

KRYLOV SPACE METHODS ON STATE-SPACE CONTROL MODELS*

*Daniel L. Boley*¹

Abstract. We give an overview of various Lanczos/Krylov space methods and the way in which they are being used for solving certain problems in Control Systems Theory based on state-space models. The matrix methods used are based on Krylov sequences and are closely related to modern iterative methods for standard matrix problems such as sets of linear equations and eigenvalue calculations. We show how these methods can be applied to problems in Control Theory such as controllability, observability, and model reduction. All the methods are based on the use of state-space models, which may be very sparse and of high dimensionality. For example, we show how one may compute an approximate solution to a Lyapunov equation arising from a discrete-time linear dynamic system with a large sparse system matrix by the use of the Arnoldi algorithm, and so obtain an approximate Gramian matrix. This has applications in model reduction. The close relation between the matrix Lanczos algorithm and the algebraic structure of linear control systems is also explored.

1. Introduction

Most computational methods currently used for control problems cannot handle very effectively very large problems with hundreds or thousand of states. In many instances, very large state-space models are constructed from physical considerations. In order to manipulate these models effectively, it is usually necessary to either form reduced order state-space models or to compute a frequency domain transfer function description of the system. Both of these methods entail some loss of information. This paper is an attempt to show how techniques designed for classic problems of very large dimensionality in linear algebra may be applicable to certain problems in control. These techniques allow one to operate directly on state-space models without needing to construct frequency domain transfer function descriptions. For example, we may compute a reduced order model for a system in a computationally efficient manner without computing a transfer function

* Received April 2, 1992; revised June 1, 1992; accepted July 20, 1992.

¹ Department of Computer Science, University of Minnesota, Minneapolis, Minnesota 55455.

description, but operating directly on the original state-space model of possibly very large dimensionality.

Many large Control problems involve state-space representations of large dimensionality. A major problem in Control is to reduce the dimensionality of the models (“model reduction”) or just to analyze the stability or other properties of such large systems [21]. In the literature of computational methods in linear algebra, much study has been done on methods for solving standard linear algebra problems involving large sparse matrix operators: linear equations using conjugate gradients, least squares also with conjugate gradients, eigenvalues using Lanczos methods, singular values, etc. Many of the techniques are mentioned in [29].

It is only recently that the “technology” developed for large sparse linear algebra problems has been applied to control problems of large dimensionality [20]. Most of the methods are based on the recursive generation of Krylov spaces. These methods are based on the idea of projecting the problem onto ever-expanding subspaces generated by the matrices occurring in the problem itself. It has been shown that for many linear algebra problems, the convergence rate is much faster than what would be expected from arbitrary projections. Krylov space methods also allow iteration by expanding the dimension of the target space in a very natural way. The algorithms based on recursive generation of Krylov spaces can be referred to collectively as Lanczos-type algorithms, and are the basis for this paper.

In this paper, we show some existing examples where Lanczos-type algorithms have been applied to some standard problems in Control, in particular model reduction and controllability/observability. The rest of this paper is organized as follows: In Section 2 we describe the basic Lanczos algorithms, the Arnoldi algorithm and the block variants. In Sections 3–5 we show how these algorithms may be used to solve certain model reduction problems in Control. We also show how certain theories of Linear Control Systems actually contribute to our knowledge of the behavior of the nonsymmetric Lanczos algorithm.

2. Krylov sequence methods

Most iterative methods for large linear algebra problems are based on the recursive generation of so-called Krylov spaces of ever-increasing dimensions, and the projection of the original matrix operator(s) onto these Krylov spaces. The success of these algorithms comes from the fact that as the Krylov space dimension increases, the new projection can be obtained very quickly from the previous projection onto the next lower dimensional Krylov space. In the domain of Control, the system matrix for a model can be projected in this way. In fact this is the idea behind most of the methods in this paper.

In this section, we define a Krylov space and describe some of the basic algorithms for computing bases for the Krylov space and the projection of a given

matrix operator onto these spaces. In the next section, we apply these algorithms to some standard control problems.

A Krylov sequence is a sequence of vectors generated by a matrix as follows. Given an $(n \times n)$ -matrix A and a vector \mathbf{x} , the k th Krylov sequence $K(A, \mathbf{x}, k)$ is a sequence of k column vectors:

$$K(A, \mathbf{x}, k) \equiv (\mathbf{x}, A\mathbf{x}, A^2\mathbf{x} \cdots A^{k-1}\mathbf{x}).$$

As k increases, we get a sequence of sequences $K(A, \mathbf{x}, 0) \equiv \mathbf{x}$, $K(A, \mathbf{x}, 1) = (\mathbf{x}, A\mathbf{x}), \dots$. A block Krylov sequence is generated by a matrix A and an $(n \times p)$ -matrix X as follows

$$K(A, X, k) \equiv (X, AX, A^2X, \dots, A^{k-1}X),$$

and the corresponding column space is called the k th Krylov space and is denoted by $\mathcal{K}(A, X, k)$.

The basic theorem for Krylov sequences shows how the rank increases as the vectors are appended.

Theorem 1. *Given any $(n \times n)$ -matrix A and $(n \times p)$ -matrix X ,*

$$\text{rank } K(A, X, k) = \text{rank } K(A, X, k + 1)$$

if and only if

$$\begin{aligned} \text{rank } K(A, X, k) &= \max_j \text{rank } K(A, X, j) \equiv \text{rank } K(A, X, n) \\ &\equiv \text{rank } K(A, X, \infty). \end{aligned}$$

Furthermore, the space $\mathcal{K}(A, X, \infty)$ is invariant under A .

Proof. If $\text{rank } K(A, X, k) = \text{rank } K(A, X, k + 1)$, then $A^k X$ is a linear combination of $X, \dots, A^{k-1} X$. But then $A^{k+j} X$ is the same linear combination of $A^j X, \dots, A^{k+j-1} X$ for every $j \geq 0$. Hence by induction, $A^{k+j} X$ is a linear combination of $X, \dots, A^{k-1} X$ for all $j \geq 0$. The maximal rank is achieved in at most n steps. \square

As the vectors X, AX, A^2X, \dots are generated and appended to form ever larger and larger Krylov sequences, the theorem states that the rank increases at every step until it reaches a maximal value.

The Lanczos-type algorithms are algorithms for recursively generating bases for the Krylov Spaces, together with a matrix of recurrence coefficients. This latter matrix will be in reduced form, such as tridiagonal, and represents the projection of the original matrix operator A onto the Krylov space, typically of smaller dimension.

The Lanczos algorithm was originally proposed by Lanczos [46] as a method for the computation of eigenvalues of symmetric and nonsymmetric matrices. The idea was to reduce a general matrix to tridiagonal form, from which the

eigenvalues could be easily determined. For symmetric matrices, the Lanczos algorithm has been studied extensively [17], [52], and the convergence of the algorithm, when used to compute eigenvalues, has been extensively analyzed in [42], [51], [56], [60], and [68, p. 270ff]. This algorithm is particularly suited for large sparse matrix problems. A block Lanczos analog has been studied and analyzed in [28], [17], [52]. However, until recently, the nonsymmetric Lanczos algorithm has received much less attention. Some recent computational experience with this algorithm can be found in [16]. In addition to some numerical stability problems, the method suffered from the possibility of an incurable breakdown from which the only way to “recover” was to restart the whole process from the beginning with different starting vectors [68, p. 388ff]. More recently, several modifications allowing the Lanczos process to continue after such breakdowns have been proposed in [33], [34], [54], and a numerical implementation has been developed in [23]. The close connection between the modified non-symmetric Lanczos algorithm and orthogonal polynomials with respect to indefinite inner products is discussed in [7], [8], [27]. Recently, in [10], [53] the close relation was observed independently between the Lanczos algorithm and the controllability-observability structure of dynamical systems.

All of the above papers address the Lanczos algorithm with single starting vectors. Recently, a block nonsymmetric Lanczos algorithm was proposed in [43], [44] that is capable of starting with several starting vectors, analogous to the block Arnoldi algorithm, but so far the breakdown situation has not been addressed for the block case.

The Lanczos algorithm [46] is an example of a method that generates bases for Krylov subspaces starting with a given vector. The Arnoldi algorithm [4] can be thought of as a “one-sided” method, which generates one sequence of vectors that span the reachable space.

2.1. Arnoldi algorithm

The first algorithm we will describe is the Arnoldi algorithm [4], [9], [68], which is a recursive way to generate an orthonormal basis for the Krylov space generated by a given matrix A and vectors X . It will be seen that it is also a way to reduce a given matrix to block upper Hessenberg form H . The algorithm proceeds by recursively filling in the first few columns in the relation

$$AX = XH \text{ subject to } X^T X = I, \quad (1)$$

where X is an $(n \times r_{\max})$ -matrix of orthonormal columns spanning the Krylov space of maximal rank r_{\max} and H is an $(r_{\max} \times r_{\max})$ block upper Hessenberg matrix. Unless the starting vector is deficient in certain eigendirections (normally an unusual circumstance), $r_{\max} = n$. The following description is taken from [9], suitably modified to use modified Gram-Schmidt orthogonalization [29, p. 218] as suggested in [59].

Block Arnoldi Algorithm [9].

0. Start with $(n \times n)$ -matrix A and $(n \times p)$ -matrix X .
 {Generate orthonormal vectors \mathbf{X} and block upper Hessenberg matrix \mathbf{H} }
1. Factor $X_0 R = \text{QR}$ factorization of X . {normalization}
2. For $k = 1, 2, 3, \dots$, while X_k has nonzero rank,
3. Set $X_k^{(0)} = AX_{k-1}$; {expand Krylov space}
4. For $j = 1, \dots, k$, {modified Gram-Schmidt orthogonalization}
5. Set $H_{k-j,k-1} = X_{k-j}^T X_k^{(j-1)}$; {get coeffs. to enforce ortho. cond.}
6. Set $X_k^{(j)} = X_k^{(j-1)} - X_{k-j} H_{k-j,k-1}$, {enforce ortho. cond.}
7. Factor $X_k H_{k,k-1} = \text{QR}$ factorization of $X_k^{(k)}$; {normalize next columns}
8. Collect generated vectors $\mathbf{X} = (X_0, X_1, \dots)$ and coefficients $\mathbf{H} = (H_{ij})$.

