

---

Lecture Notes  
**“Model Reduction”**

University of Hamburg  
(summer term 2019)

Dr. Matthias Voigt  
`matthias.voigt@uni-hamburg.de`

---



---

## Preface

---

This document is based in large parts on the hand-written lecture notes of Christian Schröder who gave this course on model reduction at the TU Berlin in the winter term 2016/17. Special thanks go to Martijn Nagtegaal, Nora Heinrich, and Ines Ahrens for finding so many typos in the initial version of these lecture notes from winter term 2017/18. I believe that there are more errors and typos in this document, please send an email to [matthias.voigt@uni-hamburg.de](mailto:matthias.voigt@uni-hamburg.de) if you find any.

There are not many textbooks on model reduction, the most commonly known one has been written by A. C. Antoulas:

A. C. Antoulas. *Approximation of Large-Scale Dynamical Systems*, volume 6 of *Adv. Des. Control*. SIAM Publications, Philadelphia, PA, 2005. doi:10.1137/1.9780898718713.

Most aspects discussed in this course have also been covered by Peter Benner on the Gene Golub SIAM Summer School 2013 at Fudan University in Shanghai, China. Some more applications and illustrative examples on model reduction can be found in his slides that you can download from the summer school's website<sup>1</sup>. Since this course is strongly based on control-theoretic basics, I recommend to read Chapters 3 and 4 of

K. Zhou, J. C. Doyle, and K. Glover. *Robust and Optimal Control*, Prentice-Hall, Englewood Cliffs, NJ, 1996.

to look up these concepts. Further, more recent results discussed here will be cited throughout the lecture notes, so that you can read the original sources.

---

<sup>1</sup><http://g2s3.cs.ucdavis.edu/lecturers/Benner/Benner-lectures-online.pdf>



---

## Contents

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	What is Model Reduction? . . . . .	1
1.2	Examples of Large-Scale Dynamical Systems . . . . .	3
1.2.1	A Controlled Discretized Heat Equation . . . . .	3
1.2.2	Further Examples . . . . .	4
<b>2</b>	<b>Basics of Systems and Control Theory</b>	<b>7</b>
2.1	Properties of LTI Systems . . . . .	7
2.2	Laplace Transformation and Transfer Functions . . . . .	11
2.3	Realizations . . . . .	13
2.4	Hardy Spaces . . . . .	14
2.4.1	The Hilbert Space $\mathcal{H}_2^{p \times m}$ . . . . .	15
2.4.2	The Banach Space $\mathcal{H}_\infty^{p \times m}$ . . . . .	18
<b>3</b>	<b>Eigenvalue-Based Approaches</b>	<b>21</b>
3.1	Modal Truncation . . . . .	23
3.2	The Dominant Pole Algorithm . . . . .	25
<b>4</b>	<b>Balancing-Based Approaches</b>	<b>31</b>
4.1	Input and Output Energy . . . . .	31
4.2	Balancing Transformations and Balanced Truncation . . . . .	33
4.3	Hankel Operator and Hankel Singular Values . . . . .	37
4.4	Properties of Balanced Truncation . . . . .	40
4.5	Numerical Solution of Large-Scale Lyapunov Equations . . . . .	46
4.5.1	Derivation of the ADI Iteration . . . . .	46
4.5.2	The ADI Shift Parameter Problem . . . . .	48
4.5.3	The Low-Rank Phenomenon . . . . .	48

---

4.5.4	The Low-Rank Cholesky Factor ADI Iteration . . . . .	49
4.5.5	Balanced Truncation Using the LRCF-ADI Method . . . . .	52
<b>5</b>	<b>Passivity-Preserving Balancing-Based Model Reduction</b>	<b>55</b>
5.1	Passivity and Positive Real Transfer Functions . . . . .	55
5.2	Positive Real Balanced Truncation . . . . .	60
5.3	Analysis of the Method . . . . .	63
5.4	Numerical Solution of Large-Scale Algebraic Riccati Equations . . . . .	66
5.4.1	Newton's Method for Solving Algebraic Riccati Equations . . . . .	66
5.4.2	The Low-Rank Newton-Kleinman Method . . . . .	70
<b>6</b>	<b>Interpolatory Model Reduction</b>	<b>73</b>
6.1	Moment Matching . . . . .	73
6.1.1	Moments . . . . .	73
6.1.2	One-Sided Moment Matching . . . . .	75
6.1.3	Two-Sided Moment Matching . . . . .	79
6.2	$\mathcal{H}_2$ -Optimal Interpolation: The Iterative Rational Krylov Algorithm . . . . .	84
6.3	Interpolation from Data: The Loewner Framework . . . . .	94
<b>7</b>	<b>Outlook</b>	<b>101</b>
7.1	Parametric Model Reduction . . . . .	101
7.2	Sampling-Based Methods . . . . .	103

---

### 1.1 What is Model Reduction?

Today, for the study of real-world processes, one usually sets up mathematical models usually consisting of differential (or differential-algebraic) equations that describe the behavior of the system under consideration. However, there is an ever-increasing need for higher accuracy which means that these models get more and more complex. The simulation, optimization, and control using such models then often leads to a very high demand in computational resources, both in terms of consumed time and memory – often even forbidding performing the desired task at all. Therefore, there is need for replacing the complex mathematical model by a much simpler model, that approximately behaves like the original model but which is computationally much less demanding. The process of finding this simpler representation is called *model reduction*. The typical set-up is depicted in Figure 1.1.

