

STABILITY CRITERIA FOR ARNOLDI-BASED MODEL-ORDER REDUCTION

I. M. Elfadel L. Miguel Silveira J. White

Research Laboratory of Electronics
Massachusetts Institute of Technology
Cambridge, MA 02139

ABSTRACT

Padé approximation is an often-used method for reducing the order of a finite-dimensional, linear, time invariant, signal model. It is known to suffer from two problems: numerical instability during the computation of the Padé coefficients and lack of guaranteed stability for the resulting reduced model *even* when the original system is stable. In this paper, we show how the numerical instability problem can be avoided using the Arnoldi algorithm applied to an appropriately chosen Krylov subspace. Moreover, we give an easily computable sufficient condition on the system matrix that guarantees the stability of the reduced model at any approximation order.

1. INTRODUCTION

Consider the discrete-time, linear, time-invariant (LTI) system

$$\begin{aligned} \mathbf{x}(t) &= \mathbf{A}\mathbf{x}(t-1) + \mathbf{b}u(t-1) \\ y(t) &= \mathbf{c}\mathbf{x}(t) \end{aligned} \quad (1)$$

where $t \in \{1, 2, \dots\}$, $\mathbf{x}(t), \mathbf{b}, \mathbf{c}^T \in \mathbb{R}^n$, $u(t), y(t) \in \mathbb{R}$, and $\mathbf{A} \in \mathbb{R}^{n \times n}$. We will assume this system both reachable and observable, which means that the ranks of the reachability matrix

$$(\mathbf{b}, \mathbf{A}\mathbf{b}, \dots, \mathbf{A}^{n-1}\mathbf{b})$$

and observability matrix

$$(\mathbf{c}^T, \mathbf{c}^T \mathbf{A}^T, \dots, \mathbf{c}^T (\mathbf{A}^{n-1})^T)$$

are both equal to n . The transfer function of this linear system is $G(z) = z^{-1}\mathbf{c}(\mathbf{I}_n - z^{-1}\mathbf{A})^{-1}\mathbf{b}$, where

\mathbf{I}_p is the identity matrix of order p . The rational function $G(z)$ can also be written as the power series

$$G(z) = \sum_{k=0}^{\infty} \mathbf{c}\mathbf{A}^k \mathbf{b} z^{-k-1} = \sum_{k=0}^{\infty} g_k z^{-k-1}, \quad (2)$$

where the scalars $g_k \triangleq \mathbf{c}\mathbf{A}^k \mathbf{b}$ are called the Markov coefficients of (1). When $u(t)$ is an impulse, $y(t) = g(t)$, $t \geq 0$, and the absolute summability of the sequence $g(t)$ is equivalent to the stability of the transfer function $G(z)$ in the bounded-input/bounded-output sense. We denote by λ_i the i -th eigenvalue of \mathbf{A} and define the spectral radius of \mathbf{A} as $\rho(\mathbf{A}) = \max\{|\lambda_i|, 1 \leq i \leq n\}$. The stability of the LTI system (1) is then equivalent to $\rho(\mathbf{A}) < 1$. In this paper, we will speak interchangeably of the stability of the transfer function or the stability of its system matrix \mathbf{A} .

Assuming that the dimension of the state space n is very large and that we are given an integer $q < n$, the Padé approximation method aims at finding a triplet $(\tilde{\mathbf{A}}, \tilde{\mathbf{b}}, \tilde{\mathbf{c}}) \in \mathbb{R}^{q \times q} \times \mathbb{R}^{q \times 1} \times \mathbb{R}^{1 \times q}$, such that the transfer function $\tilde{G}(z) = z^{-1}\tilde{\mathbf{c}}(\mathbf{I}_q - z^{-1}\tilde{\mathbf{A}})^{-1}\tilde{\mathbf{b}}$ of the reduced system approximates the transfer function $G(z)$ in the sense that there is perfect matching between the coefficients of the Markov coefficients of the original and reduced models up to a certain order m , i.e.,

$$\tilde{g}_k = \tilde{\mathbf{c}}\tilde{\mathbf{A}}^k \tilde{\mathbf{b}} = \mathbf{c}\mathbf{A}^k \mathbf{b} = g_k, \quad 0 \leq k \leq m-1.$$