The QR factorization used in steps 1 and 7 is used to compute an orthonormal basis for the column space of a given matrix. The particular formulation we use is, given an $n \times m$, rank r matrix M with $r \leq m \leq n$, we compute an $(n \times r)$ -matrix Q of orthonormal columns and an upper triangular $(r \times m)$ -matrix R such that $M = QR$. The method used can be either based on a Gram-Schmidt process or can be one that applies a series of Givens rotations or Householder transformations [29]. With this formulation, the diagonal blocks H_{kk} will be square of progressively smaller (or nonincreasing) dimensions as k increases, and the off-diagonal blocks will be rectangular.

It is useful to explain the steps individually, since the same process is used in all the Krylov space algorithms. The purpose of step 3 is to generate the next set of vectors expanding the Krylov space. The purpose of steps 4–6 is to orthogonalize the result against all the previous vectors generated. The process used is a modified Gram-Schmidt orthogonalization, which is mathematically equivalent to the ordinary Gram-Schmidt process, but behaves better numerically [29]. Numerical experience [59], [57] has also shown that it is useful to repeat steps 4–6 to ensure that the orthogonality condition is satisfied to the precision of the computer. The rank of $X_k^{(k)}$ may be less than that of $X_k^{(0)}$, in which case the matrix block H_{kk} generated by the method will not be square. In fact the method proceeds by reducing the dimension of the blocks progressively to zero. Step 7 is used to orthonormalize the vectors once they are made orthogonal to the preceding vectors.

Let us examine the situation during intermediate steps of the algorithm. We denote the items generated after k steps by

$$\mathbf{X}_k \equiv (X_0, X_1, \dots, X_{k-1}) \text{ and } \mathbf{H}_k \equiv \begin{pmatrix} H_{00} & H_{01} & \cdots & H_{0,k-1} \\ H_{10} & H_{11} & \ddots & H_{1,k-1} \\ & \ddots & \ddots & \vdots \\ 0 & & H_{k-1,k-2} & H_{k-1,k-1} \end{pmatrix}.$$

A simple induction argument then demonstrates that at each intermediate stage k , the items generated by that stage satisfy

$$\text{COLSP } \mathbf{X}_k = \mathcal{K}(A, X, k), \quad (2)$$

$$\mathbf{X}_k^T A \mathbf{X}_k = \mathbf{H}_k, \quad (3)$$

and

$$A \mathbf{X}_k = \mathbf{X}_k \mathbf{H}_k + X_k H_{k,k-1} (0, \dots, 0, I), \quad (3)$$

for every k , where the I in (4) is a square identity matrix of dimension equal to the number of columns in X_k . (Here $\text{COLSP } M$ denotes the *column space* of M .) When the maximal rank is reached, the term in (4) involving $X_k H_{k,k-1}$ disappears.

Two special cases for this algorithm are worthy of note. The first is the special case when the starting vectors X consist of just a single vector. In this case, the QR factorization (used in steps 1 and 7) reduces to a simple normalization of the vector to unit norm, and the resulting matrix \mathbf{H} will be “scalar” upper Hessenberg. The algorithm is considerably simpler, since all the “ H -blocks” are just scalars. The algorithm then proceeds as long as the result from step 6 is nonzero. This algorithm is called simply the *Arnoldi algorithm*, to distinguish it from the *block* version above.

The second special case is as follows. If the matrix A is symmetric, then so will be the generated \mathbf{H} . In other words, \mathbf{H} will be (block) tridiagonal. This implies that in step 5, $H_{k-j,k-1}$ will be zero for $j \geq 3$, so that the loop in steps 4–6 need be carried out only for $j = 1, 2$. However, numerical experience has shown that it is often necessary to reorthogonalize the vectors against all the previous vectors when computing in approximate floating point arithmetic (see the extensive literature reported in [29, Ch. 9]). The resulting algorithm is the *Lanczos algorithm*, and indeed it is historically the first recursive Krylov sequence method proposed [46].

2.2. Nonsymmetric Lanczos algorithm

If one starts with the symmetric Lanczos algorithm and relaxes the condition that the generated vectors be orthonormal, but maintains the condition that the generated matrix of coefficients be tridiagonal, one obtains the *nonsymmetric Lanczos algorithm*, capable of reducing most any matrix to tridiagonal form.

We present three variations of the nonsymmetric Lanczos algorithm. To aid in the exposition, we first present the version assuming no breakdown occurs. Then we present the modifications necessary to handle breakdowns. Finally we present the generalizations to obtain the “block” algorithm. The simple scalar algorithm starts with a general matrix A and two starting vectors $\mathbf{x}_0, \mathbf{y}_0$ such that $\mathbf{y}_0^T \mathbf{x}_0 \neq 0$. The method then proceeds to generate two sequences of vectors

$$\mathbf{X}_k = (\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{k-1}) \quad \text{and} \quad \mathbf{Y}_k = (\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{k-1}) \quad (5)$$

such that for any k we have the following bases for the Krylov spaces

$$\text{COLSP } \mathbf{X}_k = \mathcal{K}(A, \mathbf{x}_0, k), \quad \text{COLSP } \mathbf{Y}_k = \mathcal{K}(A^T, \mathbf{y}_0, k), \quad (6)$$

and the matrix

$$\mathbf{Y}^T \mathbf{X} = \mathbf{D}_K = \begin{pmatrix} d_0 & & & 0 \\ & \ddots & & \\ 0 & & & d_{k-1} \end{pmatrix} \quad (7)$$

is nonsingular and diagonal. Condition (7) is called the *bi-orthogonality* condition. The method is based on the recursive application of the identities

$$\mathbf{A}\mathbf{X}_k = \mathbf{X}_k\mathbf{H}_k + \mathbf{x}_k(0, \dots, 0, 1)\beta_k, \quad \mathbf{A}^T\mathbf{Y}_k = \mathbf{Y}_k\mathbf{G}_k + \mathbf{y}_k(0, \dots, 0, 1)\gamma_k,$$

and hence

$$\mathbf{G}_k^T \mathbf{D}_k = \mathbf{Y}_k^T \mathbf{A} \mathbf{X}_k = \mathbf{D}_k \mathbf{H}_k, \quad (8)$$

where $\mathbf{H}_k, \mathbf{G}_k$ are upper Hessenberg matrices of coefficients

$$\mathbf{H}_k = \begin{pmatrix} h_{00} & h_{01} & & & 0 \\ \beta_1 & h_{11} & \ddots & & \\ & \ddots & \ddots & & h_{k-2,k-1} \\ 0 & & \beta_{k-1} & & h_{k-1,k-1} \end{pmatrix},$$

$$\mathbf{G}_k = \begin{pmatrix} g_{00} & g_{01} & & & 0 \\ \gamma_1 & g_{11} & \ddots & & \\ & \ddots & \ddots & & g_{k-2,k-1} \\ 0 & & \gamma_{k-1} & & g_{k-1,k-1} \end{pmatrix},$$

generated in a manner analogous to that of the Arnoldi algorithm. Identity (8) implies that the resulting matrices \mathbf{H}, \mathbf{G} will actually take on the above tridiagonal form. The resulting algorithm is then taken from [68, pp. 388ff], but uses a modified Gram-Schmidt bi-orthogonalization process analogous to the ordinary process used above in the Arnoldi algorithm.

Simple Nonsymmetric Lanczos Algorithm [68].

0. Start with matrix A and two vectors $\mathbf{x}_0, \mathbf{y}_0$ such that $\mathbf{y}_0^T \mathbf{x}_0 = d_0 \neq 0$
1. For $k = 1, 2, \dots$, while $\mathbf{x}_k \neq 0$ and/or $\mathbf{y}_k \neq 0$
2. Set $\mathbf{x}_k^{(0)} = \mathbf{A}\mathbf{x}_{k-1}$ {expand Krylov spaces} and $\mathbf{y}_k^{(0)} = \mathbf{A}^T\mathbf{y}_{k-1}$.
3. For $j = 1, 2$ {enforce bi-ortho. cond.}
4. Set $h_{k-j,k-1} = d_{k-j}^{-1} \mathbf{y}_{k-j}^T \mathbf{x}_k^{(j-1)}$ and $\mathbf{x}_k^{(j)} = \mathbf{x}_k^{(j-1)} - \mathbf{x}_{k-j} h_{k-j,k-1}$.
Set $g_{k-j,k-1} = d_{k-j}^{-1} \mathbf{x}_{k-j}^T \mathbf{y}_k^{(j-1)}$ and $\mathbf{y}_k^{(j)} = \mathbf{y}_k^{(j-1)} - \mathbf{y}_{k-j} g_{k-j,k-1}$.
5. Set $\beta_k \mathbf{x}_k = \mathbf{x}_k^{(2)}$ and $\gamma_k \mathbf{y}_k = \mathbf{y}_k^{(2)}$,
where $h_{k,k-1} \equiv \beta_k, g_{k,k-1} \equiv \gamma_k$ are scale factors to be chosen.
6. If $(\mathbf{y}_k)^T \mathbf{x}_k = 0$ then we have a *breakdown error*.
7. Set $d_k = \mathbf{y}_k^T \mathbf{x}_k$.

There are two typical choices for the scale factors β_k, γ_k . Choice I is to choose the factors so that $d_k = 1$, i.e., so that $\mathbf{D}_k = I$. In this case (8) implies that $\mathbf{G}_k^T = \mathbf{H}_k$. Choice II is to set $\beta_k = 1$ and $\gamma_k = 1$. In this case, it is easily verified from (8) that $\mathbf{G}_k = \mathbf{H}_k$. Thus with either choice, the computation of these coefficients in step 4 can be cut in half.

In order to handle the possibility of a breakdown, the method must be modified. Originally a limited recovery method was proposed in [54], but a full recovery method was not proposed until more recently in [8], [53]. Numerical implementations have been described in [23]. The modification is based on the idea of generating the individual vectors to satisfy (6), but grouping the generated vectors into clusters, and enforcing the bi-orthogonality condition only between different clusters. The clusters are denoted $X_0 = (\mathbf{x}_0, \dots, \mathbf{x}_{k_1-1})$, $X_1 = (\mathbf{x}_{k_1}, \dots, \mathbf{x}_{k_2-1})$, \dots and $Y_0 = (\mathbf{y}_0, \dots, \mathbf{y}_{k_1-1})$, $Y_1 = (\mathbf{y}_{k_1}, \dots, \mathbf{y}_{k_2-1})$, \dots . The resulting method is the following, taken from [8], but using the modified Gram-Schmidt algorithm as before:

Clustered Nonsymmetric Lanczos Algorithm [8].