In this course we mainly consider *control systems* of the general form

$$\begin{aligned}\dot{x}(t) &= f(t, x(t), u(t)), & x(t_0) &= x_0 \\ y(t) &= g(t, x(t), u(t)),\end{aligned}\tag{1.1}$$

where  $\mathbb{I} = [t_0, t_f]$  is a time interval of interest,  $x : \mathbb{I} \rightarrow \mathbb{R}^n$  is the *state function* with initial value  $x_0 \in \mathbb{R}^n$ ,  $u : \mathbb{I} \rightarrow \mathbb{R}^m$  is the *input function*,  $y : \mathbb{I} \rightarrow \mathbb{R}^p$  is the *output function*, and  $f : \mathbb{I} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$  and  $g : \mathbb{I} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^p$ . Usually, the input is a function that can be used to control the state of the system to

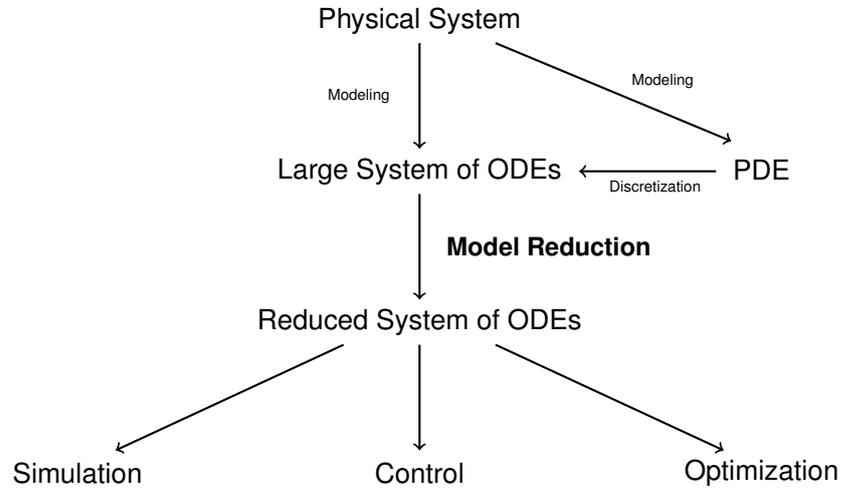


Figure 1.1: The broad setup of model reduction.

achieve a desired behavior. The output consists of “quantities of interest” that can often be measured in the real physical process.

The goal of model reduction is to replace the functions  $f$  and  $g$  in (1.1) by a *reduced-order model*

$$\begin{aligned}\dot{\tilde{x}}(t) &= \tilde{f}(t, \tilde{x}(t), u(t)), & \tilde{x}(t_0) &= \tilde{x}_0 \\ \tilde{y}(t) &= \tilde{g}(t, \tilde{x}(t), u(t)),\end{aligned}\tag{1.2}$$

where  $\tilde{x} : \mathbb{I} \rightarrow \mathbb{R}^r$  with  $r \ll n$  is the reduced state function with initial value  $\tilde{x}_0 \in \mathbb{R}^r$ , and  $\tilde{f} : \mathbb{I} \times \mathbb{R}^r \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ ,  $\tilde{g} : \mathbb{I} \times \mathbb{R}^r \times \mathbb{R}^m \rightarrow \mathbb{R}^p$ . This model should be constructed such that  $\|y - \tilde{y}\|$  is “small” for all admissible inputs  $u$ . This will be made precise later. Note that we are only interested in the map from the input to the output, not in the evolution of the state itself.

In this course we focus on *linear time-invariant systems*, which are of the simpler form

$$\begin{aligned}\dot{x}(t) &= Ax(t) + Bu(t), \\ y(t) &= Cx(t) + Du(t),\end{aligned}\tag{1.3}$$

for some matrices  $A \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times m}$ ,  $C \in \mathbb{R}^{p \times n}$ , and  $D \in \mathbb{R}^{p \times m}$ . Normally, we will also assume that  $x(t_0) = x(0) = 0$ .

We will discuss a rigorous mathematical theory for model reduction. We will discuss efficient numerical algorithms as well as theorems on the approximation quality, e. g., we state and prove error bounds. We will also touch on aspects of structure-preservation. This means, that if the original model has a certain structure, then also the reduced model should have this structure to account for physical properties that are encoded in the model.

## 1.2 Examples of Large-Scale Dynamical Systems

### 1.2.1 A Controlled Discretized Heat Equation

Consider the temperature distribution  $T(t, \xi)$  of a one-dimensional beam of length  $\ell = 1$ . Here,  $\xi \in [0, 1]$  is the space variable and  $t \geq 0$  denotes the time. On the right end of the beam we impose the boundary condition  $T(t, 1) = 0$  for all  $t \geq 0$ . On the left end we have a heat source that results in a controllable temperature flux

$$-\frac{\partial}{\partial \xi} T(t, 0) = u(t).$$

The heat diffusion inside the beam is described by the heat equation

$$\frac{\partial}{\partial t} T(t, \xi) = k \cdot \frac{\partial^2}{\partial \xi^2} T(t, \xi) \quad \text{for all } \xi \in (0, 1), t > 0.$$

Moreover, we are interested in the average beam temperature, i. e., our output is

$$y(t) = \int_0^1 T(t, \xi) d\xi.$$

Finally, we need an initial condition which is given by

$$T(0, \xi) = 0 \quad \text{for all } \xi \in [0, 1].$$

Now we discretize in space at  $n$  equidistant points and obtain

$$x(t) = \begin{bmatrix} x_1(t) \\ \vdots \\ x_n(t) \end{bmatrix} := \begin{bmatrix} T(t, 0) \\ T(t, \frac{1}{n}) \\ \vdots \\ T(t, \frac{n-1}{n}) \end{bmatrix}.$$

For  $i = 2, 3, \dots, n-1$  we find

$$\begin{aligned} \dot{x}_i(t) &= \frac{\partial}{\partial t} T(t, \frac{i-1}{n}) = k \cdot \frac{\partial^2}{\partial \xi^2} T(t, \frac{i-1}{n}) \\ &\approx k \cdot n^2 (T(t, \frac{i-2}{n}) - 2T(t, \frac{i-1}{n}) + T(t, \frac{i}{n})) \\ &= k \cdot n^2 (x_{i-1}(t) - 2x_i(t) + x_{i+1}(t)). \end{aligned}$$

Analogously, we find

$$\dot{x}_n(t) \approx k \cdot n^2 (x_{n-1}(t) - 2x_n(t)),$$

since  $x_{n+1}(t) := T(t, \frac{n}{n}) = 0$ . Moreover, we have

$$\begin{aligned} \dot{x}_1(t) &= \frac{\partial}{\partial t} T(t, 0) = k \cdot \frac{\partial^2}{\partial \xi^2} T(t, 0) \\ &\approx k \cdot n \left( \frac{\partial}{\partial \xi} T(t, \frac{1}{n}) - \frac{\partial}{\partial \xi} T(t, 0) \right) \\ &\approx k \cdot n (n(x_2(t) - x_1(t)) + u(t)). \end{aligned}$$

