now verified by equating the right sides of the above equations with  $F_2$  and  $F_{2k}$ . So

$$\alpha_1 = -\beta_2; \quad \alpha_k = -(\gamma_k + \beta_{k+1}), \quad k > 1$$
  
$$\Rightarrow \beta_2 \equiv -\alpha_1; \quad \beta_{k+1} \equiv -(\alpha_k + \gamma_k), \quad k > 1.$$
(131)

This was the choice of  $\beta_{k+1}$  that we made in § 4.6. As long as it is nonzero, it is a satisfactory normalizer for the process, and  $\bar{U}_{2k+2}$  exists. For ease of calculation, we make the following definitions.

$$\bar{W}_{2k} \equiv \begin{bmatrix} \bar{w}_1 & \bar{w}_2 & \cdots & \bar{w}_{2k-1} & \bar{w}_{2k} \end{bmatrix} \equiv \bar{V}_{2k} G_{2k} \bar{U}_{2k}^{-1}$$
and
$$\bar{z}_{2k} \equiv \begin{bmatrix} \bar{\zeta}_1 \\ \bar{\zeta}_2 \\ \vdots \\ \bar{\zeta}_{2k} \end{bmatrix} \equiv \bar{U}_{2k} \bar{y}_{2k}$$

$$\Rightarrow \bar{\zeta}_{2k} = \bar{y}_{2k,2k}.$$
(132)

So our system is transformed to

$$L_{2k}\bar{D}_{2k}\bar{z}_{2k}=e_1, \quad \bar{x}_k=\bar{W}_{2k}\bar{z}_{2k}.$$

Solving

$$\begin{split} \mathcal{I}_{2k}\bar{D}_{2k}\bar{z}_{2k} &= \begin{pmatrix} -\beta_2\bar{\zeta}_1 \\ \beta_2\bar{\zeta}_1 + \frac{1}{\omega_1}\bar{\zeta}_2 \\ -\frac{1}{\omega_1}\bar{\zeta}_2 - \beta_3\bar{\zeta}_3 \\ \vdots \\ -\frac{1}{\omega_1}\bar{\zeta}_{2k-2} - \beta_{k+1}\bar{\zeta}_{2k-1} \\ \beta_{k+1}\bar{\zeta}_{2k-1} + \frac{1}{\omega_k}\bar{\zeta}_{2k} \end{bmatrix} &= c_1 \\ \Rightarrow \bar{\zeta}_1 &= -1/\beta_2, \quad \bar{\zeta}_2 = \omega_1 \\ \bar{\zeta}_{2k-1} &= -\frac{1}{\beta_{k+1}\omega_{k-1}}\bar{\zeta}_{2k-2}, \quad k > 1 \\ \bar{\zeta}_{2k} &= -\beta_{k+1}\omega_k\bar{\zeta}_{2k-1} = \frac{\omega_k}{\omega_{k-1}}\bar{\zeta}_{2k-2} = \omega_k, \quad k > 1 \\ \bar{\zeta}_{2k-1} &= -\frac{1}{\beta_{k+1}}, \quad k > 1. \end{split}$$

Consider the odd and even columns of  $\bar{W}_{2k}\bar{U}_{2k} = \bar{V}_{2k}G_{2k}$ .

$$\bar{w}_1 = v_1 \tag{133}$$

$$\bar{w}_{2k-1} = \omega_{k-1} \gamma_k (v_{k-1} - s_{k-1}) + v_k + \frac{\kappa}{\beta_k} \bar{w}_{2k-3}, \quad k > 1$$
(134)  
$$\bar{w}_{2k} = s_k, \quad k \ge 1.$$

This causes us to redefine  $\bar{W}_{2k} \equiv \begin{bmatrix} p_1 & s_1 & \cdots & p_k & s_k \end{bmatrix}$  and rewrite (133) and (134) as

$$p_{1} = v_{1}$$

$$p_{k} = v_{k} + \omega_{k-1}\gamma_{k}(v_{k-1} - s_{k-1}) + \frac{\gamma_{k}}{\beta_{k}}p_{k-1}, \quad k > 1.$$
Now  $r_{k} = \frac{\bar{y}_{2k,2k}}{\bar{y}_{2k,2k}} = \frac{\bar{\zeta}_{2k}}{\bar{\zeta}_{2k}}v_{k+1} = v_{k+1}$ 
(135)