0. Start with matrix A , two vectors $\mathbf{x}_0, \mathbf{y}_0$
1. Initialize clusters $X_0 = (\mathbf{x}_0)$, $Y_0 = (\mathbf{y}_0)$.
2. Initialize cluster counter $p = 0$ and marker $k_0 = 0$.
3. For $k = 1, 2, \dots$ while $\mathbf{x}_k \neq 0$ and/or $\mathbf{y}_k \neq 0$
4. If $Y_p^T X_p$ is nonsingular then {close current cluster and increment counter}
5. Set $D_p = Y_p^T X_p$ {next diagonal block}
6. Increment $p = p + 1$.
7. Set $k_p = k$ {save the index k corresponding to the beginning of cluster p }
8. Set $\mathbf{x}_k^{(0)} = A\mathbf{x}_{k-1}$ {expand Krylov spaces}
and $\mathbf{y}_k^{(0)} = A^T\mathbf{y}_{k-1}$.
9. For $j = 1, 2, \dots, p$ {enforce ortho. cond.}
10. Set $\mathbf{h}_{p-j,k-1} = D_{p-j}^{-1} Y_{p-j}^T \mathbf{x}_k^{(j-1)}$ and $\mathbf{x}_k^{(j)} = \mathbf{x}_k^{(j-1)} - X_{p-j} \mathbf{h}_{p-j,k-1}$.
Set $\mathbf{g}_{p-j,k-1} = D_{p-j}^{-T} X_{p-j}^T \mathbf{y}_k^{(j-1)}$ and $\mathbf{y}_k^{(j)} = \mathbf{y}_k^{(j-1)} - Y_{p-j} \mathbf{g}_{p-j,k-1}$.
11. Set $\beta_k \mathbf{x}_k = \mathbf{x}_k^{(p)} - \{X_p \mathbf{h}_{p,k-1}\}$,
and $\gamma_k \mathbf{y}_k = \mathbf{y}_k^{(p)} - \{Y_p \mathbf{g}_{p,k-1}\}$,
where β_k, γ_k are scale factors to be chosen, and the braces enclose optional terms (see text).
12. Append to current clusters:
 $X_p = (X_p, \mathbf{x}_k)$,
 $Y_p = (Y_p, \mathbf{y}_k)$.

In step 10, the vectors of coefficients $\mathbf{h}_{p-j,k-1}$, $\mathbf{g}_{p-j,k-1}$ are chosen so that \mathbf{x}_k , \mathbf{y}_k satisfy the bi-orthogonality conditions

$$\mathbf{x}_k^T (Y_0, \dots, Y_{p-1}) = 0, \quad \mathbf{y}_k^T (X_0, \dots, X_{p-1}) = 0.$$

The coefficients $\mathbf{h}_{p,k-1}$, $\mathbf{g}_{p,k-1}$ in step 11 can be chosen to enforce "internal orthogonality" conditions $X_p^T \mathbf{x}_k = 0$, $Y_p^T \mathbf{y}_k = 0$ with the result that when the algorithm finishes, the individual clusters will satisfy $X_k^T X_k = I = Y_k^T Y_k$ for each k [10]. They may also be chosen to be zero [53], in which case each diagonal block D_p will be lower *anti*-triangular with Hankel structure. In the latter case, the terms enclosed in braces "{ }" are omitted.

At each stage $k = k_p$ we group the vectors generated into matrices, partitioned according to the clusters:

$$\mathbf{X}_k = (\mathbf{x}_0, \dots, \mathbf{x}_{k-1}) = (X_0, \dots, X_p), \quad \mathbf{Y}_k = (\mathbf{y}_0, \dots, \mathbf{y}_{k-1}) = (Y_0, \dots, Y_p), \quad (9)$$

where the vectors satisfy the clustered bi-orthogonality conditions $Y_i^T X_j = 0$ for $i \neq j$; i.e.,

$$\mathbf{Y}_k^T \mathbf{X}_k = \mathbf{D}_k = \text{diag}(D_0, \dots, D_p), \quad (10)$$

where $D_i = Y_i^T X_i$. Since each cluster i is filled until the corresponding matrix D_i is nonsingular, we see that upon termination of the algorithm, all the diagonal blocks of \mathbf{D}_k will be nonsingular, except the last one, D_p . It will be seen in the next section, from considerations in Control Theory, that in fact D_p will be either entirely zero, or empty. Hence the rank of \mathbf{D}_k will be simply the dimension of $\text{diag}(D_0, \dots, D_{p-1})$. We also partition the upper Hessenberg matrices of coefficients

$$\mathbf{H}_k = ((H_{ij})_{i,j=0}^p), \quad \text{with } H_{ij} \text{ satisfying } AX_j = (X_0, \dots, X_{j+1}) \begin{pmatrix} H_{0j} \\ \vdots \\ H_{j+1,j} \end{pmatrix}, \quad (11)$$

and

$$\mathbf{G}_k = ((G_{ij})_{i,j=0}^p), \quad \text{with } G_{ij} \text{ satisfying } A^T Y_j = (Y_0, \dots, Y_{j+1}) \begin{pmatrix} G_{0j} \\ \vdots \\ G_{j+1,j} \end{pmatrix}, \quad (12)$$

where the partitioning of \mathbf{H}_k , \mathbf{G}_k into $((H_{ij})_{i,j=0}^p)$, $((G_{ij})_{i,j=0}^p)$ is consistent with the partitioning (9). With this partitioning, we have the formula analogous to (8):

$$\mathbf{G}_k^T \mathbf{D}_k = \mathbf{Y}_k^T \mathbf{A} \mathbf{X}_k = \mathbf{D}_k \mathbf{H}_k \quad (13)$$

from which it follows that \mathbf{H}_k , \mathbf{G}_k are block tridiagonal. That is, $H_{ij} = G_{ij} = 0$ for $|i - j| > 1$, and this means that (at least in exact arithmetic) in steps 9–11 j need take on only the values 1,2; the remaining coefficients are all zero.

We remark that in floating point arithmetic, the condition of "nonsingularity" in step 4 is necessarily replaced with "not too ill-conditioned," where the amount of acceptable ill conditioning depends upon a user-supplied parameter. In this case,

the last block D_p will no longer be exactly zero. The block tridiagonal structure of \mathbf{H}, \mathbf{G} will remain, but will have to be especially enforced via complete bi-orthogonalization. A detailed numerical implementation that addresses most of these issues can be found in [23].

We note that also for this algorithm, the matrices of generated vectors $\mathbf{X}_k, \mathbf{Y}_k$ satisfy (6) for every k . In fact, the clustered Lanczos algorithm is based on the idea of generating the vectors one by one so that they satisfy (6) at every stage, enforcing (10) at every stage where possible. This is accomplished by filling in the identities (11) and (12).

The block nonsymmetric Lanczos algorithm is a generalization of the simple nonsymmetric Lanczos algorithm to take several starting vectors on the left and right. No one has yet proposed a clustering scheme to take care of the breakdown, so we describe the algorithm assuming breakdowns do not occur. The matrices generated will be just block analogs of those generated by the simple nonsymmetric Lanczos algorithm, with uniform block sizes:

$$\mathbf{X}_k = (X_0, X_1, \dots, X_{k-1}) \quad \text{and} \quad \mathbf{Y}_k = (Y_0, Y_1, \dots, Y_{k-1}) \quad (14)$$

$$\mathbf{D}_k = \begin{pmatrix} D_0 & & 0 \\ & \ddots & \\ 0 & & D_{k-1} \end{pmatrix},$$

$$\mathbf{H}_k = \begin{pmatrix} H_{00} & H_{01} & & 0 \\ \mathcal{B}_1 & H_{11} & \ddots & \\ & \ddots & \ddots & H_{k-2,k-1} \\ 0 & & \mathcal{B}_{k-1} & H_{k-1,k-1} \end{pmatrix},$$

$$\mathbf{G}_k = \begin{pmatrix} G_{00} & G_{01} & & 0 \\ \Gamma_1 & G_{11} & \ddots & \\ & \ddots & \ddots & G_{k-2,k-1} \\ 0 & & \Gamma_{k-1} & G_{k-1,k-1} \end{pmatrix},$$

which satisfy

$$A\mathbf{X}_k = \mathbf{X}_k\mathbf{H}_k + X_k\mathcal{B}_k(0, \dots, 0, I), \quad (15)$$

$$A^T\mathbf{Y}_k = \mathbf{Y}_k\mathbf{G}_k + Y_k\Gamma_k(0, \dots, 0, I), \quad (16)$$

$$\mathbf{Y}_k^T\mathbf{X}_k = \mathbf{D}_k.$$

The algorithm is a block analog of the simple algorithm. The following is the form of the algorithm taken from [43], [44], but using the modified Gram-Schmidt bi-orthogonalization process.

Block Nonsymmetric Lanczos Algorithm [43], [44].

0. Start with matrix A and two sets of vectors X, Y such that $Y^T X$ is nonsingular
1. Compute factorizations $X_0\mathcal{B}_0 = X$ and $Y_0\Gamma_0 = Y$.

2. For $k = 1, 2, \dots$ while $X_k \neq 0$ and/or $Y_k \neq 0$
3. Set $X_k^{(0)} = AX_{k-1}$ {expand Krylov spaces}
and $Y_k^{(0)} = A^T Y_{k-1}$.
4. For $j = 1, 2$ {enforce bi-ortho. cond.}
5. Set $H_{k-j,k-1} = D_{k-j}^{-1} Y_{k-j}^T X_k^{(j-1)}$ and $X_k^{(j)} = X_k^{(j-1)} - X_{k-j} H_{k-j,k-1}$.
Set $G_{k-j,k-1} = D_{k-j}^{-1} X_{k-j}^T Y_k^{(j-1)}$ and $Y_k^{(j)} = Y_k^{(j-1)} - Y_k G_{k-j,k-1}$.
6. Set $X_k \mathcal{B}_k = X_k^{(2)}$ and $Y_k \Gamma_k = Y_k^{(2)}$,
where \mathcal{B}_k, Γ_k are matrices chosen to normalize the vector,
as indicated below.
7. If $Y_k^T X_k$ is singular then we have a *breakdown error*, unless $X_k \neq 0$
or $Y_k \neq 0$.
8. Set $D_k = Y_k^T X_k$.