It is important to point out that this matching condition corresponds to the requirement that the *transient* behavior of the original and the reduced-order models be the same. This is because the expansion of the transfer function $G(z)$ given in Equation 2 is accomplished in the neighborhood of $z^{-1} = 0$. If we were to require the matching to be at *steady-state*, then we will have to formulate the matching condition in terms of Markov coefficients obtained from an expansion of $G(z)$ in the neighborhood of

This work was supported by ARPA contracts N00014-91-J-1698, N00174-93-C-0035, and grants from the Semiconductor Research Corporation (SJ-558), IBM and Digital Equipment Corporation.

$z^{-1} = 1$. This aspect of the problem will be addressed elsewhere.

The classical procedure ([1], Chapter 3) to find the approximate transfer function is to solve a Hankel linear system in which the Hankel matrix is based on the Markov coefficients of (1) and the unknowns are the coefficients of the rational function $\tilde{G}(z)$. This procedure suffers from two main problems. The first is computational and is related to the computation of the Markov coefficients of the original system. This computation involves the power iterations of the large system matrix \mathbf{A} . It is well known ([2], Chapter 7), that for $\mathbf{r} \in \mathbb{R}^n$ the iterates $\mathbf{A}^k \mathbf{r}$ converges generically to the eigenvector of \mathbf{A} corresponding to the eigenvalue of the largest magnitude. In other words, the Markov coefficients become close to each other, thus making the Hankel system very ill-conditioned [3].

Another more fundamental problem is that even when the LTI system (1) is stable and the computations are well-conditioned, there is no guarantee that the resulting reduced order system $\tilde{G}(s)$ will be stable. This instability could occur *even* when the original system is an FIR filter, i.e, $H(z) = \sum_{k=0}^N h_k z^k$ ([1], Example 3.9).

In this paper, we address these two problems and show that the Arnoldi algorithm (see [4] and the references therein) provides a numerically stable way for obtaining a reduced-order model. We also show that when the system matrix \mathbf{A} is normal and stable, the Arnoldi reduced-order matrix $\tilde{\mathbf{A}}$ is guaranteed stable at any approximation order q . When \mathbf{A} is not normal, we provide an easily computable sufficient condition for guaranteeing the stability of $\tilde{\mathbf{A}}$ at any approximation order q .

2. REDUCED-ORDER MODEL

The computation of the Markov coefficients involves the power iterates $\mathbf{A}^k \mathbf{b}$. It is therefore natural to consider the Krylov subspace $\mathcal{K}_q(\mathbf{A}, \mathbf{b}) = \text{span}\{\mathbf{b}, \mathbf{A}\mathbf{b}, \mathbf{A}^2\mathbf{b}, \dots, \mathbf{A}^{q-1}\mathbf{b}\}$. Because of the reachability assumption, this subspace is of dimension q . The essence of the Arnoldi algorithm is to use the Gram-Schmidt procedure to build an orthonormal basis $\mathbf{V}_q = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_q\}$ of the Krylov subspace $\mathcal{K}_q(\mathbf{A}, \mathbf{b})$. At the k -th step of the algorithm as shown in (2.1), a unit-norm vector \mathbf{v}_k is constructed such that \mathbf{v}_k is orthogonal to $\mathcal{K}_{k-1}(\mathbf{A}, \mathbf{b})$.