For the output we take a piecewise constant approximation, i. e.,

$$y(t) = \int_0^1 T(t, \xi) d\xi \approx \frac{1}{n} \sum_{i=0}^{n-1} T(t, \frac{i}{n}) = \frac{1}{n} \sum_{i=1}^n x_i(t).$$

The zero initial conditions imply  $x(0) = 0$ . Our final controlled discretized heat equation now attains the form

$$\begin{aligned} \dot{x}(t) &= Ax(t) + Bu(t), \quad x(0) = 0, \\ y(t) &= Cx(t), \end{aligned} \tag{1.4}$$

where

$$A = kn^2 \begin{bmatrix} -1 & 1 & & & \\ 1 & -2 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & 1 & -2 & 1 \\ & & & 1 & -2 \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad B = kn \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{n \times 1}, \tag{1.5}$$

$$C = \frac{1}{n} [1 \quad \dots \quad 1] \in \mathbb{R}^{1 \times n}.$$

The larger  $n$  the better the solution of the PDE will be approximated, but the size of the system of ODEs in (1.4) and (1.5) will also grow and thus its evaluation will be more expensive.

Let  $k = 1$  and  $n = 1000$ . Using the method of balanced truncation (discussed later), we can approximate the system by

$$\begin{aligned} \dot{\tilde{x}}(t) &= \tilde{A}\tilde{x}(t) + \tilde{B}u(t), \quad \tilde{x}(0) = 0, \\ \tilde{y}(t) &= \tilde{C}\tilde{x}(t), \end{aligned}$$

with

$$\begin{aligned} \tilde{A} &= \begin{bmatrix} -2.256 & 1.775 & -0.6057 \\ -1.775 & -16.63 & 12.21 \\ -0.6057 & -12.21 & -40.66 \end{bmatrix}, \quad \tilde{B} = \begin{bmatrix} -1.074 \\ -0.4136 \\ -0.1442 \end{bmatrix}, \\ \tilde{C} &= [-1.074 \quad 0.4136 \quad -0.1442]. \end{aligned}$$

Simulation with various inputs shows that the outputs  $y$  and  $\tilde{y}$  are almost the same.

## 1.2.2 Further Examples

Here will briefly mention a few more examples to illustrate the importance of model reduction in practice.

---

**Electrical Circuits.** Electrical circuits containing only inductors, capacitors, and resistors can be modeled using modified modal analysis. This results in a linear system of the form

$$\begin{aligned} E\dot{x}(t) &= Ax(t) + Bu(t), \\ y(t) &= B^T x(t), \end{aligned} \quad (1.6)$$

where the state  $x(\cdot)$  contains the node potentials and currents through inductors and voltage sources. The input  $u(\cdot)$  contains the currents of the current sources as well as the voltages of the voltage sources. The output  $y(\cdot)$  contains the negative of the voltages of the current sources and the currents of the voltage sources. Here the matrices  $E$ ,  $A$ , and  $B$  have the form

$$\begin{aligned} E &= \begin{bmatrix} A_C C A_C^T & 0 & 0 \\ 0 & \mathcal{L} & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad A = \begin{bmatrix} -A_{\mathcal{R}} \mathcal{G} A_{\mathcal{R}}^T & -A_{\mathcal{L}} & -A_{\mathcal{V}} \\ A_{\mathcal{L}}^T & 0 & 0 \\ A_{\mathcal{V}}^T & 0 & 0 \end{bmatrix}, \\ B &= \begin{bmatrix} -A_{\mathcal{I}} & 0 \\ 0 & 0 \\ 0 & -I \end{bmatrix}, \end{aligned} \quad (1.7)$$

where  $\mathcal{G}$ ,  $\mathcal{L}$ ,  $\mathcal{C}$  are positive definite matrices containing the conductances, inductances, and capacities of the resistors, inductors, and capacitors, respectively. The matrices  $A_C$ ,  $A_{\mathcal{R}}$ ,  $A_{\mathcal{L}}$ ,  $A_{\mathcal{V}}$ , and  $A_{\mathcal{I}}$  are incidence matrices that describe the network topology of the circuit. This model differs from (1.3), namely an additional matrix  $E$  is in front of  $\dot{x}$  and moreover,  $E$  is singular. This means, that not all of the equations in (1.6) are differential equations, but there are also algebraic equations that result from Kirchhoff's laws. Therefore, such a system is called a *differential-algebraic system*. Moreover, the system (1.6) with (1.7) has certain symmetries that account for the physical properties of the circuit. For example, (1.6) with (1.7) is a *passive system*, meaning that

$$\int_0^T y(t)^T u(t) dt \geq 0$$

for all  $T \geq 0$  and all smooth solution trajectories with  $Ex(0) = 0$ . This property must be reflected in the reduced-order model in order to get meaningful simulation results. In other words, structure-preserving methods are of great importance in applications.

**Structural Mechanics.** The goal of structural mechanics is the computation of mechanical deformations and internal forces and stresses within mechanical structures, such as buildings, bridges, machines, etc. Using the finite element method, the mechanical structure is decomposed into masses that are stiffly

connected. This leads to a large ordinary differential equation of second order of the form

$$\begin{aligned} M\ddot{x}(t) + D\dot{x}(t) + Kx(t) &= Bu(t), \\ y(t) &= C_1x(t) + C_2\dot{x}(t), \end{aligned} \tag{1.8}$$

where the state  $x(\cdot)$  is the displacement of the masses from the equilibrium position and the input  $u(\cdot)$  is an external force. Moreover,  $M$  and  $K$  are the positive definite mass and stiffness matrices and  $D$  is a positive definite damping matrix. Using a linearization, one can in principal rewrite (1.8) as a first-order system as in (1.3) and do model reduction on the first-order system. However, often it is important to have a reduced-order model of the form (1.8). It is often not possible to gain such a system when using methods for first-order systems. There are methods that work directly on (1.8), but there are still many open research problems.

