



ELSEVIER

Available at

www.ElsevierMathematics.com

POWERED BY SCIENCE @ DIRECT®

JOURNAL OF
COMPUTATIONAL AND
APPLIED MATHEMATICS

Journal of Computational and Applied Mathematics 162 (2004) 213–229

www.elsevier.com/locate/cam

Sylvester equations and projection-based model reduction^{☆,☆☆}

K. Gallivan^a, A. Vandendorpe^b, P. Van Dooren^{b,*}

^aFlorida State University, FL, USA

^bUniversité catholique de Louvain, CESAME, Avenue G. Lemaître 4, 1348 Louvain-la-Neuve, Belgium

Received 10 December 2001; received in revised form 15 December 2002

Abstract

In this paper, we establish a connection between Krylov subspace techniques for Multipoint Padé interpolation, and the use of Sylvester equations for constructing reduced-order models. We also briefly point out that this connection partly extends to ADI-type techniques and to the Smith iteration for computing approximate solutions of Lyapunov equations.

© 2003 Elsevier B.V. All rights reserved.

Keywords: Model reduction; Time-invariant linear systems; large-scale system; Rational interpolation; Multipoint Padé

1. Introduction

Every linear time-invariant continuous-time system can be represented by a generalized state-space model of the type

$$\begin{aligned} E\dot{x} &= Ax + Bu, \\ y &= Cx + Du, \end{aligned} \tag{1}$$

[☆] A research fellowship from the Belgian National Fund for Scientific Research is gratefully acknowledged by the second author. This paper presents research supported by the Belgian Programme on Inter-university Poles of Attraction, initiated by the Belgian State, Prime Minister's Office for Science, Technology and Culture. This work was also supported by the National Science Foundation under Grant No. CCR-9912415.

^{☆☆} Although most of the theory presented in this paper holds for both continuous-time and discrete-time systems, we will restrict ourselves for the largest part of the paper to continuous-time systems. We will point out how to extend these results to discrete-time systems in the last section.

* Corresponding author.

E-mail address: vdooren@csam.ucl.ac.be (P. Van Dooren).

where the input $u(t) \in \mathbb{R}^m$, the state variable of the system $x(t) \in \mathbb{R}^N$ and the output $y(t) \in \mathbb{R}^p$ (where typically $m, p \ll N$). Without loss of generality, we can assume that the system is minimal, i.e., controllable and observable, since otherwise we can always find a smaller dimensional model that is controllable and observable, and that has exactly the same transfer function. In addition to this, we will assume that the system is asymptotically stable, i.e., the generalized eigenvalues of the pencil $sE - A$ lie in the open left-half plane (this also implies that E is nonsingular).

Suppose that the system order N is too big for solving various control problems within a reasonable amount of computing time. The goal of model reduction is to construct a reduced-order system

$$\begin{aligned} \hat{E}\dot{\hat{x}} &= \hat{A}\hat{x} + \hat{B}u, \\ \hat{y} &= \hat{C}\hat{x} + \hat{D}u, \end{aligned} \tag{2}$$

where, $u(t) \in \mathbb{R}^m$, $\hat{x}(t) \in \mathbb{R}^n$ and $\hat{y}(t) \in \mathbb{R}^p$ with $n \ll N$, that satisfies the constraints on computing time. For the same reasons as above, we will assume that the reduced-order model is minimal. Minimality of the reduced-order system can be achieved for most problems although it may require some postprocessing.

The basic idea of the projection-based model reduction is to project the system’s state-space of dimension N onto a space of lower dimension, $n \ll N$, in such a way that the behavior of the reduced-order model is sufficiently close to the original model, i.e., given the input $u(t)$ to both systems, we want $\hat{y}(t)$ to be close to $y(t)$. This implies that the reduced-order system must be stable, otherwise the system responses cannot be close to each other.

This approximation in the time domain is equivalent to imposing conditions on the frequency responses of both systems [12]: we want a reduced-order model such that the transfer functions of the models

$$T(s) \doteq C(sE - A)^{-1}B + D \quad \text{and} \quad \hat{T}(s) \doteq \hat{C}(s\hat{E} - \hat{A})^{-1}\hat{B} + \hat{D}$$

are “close” to each other, i.e., such that the error $\|T(\cdot) - \hat{T}(\cdot)\|$ is small in the H_∞ norm.

The reduced-order models we will consider in this paper are built using projections. In order to construct an n th-order reduced system, we project the matrices of the original system using $N \times n$ matrices Z and V as follows:

$$\hat{A} = Z^T A V, \quad \hat{E} = Z^T E V, \quad \hat{B} = Z^T B, \quad \hat{C} = C V, \quad \hat{D} = D. \tag{3}$$

Only the spaces onto which we are projecting influence the relationship between the transfer functions, i.e., the particular projection matrices are not important in this regard (although the choices may have computational significance). This is seen in the following lemma.

Lemma 1.1. *The projected transfer function $\hat{T}(s)$ defined in (3) is unchanged if we replace V and Z by other matrices $\tilde{V} \doteq VR$ and $\tilde{Z} \doteq ZL$ which span the same respective spaces, i.e., where R and L are invertible.*

Proof. The factors L, R drop out in the transfer function $\hat{T}(s) = \hat{C}(s\hat{E} - \hat{A})^{-1}\hat{B} + \hat{D}$. \square

A geometric interpretation of this lemma is that the spaces $\mathcal{X} \doteq \text{Im}(Z)$ and $\mathcal{V} \doteq \text{Im}(V)$ can be viewed as, respectively, left and right, *projection spaces*. Notice also that the D matrix is maintained in the approximation. Since it does not play a role in our analysis, we will assume $D = \hat{D} = 0$ below.