In steps 1 and 6, the normalization matrices \mathcal{B}_k, Γ_k are normally chosen to make $D_k = Y_k^T X_k = I$. To derive them, we write

$$I = D_k = Y_k^T X_k = \Gamma_k^{-T} (Y_k^{(1)})^T X_k^{(1)} \mathcal{B}_k^{-1}. \quad (17)$$

Among all the possible solutions to this equation, the authors of [43] have proposed using the particular choice

$$\Gamma_k^T \mathcal{B}_k = \text{LU decomposition of } (Y_k^{(1)})^T X_k^{(1)},$$

and they have further suggested the use of pivoting to enhance numerical stability [44]. However, if $(Y_k^{(1)})^T X_k^{(1)}$ should be exactly singular, then we have a *breakdown condition*, and the only “remedy” that has been proposed for the block algorithm is to start over with different starting vectors.

3. Controllability and observability

In this section we discuss the concepts of controllability and observability and the computation of the minimal realization via the Kalman decomposition. We show how the Lanczos algorithms of the previous section can be used to compute the Kalman decomposition and the minimal realization. We will also see how the controllability and observability theory provides understanding about the behavior of the Lanczos algorithm.

Our discussion is based on time invariant linear systems in continuous time:

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \quad \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t), \quad (18)$$

and in discrete time:

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k, \quad \mathbf{y}_k = \mathbf{C}\mathbf{x}_k, \quad (19)$$

where \mathbf{x} , \mathbf{u} , \mathbf{y} are, respectively, an n -vector of internal states, an m -vector of inputs, and a p -vector of outputs. The systems (18) and (19) are said to be *stable* if all the eigenvalues of A are, respectively, in the left half plane and within the open unit circle. We define the *controllable subspace* \mathcal{S}_c as the subspace of the \mathbf{x} -space from which one can find an input that will drive the state to zero, and the *unobservable space* \mathcal{S}_o as the subspace of the \mathbf{x} -space which one cannot detect using the outputs.

A classical algebraic characterization of these spaces is given by

Theorem 2 [19].

$$\begin{aligned} \mathcal{S}_c &\equiv \mathcal{K}(A, B, \infty), \\ \mathcal{S}_o &\equiv \text{NULLSPACE}\{(K(A^T, C^T, \infty))^T\} \equiv \mathcal{K}(A^T, C^T, \infty)^\perp, \end{aligned} \tag{20}$$

where S^\perp denotes the orthogonal complement of the set S .

One may also define the *uncontrollable space* \mathcal{S}_c and the *observable space* \mathcal{S}_o as the complement of \mathcal{S}_c and \mathcal{S}_o , respectively, as well as their mutual intersections:

$$\mathcal{S}_{c\bar{o}} = \mathcal{S}_c \cap \mathcal{S}_{\bar{o}}, \quad \mathcal{S}_{c\bar{o}} = \mathcal{S}_c \cap \mathcal{S}_o, \quad \mathcal{S}_{\bar{c}\bar{o}} = \mathcal{S}_{\bar{c}} \cap \mathcal{S}_{\bar{o}}, \quad \mathcal{S}_{\bar{c}o} = \mathcal{S}_{\bar{c}} \cap \mathcal{S}_o. \tag{21}$$

The *Kalman decomposition* [25], [40] or (18), (19) is obtained by applying a similarity transformation $T = (T_{\bar{c}o}, T_{\bar{c}\bar{o}}, T_{co}, T_{c\bar{o}})$ where each T_{ab} is a basis for the corresponding \mathcal{S}_{ab} . When T is applied, we obtain the new system (for continuous time)

$$\dot{\hat{\mathbf{x}}}(t) = \hat{A}\hat{\mathbf{x}}(t) + \hat{B}\mathbf{u}(t), \quad \mathbf{y}(t) = \hat{C}\hat{\mathbf{x}}(t), \tag{22}$$

where

$$\begin{aligned} \hat{A} = T^{-1}AT &= \begin{pmatrix} \hat{A}_{\bar{c}o} & 0 & 0 & 0 \\ \hat{A}_{21} & \hat{A}_{\bar{c}\bar{o}} & 0 & 0 \\ \hat{A}_{31} & 0 & \hat{A}_{co} & 0 \\ \hat{A}_{41} & \hat{A}_{42} & \hat{A}_{43} & \hat{A}_{c\bar{o}} \end{pmatrix}, \\ \hat{B} = T^{-1}B &= \begin{pmatrix} 0 \\ 0 \\ \hat{B}_{co} \\ \hat{B}_{c\bar{o}} \end{pmatrix}, \quad \hat{C} = CT = (\hat{C}_{\bar{c}\bar{o}} \ 0 \ \hat{C}_{co} \ 0). \end{aligned} \tag{23}$$

The minimal realization for (18) is then

$$\dot{\mathbf{z}} = \hat{A}_{co}\mathbf{z} + \hat{B}_{co}\mathbf{u}, \quad \mathbf{y} = \hat{C}_{co}\mathbf{z}, \tag{24}$$

where $\mathbf{x} = T_{co}\mathbf{z}$. Analogous statements apply for discrete time systems.

If we define the spaces $\mathcal{S}_{\bar{c}}$ and \mathcal{S}_o as the *orthogonal* complements of \mathcal{S}_c and $\mathcal{S}_{\bar{o}}$, respectively, and we use *orthonormal* bases of resulting intersection spaces in the transformation T , then the resulting T is the best conditioned (in the 2-norm) of all transformations yielding the Kalman decomposition [6].

As shown in [9] it is evident that the Arnoldi algorithm starting with A and B computes an orthonormal basis for \mathcal{S}_c . If the Arnoldi algorithm is applied to the matrix A^T and starting vectors C^T , we obtain an orthonormal basis for the

orthogonal complement of \mathcal{S}_0 . Then the individual spaces (21) may be obtained by numerically computing bases for the intersections and orthogonal complements. One method to compute the intersection and orthogonal complements of two spaces $\mathcal{S}_1, \mathcal{S}_2$ given two respective orthonormal bases T_1, T_2 is described in [30, §12.4.4]. Briefly, one computes the *singular value decomposition* (SVD) of $T_1^T T_2$ (see the next section for further discussion on the SVD of a matrix product). If $\mathbf{v}_1, \dots, \mathbf{v}_r$ are the right singular vectors corresponding to the multiple singular value 1, and $\mathbf{v}_{r+1}, \dots, \mathbf{v}_s$ are all the remaining right singular vectors, then $\{T_2 \mathbf{v}_1, \dots, T_2 \mathbf{v}_r\}$ is an orthonormal basis for the intersection space $\mathcal{S}_1 \cap \mathcal{S}_2$, and $\{T_2 \mathbf{v}_{r+1}, \dots, T_2 \mathbf{v}_s\}$ is an orthonormal basis for the subspace within \mathcal{S}_2 orthogonal to the intersection.

For *single input single output* (SISO) systems, a very special structure arises from the clustered (or look-ahead) nonsymmetric Lanczos algorithm. The rest of this section is devoted to a discussion of this case. To the author's knowledge, no one has proposed a *block clustered (or look-ahead) Lanczos algorithm* that would correspond to the *multiple input multiple output* (MIMO) case. For SISO systems, we may apply the clustered nonsymmetric Lanczos algorithm with matrix A and starting vectors $\mathbf{x}_0 = B$ and $\mathbf{y}_0 = C^T$. We will see that not only does this algorithm yield bases for the various spaces, but the theory behind the controllable and observable spaces contributes to our understanding of the behavior of the algorithm. If the algorithm is carried out until both \mathbf{x}_{k^*} and \mathbf{y}_{k^*} are zero, appending zero vectors if necessary, then at the termination of the algorithm we will have the generated vectors grouped into clusters, which are further grouped together as follows:

$$\begin{aligned} \mathbf{X}_{k^*} &= (X_0 \ X_1 \ \cdots \ X_{p-1} \ | \ X_{p^*}) = (X_{\hat{k}} | X_{p^*}), \\ \mathbf{Y}_{k^*} &= (Y_0 \ Y_1 \ \cdots \ Y_{p-1} \ | \ Y_{p^*}) = (\mathbf{Y}_{\hat{k}} | Y_{p^*}), \end{aligned} \quad (25)$$

where \hat{k} is the index k corresponding to the last *complete* cluster, i.e., \hat{k} is the rank of \mathbf{D}_{k^*} . It was shown in [10] that indeed at termination of the algorithm the matrix \mathbf{D}_{k^*} (10) will have the form $\mathbf{D}_{k^*} = \text{diag}(D_0, \dots, D_{p-1}, 0)$. This follows from

Theorem 3 [10], [53]. *Let the vectors $\mathbf{X}_{k^*} \equiv (\mathbf{x}_0 \cdots \mathbf{x}_{k^*-1})$ and $\mathbf{Y}_{k^*} \equiv (\mathbf{y}_0 \cdots \mathbf{y}_{k^*-1})$ be the vectors generated by the clustered nonsymmetric Lanczos algorithm with matrix A and starting vectors $\mathbf{x}_0, \mathbf{y}_0$. Then there exists an index r such that the generated vectors can be partitioned in the form*

$$\begin{aligned} \mathbf{X}_{k^*} &= (\mathbf{X}_r \ | \ \Xi_r) = (\mathbf{x}_0 \cdots \mathbf{x}_{r-1} \ | \ \mathbf{x}_r \cdots \mathbf{x}_{k^*-1}) \\ \mathbf{Y}_{k^*} &= (\mathbf{Y}_r \ | \ \Upsilon_r) = (\mathbf{y}_0 \cdots \mathbf{y}_{r-1} \ | \ \mathbf{y}_r \cdots \mathbf{y}_{k^*-1}) \end{aligned}$$

such that $\mathbf{D}_r = \mathbf{Y}_r^T \mathbf{X}_r$ is nonsingular and both $\Xi_r^T \mathbf{Y}_{k^*}$ and $\Upsilon_r^T \mathbf{X}_{k^*}$ are identically zero or empty.