---

---

## Basics of Systems and Control Theory

---

In this chapter we consider linear time-invariant (LTI) control systems

$$\begin{aligned} \dot{x}(t) &= Ax(t) + Bu(t), & x(t_0) &= x_0, \\ y(t) &= Cx(t) + Du(t), \end{aligned} \tag{2.1}$$

where  $A \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times m}$ ,  $C \in \mathbb{R}^{p \times n}$ ,  $D \in \mathbb{R}^{p \times m}$ ,  $x : [t_0, t_f] \rightarrow \mathbb{R}^n$  is the *state* of the system,  $u : [t_0, t_f] \rightarrow \mathbb{R}^m$  denotes a *control input* and  $y : [t_0, t_f] \rightarrow \mathbb{R}^p$  is a *measurable output*. The set of LTI systems with state-space dimension  $n$ ,  $m$  inputs, and  $p$  outputs is denoted by  $\Sigma_{n,m,p}$  and we write  $[A, B, C, D] \in \Sigma_{n,m,p}$ . The goal of this chapter is to give a basic analysis and discussion of such systems in order to set the foundations for the model reduction methods we discuss later. Here we will rather skip the proofs or keep them short since this will mainly be the topic of the course on control theory. A more detailed introduction to the concepts presented here can be found in the textbook [ZDG96, Chapters 3 & 4].

### 2.1 Properties of LTI Systems

Next we discuss some fundamental properties of LTI dynamical systems. In the next definition we assume for simplicity that  $t_f = \infty$  and that  $\mathcal{U}_{\text{ad}} := \mathcal{PC}([t_0, t_f], \mathbb{R}^m)$  is the set of admissible inputs, i. e., the set of all piecewise continuous functions mapping from  $[t_0, t_f]$  to  $\mathbb{R}^m$ , but in principal we could also take  $\mathcal{U}_{\text{ad}} = \mathcal{L}_2([t_0, t_f], \mathbb{R}^m)$ .

**Definition 2.1:** The LTI system  $[A, B, C, D] \in \Sigma_{n,m,p}$  is called

- a) *asymptotically stable*, if all solutions of the linear homogeneous ODE  $\dot{x}(t) = Ax(t)$  satisfy  $\lim_{t \rightarrow \infty} x(t) = 0$  for all initial conditions  $x(t_0) = x_0$ .
- b) *controllable*, if for all initial conditions  $x(t_0) = x_0$  and all  $x_1 \in \mathbb{R}^n$ , there exists a  $t_1 > t_0$  and a control function  $u \in \mathcal{U}_{\text{ad}}$  such that  $x(t_1) = x_1$ .
- c) *stabilizable*, if for all initial conditions  $x(t_0) = x_0$  there exists a control function  $u \in \mathcal{U}_{\text{ad}}$  such that  $\lim_{t \rightarrow \infty} x(t) = 0$ .
- d) *observable*, if for two solution trajectories (obtained with the same input  $u \in \mathcal{U}_{\text{ad}}$ )  $x(\cdot)$  and  $\tilde{x}(\cdot)$  it holds

$$Cx(t) = C\tilde{x}(t) \quad \forall t \geq t_0 \Rightarrow x(t) = \tilde{x}(t) \quad \forall t \geq t_0.$$

- e) *detectable*, if for any solution  $x(\cdot)$  of  $\dot{x}(t) = Ax(t)$  with  $Cx(t) \equiv 0$  it follows that  $\lim_{t \rightarrow \infty} x(t) = 0$ .

The following lemma characterizes these properties algebraically.

**Lemma 2.2:** The LTI system  $[A, B, C, D] \in \Sigma_{n,m,p}$  is

- a) asymptotically stable  $\Leftrightarrow \Lambda(A) \subset \mathbb{C}^- := \{\lambda \in \mathbb{C} : \text{Re}(\lambda) < 0\}$ ,
- b) controllable  $\Leftrightarrow \text{rank} \begin{bmatrix} \lambda I_n - A & B \end{bmatrix} = n \quad \forall \lambda \in \mathbb{C}$   
 $\Leftrightarrow \text{rank} \begin{bmatrix} B & AB & \dots & A^{n-1}B \end{bmatrix} = n$ ,
- c) stabilizable  $\Leftrightarrow \text{rank} \begin{bmatrix} \lambda I_n - A & B \end{bmatrix} = n \quad \forall \lambda \in \overline{\mathbb{C}^+} := \{\lambda \in \mathbb{C} : \text{Re}(\lambda) \geq 0\}$   
 $\Leftrightarrow \exists F \in \mathbb{R}^{m \times n}$  such that  $\Lambda(A + BF) \subset \mathbb{C}^-$ ,
- d) observable  $\Leftrightarrow \text{rank} \begin{bmatrix} \lambda I_n - A \\ C \end{bmatrix} = n \quad \forall \lambda \in \mathbb{C}$   
 $\Leftrightarrow \text{rank} \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{bmatrix} = n$ ,
- e) detectable  $\Leftrightarrow \text{rank} \begin{bmatrix} \lambda I_n - A \\ C \end{bmatrix} = n \quad \forall \lambda \in \overline{\mathbb{C}^+}$   
 $\Leftrightarrow \exists G \in \mathbb{R}^{n \times p}$  such that  $\Lambda(A + GC) \subset \mathbb{C}^-$ .

**Remark 2.3:** a) Stabilizability weakens the concept of controllability in the sense that not all possible states are reachable, but uncontrollable parts tend to zero.

- b) Detectability weakens observability in the same sense as stabilizability weakens controllability: not all of  $x$  can be observed but unobserved parts are asymptotically stable, i. e., deviations vanish over time.
- c) The above concepts are *dual* in the sense that an LTI system is observable (detectable) if and only if the dual system

$$\dot{z}(t) = A^T z(t) + C^T v(t)$$

is controllable (stabilizable).