In the remainder of this paper, we explore the connections between several methods of model reduction and unify them via a projection-based point of view. We consider ADI-type techniques, Krylov subspace techniques for multipoint Padé interpolation, and projected Lyapunov techniques. We show a connection between these methods and the general framework of rational interpolation.

In Section 2, we present a general theorem linking the spaces \mathcal{L} and \mathcal{V} with interpolation. In Section 3, we analyze a particular step of an algorithm recently proposed by Antoulas and Sorensen ([1,11]) that uses Sylvester equations to construct projectors that define the reduced-order model. We show a connection between this approach and the *moment matching* approach. In Section 4, we generalize the results of the earlier sections and discuss some particular important cases. In Section 5, we conclude by presenting links between several existing techniques for model reduction of large-scale dynamical systems. Finally, we discuss how to extend these results to more general cases.

2. Multipoint Padé interpolation

We first review some results about *moment matching* methods. The objective of these methods is to construct a reduced-order model such that $\hat{T}(s)$ interpolates the frequency response of $T(s)$ and its derivatives at multiple points $\sigma_1, \sigma_2, \dots, \sigma_K$. We suppose that these points are not poles of the original model.

Let us consider an expansion of $T(s)$ about a point σ that is *not* a pole of $T(s)$. It then follows that $\sigma E - A$ is invertible and one obtains the following formal series expansion:

$$T(s) = C(\sigma E - A - (\sigma - s)E)^{-1}B \tag{4}$$

$$= C(I - (\sigma E - A)^{-1}E(\sigma - s))^{-1}(\sigma E - A)^{-1}B \tag{5}$$

$$= \sum_{j=0}^{+\infty} C((\sigma E - A)^{-1}E)^j(\sigma E - A)^{-1}B \cdot (\sigma - s)^j \tag{6}$$

$$\doteq \sum_{j=0}^{+\infty} T_\sigma^{(j)} \cdot (\sigma - s)^j, \tag{7}$$

which defines the so-called *moments*

$$T_\sigma^{(j)} \doteq C((\sigma E - A)^{-1}E)^j(\sigma E - A)^{-1}B \tag{8}$$

about an expansion point σ . These moments exist for every σ for which $(\sigma E - A)$ is nonsingular. The same expansion can be made of a lower-order model $\hat{T}(s)$ about the same point σ provided $(\sigma \hat{E} - \hat{A})$ is invertible. This can be expressed in an alternative and sometimes more convenient form.

Definition 2.1. We say that the reduced-order system $\hat{T}(s)$ interpolates the original system $T(s)$ at points $\sigma_k, 1 \leq k \leq K$, with an order $\alpha_k \in \mathbb{N}$, when

$$T(s) - \hat{T}(s) = O(s - \sigma_k)^{\alpha_k} \quad \forall 1 \leq k \leq K, \tag{9}$$

and that it interpolates $T(s)$ at ∞ with an order α , when

$$T(s) - \hat{T}(s) = O(s^{-\alpha}). \tag{10}$$

From the preceding definition, any two strictly proper transfer function $T(s)$ and $\hat{T}(s)$ interpolate each other at infinity.

We will show in this section that such expansions must match for a certain number of moments when the lower-order model is constructed via a projection on particular Krylov spaces. For matrices $G \in \mathbb{C}^{N \times N}$, $H \in \mathbb{C}^{N \times m}$, we define the Krylov space of index j as follows:

$$\mathcal{K}_j(G, H) = \text{Im}\{H, GH, G^2H, \dots, G^{j-1}H\}. \tag{11}$$

The following lemma provides a useful property of Krylov subspaces.

Lemma 2.1. *Let $V \in \mathbb{C}^{N \times n}$ be a full-rank matrix such that*

$$\mathcal{K}_j(G, H) \subseteq \mathcal{V} \doteq \text{Im}(V)$$

and let W be an arbitrary $N \times n$ matrix such that $W^T V = I_n$. Then the projected matrices

$$\hat{G} \doteq W^T G V, \quad \hat{H} \doteq W^T H$$

satisfy the equalities

$$G^i H = V \hat{G}^i \hat{H}, \quad i = 0, \dots, j - 1.$$

Proof. Since V is full rank, there exists a matrix W such that $W^T V = I_n$. Since the image of each $G^i H$ is spanned by the columns of V there exist for all $i = 0, \dots, j - 1$ a matrix Y_i such that

$$G^i H = V Y_i \quad \text{and hence} \quad Y_i = W^T G^i H.$$

The proof proceeds by induction. For $i = 0$ clearly $Y_0 = W^T H = \hat{H}$. If $Y_i = \hat{G}^i \hat{H}$ then it follows that $Y_{i+1} = W^T G G^i H = W^T G V \hat{G}^i \hat{H} = \hat{G}^{i+1} \hat{H}$, which proves the result. \square

The following theorem is shown in [5,7] for a single-input–single-output (SISO) system, and is extended here to the multiple-input–multiple-output (MIMO) case.

Theorem 2.1. *If*

$$\bigcup_{k=1}^K \mathcal{K}_{J_{b_k}}((\sigma_k E - A)^{-1} E, (\sigma_k E - A)^{-1} B) \subseteq \mathcal{V} = \text{Im}(V) \tag{12}$$

and

$$\bigcup_{k=1}^K \mathcal{K}_{J_{c_k}}((\sigma_k E - A)^{-T} E^T, (\sigma_k E - A)^{-T} C^T) \subseteq \mathcal{Z} = \text{Im}(Z), \tag{13}$$

where the interpolation points σ_k are chosen such that the matrices $\sigma_k E - A$ are invertible $\forall k \in \{1, \dots, K\}$ then the moments of systems (1) and (2) at the points σ_k satisfy

$$T_{\sigma_k}^{(j_k)} = \hat{T}_{\sigma_k}^{(j_k)} \tag{14}$$

for $j_k = 1, 2, \dots, J_{b_k} + J_{c_k}$ and $k = 1, 2, \dots, K$, provided these moments exist, i.e., provided the matrices $\sigma_k \hat{E} - \hat{A}$ are invertible.