Proof. Let $K \equiv K(A, \mathbf{x}_0, k^*)$ and $L = K(A^T, \mathbf{y}_0, k^*)$ be the right and left Krylov sequences, respectively, corresponding to this application of the Lanczos algorithm. Then (6) can be expressed as

$$\mathbf{X}_{k^*} = KU, \quad \mathbf{Y}_{k^*} = LV, \quad (26)$$

where U, V are some upper triangular nonsingular $(k^* \times k^*)$ -matrices. Then the block diagonal matrix (10) will be

$$D_{k^*} = Y_{k^*}^T X_{k^*} = V^T L^T K U.$$

The matrix $L^T K$ is a leading principal submatrix of the infinite Hankel matrix $[K(A^T, y_0, \infty)]^T K(A, x_0, \infty)$ of finite rank r . By the Corollary in [24, p. 206], the leading $(r \times r)$ -submatrix of this infinite Hankel matrix must be nonsingular. Hence $k^* \geq r$, and D_r , the leading $(r \times r)$ -submatrix of D_{k^*} , must also be nonsingular.

The condition in step 4 of the clustered nonsymmetric Lanczos algorithm is equivalent to the condition that $D_k = Y_k^T X_k$ be nonsingular, which is in turn equivalent to the condition that $[K(A^T, y_0, k)]^T K(A, x_0, k)$ be nonsingular. Hence for every k for which this condition holds, we have a diagonal partitioning of D_{k^*} as

$$D_{k^*} = \begin{pmatrix} D_k & 0 \\ 0 & \Delta_k \end{pmatrix}$$

with a leading $(k \times k)$ part D_k and the remaining part denoted Δ_k . In particular, we have this partitioning for $k = r$.

Since $\text{rank } D_{k^*} = \text{rank } D_r$, it follows that $\Delta_r = 0$. □

This theorem, together with Theorem 2, implies that at termination of the Lanczos algorithm, $X_r = T_{co}$, $\Xi_r = T_{c\bar{o}}$, $\Upsilon_r = T_{\bar{c}o}$ will be, respectively, bases for the spaces S_{co} , $S_{c\bar{o}}$, $S_{\bar{c}o}$. The remaining space $S_{\bar{c}\bar{o}}$ has a basis $T_{\bar{c}\bar{o}}$ which may be computed as a basis for the orthogonal complement of the other three spaces. Just as the transformation T is constructed from the vectors obtained directly from the Lanczos algorithm, the resulting transformed system matrix coefficients are obtained directly from the coefficients generated during the course of the Lanczos algorithm. As shown in detail in [10], one can substitute $x = T_{co}z = X_r z$ into (18), apply $Y_{k^*}^T$ to extract \hat{A}_{co} , \hat{B}_{co} , \hat{C}_{co} , and use the identities (10), (13) from the clustered Lanczos algorithm to obtain the minimal realization (24) of the system (18):

$$\dot{z} = H_r z + e_1 u, \quad y = (e_1^T D_0 \ 0 \ \cdots \ 0)z,$$

where H_r is the leading $(r \times r)$ upper Hessenberg matrix of coefficients from the clustered nonsymmetric Lanczos algorithm, D_0 is the leading diagonal block of D_r , and $e_1 = (1 \ 0 \ \cdots \ 0)^T$ denotes the initial coordinate unit vector of appropriate dimensions.

4. Model reduction via balanced realization

The *balanced realization* [50] is a method for balancing the gains between inputs and states with those between states and outputs, and isolating the states with small gains. These small gain states can be truncated away without disturbing the system

by much [2], [22], [26]. The Lanczos-type algorithms can be used to compute approximate balanced realizations for systems of very large dimensionality, whereas most other methods are restricted to systems of only modest dimensionality.

The balanced realization for stable systems (18) is defined in terms of the *controllability Gramian*

$$W_c = \int_0^{\infty} e^{At} B B^T e^{A^T t} dt$$

and the *observability Gramian*

$$W_o = \int_0^{\infty} e^{A^T t} C^T C e^{At} dt.$$

These Gramians are, respectively, the solutions to the Lyapunov equations

$$A W_c + W_c A^T + B B^T = 0 \quad (27)$$

and

$$A^T W_o + W_o A + C^T C = 0. \quad (28)$$

If the transformation T is applied to (18) to obtain the method

$$\dot{\hat{\mathbf{x}}}(t) = \hat{A} \hat{\mathbf{x}}(t) + \hat{B} \mathbf{u}(t), \quad \mathbf{y}(t) = \hat{C} \hat{\mathbf{x}}(t), \quad (29)$$

where

$$\hat{A} = T^{-1} A T, \quad \hat{B} = T^{-1} B, \quad \hat{C} = C T \quad (30)$$

then the Gramians will be transformed by contragredient (or congruence) transformations into

$$\hat{W}_c = T^{-1} W_c T^{-T}, \quad \hat{W}_o = T^T W_o T. \quad (31)$$

The balanced realization may be computed by the following prescription [47], [49]. Define $L_c L_c^T$, $L_o L_o^T$ as the Cholesky factorizations of W_c , W_o , respectively. Compute the SVD of the product $L_o^T L_c$ to obtain the factorization $U \Sigma V^T = L_o^T L_c$, where U , V are orthogonal matrices, and $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$. Then the particular choice $T = L_c V \Sigma^{-1/2}$ will reduce the Gramians to the same diagonal form $\hat{W}_c = \hat{W}_o = \Sigma$ [49]. The resulting system (29) with this particular choice for T is the *balanced realization* for the system (18) [50]. Effective methods for solving (27), (28) directly for the Cholesky factors without forming W_c , W_o and for finding the SVD of a matrix product without forming the product have been given in [36] and [38], [5], respectively, thus enhancing the numerical accuracy of the results.

Partition the matrices in (30) conformally as

$$\hat{A} = \begin{pmatrix} \hat{A}_{11} & \hat{A}_{12} \\ \hat{A}_{21} & \hat{A}_{22} \end{pmatrix}, \quad \hat{B} = \begin{pmatrix} \hat{B}_1 \\ \hat{B}_2 \end{pmatrix}, \quad \hat{C} = (\hat{C}_1 \quad \hat{C}_2),$$

as well as the Gramians

$$\hat{W}_c = \hat{W}_o = \begin{pmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{pmatrix},$$

where \hat{A}_{11} , Σ_1 are the leading $(r \times r)$ part of the respective matrices, for some arbitrary r . Define $F(s) = C(sI - A)^{-1}B = \hat{C}(sI - \hat{A})^{-1}\hat{B}$ as the transfer function for (18) and (29), and $G_1(s) = \hat{C}_1(sI - \hat{A}_{11})^{-1}\hat{B}_1$ as the transfer function of the system obtained when the last $(n - r)$ states in the balanced realization are truncated away. Then it has been shown (e.g., [22]) that

$$\|G(s) - G_1(s)\|_\infty \leq 2 \sum_{r+1}^n \sigma_i,$$

where the norm $\|\cdot\|_\infty$ on functions analytic everywhere except at some poles in the left half plane is defined by

$$\|G(s)\|_\infty = \sup_{\omega} \sigma_{\max} G(j\omega),$$

where $j\omega$ varies over the entire imaginary axis. Thus the truncation of the balanced realization is a reduced order model whose behavior does not differ from that of the original system by very much.

A similar development may be carried out for stable discrete time systems of the form (19). In this case, the Gramians are defined as

$$W_c = \sum_{i=0}^{\infty} A^i B B^T (A^T)^i, \quad W_o = \sum_{i=0}^{\infty} (A^T)^i C^T C A^i,$$

and satisfy the Lyapunov equations

$$A W_c A^T + B B^T = W_c \tag{32}$$

$$A^T W_o A + C^T C = W_o. \tag{33}$$

When a transformation T is applied to (19) to obtain the transformed system

$$\hat{\mathbf{x}}_{k+1} = \hat{A} \hat{\mathbf{x}}_k + \hat{B} \mathbf{u}_k, \quad \mathbf{y}_k = \hat{C} \hat{\mathbf{x}}_k, \tag{34}$$

where the matrices are defined by (30), the Gramians are again transformed by the same congruence transformations (31). Thus the balanced realization can be computed using the same prescription as for the continuous time case [37]. A similar error bound for truncated balanced realizations was obtained in [2].

It is difficult to solve the Lyapunov equations for very large sparse systems. The Arnoldi algorithm may be used to compute approximate solutions to such Lyapunov equations (27), (28), (32), (33), which will satisfy certain Galerkin conditions. We will describe these conditions for the discrete time case. The continuous time case is analogous and is treated in [58]. In addition, one can obtain an error bound for the discrete time solution. Similar techniques have been proposed for the more general Sylvester equation $AX - XB = C$ in [55], in which the approximate solutions can be recursively generated as the Arnoldi process advances.

We show how the Arnoldi algorithm yields an approximate solution to (32), then we show that it satisfies the Galerkin property and an error bound. If the Arnoldi algorithm is carried out with a matrix A and starting vectors $X = B$, then after

k steps, we have the relation (4). We also have that $R = X_0^T B$ is the nonsingular upper triangular matrix constructed in step 1 of the Arnoldi algorithm. Define

$$G = \sum_{i=0}^{\infty} (\mathbf{H}_k)^i \begin{pmatrix} R \\ 0 \end{pmatrix} (R^T \ 0) (\mathbf{H}_k^T)^i \quad (35)$$

so that G satisfies the Lyapunov equation

$$\mathbf{H}_k G \mathbf{H}_k^T + \begin{pmatrix} R \\ 0 \end{pmatrix} (R^T \ 0) = G. \quad (36)$$

Then the approximate solution to (32) is $\tilde{W} = \mathbf{X}_k G \mathbf{X}_k^T$. One can take advantage of the recursive nature of the Arnoldi algorithm to use the approximate solution obtained after each step of the Arnoldi algorithm to generate subsequent approximate solutions rapidly, but for the sake of brevity and clarity we do not describe that process here. This is the basis of the methods in, for example, [55], [59]. This solution satisfies the following theorem, proved in [58] for the continuous time case:

Theorem 4. Define the inner product on $n \times n$ matrices as $\langle X, Y \rangle \equiv \text{tr}(X^T Y)$. Define the space $\mathcal{G} = \{\mathbf{X}_k M \mathbf{X}_k^T : \forall M \in \mathcal{R}^{k \times k}\}$. Then the solution $\tilde{W} = \mathbf{X}_k G \mathbf{X}_k^T$ lies in \mathcal{G} and the residual:

$$\text{RES}(\tilde{W}) \equiv A \tilde{W} A^T + B B^T - \tilde{W}$$

is orthogonal to \mathcal{G} with respect to this inner product. Hence \tilde{W} satisfies a Galerkin condition.