The following considerations motivate the Gramians that we define next. First we consider the input-to-state map

$$\zeta(t) = e^{At} B,$$

which is motivated by the fact that for  $x(0) = x_0 = 0$  and impulsive inputs  $u = u_0 \cdot \delta$  (where  $u_0 \in \mathbb{R}^m$  and  $\delta$  denotes the Dirac delta distribution), we obtain

$$\begin{aligned} x(t) &= e^{At} x_0 + \int_0^t e^{A(t-\tau)} B u(\tau) d\tau \\ &= \int_0^t e^{A(t-\tau)} B u_0 \delta(\tau) d\tau \\ &= e^{At} B u_0 = \zeta(t) u_0. \end{aligned}$$

Note that the above is formally *not correct*, since  $\delta$  is not a function mapping from  $\mathbb{R}$  to  $\mathbb{R}$ , but a distribution (often called generalized function) that is defined by

$$\delta : \mathcal{C}^\infty(\mathbb{R}, \mathbb{R}^n) \rightarrow \mathbb{R}^n, \quad f \mapsto f(0).$$

So actually we have more correctly

$$x(t) := \delta \left( e^{A(t-\cdot)} B u_0 \right) = e^{At} B u_0.$$

Consider on the other hand the state-to-output map

$$\eta(t) = C e^{At},$$

which is motivated by the fact that for  $x(0) = x_0$  and  $u(t) \equiv 0$ , we obtain

$$\begin{aligned} y(t) &= C e^{At} x_0 + C \int_0^t e^{A(t-\tau)} B u(\tau) d\tau \\ &= C e^{At} x_0 = \eta(t) x_0. \end{aligned}$$

For the analysis of LTI control systems we now make use of the following Gramians.

---

**Definition 2.4:** The matrix

$$P(T) = \int_0^T e^{At} B B^T e^{A^T t} dt$$

is called the  $(0, T)$ -controllability Gramian of the system (2.1).

The matrix

$$Q(T) = \int_0^T e^{A^T t} C^T C e^{At} dt$$

is called the  $(0, T)$ -observability Gramian of the system (2.1).

**Remark 2.5:** For a system (2.1),  $P(T)$  and  $Q(T)$  can be used to identify states of the system that are easily reachable and easily observable in the interval  $(0, T)$  in the following sense:

a) For a reachable state  $x_*$  of the system (2.1), one can show that  $\hat{u}(t) = B^T e^{A^T(t_*-t)} P(t_*)^\dagger x_*$ , where  $P(t_*)^\dagger$  is the Moore-Penrose inverse of  $P(t_*)$ , controls the system from  $x(0) = 0$  to  $x(t_*) = x_*$ . Moreover, among all such controls,  $\hat{u}(t)$  is the one with minimal  $\mathcal{L}_2$ -norm.

b) For any  $t_* > 0$  and  $x_0 \in \mathbb{R}^n$ , we have

$$x_0^T Q(t_*) x_0 = \int_0^{t_*} x_0^T e^{A^T t} C^T C e^{At} x_0 dt = \int_0^{t_*} \|C e^{At} x_0\|_2^2 dt = \|y_{x_0}(\cdot)\|_{\mathcal{L}_2}^2.$$

Now we consider the above Gramians for  $T \rightarrow \infty$ .

**Lemma 2.6:** If  $A$  in (2.1) is asymptotically stable, then

a) the infinite controllability and observability Gramians

$$P = \lim_{T \rightarrow \infty} P(T) = \int_0^\infty e^{At} B B^T e^{A^T t} dt$$

and

$$Q = \lim_{T \rightarrow \infty} Q(T) = \int_0^\infty e^{A^T t} C^T C e^{At} dt$$

exist,

b) they solve the two *Lyapunov equations*

$$\begin{aligned} AP + PA^T &= -BB^T, \\ A^T Q + QA &= -C^T C. \end{aligned}$$

c) If  $(A, B)$  is controllable and  $(A, C)$  is observable, it moreover holds that  $P = P^T > 0$  and  $Q = Q^T > 0$ . (Otherwise we just have  $P = P^T \geq 0$  and  $Q = Q^T \geq 0$ .)

*Proof.* Exercise or lecture “Control Theory”. □

## 2.2 Laplace Transformation and Transfer Functions

In applications it is often useful to consider a dynamical system in the frequency domain. When doing so, the system can be treated using tools from linear algebra instead of from differential equations. A function  $f : [0, \infty) \rightarrow \mathbb{R}^n$  is called *exponentially bounded*, if there exist numbers  $M$  and  $\alpha$  such that  $\|f(t)\|_2 \leq M e^{\alpha t}$  for all  $t \geq 0$ . The value  $\alpha$  is called a *bounding exponent*.

**Definition 2.7:** Let  $f : [0, \infty) \rightarrow \mathbb{R}^n$  be exponentially bounded with bounding exponent  $\alpha$ . Then

$$\mathcal{L}\{f\}(s) := \int_0^{\infty} f(\tau) e^{-s\tau} d\tau$$

for  $\operatorname{Re}(s) > \alpha$  is called the *Laplace transform* of  $f$ . The process of forming the Laplace transform is called *Laplace transformation*.

It can be shown that the integral converges uniformly in a domain of the form  $\operatorname{Re}(s) \geq \beta$  for all  $\beta > \alpha$ .

Moreover, the following two fundamental properties hold.

**Theorem 2.8:** Let  $f, g, h : [0, \infty) \rightarrow \mathbb{R}^n$  be given. Then the following two statements hold true:

a) The Laplace transformation is linear, i. e., if  $f$  and  $g$  are exponentially bounded, then  $h := \gamma f + \delta g$  is also exponentially bounded and

$$\mathcal{L}\{h\} = \gamma \mathcal{L}\{f\} + \delta \mathcal{L}\{g\}$$

holds for all  $\gamma, \delta \in \mathbb{C}$ .

b) If  $f \in \mathcal{PC}^1([0, \infty), \mathbb{R}^n)$  and  $\dot{f}$  is exponentially bounded, then  $f$  is exponentially bounded and

$$\mathcal{L}\{\dot{f}\}(s) = s \mathcal{L}\{f\}(s) - f(0).$$

Now we apply the Laplace transformation to the system  $[A, B, C, D] \in \Sigma_{n,m,p}$ . Assume that each of the Laplace transforms  $X(s) := \mathcal{L}\{x\}(s)$ ,  $U(s) := \mathcal{L}\{u\}(s)$ , and  $Y(s) := \mathcal{L}\{y\}(s)$  exist. By using Theorem 2.8, we obtain the Laplace transformed system

$$\begin{aligned} sX(s) - x(0) &= AX(s) + BU(s), \\ Y(s) &= CX(s) + DU(s). \end{aligned}$$

Under the assumption that  $x(0) = 0$ , we obtain the relation

$$Y(s) = (C(sI_n - A)^{-1}B + D)U(s).$$

This leads to the following definition.