Proof. We prove the result for an arbitrary expansion point σ_k . Since

$$\mathcal{H}_{J_{b_k}}((\sigma_k E - A)^{-1}E, (\sigma_k E - A)^{-1}B) \subseteq \text{Im}(V), \tag{15}$$

it follows from Lemma 2.1 that for $0 \leq j \leq J_{b_k} - 1$,

$$V((\sigma_k \hat{E} - \hat{A})^{-1} \hat{E})^j (\sigma_k \hat{E} - \hat{A})^{-1} \hat{B} = ((\sigma_k E - A)^{-1} E)^j (\sigma_k E - A)^{-1} B. \tag{16}$$

In order to apply this lemma, it suffices to choose

$$W^T = (\sigma_k \hat{E} - \hat{A})^{-1} Z^T (\sigma_k E - A), \quad G = (\sigma_k E - A)^{-1} E, \quad H = (\sigma_k E - A)^{-1} B, \tag{17}$$

from which it follows that

$$\hat{G} = (\sigma_k \hat{E} - \hat{A})^{-1} \hat{E}, \quad \hat{H} = (\sigma_k \hat{E} - \hat{A})^{-1} \hat{B}.$$

Similarly, since

$$\mathcal{H}_{J_{c_k}}((\sigma_k E - A)^{-T} E^T, (\sigma_k E - A)^{-T} C^T) \subseteq \text{Im}(Z), \tag{18}$$

it follows from Lemma 2.1 that for $0 \leq j \leq J_{c_k} - 1$,

$$\hat{C}(\sigma_k \hat{E} - \hat{A})^{-1} (\hat{E}(\sigma_k \hat{E} - \hat{A})^{-1})^j Z^T = C(\sigma_k E - A)^{-1} (E(\sigma_k E - A)^{-1})^j. \tag{19}$$

Multiplying E from the left with (19) and from the right with (16), it then follows that the first $J_{b_k} + J_{c_k}$ moments of $T(s)$ and $\hat{T}(s)$ at σ_k are equal. \square

Remark 2.1. In the proof of Theorem 2.1, the conditions requiring the nonsingularity of the matrices $\sigma_k E - A$ and $\sigma_k \hat{E} - \hat{A}$ at all interpolation points are crucial to the theory but not always guaranteed. The following observations govern the possible scenarios:

- (1) All the matrices $\sigma_k E - A$ of the original model *must* be invertible. If this is not the case, the Krylov subspaces defined in Theorem 2.1 are no longer well defined. Fortunately, if the original system is stable and the interpolation points are restricted to the imaginary axis and the right-half of the complex plane this is easily guaranteed.
- (2) All the matrices $\sigma_k \hat{E} - \hat{A}$ of the reduced-order model must be invertible, otherwise W defined in (17) would not exist and Theorem 2.1 cannot be used to guarantee moment matching. This, of course, does not mean that moment matching may not occur.
- (3) If all the matrices $\sigma_k E - A$ are invertible, the Krylov subspaces that must be contained in the images of the projectors V and Z are then well defined and the singularity of some or all $\sigma_k \hat{E} - \hat{A} = Z^T(\sigma_k E - A)V$ may have different causes:
 - (a) Either Z or V is not full rank. In such a situation, \hat{E} and \hat{A} are singular, and the pencil (\hat{E}, \hat{A}) is singular. To avoid these difficulties, we construct bases Z_{new} for Z and V_{new} for V ; at least one of which must have a rank lower than, respectively, Z or V . If the dimension of $\text{Im}(Z_{\text{new}})$ is equal to the dimension of $\text{Im}(V_{\text{new}})$ and if the reduced-order system constructed via these new projectors is nonsingular at the interpolation points, it satisfies the moments matching properties of Theorem 2.1. If the dimension of $\text{Im}(Z_{\text{new}})$ is not equal to the dimension of $\text{Im}(V_{\text{new}})$, we can augment the projector of smaller dimension by simply choosing random vectors until the two subspaces spanned by Z and V are of

equal dimension. In either case, minimality of the reduced-order system is guaranteed by postprocessing. In the SISO case, the fact that Z or V is not full rank can only occur when $T(s)$ is nonminimal.

- (b) It is possible that Z , V and the $\sigma_k E - A$ are all full-rank matrices and that their product is singular. This case requires more investigation and is not discussed in this paper.

3. Sylvester equations and moment matching

The system's observability gramian P and controllability gramian Q satisfy the equations

$$APE^T + EPA^T + BB^T = 0, \quad (20)$$

$$A^TQE + E^TQA + C^TC = 0. \quad (21)$$

These gramians are very important since they describe how the energy of the system is “distributed” over the coordinates of the state-space. There exist, e.g., coordinate transformations X and Y such that

$$sI_N - \tilde{A} \doteq Y^T(sE - A)X, \quad \tilde{B} \doteq Y^TB, \quad \tilde{C} \doteq CX,$$

where the transformed state-space system $\{\tilde{A}, \tilde{B}, \tilde{C}\}$ is “balanced”. This means that the transformed Lyapunov equations

$$\tilde{A}\tilde{P} + \tilde{P}\tilde{A}^T + \tilde{B}\tilde{B}^T = 0, \quad (22)$$

$$\tilde{A}^T\tilde{Q} + \tilde{Q}\tilde{A} + \tilde{C}^T\tilde{C} = 0, \quad (23)$$

yield transformed gramians that are equal and diagonal

$$\tilde{P} = \tilde{Q} = \Sigma.$$

When the diagonal matrix

$$\Sigma \doteq \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix}$$

has n “large” elements in Σ_1 and $N - n$ “small” elements in Σ_2 , one can dismiss the small ones and keep only the n dominant ones. It is shown in [11] that this amounts to choosing the projection spaces $\mathcal{V} = \text{Im}(X_1)$ and $\mathcal{Z} = \text{Im}(Y_1)$, where X_1 and Y_1 are the first n columns of the transformation matrices X and Y , respectively. Multiplying (20) with Y^T from the left and with Y from the right and multiplying (21) with X^T from the left and X from the right yields

$$P = X\Sigma X^T \approx X_1\Sigma_1 X_1^T, \quad Q = Y\Sigma Y^T \approx Y_1\Sigma_1 Y_1^T,$$

which indicates that X_1 and Y_1 contain the dominant parts of P and Q .