Now 
$$r_k = \frac{g_{2k,2k}}{\omega_k} = \frac{\zeta_{2k}}{\omega_k} v_{k+1} = v_{k+1}$$
 (135)  
and  $x_k = \bar{W}_{2k-1} \bar{z}_{2k-2} + \bar{\zeta}_{2k-1} \bar{w}_{2k-1} + \bar{\zeta}_{2k} \bar{w}_{2k}$   
 $= x_{k-1} - \frac{1}{\beta_{k+1}} p_k + \omega_k s_k$  (136)

so this BiCGSTAB process produces vectors  $v_{j+1}$  which are the residual vectors  $r_j$  of the method.

It remains to examine under what conditions  $\beta_{j+1}$  can be zero. Supose  $\beta_2, \ldots, \beta_j$  obtained from (131) are nonzero. The with  $e^T \equiv [1, \ldots, 1]$ , we have

$$e^{T}F_{2j-1} = [0, \dots, \frac{1}{\omega_{j-1}} - \frac{1}{\omega_{j-1}}, \alpha_{j} + \gamma_{j}]$$
  
=  $[0, \dots, 0, \alpha_{j} + \gamma_{j}],$ 

which will be nonzero if  $F_{2j-1}$  is nonsingular. However if  $\alpha_j + \gamma_j = 0$ then  $F_{2j-1}$  is singular, and both this and the standard CGS method fails. Otherwise it is clear that the  $L_{2k}\bar{D}_{2k}\bar{U}_{2k}$  'factorization' is possible. However small  $\beta_k$  could cause large element growth in  $\overline{U}_{2k}$ .

Equivalence of this method with the vector form: The BiCGSTAB method - slightly altered from 'Templates' - is as follows.

$$r_{0} = b$$
  

$$\tilde{r} = r_{0}$$
  
for  $i = 1, 2, ...$   

$$\rho_{i-1} = \tilde{r}^{T} r_{i-1}$$
  
if  $\rho_{i-1} = 0$ , method fails  
if  $i = 1$   

$$p_{i} = r_{i-1}$$
  
else  

$$\beta_{i-1} = (\rho_{i-1}/\rho_{i-2})(\alpha_{i-1}/\omega_{i-1})$$
  

$$\rho_{i} = r_{i-1} + \beta_{i-1}(p_{i-1} - \omega_{i-1}w_{i-1})$$
  

$$w_{i} = Ap_{i}$$
  

$$\alpha_{i} = \rho_{i-1}/\tilde{r}^{T}w_{i}$$
  

$$s_{i} = r_{i-1} - \alpha_{i}w_{i}$$
  
check norm of s; if small enough: set

t  $x_i = x_{i-1} + \alpha_i p_i$  and stop

$$t_i = As_i$$
$$\omega_i = t_i^T s_i / t_i^T t_i$$

 $\begin{aligned} x_i &= x_{i+1} + \alpha_i p_i + \omega_i s_i \\ r_i &= s_i - \omega_i t_i \\ \text{check for convergence; continue if necessary} \\ \text{for continuation it is necessary that } \omega_i \neq 0 \end{aligned}$ 

end

To avoid confusion we will use the superscript  $^{B}$  to distinguish 'Templates' terms from those here, where it will also make life easier if we re-index two of the 'Templates' terms as follows:

$$\rho_i^B \equiv \text{`Templates'}\rho_{i-1},$$
  
 $\beta_i^B \equiv \text{`Templates'}\beta_{i-1}.$ 

We now write the algorithms side by side for comparison, where in this new BiCGSTAB algorithm we replace  $v_{i+1}$  by  $r_i$ , see (135). To handle the i = 1 case we define

BiCGSTABnew BiCGSTAB $\rho_0^B = 1, \alpha_0^B = 0, \omega_0^B = 0$  $\rho_0 = 1, \beta_1 = 1, \omega_0 = 0$  $p_0^B = 0$  $v_0 = 0, s_0 = 0, p_0 = 0, u_0 = 0$ 