**Definition 2.9:** The function

$$G(s) := C(sI_n - A)^{-1}B + D \in \mathbb{R}(s)^{p \times m}$$

is called the *transfer function* of the system  $[A, B, C, D] \in \Sigma_{n,m,p}$ . Here,  $\mathbb{R}(s)^{p \times m}$  denotes the set of all  $p \times m$  matrices that have real-rational functions as entries.

The following properties of rational functions will play an important role in the characterization of transfer functions.

**Definition 2.10 (Properness):** Let  $G(s) \in \mathbb{R}(s)^{p \times m}$  be given. We call  $G(s)$

- a) *strictly proper*, if  $\lim_{\omega \rightarrow \infty} \|G(i\omega)\|_2 = 0$ ;
- b) *proper*, if  $\lim_{\omega \rightarrow \infty} \|G(i\omega)\|_2 < \infty$ ;
- c) *improper*, if  $\lim_{\omega \rightarrow \infty} \|G(i\omega)\|_2 = \infty$ .

Since  $\lim_{\omega \rightarrow \infty} (i\omega I_n - A)^{-1} = 0$ , it is easy to see that transfer functions of systems  $[A, B, C, D] \in \Sigma_{n,m,p}$  are always proper, improper transfer functions can only be realized by DAE systems. Furthermore, the transfer function of a system  $[A, B, C, D] \in \Sigma_{n,m,p}$  is strictly proper if and only if  $D = 0$ . Now we define the notions of poles and zeros of rational matrices. For this we need the following terms.

**Definition 2.11 (Unimodular matrix, monic/copprime polynomials):** Let  $\mathbb{R}[s]$  denote the set of polynomials with real coefficients.

- a) A polynomial matrix  $U(s) \in \mathbb{R}[s]^{n \times n}$  is called *unimodular*, if its determinant is a nonzero constant in  $\mathbb{R}$ .
- b) A polynomial  $p(s) \in \mathbb{R}[s]$  is called *monic*, if its leading coefficient is one.

c) Two polynomials  $p(s), q(s) \in \mathbb{R}[s]$  are called *coprime*, if their greatest common divisor is 1.

Matrices with rational entries can, via multiplication with suitable unimodular matrices, be transformed to *Smith-McMillan* form, described in the next theorem.

**Theorem 2.12** (Smith-McMillan form): For  $G(s) \in \mathbb{R}(s)^{p \times m}$  there exist unimodular matrices  $U(s) \in \mathbb{R}[s]^{p \times p}$  and  $V(s) \in \mathbb{R}[s]^{m \times m}$ , such that

$$U^{-1}(s)G(s)V^{-1}(s) = \begin{bmatrix} \tilde{G}(s) & 0 \\ 0 & 0 \end{bmatrix} \quad \text{with} \quad \tilde{G}(s) = \text{diag} \left( \frac{\varepsilon_1(s)}{\psi_1(s)}, \dots, \frac{\varepsilon_r(s)}{\psi_r(s)} \right) \quad (2.2)$$

for some monic and coprime polynomials  $\varepsilon_j(s), \psi_j(s) \in \mathbb{R}[s]$  such that  $\varepsilon_j(s)$  divides  $\varepsilon_{j+1}(s)$  and  $\psi_{j+1}(s)$  divides  $\psi_j(s)$  for  $j = 1, \dots, r-1$ .

The Smith-McMillan form can now be utilized to define poles and zeros of rational matrices.

**Definition 2.13** (Poles and zeros): Let  $G(s) \in \mathbb{R}(s)^{p \times m}$  with Smith-McMillan form (2.2) be given. Then  $\lambda \in \mathbb{C}$  is called

- a) a *zero* of  $G(s)$  if  $\varepsilon_r(\lambda) = 0$ ;
- b) a *pole* of  $G(s)$  if  $\psi_1(\lambda) = 0$ .

Roughly speaking, the poles of  $G(s)$  are the points  $\lambda_0 \in \mathbb{C}$  where we have  $\lim_{\lambda \rightarrow \lambda_0} \|G(\lambda)\| = \infty$ . The zeros are the points  $\lambda_0 \in \mathbb{C}$  where a rank drop occurs, i. e., those points where the rank of  $G(\lambda_0)$  is strictly less than the rank for all other matrices  $G(\lambda)$ , where  $\lambda$  is in some neighborhood of  $\lambda_0$ .

## 2.3 Realizations

It is also possible to assign a dynamical system  $[A, B, C, D] \in \Sigma_{n,m,p}$  to a given proper transfer function  $G(s) \in \mathbb{R}(s)^{p \times m}$  which is, however, not unique. This leads to the following definitions.

**Definition 2.14:** Assume that the system  $[A, B, C, D] \in \Sigma_{n,m,p}$  has the proper transfer function  $G(s) \in \mathbb{R}(s)^{p \times m}$ . Then we say that  $[A, B, C, D]$  is a *realization* of  $G(s)$ . The smallest  $n \geq 0$  such that  $[A, B, C, D] \in \Sigma_{n,m,p}$  is a realization of  $G(s)$  is called the *McMillan degree* of  $G(s)$ . A realization  $[A, B, C, D] \in \Sigma_{n,m,p}$

of  $G(s)$  is called *minimal*, if  $n$  is the McMillan degree of  $G(s)$ .

**Remark 2.15:** a) Realizations are not unique. If  $[A, B, C, D] \in \Sigma_{n,m,p}$  is a realization of  $G(s)$ , then for any nonsingular matrix  $T \in \mathbb{R}^{n \times n}$ , the system

$$[T^{-1}AT, T^{-1}B, CT, D] \in \Sigma_{n,m,p}$$

is also a realization of  $G(s)$ . Transformations of the above kind are also called *state-space transformations*.

b) A realization is minimal, if and only if it is both controllable and observable.