(Here the notation $\text{tr}(M)$ stands for the trace of M .)

Proof. All we must verify is that $\langle \mathbf{X}_k M \mathbf{X}_k^T, \text{RES}(\tilde{W}) \rangle = 0$ for any M . We make use of the identity $\text{tr}(MN) = \text{tr}(NM)$ for any matrices M, N and the identity (3). We have

$$\begin{aligned} \langle \mathbf{X}_k M \mathbf{X}_k^T, \text{RES}(\tilde{W}) \rangle &= \text{tr}((\mathbf{X}_k M \mathbf{X}_k^T) A (\mathbf{X}_k G \mathbf{X}_k^T) A^T + (\mathbf{X}_k M \mathbf{X}_k^T) B B^T \\ &\quad - (\mathbf{X}_k M \mathbf{X}_k^T) (\mathbf{X}_k G \mathbf{X}_k^T)) \\ &= \text{tr}(M (\mathbf{X}_k^T A \mathbf{X}_k) G (\mathbf{X}_k^T A^T \mathbf{X}_k) + M \mathbf{X}_k^T B B^T \mathbf{X}_k \\ &\quad - M \mathbf{X}_k^T \mathbf{X}_k G \mathbf{X}_k^T \mathbf{X}_k) \\ &= \text{tr} \left(M \mathbf{H}_k G \mathbf{H}_k^T + M \begin{pmatrix} R \\ 0 \end{pmatrix} (R^T \ 0) - M G \right) \\ &= \text{tr} \left(M \left(\mathbf{H}_k G \mathbf{H}_k^T + \begin{pmatrix} R \\ 0 \end{pmatrix} (R^T \ 0) - G \right) \right) \\ &= \text{tr}(M \cdot 0) \\ &= 0. \end{aligned}$$

□

In addition to the above Galerkin condition, we can also bound the error $\|W_c - \tilde{W}\|$. In order to do this, we need the following theorem from linear algebra.

Theorem 5. *Let $\rho(A)$ denote the spectral radius of a matrix A . If $\rho(A) < \delta < 1$, then there is a constant θ depending only on A and the choice of δ such that for any $k = 1, 2, 3, \dots$*

$$\|A^k\|_2 \leq \theta \delta^k.$$

Furthermore, if A is diagonalizable, then the above also holds for $\delta = \rho(A) < 1$.

Proof. One can find a consistent matrix norm such that $\|A\|_* \leq \delta$ [62]. Since all norms on a finite dimensional space are equivalent, one can find a constant θ such that $\|M\|_2 \leq \theta \|M\|_*$ for all matrices M . The result follows from

$$\|A^k\|_2 \leq \theta \|A^k\|_* \leq \theta \|A\|_*^k \leq \theta \delta^k.$$

If A is diagonalizable, then we have $A^k = P D^k P^{-1}$ for diagonal D , which yields equality above with $\theta = \|P\|_2 \cdot \|P^{-1}\|_2$ and $\delta = \rho$. □

We also need the following lemma.

Lemma 6. *When the Arnoldi algorithm is run with a matrix A and starting vectors B , then for $0 \leq i < k$*

$$A^i B = \mathbf{X}_k \mathbf{H}_k^i \begin{pmatrix} R \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

where $R = X_0^T B$ is the upper triangular matrix computed in step 1 of the Arnoldi algorithm.

Proof. Choose a k and assume $0 \leq i < k$. Then we have the identity from (4)

$$\begin{aligned} A \mathbf{X}_k \mathbf{H}_k^{i-1} \begin{pmatrix} R \\ 0 \\ \vdots \\ 0 \end{pmatrix} &= \mathbf{X}_k \mathbf{H}_k^i \begin{pmatrix} R \\ 0 \\ \vdots \\ 0 \end{pmatrix} \\ &+ X_k H_{k,k-1}(0, \dots, 0, I) \mathbf{H}_k^{i-1} \begin{pmatrix} R \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{X}_k \mathbf{H}_k^i \begin{pmatrix} R \\ 0 \\ \vdots \\ 0 \end{pmatrix} \end{aligned}$$

since

$$(0, \dots, 0, I) \mathbf{H}_k^{i-1} \begin{pmatrix} R \\ 0 \\ \vdots \\ 0 \end{pmatrix} = 0.$$

The lemma then follows from

$$A^i B = A^i \mathbf{X}_k \begin{pmatrix} R \\ 0 \\ \vdots \\ 0 \end{pmatrix} = A^{i-1} \mathbf{X}_k \mathbf{H}_k \begin{pmatrix} R \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \cdots = A^0 \mathbf{X}_k \mathbf{H}_k^i \begin{pmatrix} R \\ 0 \\ \vdots \\ 0 \end{pmatrix}. \quad \square$$

Theorem 7. Suppose the Arnoldi algorithm is run with matrix A and starting vectors B , where $\rho(A) < \delta < 1$ for some δ . Suppose also that $\rho(\mathbf{H}_k) < \delta < 1$. Then

$$\|W_c - \tilde{W}\|_2 \leq \frac{2\|B\|_2^2 \theta^2 \delta^{2k}}{1 - \delta^2},$$

where θ is the constant from Theorem 5.

Proof. By Lemma 6, for $0 \leq i < k$,

$$A^i B B^T (A^T)^i = \mathbf{X}_k \mathbf{H}_k^i \begin{pmatrix} R \\ 0 \end{pmatrix} (R^T \ 0) (\mathbf{H}_k^T)^i \mathbf{X}_k^T.$$

Therefore, we have the identity

$$\begin{aligned} W_c - \tilde{W} &= \sum_{i=0}^{\infty} A^i B B^T (A^T)^i - \mathbf{X}_k \sum_{i=0}^{\infty} \left(\mathbf{H}_k^i \begin{pmatrix} R \\ 0 \end{pmatrix} (R^T \ 0) (\mathbf{H}_k^T)^i \right) \mathbf{X}_k^T \\ &= \sum_{i=k}^{\infty} A^i B B^T (A^T)^i - \mathbf{X}_k \sum_{i=k}^{\infty} \left(\mathbf{H}_k^i \begin{pmatrix} R \\ 0 \end{pmatrix} (R^T \ 0) (\mathbf{H}_k^T)^i \right) \mathbf{X}_k^T. \end{aligned}$$

It remains to bound the last sums.

$$\begin{aligned} \left\| \sum_{i=k}^{\infty} A^i B B^T (A^T)^i \right\| &\leq \|B\|^2 \sum_{i=k}^{\infty} \|A^i\|^2 \\ &\leq \|B\|^2 \sum_{i=k}^{\infty} \theta^2 \delta^{2i} \\ &\leq \|B\|^2 \theta^2 \delta^{2k} \frac{1}{1 - \delta^2}. \end{aligned}$$

We have the same bound for the sum in \mathbf{H}_k^i , where $\|R\| = \|B\|$. So the final bound on the above sums is as given in the theorem. \square

Normally, as k increases, i.e., as we carry out the Arnoldi algorithm to more steps, the eigenvalues of \mathbf{H}_k should become ever better approximations to those of A , and hence should be bounded by some $\delta < 1$. So as k increases, the bound in Theorem 7 will go to zero, and we obtain ever better approximations to the Gramian.

5. Model reduction via parameter matching

An alternative approach to model reduction is to expand the transfer function in a Laurent series, and then to create a lower dimensional model for which the first few terms of its Laurent series match those of the original model. If

$$F(s) = \sum_{i=0}^{\infty} \frac{F_i}{s^{i+1}} \quad (37)$$

is a transfer function, then the idea is to construct a smaller model that yields a transfer function

$$\hat{F}(s) = \sum_{i=0}^{\infty} \frac{\hat{F}_i}{s^{i+1}} \quad (38)$$

for which $\hat{F}_i = F_i$ for $i = 0, \dots, k$ for some given k . The parameters F_i are called the *high frequency moments* [67] or *Markov parameters* [39], [61]. It is a simple consequence of the matrix identity (Neumann series)

$$(I - M)^{-1} = \sum_{i=0}^{\infty} M^i \text{ whenever } \rho(M) < 1$$

that the transfer function for (18) and (19) satisfies the identity

$$F(s) = C(sI - A)^{-1}B = \sum_{i=0}^{\infty} \frac{CA^iB}{s^{i+1}} \text{ for } |s| > \rho(A). \quad (39)$$

With respect to (19), if $\mathbf{u}_i = 0$, $i = 1, 2, \dots$, and $\mathbf{x}_0 = 0$, then the sequence $\{\mathbf{y}_{i+1} = CA^iB\mathbf{u}_0\}$, $i = 0, 1, 2, \dots$ is the *impulse response* from the impulse \mathbf{u}_0 . The systems (18) and (19) are *state-space realizations* of the transfer function (39).

In applications, $F(s)$ may be either the usual case of a finite dimensional system of large order or an infinite dimensional system [32]. In either case, the model reduction problem is "solved" by constructing (*realizing*) a lower order model of the form (18) such that $\hat{C}\hat{A}^i\hat{B} = F_i$ for $i = 0, \dots, k$, for a given k . In the SISO case, such a reduced order realization may be found via computation of a transfer function $\hat{F}(s)$, expressed as a rational function of polynomials (see e.g., [24, Ch. 15, §10], [18], [31], [41]). Several authors have also looked at the MIMO realization problem (see e.g., [63] and refs. therein, such as [3, 14, 45, 69]). Many of the algorithms involved are recursive algorithms that are intimately related to the clustered Lanczos algorithm (see e.g., [11], [31], [53]). We describe two approaches for the MIMO case proposed recently for constructing such a smaller state-space model, both based on the use of a large state-space model.