In the well-known balanced truncation model reduction technique [9], such dominant spaces $\text{Im}(X_1)$ and $\text{Im}(Y_1)$ play an important role. Unfortunately, for large-scale models (when N is very large), solving (20) and (21) and finding their dominant eigenspaces requires too many computations. Hence, we must approximate the dominant subspaces via another scheme. In [11], this idea is used to justify

the following projected Lyapunov approach related to subspace iteration, which is used as part of the algorithm’s basic step.

Suppose one has orthonormal bases $V_{in} \in \mathbb{C}^{N \times k}$ and $Z_{in} \in \mathbb{C}^{N \times k}$ of initial estimates of \mathcal{V} and \mathcal{Z} . In subspace iteration one updates these bases via $V_{up} = PV_{in}$ and $Z_{up} = QZ_{in}$. Consider the projectors $V_{in}V_{in}^T$ and $\tilde{V}_{in}\tilde{V}_{in}^T = I_N - V_{in}V_{in}^T$, and multiply (20) with a matrix $R \in \mathbb{C}^{N \times n}$ on the right. This yields the following equation:

$$AV_{up}F_r^T + EV_{up}G_r^T + BH_r^T = -AP\tilde{V}_{in}\tilde{V}_{in}^TE^TR - EP\tilde{V}_{in}\tilde{V}_{in}^TA^TR \tag{24}$$

where $F_r = R^TEV_{in}$, $G_r = R^TAV_{in}$ and $H_r = R^TB$. If V_{in} is a good approximation of the dominant n -dimensional space X_1 of P , then $P\tilde{V}_{in}$ is small and the right-hand side will be small compared to the terms of the left-hand side.

The basic idea of the algorithm is to construct a basis for \mathcal{V} by solving the Sylvester equation that results from setting the right-hand side of the projected Lyapunov to 0

$$AV_{up}F_r^T + EV_{up}G_r^T + BH_r^T = 0 \tag{25}$$

rather than (24). A similar argument can be used for \mathcal{Z} by considering the projector $Z_{in}Z_{in}^T$ and by multiplying (21) with a matrix L to construct a basis for \mathcal{Z} . The corresponding projected equation is then

$$A^TZ_{up}F_1 + E^TZ_{up}G_1 + C^TH_1 = 0, \tag{26}$$

where $F_1 = Z_{in}^TE L$, $G_1 = Z_{in}^T A L$ and $H_1 = C L$.

These two equations are the basis of the iterative updating procedure for V and Z , described in [11]. We show below that this is, in fact, equivalent to multipoint interpolation according to our earlier definitions.

The N th-order transfer function from which we start is $T(s)$, and we want to approximate it with $\hat{T}(s)$ based on projection spaces \mathcal{V} and \mathcal{Z} as defined in the introduction. We choose the matrix V by solving the following Sylvester equation—dropping the indices to simplify the notation

$$AVF^T + EVG^T + BH^T = 0, \tag{27}$$

where $A, E \in \mathbb{C}^{N \times N}$, $B \in \mathbb{C}^{N \times m}$, $F, G \in \mathbb{C}^{n \times n}$ and $H \in \mathbb{C}^{n \times m}$. We call (27) a Projected Lyapunov equation because the matrices F, G and H may be seen as projections of E, A and B . Indeed, if we replace F, G and H by E, A and B , we recover (20). All the results obtained for the equation in V can be transposed to the corresponding equation for Z , so we focus on only one of the equations.

In this paper, we suppose that n is an integer multiple of m , i.e., $n = rm$. We further assume that the spectra of the generalized eigenvalue problems $sE - A$ and $sF - G$ are well defined and such that

$$\lambda_i(E, A) + \lambda_j(F, G) \neq 0,$$

in order to ensure the existence and uniqueness of V . We then show that under certain conditions, any reduced-order model $\hat{T}(s)$ built with $\mathcal{V} \supseteq \text{Im}(V)$ interpolates the transfer function $T(s)$ at the eigenvalues of $sF - G$ up to a particular order.

We first suppose for simplicity that F is invertible and that $sF - G$ has only one eigenvalue σ with m Jordan blocks of length r (the general case will be developed in later sections). Under these

conditions, there exist invertible matrices S and T such that

$$S(G - sF)T = sI - G_J, \tag{28}$$

where G_J has the particular form

$$G_J = \begin{bmatrix} \sigma I_m & -I_m & & & \\ & \ddots & \ddots & & \\ & & \ddots & -I_m & \\ & & & \ddots & \sigma I_m \end{bmatrix}, \tag{29}$$

which is a “modified” Jordan canonical form that will simplify further calculations. Accordingly, we define $H_J \doteq SH$ where H_J is partitioned in r square blocks $H_i \in \mathbb{C}^{m \times m}$

$$H_J = \begin{bmatrix} H_1 \\ \vdots \\ H_r \end{bmatrix}. \tag{30}$$

Multiplying (27) by S^T on the right and defining $\bar{V} = VT^{-T}$, we obtain

$$A\bar{V} + E\bar{V}G_J^T - BH_J^T = 0. \tag{31}$$

Notice that $\text{Im}(\bar{V}) = \text{Im}(V)$ since T is invertible. This equation yields