If a realization is not minimal, we can obtain a minimal realization by using *Kalman decompositions*. There is a *controllability Kalman decomposition*, meaning that for  $[A, B, C, D] \in \Sigma_{n,m,p}$  there exists an orthogonal matrix  $Q \in \mathbb{R}^{n \times n}$  such that

$$Q^T A Q = \begin{bmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{bmatrix}, \quad Q^T B = \begin{bmatrix} B_1 \\ 0 \end{bmatrix}, \quad C Q = [C_1 \quad C_2]$$

where the system  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,p}$  is controllable. In the above decomposition we have  $\Lambda(A) = \Lambda(A_{11}) \cup \Lambda(A_{22})$ . Here, the eigenvalues  $\lambda \in \Lambda(A_{22})$  are called *uncontrollable modes* of the system  $[A, B, C, D]$  since  $B^T v = 0$  holds for all eigenvectors  $v \in \mathbb{C}^n \setminus \{0\}$  of  $A^T$  associated with eigenvalues in  $\Lambda(A_{22})$ .

On the other hand, there is the *observability Kalman decomposition*, i. e., there exists an orthogonal matrix  $\tilde{Q} \in \mathbb{R}^{n \times n}$  such that

$$\tilde{Q}^T A \tilde{Q} = \begin{bmatrix} \tilde{A}_{11} & 0 \\ \tilde{A}_{21} & \tilde{A}_{22} \end{bmatrix}, \quad \tilde{Q}^T B = \begin{bmatrix} \tilde{B}_1 \\ \tilde{B}_2 \end{bmatrix}, \quad C \tilde{Q} = [\tilde{C}_1 \quad 0],$$

where the system  $[\tilde{A}_{11}, \tilde{B}_1, \tilde{C}_1, D] \in \Sigma_{\tilde{r},m,p}$  is observable. Similarly as above, eigenvalues  $\lambda \in \Lambda(\tilde{A}_{22})$  are called *unobservable modes*, since it holds  $C v = 0$  for all eigenvectors  $v \in \mathbb{C}^n \setminus \{0\}$  of  $A$  associated with eigenvalues in  $\Lambda(\tilde{A}_{22})$ .

A minimal realization is then obtained by first computing a controllability Kalman decomposition and applying an observability Kalman decomposition to the resulting controllable subsystem.

## 2.4 Hardy Spaces

In this section we consider linear spaces of rational functions in  $\mathbb{R}(s)^{p \times m}$ . These spaces are normed spaces or even inner product spaces that allow for

geometric concepts such as length of transfer functions or distances and angles between them. Later this will be useful to measure the approximation quality of reduced-order models in terms of distances between the transfer functions of the original model and the reduced one.

### 2.4.1 The Hilbert Space $\mathcal{H}_2^{p \times m}$

The space  $\mathcal{H}_2^{p \times m}$  is defined by

$$\mathcal{H}_2^{p \times m} := \left\{ G : \mathbb{C}^+ \rightarrow \mathbb{C}^{p \times m} : G \text{ is analytic in } \mathbb{C}^+ \text{ and } \int_{-\infty}^{\infty} \|G(i\omega)\|_{\mathbb{F}}^2 d\omega < \infty \right\}.$$

Since every  $G \in \mathcal{H}_2^{p \times m}$  is analytic, there exists a unique continuation to the imaginary axis. The space  $\mathcal{H}_2^{p \times m}$  is a Hilbert space with the inner product

$$\langle F, G \rangle_{\mathcal{H}_2} := \frac{1}{2\pi} \int_{-\infty}^{\infty} \text{tr} \left( F(i\omega)^H G(i\omega) \right) d\omega.$$

This inner product induces the  $\mathcal{H}_2$ -norm

$$\|G\|_{\mathcal{H}_2} := \langle G, G \rangle_{\mathcal{H}_2}^{1/2} = \left( \frac{1}{2\pi} \int_{-\infty}^{\infty} \|G(i\omega)\|_{\mathbb{F}}^2 d\omega \right)^{1/2}.$$

We are now interested in *rational* functions, i. e., in functions that are in  $\mathcal{RH}_2^{p \times m} := \mathcal{H}_2^{p \times m} \cap \mathbb{R}(s)^{p \times m}$ . First we have the following.

**Lemma 2.16:** The following statements are equivalent:

- The function  $G$  is an element of  $\mathcal{RH}_2^{p \times m}$ .
- The function  $G$  is strictly proper and all its poles are in  $\mathbb{C}^-$ .
- The function  $G$  can be realized by a system  $[A, B, C, D]$  with  $\Lambda(A) \subset \mathbb{C}^-$  and  $D = 0$ .

The  $\mathcal{H}_2$ -norm of a transfer function can be utilized to bound the norm of the output by the norm of the input as follows. For this we will make use of the following result. It basically says that the  $\mathcal{L}_2$ -norm of a function on  $\mathbb{R}$  is equal to the  $\mathcal{L}_2$ -norm of its Fourier transform on  $i\mathbb{R}$  (scaled by a constant).

**Theorem 2.17** (Plancherel's Theorem): Let  $f \in \mathcal{L}_1(\mathbb{R}, \mathbb{R}^n) \cap \mathcal{L}_2(\mathbb{R}, \mathbb{R}^n)$ . Then the *Fourier transform* of  $f$ , given by

$$F(i\omega) := \mathcal{F}\{f\}(i\omega) := \int_{-\infty}^{\infty} f(t) e^{-i\omega t} dt$$

exists, it satisfies  $F \in \mathcal{L}_2(i\mathbb{R}, \mathbb{C}^n)$  and, moreover, it holds

$$\|f\|_{\mathcal{L}_2}^2 = \frac{1}{2\pi} \int_{-\infty}^{\infty} \|F(i\omega)\|_2^2 d\omega.$$

When we consider functions  $f$  with  $f(t) = 0$  for all  $t < 0$ , then the Fourier transform of  $f$  coincides with the Laplace transform of  $f$  restricted to the imaginary axis.

In fact, it can even be shown that the Laplace transform of  $f \in \mathcal{L}_2([0, \infty), \mathbb{C}^n)$  will always give a result that is in  $\mathcal{H}_2^n$ . Conversely, applying the inverse Laplace transform to  $F \in \mathcal{H}_2^n$  will return a function in  $\mathcal{L}_2([0, \infty), \mathbb{C}^n)$ . Summarizing, we can write

$$\mathcal{L} \{ \mathcal{L}_2([0, \infty), \mathbb{C}^n) \} = \mathcal{H}_2^n.$$

Now we show that the  $\mathcal{H}_2$ -norm bounds the  $\mathcal{L}_\infty$ -norm of the output by the  $\mathcal{L}_2$ -norm of the input.