After q steps, the Arnoldi algorithm returns a set of q orthonormal vectors, as the columns of the matrix $\mathbf{V}_q \in \mathbb{R}^{m \times q}$, and a $q \times q$ upper Hessenberg matrix $\mathbf{H}_q = [h_{i,j}]$. These two matrices satisfy the

Algorithm 2.1 (Arnoldi Algorithm)

```

arnoldi(input  $\mathbf{A}, \mathbf{b}, q$ ;
output  $\mathbf{V}_q, \mathbf{v}_{q+1}, \mathbf{H}_q, h_{q+1,q}$ )
{
   $\mathbf{v}_1 = \mathbf{b}/\|\mathbf{b}\|_2$ 
  for ( $j = 1$ ;  $j \leq q$ ;  $j++$ ) {
     $\mathbf{w} = \mathbf{A}\mathbf{v}_j$ 
    for ( $i = 1$ ;  $i \leq j-1$ ;  $i++$ ) {
       $h_{i,j} = \mathbf{w}^T \mathbf{v}_i$ 
       $\mathbf{w} = \mathbf{w} - h_{i,j} \mathbf{v}_i$ 
    }
     $h_{j+1,j} = \|\mathbf{w}\|_2$ 
    if ( $h_{j+1,j} \neq 0$ ) {
       $\mathbf{v}_{j+1} = \mathbf{w}/h_{j+1,j}$ 
    }
  }
   $\mathbf{V}_q = [\mathbf{v}_1 \dots \mathbf{v}_q]$ 
   $\mathbf{H}_q = (h_{i,j}), \quad i, j = 1, \dots, q$ 
}

```

following relationship:

$$\mathbf{A}\mathbf{V}_q = \mathbf{V}_q \mathbf{H}_q + h_{q+1,q} \mathbf{v}_{q+1} \mathbf{e}_q^T \quad (3)$$

where \mathbf{e}_q is the q -th unit vector in \mathbb{R}^n , and $\mathbf{v}_{q+1} \in \mathbb{R}^n$ is orthonormal to the columns of \mathbf{V}_q . Using the fact that $\mathbf{V}_q^T \mathbf{V}_q = \mathbf{I}_q$, we can write the above equation as

$$\mathbf{V}_q^T \mathbf{A}\mathbf{V}_q = \mathbf{H}_q. \quad (4)$$

As noted in [5], the above equation defines a congruence transform that allows \mathbf{H}_q to inherit the passivity of \mathbf{A} whenever the latter matrix is symmetric.

Furthermore, using the fact that $\mathbf{V}_q \mathbf{e}_1 = \mathbf{b}/\|\mathbf{b}\|_2$, it can be easily seen that after q steps of an Arnoldi process,

$$\mathbf{A}^k \mathbf{b} = \|\mathbf{b}\|_2 \mathbf{A}^k \mathbf{V}_q \mathbf{e}_1 = \|\mathbf{b}\|_2 \mathbf{V}_q \mathbf{H}_q^k \mathbf{e}_1, \quad 0 \leq k \leq q-1, \quad (5)$$

which yields

$$\mathbf{c}\mathbf{A}^k \mathbf{b} = \|\mathbf{b}\|_2 \mathbf{c}\mathbf{V}_q \mathbf{H}_q^k \mathbf{e}_1, \quad 0 \leq k \leq q-1. \quad (6)$$

Therefore, if we choose for the reduced-order model the triplet $(\tilde{\mathbf{A}}, \tilde{\mathbf{b}}, \tilde{\mathbf{c}}) = (\mathbf{H}_q, \mathbf{e}_1, \|\mathbf{b}\|_2 \mathbf{c}\mathbf{V}_q)$, the original system and the reduced-order system will have q

of their Markov coefficients matched. The transfer function of the reduced-order model is then

$$\tilde{G}(z) = z^{-1} \|b\|_2 c^T V_q (\mathbf{I}_q - z^{-1} \mathbf{H}_q)^{-1} e_1 \quad (7)$$

Note that the state-space realization of the Arnoldi reduced-order model comes naturally in a *system Hessenberg form* [6]. This algorithm has been successfully used in [4] to produce reduced-order models for a variety of very large linear circuits encountered in the analysis and simulation of VLSI interconnect.