$$A[\bar{V}_1, \dots, \bar{V}_r] + E[\bar{V}_1, \dots, \bar{V}_r] \begin{bmatrix} \sigma I \\ -I & \ddots \\ & \ddots & \sigma I \\ & & -I & \sigma I \end{bmatrix} = B[H_1^T, \dots, H_r^T], \tag{32}$$

where we partitioned \bar{V} in r blocks $\bar{V}_i \in \mathbb{C}^{N \times m}$. This last equation allows us to write

$$[\bar{V}_1, \dots, \bar{V}_r] = [(A + \sigma E)^{-1}B, \dots, ((A + \sigma E)^{-1}E)^{r-1}(A + \sigma E)^{-1}B] \begin{bmatrix} H_1^T & \dots & H_r^T \\ \vdots & & \\ H_r^T & & \ddots \end{bmatrix}. \tag{33}$$

H_r is invertible if and only if (G_J, H_J) is controllable and it follows that

$$\text{Im}(V) = \mathcal{H}_r((A + \sigma E)^{-1}E, (A + \sigma E)^{-1}B). \tag{34}$$

We therefore can show the following.

Theorem 3.1. *The reduced-order system $\hat{T}(s) = \hat{C}(s\hat{E} - \hat{A})^{-1}\hat{B}$ obtained using a matrix $Z \in \mathbb{C}^{N \times n}$ on the left and on the right the matrix V that contains the image of the solution of the Sylvester equation (27) is such that*

$$T(s) - \hat{T}(s) = O(s + \sigma)^r, \tag{35}$$

where the pencil $sF - G$ has m Jordan blocks of length r at its single eigenvalue σ , provided $\sigma E - A$ and $\sigma \hat{E} - \hat{A}$ are invertible and provided $[sF - G|H]$ has full rank for all s , i.e., the pair $(sF - G, H)$ is controllable.

Proof. Note the only condition on Z is that the reduced system satisfy the invertibility and controllability constraints. In particular, the invertibility constraints imply that Z must be full rank (cf. Remark 2.1). Since the controllability of $(sF - G, H)$ is equivalent to the controllability of the transformed pair (G_J, H_J) , this result follows directly from the discussion above and Theorem 2.1. \square

We recall that the interpolation point is the opposite of the eigenvalue of (F, G) , i.e.,

$$\sigma = -\lambda(F, G). \tag{36}$$

Remark 3.1. Notice that controllability of the pair $(sF - G, H)$ implies that the pencil $(sF - G)$ can have at most m Jordan blocks (of any length) at a given eigenvalue σ . The assumption to have exactly m Jordan blocks at σ is in a sense the extremal case. But the fact that these blocks have all equal length is clearly a severe constraint.

4. Extensions and particular cases

In this section, we first derive two generalizations of Theorem 3.1 which extend its applicability. We then consider some particular cases that admit stronger results.

4.1. The case of infinite eigenvalues

We consider here the case where $sF - G$ has a single infinite eigenvalue rather than a finite one. If we assume the pencil $(sF - G)$ is regular, i.e., that $\det(sF - G)$ is not identically 0 in the complex plane, the Jordan canonical form is then replaced by the Weierstrass canonical form. This form expresses that any regular pencil can be transformed by invertible matrices S and T to the block diagonal pencil

$$S(G - sF)T = \left[\begin{array}{c|c} sN_J - I & \\ \hline & sI - G_J \end{array} \right], \tag{37}$$

where N_J is a nilpotent matrix and where both N_J and G_J are in Jordan canonical form. The finite eigenvalues of the pencil are then the eigenvalues of G_J and the infinite eigenvalue are reflected by the nilpotency of N_J [6]. We assume again that there is only one eigenvalue, namely ∞ , and that it

has exactly m Jordan blocks of length r . The matrix block $sI - G_J$ then does not exist and N_J can then be assumed in the Jordan-like form

$$N_J = \begin{bmatrix} 0_m & & & & \\ -I_m & \ddots & & & \\ & \ddots & \ddots & & \\ & & & -I_m & 0_m \end{bmatrix}. \tag{38}$$

We again transform (27) with the transformations S and T to the particular form

$$A\bar{V}N_J^T + E\bar{V} = BH_J^T, \tag{39}$$

where $\text{Im}(\bar{V}) = \text{Im}(V)$, N_J has form (37) and $H_J \doteq SH$.

Again, we partition \bar{V} in k blocks $\bar{V}_i \in \mathbb{C}^{N \times m}$ and H_J in k blocks $H_i \in \mathbb{C}^{m \times m}$. We then find recursively

$$\bar{V}_1 = E^{-1}BH_1^T \tag{40}$$

$$\bar{V}_i = E^{-1}(BH_i^T + A\bar{V}_{i-1}) \quad \forall 2 \leq i \leq k. \tag{41}$$

For (N_J, H_J) controllable—or equivalently for $(sF - G, H)$ controllable at infinity—it then follows that

$$\mathcal{H}_k(E^{-1}A, E^{-1}B) = \text{Im}(V). \tag{42}$$

A similar theorem to Theorem 3.1 is then easily obtained, which can be interpreted as *moment matching* at infinity.

Theorem 4.1. *The reduced-order system $\hat{T}(s) = \hat{C}(s\hat{E} - \hat{A})^{-1}\hat{B}$ obtained using a matrix $Z \in \mathbb{C}^{N \times n}$ on the left and on the right a matrix V that contains the image of the solution of the Sylvester equation (27) is such that*

$$T(s) - \hat{T}(s) = O(s^{-r-1}), \tag{43}$$

where the pencil $sF - G$ has m Jordan blocks of length r at ∞ , provided E and \hat{E} are invertible and provided $[sF - G|H]$ has full rank for all s , i.e., the pair $(sF - G, H)$ is controllable.