The first approach was proposed in [32] as a method for finding a finite dimensional approximation to an infinite dimensional system, in the case where the transfer function $F(s)$ is given, but not the expansion (37) nor a state-space realization. In this case, the authors of [32] proposed applying the discrete Fourier transform to the sequence $F(1), F(\omega), F(\omega^2), \dots, F(\omega^{N-1})$ for some $N > k$,

where ω is an N th root of unity to obtain approximations to the first N terms in the series (37). Next they constructed a triple of matrices \hat{A} , \hat{B} , \hat{C} such that $\hat{C}\hat{A}^i\hat{B} = F_i$ for $i = 0, \dots, k$ and $\hat{C}\hat{A}^i\hat{B} = 0$ for $i > k$. They chose the triple \hat{A} , \hat{B} , \hat{C} to be:

$$\hat{A} = \begin{pmatrix} 0 & & & & \\ I & \ddots & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & I & 0 \end{pmatrix}, \quad \hat{B} = \begin{pmatrix} I \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad \hat{C} = (F_0 \ F_1 \ \dots \ F_k).$$

We refer the reader to [32] for the detailed derivations, including how they applied further model reduction techniques to the resulting finite dimensional model. In particular, they used the balanced realization as described in the previous section.

The second approach that we describe was proposed in [43], [44], and is based on the theory in [67]. In this case, we start with a state-space model of the form (18) (or (19)) represented by the triple A, B, C of large dimensionality and we try to find a triple $\hat{A}, \hat{B}, \hat{C}$ for which the corresponding series (38), defined by $\hat{F}_i = \hat{C}\hat{A}^i\hat{B}$, agrees with (37) up to a given number of terms. The reduced order model is constructed by the use of the block nonsymmetric Lanczos algorithm. We assume that we have an initial model (18), and we assume for simplicity in presentation that B, C , and $C^T B$ all have full rank. The block nonsymmetric Lanczos algorithm is applied with matrix A and starting vectors $X_0 B_0 = B$ and $Y_0 \Gamma_0 = C^T$, using the scaling (17), so that B_0 and Γ_0 are also constructed so that $Y_0^T X_0 = I$. After k steps, we have the generated vectors (14), which satisfy

$$\text{COLSP } \mathbf{X}_k = \mathcal{K}(A, B, k), \quad \text{COLSP } \mathbf{Y}_k = \mathcal{K}(A^T, C^T, k), \quad \text{and } \mathbf{Y}_k^T \mathbf{X}_k = I.$$

Then the reduced order model is defined by

$$\dot{\mathbf{z}}(t) = \hat{A}\mathbf{z}(t) + \hat{B}\mathbf{u}(t), \quad \mathbf{y}(t) = \hat{C}\mathbf{z}(t), \tag{40}$$

where

$$\hat{A} = \mathbf{Y}_k^T A \mathbf{X}_k, \quad \hat{B} = \mathbf{Y}_k^T B, \quad \hat{C} = C \mathbf{X}_k, \quad \mathbf{z} = \mathbf{Y}_k^T \mathbf{x}. \tag{41}$$

It remains to verify that $\hat{C}\hat{A}^i\hat{B} = F_i$ for sufficiently many values of i .

Theorem 8. *Given a system (18), together with the results from k steps of the block nonsymmetric Lanczos algorithm started with A, B, C using scaling (17), and given the reduced order model defined by (40), (41), then*

$$\hat{C}\hat{A}^i\hat{B} = C A^i B \text{ for } i = 0, \dots, 2(k - 1).$$

To prove this, we use the following lemma.

Lemma 9 [67]. $\hat{A}^i \hat{B} = \mathbf{Y}_k^T A^i B$ and $(\hat{A}^T)^i \hat{C}^T = \mathbf{X}_k^T (A^T)^i C^T$ for $0 \leq i \leq k - 1$.

Proof of Lemma 9. The properties (15), (16) are identical to (4), so the proof is analogous to that of Lemma 6.

Proof of Theorem 8. Choose any i, j such that $i < k, j < k$. Then $A^j B \in \mathcal{K}(A, B, k)$; hence we can write $A^j B = \mathbf{X}_k V$ for some matrix V . It follows that

$$\mathbf{X}_k \mathbf{Y}_k^T A^j B = \mathbf{X}_k \mathbf{Y}_k^T \mathbf{X}_k V = \mathbf{X}_k V = A^j B.$$

We then use Lemma 9 to establish

$$\hat{C} \hat{A}^i \hat{A}^j \hat{B} = C A^i \mathbf{X}_k \mathbf{Y}_k^T A^j B = C A^i A^j B. \quad \square$$

The reduced order model obtained by this method exhibits similar transient response behavior to that of the original model. In order to obtain a model that also exhibits similar steady state behavior, it is necessary to compute a model in which the *low frequency moments* $\hat{C} \hat{A}^{-i-1} \hat{B}$, $i = 0, \dots, p$, match those for the original system [67]. One such model is obtained by an oblique projection

$$\dot{\mathbf{z}} = L^T A T \mathbf{z} + L^T B \mathbf{u}, \quad \mathbf{y} = C T \mathbf{z}, \quad (42)$$

where the columns of matrices L and T include bases for the spaces $\mathcal{K}(A^{-T}, A^{-T}, C^T, p)$ and $\mathcal{K}(A^{-1}, A^{-1}B, p)$, respectively [67], where $M^{-T} \equiv (M^T)^{-1} \equiv (M^{-1})^T$. In order to create models matching both the high frequency moments and the low frequency moments, [63] proposed combining the vectors generated from the block Lanczos method using A, B, C^T with those from the Lanczos method using $A^{-1}, A^{-1}B, C^T$. This is based on the following theorem.

Theorem 10 [67]. *If L, T are bases for the spaces*

$$\begin{aligned} & \text{COLSP}((A^T)^{-p} C^T \dots (A^T)^{+q} C^T), \\ & \text{COLSP}(A^{-s} B \dots A^{+t} B), \end{aligned}$$

where $L^T T = I$ and $p + q = s + t$, then the system (42) obtained by the oblique projection satisfies

$$(C T)(L^T A T)^i (L^T B) = C A^i B \text{ for } i = -p - s - 1, \dots, +q + t + 1.$$

The proof is analogous to that of Theorem 8, which is a special case of Theorem 10.

This approach has been exploited very successfully in [15], [63], [64] to applications in large flexible space structures. One problem with this method is that the block nonsymmetric Lanczos algorithm may break down, if the matrix in step 7 of the algorithm is singular. In practice, this has not been a problem in the applications in [64]. We note that in this use of the algorithm, we do not compute the spaces $\mathcal{K}(A, B, \infty)$ or $\mathcal{K}(A^T, C^T, \infty)$. Rather, we choose a k_{\max} and run the algorithm only that many steps. If k_{\max} is much less than the dimension of the entire spaces $\mathcal{K}(A, B, \infty)$ or $\mathcal{K}(A^T, C^T, \infty)$, it is reasonable not to expect breakdowns in the first k_{\max} steps.

6. Conclusions

In this paper we have described three variants of the Lanczos algorithm for generating Krylov sequences and have given three examples of control theory applications where these Lanczos algorithms can be used. These algorithms are particularly suited to solving control problems involving state-space models with sparse matrices of very large order. The development of algorithms for such control problems is still in its infancy, especially since many algorithms for “small” dense state-space models are also of recent origin [48]. We have illustrated the use of the Lanczos algorithms on problems of controllability, observability, and model reduction. The Gramians are also useful for other computations in optimal control and \mathcal{H}_∞ theory (e.g., finding the \mathcal{H}_2 , \mathcal{H}_∞ , and Hankel norms of plants [12, and refs. therein]), hence the methods of this paper can be used to compute approximate solutions to those problems.

We remark that the reduced order systems (42), (40) have transfer functions that are rational functions whose Laurent or Macclaurin series expansions agree with the series expansions of the original transfer functions of a given number of terms. In the space of the transfer functions themselves, the problem of computing low degree rational approximations to given meromorphic functions is a classical problem in approximation theory in the complex plane. Many methods operating directly on the coefficients in the expansions of the form (37) (or more general Macclaurin expansions) exist, and many are equivalent to the Lanczos process ([11], [13], [31], [30], [53], and refs. therein), and many use related algorithms on related Hankel matrices, though based on different theory [1], [35], [65]. The methods presented in this paper are distinctive in that they are based on the use of state-space descriptions, which facilitate the generalization of many SISO methods to MIMO systems, and which may often allow one to avoid necessary conversions between state-space and frequency domain descriptions and the ensuing computational and numerical difficulties. Space does not permit us to explore here the relations and equivalences between the state-space approximations and the frequency domain methods based on optimal Hankel, \mathbf{H}_∞ , or Padé approximations (mentioned in some of the references), or the close relations with fast Hankel matrix factorization methods and orthogonal polynomials. These will be explored in a future paper.

7. Acknowledgments

The author would like to acknowledge the helpful comments from Youcef Saad, Roland Freund, and Franklin Luk, as well as the support from the Minnesota Supercomputer Institute.