Then for i = 1, 2, 3, ... we have

BICGSTAB	new BiCGSTAB
$\rho_i^B = \tilde{r}^{BT} r_{i-1}^B$	$ ho_i =  ilde{v}^T r_{i-1}$
if $i = 1$	if $i = 1$
$\beta_1^B = 0$	$\gamma_1 = 0$

else

$$\beta_{i}^{B} = (\rho_{i}^{B} / \rho_{i-1}^{B}) (\alpha_{i-1}^{B} / \omega_{i-1}^{B})$$
$$\rho_{i}^{B} = r_{i-1}^{B} + \beta_{i}^{B} (p_{i-1}^{B} - \omega_{i-1}^{B} A p_{i-1}^{B}) \qquad p_{i} =$$

$$\gamma_i = -(\rho_i/\rho_{i-1})(1/\omega_{i-1})$$
$$p_i = r_{i-1} + \frac{\gamma_i}{\beta_i} p_{i-1}$$
$$+\omega_{i-1}\gamma_i(r_{i-2} - s_{i-1})$$

else

$$\begin{split} \alpha_i^B &= \rho_i^B / \hat{r}^{BT} A p_i^B & \alpha_i &= \hat{v}^T (Ar_{i-1} - \gamma_i u_{i-1}) / \rho_i \\ \beta_{i+1} &= -(\alpha_i + \gamma_i) \\ s_i^B &= r_{i-1}^B - \alpha_i^B A p_i^B & s_i &= (Ar_{i-1} - \alpha_i r_{i-1} - \gamma_i u_{i-1}) / \beta_{i+1} \\ \omega_i^B &= (As_i^B)^T s_i^B / (As_i^B)^T (As_i^B) & \omega_i &= (As_i)^T s_i / (As_i)^T (As_i) \\ u_i &= r_{i-1} - \omega_i A r_{i-1} \\ x_i^B &= x_{i-1}^B + \alpha_i^B p_i^B + \omega_i^B s_i^B & x_i &= x_{i-1} - 1 / \beta_{i+1} p_i + \omega_i s_i \\ r_i^B &= s_i^B - \omega_i^B A s_i^B & r_i &= s_i - \omega_i A s_i \end{split}$$

There are noticeable differences between the two algorithms — the new BiCGSTAB algorithm has more complicated expressions for  $p_i$  and  $s_i$  and also forms an intermediate vector  $u_i$ , which the 'Templates' algorithm does not. However relationships exist between the two; the exact equivalences are

$$\gamma_i / \beta_i = \beta_i^B, \quad i > 1$$

$$\rho_i = \rho_i^B, \quad -1/\beta_{i+1} = \alpha_i^B, \quad \omega_i = \omega_i^B$$

$$p_i = p_i^B, \quad s_i = s_i^B, \quad u_i = r_{i-1}^B - \omega_i A r_{i-1}^B$$

$$x_i = x_i^B, \quad r_i = r_i^B.$$
(137)

A straightforward examination shows this is true for i = 1. Suppose it is true for i = 1, 2, ..., k - 1. Then it is clear that

$$\rho_k = \rho_k^B, \quad \gamma_k / \beta_k = \beta_k^B.$$

The following will be useful in proving the rest.

$$P_{k} = \tilde{W}_{2k}I_{2k,k}$$
where  $I_{2k,k} \equiv \begin{bmatrix} I_{2k-2,k-1} \\ & 1 \end{bmatrix}$ 

To show  $p_k = p_k^B$ , consider the *i*-th column of (139).

$$Ap_{i} = \beta_{i+1}(s_{i} - r_{i-1}), \quad i = 1, 2, ..., k$$
(140)  

$$p_{k} = r_{k-1} + \omega_{k-1}\gamma_{k}(r_{k-2} - s_{k-1}) + \gamma_{k}/\beta_{k}p_{k-1}$$
  

$$= r_{k-1}^{B} - \omega_{k-1}^{B}\gamma_{k}/\beta_{k}Ap_{k-1} + \beta_{k}^{B}p_{k-1}^{B}$$
  

$$= r_{k-1}^{B} - \omega_{k-1}^{B}\beta_{k}^{B}Ap_{k-1}^{B} + \beta_{k}^{B}p_{k-1}^{B} = p_{k}^{B}.$$