Proof. Note the only condition on Z is that the reduced system satisfy the invertibility and controllability constraints. Since the controllability of $(sF - G, H)$ is equivalent to the controllability of the transformed pair (G_J, H_J) , this result follows directly from the above discussion and Theorem 2.1. \square

4.2. The multiple eigenvalue case

The case where the pencil $sF - G$ has several eigenvalues can also be analyzed. For simplicity of notation we will assume they are all finite. Suppose that $(sI - G_J) \doteq S(G - sF)T$ is in Jordan

form, with

$$G_J = \begin{bmatrix} G_{J_1} & & \\ & \ddots & \\ & & G_{J_K} \end{bmatrix}, \tag{44}$$

in which each diagonal block $G_{J_k} \in \mathbb{C}^{n_k \times n_k}$ has exactly m Jordan blocks of length b_k (i.e., $n_k = b_k m$) at eigenvalue σ_k (each block G_{J_k} is then of type (29) but for a different eigenvalue σ_k). We then partition $H_J \doteq SH$ conformably with blocks $H_k \in \mathbb{C}^{n_k \times m}$:

$$H_J = \begin{bmatrix} H_1 \\ \vdots \\ H_K \end{bmatrix}. \tag{45}$$

Eq. (31) is then equivalent to the K equations

$$AV_k + EV_k G_{J_k}^T - BH_k^T = 0, \quad 1 \leq k \leq K, \tag{46}$$

where each V_k and H_k can be partitioned into subblocks with, respectively, m columns and m rows

$$V_k = [V_{k,1} \dots V_{k,b_k}], \tag{47}$$

$$H_k = \begin{bmatrix} H_{k,1} \\ \vdots \\ H_{k,b_k} \end{bmatrix}, \tag{48}$$

i.e., $H_{k,l} \in \mathbb{C}^{m \times m}$ and $V_{k,l} \in \mathbb{C}^{N \times m}$. We, therefore, have the following.

Theorem 4.2. *The reduced-order system $\hat{T}(s) = \hat{C}(s\hat{E} - \hat{A})^{-1}\hat{B}$ obtained by using a matrix $Z \in \mathbb{C}^{N \times n}$ on the left and on the right the matrix V that contains the image of the solution of Eq. (31) where G_J is defined in (44) is such that*

$$T(s) - \hat{T}(s) = O(s + \sigma_k)^{b_k} \quad \forall 1 \leq k \leq K, \tag{49}$$

where the pencil $sF - G$ has m Jordan blocks of length b_k at the eigenvalues σ_k , provided all matrices $\sigma_k E - A$ and $\sigma_k \hat{E} - \hat{A}$ are invertible and provided $[sF - G|H]$ has full rank for all s , i.e., the pair $(sF - G, H)$ is controllable.

Proof. Note the only condition on Z is that the reduced system satisfy the invertibility and controllability constraints. Applying Theorem 2.1 to each V_k defined in (47) proves the result. \square

Notice also that the case of infinite eigenvalues could have been included here without difficulty.

4.3. The two-sided case

Until now, we have only considered the image of the right projection matrix V that satisfies

$$AVF_r^T + EVG_r^T + BH_r^T = 0. \tag{50}$$

By duality, we can find a similar result for the left projection matrix Z . For this we need to construct Z such that

$$A^T Z F_1 + E^T Z G_1 + C^T H_1 = 0, \quad (51)$$

where 1 is added to F, G and H to recall that the new matrices F_1, G_1 and H_1 can be chosen independently from F_r, G_r and H_r that define V , and use the same notation as before.

By transposing all the results of Section 3 and the preceding subsection, we can state the following general theorem.

Theorem 4.3. *The reduced-order system $\hat{T}(s) = \hat{C}(s\hat{E} - \hat{A})^{-1}\hat{B}$, obtained by using the projectors Z and V such that Z has the same image as the solution of (51) and V has the same image as the solution of (50), is such that*

$$T(s) - \hat{T}(s) = O(s + \sigma_k)^{b_k + c_k} \quad \forall 1 \leq k \leq K, \quad (52)$$

where σ_k is a generalized eigenvalue of (F_r, G_r) and (F_1, G_1) , associated with the diagonal block $G_{r_{j_k}}$ that has exactly m Jordan blocks of length b_k at eigenvalue σ_k and the diagonal block $G_{1_{j_k}}$ that has p Jordan blocks of length c_k at σ_k , provided $sE - A$ and $s\hat{E} - \hat{A}$ are invertible in σ_k , the pair $(sF_r - G_r, H_r)$ is controllable, and the pair $(sF_1 - G_1, H_1)$ is observable.

Proof. This is again a direct consequence of Theorem 2.1. \square

Remark 4.1.

- (1) The spectra of (F_1, G_1) and (F_r, G_r) are in general different. Suppose for instance that σ_k is a generalized eigenvalue of (F_1, G_1) that is not a generalized eigenvalue of (F_r, G_r) , then we define the corresponding matrices $H_{r_k}, G_{r_{j_k}}$ to be of dimension zero. Hence, $c_k = 0$.
- (2) One can easily also include infinite eigenvalues as in Theorem 4.1.

4.4. The SISO case

Until now, we needed to suppose the pencils (F_1, G_1) and (F_r, G_r) to have exactly, respectively, m and p Jordan blocks for every interpolation point, as it has already been pointed out in Remark 3.1. In a SISO system, $m = p = 1$. In such a case, the above condition is not extremal anymore. In a practical point of view, it means that, when a SISO reduced-order system of McMillan degree n is built with the projected Lyapunov technique explained in this paper, it generically interpolates the original system at $2n$ different points. In general, this is not true for MIMO systems.