3. STABILITY

From now on we make the assumption that the system matrix \mathbf{A} is stable, i.e., its spectral radius $\rho(\mathbf{A}) < 1$. We would like to find under what conditions the reduced matrix $\tilde{\mathbf{A}} = \mathbf{H}_q$ is itself stable. Formula (4) will be essential in answering this question.

Case 1: \mathbf{A} is symmetric. The symmetric case occurs very often in practice especially in the context of scientific computing [3, 4]. When \mathbf{A} is symmetric, (4) implies that the upper Hessenberg matrix \mathbf{H}_q is symmetric and therefore tridiagonal. We denote the q real eigenvalues of \mathbf{H}_q by $\tilde{\lambda}_1 \geq \tilde{\lambda}_2 \geq \dots \geq \tilde{\lambda}_q$. Because both \mathbf{A} and \mathbf{H}_q are symmetric, we can use the Rayleigh-Ritz quotients ([7], p. 176) to conclude that $\lambda_{\max} \geq \tilde{\lambda}_{\max} \geq \tilde{\lambda}_{\min} \geq \lambda_{\min}$.

In other words, when \mathbf{A} is symmetric, the spectral radius of $\tilde{\mathbf{A}}$ satisfies $\rho(\tilde{\mathbf{A}}) \leq \rho(\mathbf{A}) < 1$, i.e., the stability of the reduced-order model is guaranteed at any order whenever the original matrix \mathbf{A} is stable.

In fact, using the theory of Paige-Kaniel ([2], Chapter 9), we can obtain much sharper results about the location of the eigenvalues of \mathbf{H}_q with respect to those of \mathbf{A} .

Case 2: \mathbf{A} is normal. A symmetric matrix is a special case of a normal matrix, i.e., a matrix that commutes with its transpose. However, normal matrices do not satisfy the Rayleigh-Ritz variational formulas for eigenvalues. Moreover, normality is not preserved under the projection formula (4). Notwithstanding these facts, the concept of the *numerical radius* ([8], p. 7) defined as $r(\mathbf{A}) \equiv \max\{|\mathbf{x}^* \mathbf{A} \mathbf{x}|, \mathbf{x}^* \mathbf{x} = 1\}$ can be used to prove the following

Theorem 3.1 *Assume the matrix \mathbf{A} is normal and stable. Then $\rho(\tilde{\mathbf{A}}) < 1$, i.e., the reduced-order model is stable at any order.*

Let us first establish some basic properties for the numerical radius of a matrix.

Proposition 3.2 *Let \mathbf{M} be an arbitrary, square, complex matrix of order n . Then the $\rho(\mathbf{M}) \leq r(\mathbf{M})$. Moreover, if the matrix $\mathbf{U} \in \mathbb{R}^{n \times q}$ has orthonormal columns, ($\mathbf{U}^T \mathbf{U} = \mathbf{I}_q$), then $r(\mathbf{U}^T \mathbf{M} \mathbf{U}) \leq r(\mathbf{M})$.*

Proof. To prove the first part, let \mathbf{z} be an eigenvector of unit 2-norm corresponding to the eigenvalue, λ_{\max} , of the largest magnitude of \mathbf{M} . Then

$$\begin{aligned} |\mathbf{z}^* \mathbf{M} \mathbf{z}| &= |\lambda_{\max}| \mathbf{z}^* \mathbf{z} = |\lambda_{\max}| \\ &\leq \max\{|\mathbf{x}^* \mathbf{M} \mathbf{x}|, \mathbf{x}^* \mathbf{x} = 1\} = r(\mathbf{M}). \end{aligned}$$

To prove the second part, note that under the assumption $\mathbf{U}^T \mathbf{U} = \mathbf{I}_q$, $\{\mathbf{U} \mathbf{y}, \mathbf{y} \in \mathbb{R}^q, \mathbf{y}^* \mathbf{y} = 1\} \subseteq \{\mathbf{x} \in \mathbb{R}^n, \mathbf{x}^* \mathbf{x} = 1\}$. Therefore $\{\mathbf{y}^* \mathbf{U}^T \mathbf{M} \mathbf{U} \mathbf{y}, \mathbf{y}^* \mathbf{y} = 1\} \subseteq \{\mathbf{x}^* \mathbf{M} \mathbf{x}, \mathbf{x}^* \mathbf{x} = 1\}$, and it follows that

$$r(\mathbf{U}^T \mathbf{M} \mathbf{U}) \leq r(\mathbf{M}).$$

□

We also need the following Lemma for the numerical radius of a normal matrix.