To show  $-1/\beta_{k+1} = \alpha_k^B$ , consider the scalar product of  $\tilde{v}$  with (140)

$$\tilde{v}^T A p_i = \beta_{i+1} \tilde{v}^T (s_i - r_{i-1})$$
  
=  $-\rho_i \beta_{i+1}, \quad i = 1, 2, \dots, k.$   
So  $\alpha_k^B = \rho_k^B / \tilde{r}^T A p_k^B = \rho_k / \tilde{v}^T A p_k = -1 / \beta_{k+1}$ 

(140) makes it easy to prove  $s_k = s_k^B$ .

$$s_k^B = r_{k-1}^B - \alpha_k^B A p_k^B = r_{k-1} + 1/\beta_{k+1} A p_k$$
$$= r_{k-1} + s_k - r_{k-1} = s_k.$$

Finally, it is obvious that

$$\omega_k = \omega_k^B, \quad x_k = x_k^B, \quad r_k = r_k^B.$$

The induction proof (137) is complete, showing that in theory these algorithms not only compute the same iterates, but also several of the intermediate quantities are identical. However the algorithms are clearly not identical, as the use of different vectors shows, and may have different numerical behaviours. This suggests a matrix development of BiCGSTAB would be useful, but we leave this for future work. Reference: [vdV92]

94

# 5 <u>Conclusion</u>

From the point of view of developing the research, the most interesting and satisfactory applications of the formalism stated in § 1.1 were to the methods of Conjugate Gradients and BiConjugate Gradients. In these, connections with the underlying symmetric and unsymmetric Lanczos processes respectively were exposed. In each case two methods were developed — one with unit norm vectors spanning  $\mathcal{K}^k(A; b)$  and one which produced the residual vectors of the method to span the Krylov subspace. In the case of BiCG, the Krylov process which produced the residual vectors was rather surprisingly shown to use the same coefficients in the formation of the "shadow residuals", as opposed to the familiar unsymmetric Lanczos process which "transposed" the coefficients. It was the second method in each case which was shown to be mathematically equivalent to the "vector" form, since the resulting algorithm more closely resembled the vector form. An understanding of the quantities used in the vector form of the algorithm was obtained and in general the algorithm was clarified.

The methods of SYMMLQ, MINRES, LSQR, and GMRES were originally implicitly developed using the formalism. Thus the Krylov processes and associated subproblems were easy to determine for these and there was no difference between the algorithm derived from the matrix development and the corresponding vector form.

The method of Quasi-Minimal Residuals as developed in [FN91] also implicitly used the formalism. The Krylov process and subproblem were obvious and implementation details were also given in the paper. However, the 'Templates' version of QMR differs from that in the original paper — although theoretically the two are the same, the vector form has been derived from a subproblem whose solution implementation is less strightforward, in the interests of time and storage of the algorithm. As future work, a matrix development of that 'Templates' vector form of QMR to determine the Krylov process and subproblem and hence show the resulting method is QMR would make § 4.7 more complete.

Attempts to apply the formalism to CGS and BiCGSTAB, algorithms which double the complexity of BiCG in an attempt to halve the time, required a small modification to the approach suggested in § 4. Since these two methods produce residual vectors which are elements of every other Krylov subspace, it was necessary to consider auxiliary vectors, elements of the missed subspaces, as well. Using both sets of vectors, we achieved a matrix equation of the form

$$A\tilde{V}_{K}G_{K} = \tilde{V}_{K+1}F_{K+1,K}$$

$$\Rightarrow A\tilde{V}_{K} = \tilde{V}_{K+1}H_{K+1,K} \qquad (141)$$
where  $H_{K+1,K} \equiv F_{K+1,K}G_{K}^{-1}$  is upper Hessenberg
and  $K \equiv 2k$ .

The form of (141) is exactly the same as for the general Krylov process (3) in § 2.

The next difference came in the iterates we used. The original idea was to form iterates  $x_k = V_k y_k$ , but here we chose

$$x_k = \bar{x}_K = \bar{V}_K G_K \bar{y}_K \tag{142}$$

$$= \bar{V}_K \bar{y}_K \tag{143}$$

where 
$$\tilde{y}_K = G_K \bar{y}_K$$
,

thus a simple matter of relabelling (142) gives (143) just as (47) in § 4. With this slightly different approach, attempts to reveal the process and subproblem were successful although not as elegant as might have been hoped. As was stated in those two subsections, the work done is really "work in progress" and certainly further improvements may be made.