References

- [1] V. Adamjan, D. Arov, and M. Krein, Analytic properties of Schmidt pairs for a Hankel operator and the generalized Schur-Takagi problem, *Math. USSR Sbornik* 15 (1971), pp. 31–73.
- [2] U. M. Al-Saggaf and G. F. Franklin, An error bound for a discrete reduced order model of a linear multivariable system, *IEEE Trans. Auto. Contr.* AC-32 (1987), pp. 815–819.
- [3] A. Antoulas, New results on the algebraic theory of linear systems: The solution of the cover problems, *Linear Algebra Appl.* 50 (1983), pp. 1–45.
- [4] W. E. Arnoldi, The principle of minimized iterations in the solution of the matrix eigenvalue problem, *Quart. Appl. Math.* 9 (1951), pp. 17–29.
- [5] A. W. Bojanczyk, L. W. Ewerbring, F. T. Luk, and P. Van Dooren, An accurate product SVD algorithm, *Signal Processing* 25 (1991), pp. 189–201.
- [6] D. L. Boley, Computing the Kalman decomposition, an optimal method, *IEEE Trans. Auto. Contr.* AC-29 (1984), pp. 51–53. (Correction in AC-36 (1991), p. 1341).
- [7] D. L. Boley, R. P. Brent, G. H. Golub, and F. T. Luk, Error correction via the Lanczos process, *SIAM J. Matrix Anal.* (1992), pp. 312–332.
- [8] D. L. Boley, S. Elhay, G. H. Golub, and M. H. Gutknecht, Nonsymmetric Lanczos and finding orthogonal polynomials associated with indefinite weights, *Numerical Algorithms* 1 (1991), pp. 21–44.
- [9] D. L. Boley and G. H. Golub, The Lanczos algorithm and controllability, *Systems Control Lett.* 4 (1984), pp. 317–324.
- [10] D. L. Boley and G. H. Golub, The nonsymmetric Lanczos algorithm and controllability, *Systems Control Lett.* 16 (1991), pp. 97–105.
- [11] D. L. Boley, T. J. Lee, and F. T. Luk, The Lanczos algorithm and Hankel matrix factorization, *Linear Algebra Appl.* 172 (1992), pp. 109–133.
- [12] S. P. Boyd and C. H. Barratt, *Linear Controller Design*, Prentice-Hall, Englewood Cliffs, NJ, 1991.
- [13] C. Brezinski, *Padé-Type Approximation and General Orthogonal Polynomials*, Birkhäuser, Basel and Boston, 1980.
- [14] A. Bultheel, Recursive algorithms for the matrix Padé problem, *Math. Comp.* 35 (1980), pp. 875–892.
- [15] R. R. Craig, Jr. and A. L. Hale, Block-Krylov component synthesis method for structural model reduction, *J. Guid., Contr. Dyn.* 11 (1988), pp. 562–570.
- [16] J. Cullum, W. Kerner, and R. Willoughby, A generalized nonsymmetric Lanczos procedure, *Computer Physics Communications* 53 (1989), pp. 19–48.
- [17] J. Cullum and R. Willoughby, *Lanczos Algorithms for Large Symmetric Eigenvalue Computations*, vol. I, Theory, Birkhäuser, Basel and Boston, 1985.
- [18] L. S. de Jong, Numerical aspects of recursive realization algorithms, *SIAM J. Control Optim.* 16 (1978), pp. 646–659.
- [19] C. A. Desoer, *A Second Course in Linear Systems*, Van Nostrand Reinhold, New York, 1970.
- [20] P. V. Dooren, Numerical linear algebra techniques for large scale matrix problems in systems and control, in Proc. 31st Conf. on Dec. and Contr., Tucson, AZ, Dec. 1992, Session TM6, to appear.
- [21] P. V. Dooren, *Upcoming numerical linear algebra issues in systems and control theory*, TR 1007, Univ. of Minn. IMA, August 1992.
- [22] D. F. Enns, Model reduction with balanced realizations: An error bound and a frequency weighted generalization, in Proc. 23rd IEEE Conf. Dec. Contr., 1984, pp. 127–132.
- [23] R. W. Freund, M. H. Gutknecht, and N. M. Nachtigal, An implementation of the look-ahead Lanczos algorithm for non-hermitian matrices, *SIAM J. Scientific Comp.*, 14 (1993), pp. 137–158.
- [24] F. R. Gantmacher, *Theory of Matrices*, vol. 2, Chelsea, New York, 1959.
- [25] E. G. Gilbert, Controllability and observability in multivariable control systems, *SIAM J. Control* 1 (1963), pp. 128–151.

- [26] K. Glover, All optimal Hankel norm approximations of linear multivariable systems, and their L_∞ error bounds, *Intl. J. Contr.* 39 (1984), pp. 1115–1193.
- [27] G. H. Golub and M. H. Gutknecht, Modified moments for indefinite weight functions, *Numer. Math.* 57 (1990), pp. 607–624.
- [28] G. H. Golub and R. Underwood, The block Lanczos method for computing eigenvalues, in *Mathematical Software III*, J. Rice, ed., pp. 364–377, Academic Press, New York and San Diego, 1977.
- [29] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 2d ed., Johns Hopkins Univ. Press, Baltimore, MD, 1989.
- [30] W. B. Gragg, The Padé table and its relation to certain algorithms in numerical analysis, *SIAM Rev.* 14 (1972), pp. 1–62.
- [31] W. B. Gragg and A. Lindquist, On the partial realization problem, *Linear Algebra Appl.* 50 (1985), pp. 277–319.
- [32] G. Gu, P. P. Khargonekar, and E. B. Lee, Approximation of infinite-dimensional systems, *IEEE Trans. Auto. Contr.* AC-34 (1989), pp. 610–618.
- [33] M. H. Gutknecht, The unsymmetric Lanczos algorithms and their relations to Padé approximation, continued fractions, and gd algorithm, biconjugate gradient squared algorithms, and fast Hankel solvers, preprint, 1990.
- [34] M. H. Gutknecht, A completed theory of the unsymmetric Lanczos process and related algorithms, part I, *SIAM J. Math. Anal.* 13 (1992), pp. 594–639.
- [35] M. H. Gutknecht and L. N. Trefethen, Real polynomial Chebyshev approximation by the Carathéodory-Féjer method, *SIAM J. Numer. Anal.* 19 (1982), pp. 358–371.
- [36] S. J. Hammarling, Numerical solution of stable, non-negative definite Lyapunov equations, *IMA J. Numer. Anal.* 2 (1982), pp. 303–323.
- [37] S. J. Hammarling, Numerical solutions of the discrete-time convergent, non-negative definite Lyapunov equations, *Systems Control Lett.* 17 (1991), pp. 137–139.
- [38] M. T. Heath, A. J. Laub, C. C. Paige, and R. C. Ward, Computing the SVD of a product of two matrices, *SIAM J. Sci. Statist. Comput.* 7 (1986), pp. 1147–1159.
- [39] J. Hickin and N. K. Sinha, Model reduction for linear multivariable systems, *IEEE Trans. Auto. Contr.* AC-25 (1980), p. 1121–1127.
- [40] R. E. Kalman, Mathematical description of linear systems, *SIAM J. Contr.* 1 (1963), pp. 152–192.
- [41] R. E. Kalman, On partial realizations, transfer functions and canonical forms, *Acta Polytech. Scand.* MA31 (1979), pp. 9–32.
- [42] S. Kaniel, Estimates for some computational techniques in linear algebra, *Math. Comp.* 20 (1966), pp. 369–378.
- [43] H. M. Kim and R. R. Craig, Jr., Structural dynamics analysis using an unsymmetric block Lanczos algorithm, *Internat. J. Numer. Methods Engrg.* 26 (1988), pp. 2305–2318.
- [44] H. M. Kim and R. R. Craig, Jr., Computational enhancement of an unsymmetric block Lanczos algorithm, *Internat. J. Numer. Methods Engrg.* 30 (1990), pp. 1083–1089.
- [45] S. Y. Kung, A new identification and model reduction algorithm via singular value decompositions, in Proc. 12th Asilomar Conf. Circ. Syst. Comp., November 1984, pp. 705–714.
- [46] C. Lanczos, An iteration method for the solution of the eigenvalue problem linear differential and integral operators, *J. Res. Natl. Bur. Stand.* 45 (1950), pp. 255–282.
- [47] A. J. Laub, On computing “balancing” transformations, in Proc. 1980 JACC, San Francisco, 1980, pp. FA8-E.
- [48] A. J. Laub, Numerical linear algebra aspects of control design computations, *IEEE Trans. Auto. Contr.* AC-30 (1985), pp. 97–108.
- [49] A. J. Laub, M. T. Heath, C. C. Paige, and R. C. Ward, Computation of system balancing transformations and other applications of simultaneous diagonalization algorithms, *IEEE Trans. Auto. Contr.* AC-32 (1987), pp. 115–122.
- [50] B. C. Moore, Principal component analysis in linear systems: Controllability, observability, and model reduction, *IEEE Trans. Auto. Contr.* AC-26 (1981), pp. 17–31.

- [51] C. C. Paige, The computation of eigenvalues and eigenvectors of very large sparse matrices, Ph.D. thesis, Univ. of London, 1971.
- [52] B. N. Parlett, *The Symmetric Eigenvalue Problem*, Prentice-Hall, Englewood Cliffs, NJ, 1980.
- [53] B. N. Parlett, Reduction to tridiagonal form and minimal realizations, *SIAM J. Matr. Anal.* 13 (1992), pp. 567–593.
- [54] B. N. Parlett, D. R. Taylor, and Z. A. Liu, A look-ahead Lanczos algorithm for unsymmetric matrices, *Math. Comp.* 44 (1985), pp. 105–124.
- [55] L. Reichel and D. Y. Hu, Krylov subspace methods for the Sylvester equation, *Linear Algebra Appl.* 172 (1992), pp. 283–314.
- [56] Y. Saad, On the rates of convergence of the Lanczos and the block Lanczos methods, *SIAM J. Num. Anal.* 17 (1980), pp. 687–706.
- [57] Y. Saad, Krylov subspace methods for solving large unsymmetric linear systems, *Math. Comp.* 37 (1981), pp. 105–126.
- [58] Y. Saad, Numerical solution for large Lyapunov equations, in *Signal Processing, Scattering and Operator Theory, and Numerical Methods, Proc. Intl. Symp. MTNS-89*, vol. 3, pp. 503–511, Birkhäuser, Basel and Boston, 1990.
- [59] Y. Saad and M. H. Schultz, GMRES: A generalized minimal residual algorithm for solving unsymmetric linear systems, *SIAM J. Sci. Stat. Comput.* 7 (1986), pp. 856–869.
- [60] D. Scott, Analysis of the symmetric Lanczos process, Electronic Res. Lab. Report UCB/ERL M78/40, Univ. of Calif., Berkeley, 1978.
- [61] L. S. Shieh and Y. J. Wei, A mixed method for multivariable system reduction, *IEEE Trans. Auto. Contr.* AC-20 (1975), pp. 429–432.
- [62] G. W. Stewart and J. G. Sun, *Matrix Perturbation Theory*, Academic Press, New York and San Diego, 1990.
- [63] T. J. Su, A decentralized linear quadratic control design method for flexible structures, Ph.D. thesis, Univ. of Texas, Austin, 1989.
- [64] T. J. Su and R. R. Craig, Jr., Model reduction and control of flexible structures using Krylov vectors, *J. Guid., Contr., Dyn.* 14 (1991), pp. 260–267.
- [65] L. N. Trefethen, Rational Chebyshev approximation on the unit disk, *Numer. Math.* 37 (1981), pp. 297–320.
- [66] P. Van Dooren, Numerical aspects of system and control problems, *Journal A* 30 (1987), pp. 25–32.
- [67] C. D. Vilemagne and R. E. Skelton, Model reduction using a projection formulation, *Intl. J. Contr.* 46 (1987), pp. 2141–2169.
- [68] J. H. Wilkinson, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.
- [69] H. P. Zeiger and A. J. McEwen, Approximate linear realizations of given dimensions via Ho's algorithm, *IEEE Trans. Auto. Contr.* AC-19 (1974), p. 153.