4.5. Convergence in one step

In this section, we indicate another explanation as to why this scheme can be viewed as subspace iteration. Let us suppose that the gramian P has rank n , then for almost all initial matrices V_{in} an exact subspace iteration step $V_{up} \doteq P V_{in}$ should yield $\text{Im}(V_{up}) = \text{Im}(P)$. We verify that the Sylvester equation indeed yields such a matrix V in the SISO case with E and F invertible. Without loss of

generality we can then suppose $E = I_N$ and $F = I_n$. Recall that the controllability matrix \mathcal{C} and the observability matrix \mathcal{O} are defined by

$$\mathcal{C}(A, B) = [B \ AB \ \dots \ A^{N-1}B], \tag{53}$$

$$\mathcal{O}^T(A, C) = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{N-1} \end{bmatrix}. \tag{54}$$

Lemma 4.1. *Suppose the matrices $A \in \mathbb{C}^{N \times N}$, $B \in \mathbb{C}^{N \times m}$, $G \in \mathbb{C}^{n \times n}$, and $H \in \mathbb{C}^{n \times m}$ satisfy*

$$AV + VG^T + BH^T = 0 \tag{55}$$

and we have that $\text{rank}(\mathcal{C}(A, B)) = n$, (G, H) is controllable, and $\Lambda(-G) \cap \Lambda(A) = \emptyset$, then

$$\text{Im}(V) = \text{Im}(\mathcal{C}(A, B)). \tag{56}$$

Proof. Since $\text{Im}(P)$, defined in (20), equals $\text{Im}(\mathcal{C}(A, B))$, one can always choose a coordinate system in which the standard Lyapunov equation has the form

$$\begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix} \begin{bmatrix} P_{11} & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} P_{11} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} A_{11}^T & 0 \\ A_{12}^T & A_{22}^T \end{bmatrix} + \begin{bmatrix} B_1 \\ 0 \end{bmatrix} [B_1^T \ 0] = 0, \tag{57}$$

where $P_{11} \in \mathbb{C}^{n \times n}$ is invertible, and (A_{11}, B_1) is controllable. In this coordinate system $B_2 = A_{21} = 0$ and

$$\text{Im}(P) = \text{Im} \left(\begin{bmatrix} I_n \\ 0 \end{bmatrix} \right). \tag{58}$$

Rewriting (55) in the same coordinate system yields

$$\begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} + \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} G^T + \begin{bmatrix} B_1 \\ 0 \end{bmatrix} H^T = 0. \tag{59}$$

Therefore, V_2 must satisfy

$$A_{22}V_2 + V_2G^T = 0 \tag{60}$$

and since the spectra of A_{22} and $-G$ are disjoint, $V_2 = 0$. Moreover, V_1 satisfies

$$A_{11}V_1 + V_1G^T + B_1H^T = 0. \tag{61}$$

If $V_1 \in \mathbb{C}^{n \times n}$ is nonsingular, then the matrices V and $\begin{bmatrix} I_n \\ 0 \end{bmatrix}$ have the same image and the proof would be complete. Suppose—by contradiction—that the kernel of V_1 has dimension $k > 0$. Then,

there exist matrices $X \in \mathbb{C}^{n \times k}$ and $Y \in \mathbb{C}^{k \times n}$ of rank k such that $YV_1 = 0$ and $V_1X = 0$. Multiply (61) by Y on the left and X on the right to obtain

$$YB_1H^T X = 0. \quad (62)$$

Since YB_1 and $H^T X$ are two vectors, there are now two possibilities
(1)

$$YB_1 = 0. \quad (63)$$

Then $YA_{11}V_1 = 0$ and

$$YA_{11} = TY, \quad (64)$$

because $\text{Im}(Y^T) = \ker(V_1^T)$. It follows that

$$Y[B_1 A_{11}B_1 \dots A_{11}^{n-1}B_1] = 0 \quad (65)$$

and the system (A_{11}, B_1) is not controllable. This fact contradicts our hypotheses.
(2)

$$H^T X = 0. \quad (66)$$

Then $V_1G^T X = 0$ and

$$X^T G = SX^T,$$

because $\text{Im}(X) = \ker(V_1)$. It follows that

$$X^T[H \ GH \dots G^{n-1}H] = 0 \quad (67)$$

and the system (G, H) is not controllable. This contradicts our hypotheses as well. \square

By duality, a similar result also holds for Z . It is possible to show that such a result also holds when E and F are not the identity or even when F is singular.

In practical circumstances, one often observes a fast decay of the singular values of P and Q . Recent theoretical results tend to justify these observations [3]. This means that observability and controllability gramians are often close to low-rank matrices. Lemma 4.1 shows that if the original model is not controllable and we construct a projector with dimension equal to the dimension of the controllability space, then the reduced-order model built with such a right projector is exact. By continuity, one may hope that it is not too far from the solution when gramians are close to low rank.

5. Concluding remarks

5.1. The discrete-time case

All the ideas developed for continuous-time systems can be extended to discrete-time systems. In the discrete-time case, the observability gramian defined in (20) and the controllability gramian

defined in (21) satisfy the equations

$$APA^T - EPE^T + BB^T = 0, \tag{68}$$

$$A^TQA - E^TQE + C^TC = 0. \tag{69}$$

By analogy with the continuous-time case, we construct matrices V and Z from the following Sylvester equations:

$$AVG_r^T - EVF_r^T + BH_r^T = 0, \tag{70}$$

$$A^TZG_1 - E^TZF_1 + C^TH_1 = 0. \tag{71}$$

For simplicity, we suppose that E is invertible. The result equivalent to (34) states that V constructed from (70) satisfies

$$\mathcal{K}_n \left(\left(\frac{1}{\sigma} E - A \right)^{-1} E, \left(\frac{1}{\sigma} E - A \right)^{-1} B \right) = \text{Im}(V). \tag{72}$$

From this a result similar to Theorem 3.1 easily follows. We see that we interpolate the first n moments of $T(s)$ at the point $1/\sigma$ which is the inverse of the generalized eigenvalue of (\hat{E}, \hat{A}) . The other continuous-time theorems of this paper are easily extended in a similar manner.