On the whole the formalism appears to be very useful in the analysis of algorithms and definitely makes the more mysterious "vector" forms much easier to understand and learn. Unfortunately, as was observed in the cases of CGS and BiCGSTAB, it appears somewhat difficult in the more complex algorithms to determine an optimally clear, simple and elegant division into process and problem.

.....

## References

- [Arn51] W.E. Arnoldi. The principle of minimized iterations in the solution of the matrix eigenvalue problem. Quart. Appl. Math., 9:17-29, 1951.
- [BBC+94] R. Barrett, M. Berry, T.F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. van der Vorst. Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods. SIAM Publications, Philadelphia, Pennsylvania 19104-2688, 1994. There may be a later version.
- [FN91] R.W. Freund and N.M. Nachtigal. QMR: A quasi-minimal residual method for non-hermitian linear systems. Numer. Math., 60:315-339, 1991.
- [GK65] G.H. Golub and W. Kahan. Calculating the singular values and pseudoinverse of a matrix. SIAM J. Numer. Anal., 2:205-224, 1965.
- [GL89] Gene H. Golub and Charles F. Van Loan. Matrix Computations, second edition. The Johns Hopkins University Press, Baltimore, Maryland 21211, 1989.
- [Hou74] A.S. Householder. The Theory of Matrices in Numerical Analysis. Dover Publications, New York, 1974. Originally published by Blaisdell, NY, 1964.
- [HS52] M.R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. Journal of Research of the National Bureau of Standards, 49:409-436, 1952.

- [Kry31] A.N. Krylov. O čislennom rešenii uravnenija, kotorym v techničeskih voprasah opredeljajutsja častoty malyh kolebanii material'nyh sistem. Izv. Akad. Nauk SSSR Otd. Mat. Estest., pages 491-539, 1931.
- [Lan50] C. Lanczos. An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. J. Res. Nat. Bur. Stand., 45:255-282, 1950.
- [Lan52] C. Lanczos. Solution of systems of linear equations by minimized iterations. J. Res. Nat. Bur. Stand., 49:33-53, 1952.
- [Pai71] C.C. Paige. The Computation of Eigenvalues and Eigenvectors of Vcry Large Sparse Matrices. PhD thesis, London University, London, England, 1971.
- [Pai72] C.C. Paige. Computational variants of the Lanczos method for the eigenproblem. J. Inst. Math. Appl., 10:373-381, 1972.
- [Pai94] C.C. Paige. Krylov subspace processes, Krylov subspace methods, and iteration polynomials. In M. Chu, R. Plemmons, D. Brown, and D. Ellison, editors, Proceedings of the Lanczos Centenary Conference, Raleigh, NC, Dec. 1993, pages 83-92, Philadelphia, Pa, 1994. SIAM Publications.
- [PPdV95] C.C. Paige, B.N. Parlett, and H.A. Van der Vorst. Approximate solutions and eigenvalue bounds from Krylov subspaces. Numerical Linear Algebra with Applications, 1995. To appear, May 1995.

- [PS75] C.C. Paige and M.A. Saunders. Solution of sparse indefinite systems of linear equations. SIAM J. Numer. Anal., 12:617-629, 1975.
- [PS82] C.C. Paige and M.A. Saunders. LSQR: An algorithm for sparse linear equations and sparse least squares. ACM Transactions on Mathematical Software, 8:43-71, 1982.
- [Rci71] J.K. Reid. On the method of conjugate gradients for the solution of large sparse linear equations. In J.K. Reid, editor, Large Sparse Sets of Linear Equations, pages 231-254. Academic Press, New York, 1971.
- [Son89] P. Sonneveld. CGS: A fast lanczos-type solver for nonsymmetric linear systems. SIAM J. Sci. Stat. Comput., 10:36-52, 1989.
- [SSS6] Y. Saad and M.H. Schultz. GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems. SIAM J. Sci. Stat. Comput., 7:856-869, 1986.
- [vdV92] H.A. van der Vorst. Bi-CGSTAB: A fast and smoothly converging variant of Bi-CG for the solution of nonsymmetric linear systems. SIAM J. Sci. Stat. Comput., 13:631-644, 1992.