Lemma 3.3 *Assume the matrix \mathbf{M} is normal. Then its numerical radius is equal to its spectral radius, i.e., $r(\mathbf{M}) = \rho(\mathbf{M})$.*

Proof. If the matrix is normal, then it is diagonalizable with a unitary matrix, i.e., there exists $\mathbf{U} \in \mathbb{C}^{n \times n}$ such that $\mathbf{U}^H \mathbf{U} = \mathbf{I}_n$ and $\mathbf{M} = \mathbf{U}^H \mathbf{A} \mathbf{U}$, where \mathbf{A} is the diagonal matrix of the complex eigenvalues of \mathbf{M} . Then we have

$$\begin{aligned} r(\mathbf{M}) &= \max\{|\mathbf{x}^* \mathbf{M} \mathbf{x}|, \mathbf{x}^* \mathbf{x} = 1\} \\ &= \max\{|\mathbf{x}^* \mathbf{U}^H \mathbf{A} \mathbf{U} \mathbf{x}|, \mathbf{x}^* \mathbf{x} = 1\} \\ &= \max\{|\mathbf{v}^* \mathbf{A} \mathbf{v}|, \mathbf{v}^* \mathbf{v} = 1\} \\ &\leq |\lambda_{\max}| = \rho(\mathbf{M}). \end{aligned}$$

Using the first part of Proposition 3.2, we conclude that $r(\mathbf{M}) = \rho(\mathbf{M})$. □

Now, to the proof of Theorem 3.1.

Proof. The theorem results readily from the following sequence of inequalities

$$\rho(\tilde{\mathbf{A}}) \leq r(\tilde{\mathbf{A}}) \leq r(\mathbf{A}) \leq \rho(\mathbf{A}) < 1,$$

where the first and second inequalities are due to the second and the first parts of Proposition 3.2, respectively, while the third results from the normality assumption and Lemma 3.3, and the last inequality is just the stability assumption on the original systems. □

Theorem 3.1 therefore means that for any normal, stable matrix \mathbf{A} the reduced-order system matrix $\tilde{\mathbf{A}}$ is guaranteed stable at any order $q \leq n$.

Case 3: A is not normal. When A is not normal, a sufficient condition on the size of the coefficients of the matrix A can be imposed to get guaranteed stability at any order. Indeed, we have

Theorem 3.4 Assume $\|A\|_1 + \|A\|_\infty < 2$. Then $\rho(\tilde{A}) < 1$, i.e., the reduced-order model is stable at any order. Moreover the assumption is necessary for all stable matrices of the form $A = \alpha P$, where P is a permutation matrix and $|\alpha| < 1$.

Proof. As in the proof of Theorem 3.1, we have the inequalities $\rho(\tilde{A}) \leq r(\tilde{A}) \leq r(A)$, resulting from the Proposition 3.2. On the other hand, the numerical radius of any matrix is always no greater than the average of its ℓ^1 and ℓ^∞ norms, i.e., $r(A) \leq \frac{1}{2}(\|A\|_1 + \|A\|_\infty)$. This latter fact can be derived from a generalization of Gershgorin's disk theorem applied to the field of values $\{x^*Ax \in \mathbb{C}, x^*x = 1\}$. A complete proof can be found in ([8], p. 31-33). Combining these inequalities along with the assumptions leads to the conclusion that for any order $q \leq n$, $\rho(\tilde{A}) < 1$, which means that the reduced-order model is stable at any order.