5.2. Link to other methods

Many model reduction techniques make use of Krylov spaces and are in fact particular versions of Multipoint Padé Interpolation. Arnoldi/Lanczos processes build reduced-order models that match moments at infinity

$$CA^jB = \hat{C}\hat{A}^j\hat{B}$$

and Padé approximation corresponds to moment matching at zero. More generally, interpolation methods often construct reduced-order models that match the original model at several points up to certain orders. These methods can be seen as particular cases of the general method explained in Section 2. Interpolation methods appear to work quite well in practice and are easy to construct. Their two main drawbacks are that, so far, there is no global bound of the error between the two models and stability may be lost in the reduced-order model. The latter can be fixed via post-processing thereby incurring extra computations and losing strict moment matching. The lack of a global error bound is seen in practice by the many heuristics that are possible for placing interpolation points and selecting at which point the next matched moment should be added. Having to select interpolation points is therefore at the same time a practical advantage—it allows to monitor the error locally—and a limiting disadvantage—a heuristic for choosing the interpolation points is needed.

The balanced truncation and optimal Hankel norm approximations are based on the characteristics of the controllability and observability gramians. They have the important advantages to provide a global error bound between the reduced-order model and the original model and to guarantee stability of the reduced-order model. Their drawback is their cubic complexity, which is too high for very large-scale dynamical systems (say, with an order $N > 1000$).

Recently, techniques that attempt to construct a reduced model with a global error bound for large-scale dynamical systems have appeared. In this third class are the ADI-technique of [8] and the low-rank Smith methods of [10,2]. These construct low-rank square root approximations of the controllability and observability gramians. As shown in [8], the image of these square root approximations of P and Q can be expressed as a union of Krylov subspaces. This establishes a connection between approximate balanced truncation via ADI and *moment matching* in the symmetric case. Another method in this class is proposed in [11] and constructs projectors Z and V by recursively adding images of solutions of Sylvester equations and then extracting from them some dominant subset. In Section 3, we have shown that the image of solutions of Sylvester equations can also be expressed as a union of Krylov subspaces. Even though these methods cannot directly be interpreted as interpolation techniques, they are linked to Krylov based interpolation by Theorem 2.1.

5.3. Future work

Any strictly proper scalar rational function of degree n is determined by $2n$ parameters and can, therefore, be uniquely defined by $2n$ scalar interpolating conditions with respect to an original higher degree transfer function. It can be shown that such a reduced-order transfer function is the unique minimal transfer function of order n that satisfies the $2n$ interpolating conditions if and only if this reduced-order transfer function can be built via multipoint Padé. It follows from this that, generically, imposing $2n$ interpolating conditions determines a unique minimal reduced-order transfer function, and it can be constructed via multipoint Padé. We may wonder if the following stronger statement is true: “Every transfer function of strict degree n can be constructed via multipoint Padé approximation of a higher degree transfer function”. If this is true, then every reduced-order model can be constructed via interpolation. Of course, finding the *good* interpolating conditions is not an easy task and might be more difficult than using other model reduction techniques.

In the MIMO case, it is possible to extend Multipoint Padé Interpolation to two-sided tangential interpolation in multiple points (see [4] for a definition) and the same question about generality of the approach can be asked here. These issues will be addressed in a subsequent paper.

References

- [1] A.C. Antoulas, D.C. Sorensen, Approximation of large scale dynamical systems: an overview, Technical Report 01-01, Department of Computational and Applied Mathematics, Rice University, 2001.
- [2] A.C. Antoulas, D.C. Sorensen, S. Gugercin, A modified low-rank Smith method for large-scale Lyapunov equations, Technical Report 01-10, Department of Computational and Applied Mathematics, Rice University, May 2001.
- [3] A.C. Antoulas, D.C. Sorensen, Y. Zhou, On the decay rate of Hankel singular values and related issues, *Systems and Control Letters* 46 (2002) 323–342.
- [4] J.A. Ball, I. Gohberg, L. Rodman, *Interpolation of Rational Matrix Functions*, Birkhäuser Verlag, Basel, 1990.
- [5] K. Gallivan, E. Grimme, P. Van Dooren, Model reduction of large-scale systems, rational Krylov versus balancing techniques in: H. Bulgak, C. Zenger (Eds.), *Error Control and adaptivity in Scientific Computing*, Kluwer Academic Publishers, Dordrecht, 1999, pp. 177–190.
- [6] F.R. Gantmacher, *Theory of Matrices*, Vol. 2, Chelsea Publishing Company, Chelsea, New York, 1959.
- [7] E. Grimme, Krylov projection methods for model reduction, Ph.D. Thesis, Department of Electrical Engineering, University of Illinois at Urbana-Champaign, 1997.

- [8] J.-R. Li, Model reduction of large linear systems via low rank system gramians, Ph.D. Thesis, Department of Mathematics, Massachusetts Institute of Technology, 2000.
- [9] B.C. Moore, Principal component analysis in linear systems: controllability, observability, and model reduction, *IEEE Trans. Automat. Control* 26 (1981) 17–31.
- [10] T. Penzl, A cyclic low rank Smith method for large sparse Lyapunov equations, *SIAM J. Sci. Comput.* 21 (4) (1999/2000) 1401–1418 (electronic).
- [11] D.C. Sorensen, A.C. Antoulas, Projection methods for balanced model reduction, Technical Report 01-03, Department of Computational and Applied Mathematics, Rice University, March 2001. *Linear Algebra Appl.*, to appear.
- [12] P. Van Dooren, Gramian based model reduction of large-scale dynamical systems, *Numer. Anal.* 1999 (2000) 231–247.