To show the necessity of the assumption for $A = \alpha P$, $A = \alpha P$, where P is a permutation matrix and $|\alpha| < 1$, note that for a permutation matrix, we have $\|P\|_1 = \|P\|_\infty = 1$, which implies that $\frac{1}{2}(\|A\|_1 + \|A\|_\infty) = |\alpha| < 1$. Note that a permutation matrix is normal, and therefore $r(P) = \rho(P) = 1$. \square

Another instance where the condition of the above theorem is necessarily satisfied is when the matrix $A = \alpha S$, where $|\alpha| < 1$ and S is a doubly stochastic matrix, i.e., both its row and column sums are equal to 1. A celebrated theorem due to Birkhoff ([7], Theorem 8.7.1) states that a doubly stochastic matrix is a convex combination of permutation matrices. This fact can be used to prove that $r(S) \leq 1 = \frac{1}{2}(\|S\|_1 + \|S\|_\infty)$. In other words, $r(A) < 1$, and the Arnoldi-based reduced-order models will be guaranteed stable at any order.

The importance of the numerical radius stems from the fact that for any matrix $A \in \mathbb{R}^{n \times n}$ and any isometry $V \in \mathbb{R}^{n \times q}$, ($V^T V = I_q$), we have $r(V^T A V) \leq r(A)$. This latter inequality is in general not satisfied by the spectral radius. It is also worthwhile noting that if a Lanczos-type algorithm [3, 9] is used to derive the reduced-order matrix \tilde{A} , a guaranteed-stability result similar to Theorem 3.4 is in general not possible. This is because the matrices V_q and W_q produced by the Lanczos algorithm such that $V_q^T A W_q = H_q$ do not allow us to conclude that $r(H_q) \leq r(A)$.

4. CONCLUSION

The contributions of this paper are twofold. First, we have used the Arnoldi algorithm to show how a reduced-order model can be obtained without explicitly solving the ill-conditioned Hankel linear system of the classical Padé approximation procedure. Then we have established some results regarding the guaranteed stability of the reduced model. Although this paper has addressed only single-input/single-output systems, all our conclusions remain valid for the multiple-input/multiple-output case. Indeed, the block Arnoldi algorithm [4] could be used to obtain the reduced-order model, while the stability results continue to hold since they depend solely on Equation (4) and the system matrix A .

5. REFERENCES

- [1] P. A. Regalia. *Adaptive IIR Filtering in Signal Processing and Control*. Dekker, 1995.
- [2] G. H. Golub and C. F. Van Loan. *Matrix Computation*. The Johns Hopkins University Press, second edition, 1989.
- [3] Peter Feldmann and Roland W. Freund. Efficient linear circuit analysis by Padé approximations via the Lanczos process. In *EURO-DAC'94 with EURO-VHDL'94*, September 1994.
- [4] L. M. Silveira, M. Kamon, I. M. Elfadel, and J. White. Coupled circuit-interconnect analysis using Arnoldi-based model order reduction. *IEEE Trans. on CAD*, 1995. Submitted.
- [5] Kevin J. Kerns, Ivan L. Wemple, and Andrew T. Yang. Stable and efficient reduction of substrate model networks using congruence transforms. In *IEEE/ACM International Conference on Computer Aided Design*, pages 207 - 214, San Jose, CA, November 1995.
- [6] A. J. Laub and A. Linnemann. Hessenberg and Hessenberg/triangular form in linear system theory. *International Journal of Control*, 44:1523-1547, 1986.
- [7] R. A. Horn and C. R. Johnson. *Matrix Theory*. Cambridge University Press, 1985.
- [8] R. A. Horn and C. R. Johnson. *Topics in Matrix Theory*. Cambridge University Press, 1991.
- [9] W. B. Gragg. Matrix interpretations and applications of the continued fraction algorithm. *Rocky Mountain Journal of Mathematics*, 4:213 - 225, 1974.