**Theorem 2.18:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  with a transfer function  $G \in \mathcal{RH}_2^{p \times m}$  be given. Then it holds

$$\|G\|_{\mathcal{H}_2} \geq \sup_{\substack{u \in \mathcal{L}_2([0, \infty), \mathbb{R}^m) \\ u \neq 0}} \frac{\|y\|_{\mathcal{L}_\infty}}{\|u\|_{\mathcal{L}_2}}.$$

*Proof.* Since  $G \in \mathcal{RH}_2^{p \times m}$ , we have  $D = 0$  and therefore, it holds

$$y(t) = \int_0^t C e^{A(t-\tau)} B u(\tau) d\tau.$$

Set

$$g(t) := \begin{cases} C e^{At} B, & t \geq 0, \\ 0, & t < 0. \end{cases}, \quad \tilde{u}(t) := \begin{cases} u(t), & t \geq 0, \\ 0, & t < 0. \end{cases}$$

Taking norms, we obtain

$$\begin{aligned}
\|y(t)\|_2 &= \left\| \int_{-\infty}^t g(t-\tau)\tilde{u}(\tau)d\tau \right\|_2 \\
&\leq \int_{-\infty}^t \|g(t-\tau)\|_{\mathbb{F}} \|\tilde{u}(\tau)\|_2 d\tau \\
&\leq \left( \int_{-\infty}^t \|g(t-\tau)\|_{\mathbb{F}}^2 d\tau \right)^{1/2} \left( \int_{-\infty}^t \|\tilde{u}(\tau)\|_2^2 d\tau \right)^{1/2} \\
&\leq \left( \int_{-\infty}^{\infty} \|g(t-\tau)\|_{\mathbb{F}}^2 d\tau \right)^{1/2} \left( \int_{-\infty}^{\infty} \|\tilde{u}(\tau)\|_2^2 d\tau \right)^{1/2},
\end{aligned}$$

where the last inequality follows from the Cauchy-Schwarz inequality. It can be shown that (exercise!)

$$\mathcal{F}\{g\}(i\omega) = C(i\omega I_n - A)^{-1}B.$$

Using Plancherel's Theorem we obtain

$$\int_{-\infty}^{\infty} \|g(t-\tau)\|_{\mathbb{F}}^2 d\tau = \frac{1}{2\pi} \int_{-\infty}^{\infty} \|C(i\omega I_n - A)^{-1}B\|_{\mathbb{F}}^2 d\omega = \|G\|_{\mathcal{H}_2}^2.$$

Therefore, we obtain  $\|y(t)\|_2 \leq \|G\|_{\mathcal{H}_2} \|u\|_{\mathcal{L}_2}$ . Since this inequality holds for all  $t \geq 0$ , we can take the supremum on the left-hand side and obtain the result.  $\square$

For SISO (single-input single-output) systems, it even holds

$$\|G\|_{\mathcal{H}_2} = \sup_{\substack{u \in \mathcal{L}_2([0, \infty), \mathbb{R}^m) \\ u \neq 0}} \frac{\|y\|_{\mathcal{L}_\infty}}{\|u\|_{\mathcal{L}_2}},$$

i. e., the  $\mathcal{H}_2$ -norm is the  $\mathcal{L}_2$ - $\mathcal{L}_\infty$ -induced norm of the system. For general MIMO (multi-input multi-output) systems, the interpretation of the  $\mathcal{H}_2$ -norm is more involved. The  $\mathcal{H}_2$ -norm can be computed by using Plancherel's Theorem noticing that

$$\begin{aligned}
\|G\|_{\mathcal{H}_2} &= \left( \frac{1}{2\pi} \int_{-\infty}^{\infty} \|G(i\omega)\|_{\mathbb{F}}^2 d\omega \right)^{1/2} \\
&= \left( \int_0^{\infty} \|Ce^{At}B\|_{\mathbb{F}}^2 dt \right)^{1/2} \\
&= \left( \int_0^{\infty} \text{tr} \left( Ce^{At}BB^Te^{A^T t}C^T \right) dt \right)^{1/2} \\
&= \text{tr} \left( CPC^T \right)^{1/2},
\end{aligned}$$

where  $P$  is the controllability Gramian of the system. A similar expression can be obtained using the observability Gramian.

---

### 2.4.2 The Banach Space $\mathcal{H}_\infty^{p \times m}$

The space  $\mathcal{H}_\infty^{p \times m}$  is defined by

$$\mathcal{H}_\infty^{p \times m} := \left\{ G : \mathbb{C}^+ \rightarrow \mathbb{C}^{p \times m} : G \text{ is analytic in } \mathbb{C}^+ \text{ and } \sup_{\omega \in \mathbb{R}} \|G(i\omega)\|_2 < \infty \right\}.$$

Again, since every  $G \in \mathcal{H}_\infty^{p \times m}$  is analytic, there exists a unique continuation to the imaginary axis. The space  $\mathcal{H}_\infty^{p \times m}$  is a Banach space equipped with the  $\mathcal{H}_\infty$ -norm

$$\|G\|_{\mathcal{H}_\infty} := \sup_{\omega \in \mathbb{R}} \|G(i\omega)\|_2.$$

Again, we focus on *rational* functions, i. e., in functions that are in  $\mathcal{RH}_\infty^{p \times m} := \mathcal{H}_\infty^{p \times m} \cap \mathbb{R}(s)^{p \times m}$ . First we have the following.

**Lemma 2.19:** The following statements are equivalent:

- The function  $G$  is an element of  $\mathcal{RH}_\infty^{p \times m}$ .
- The function  $G$  is proper and all its poles are in  $\mathbb{C}^-$ .
- The function  $G$  can be realized by a system  $[A, B, C, D]$  with  $\Lambda(A) \subset \mathbb{C}^-$ .

Now we show that the  $\mathcal{H}_\infty$ -norm bounds the  $\mathcal{L}_2$ -norm of the output by the  $\mathcal{L}_2$ -norm of the input..

**Theorem 2.20:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  with a transfer function  $G \in \mathcal{RH}_\infty^{p \times m}$  be given. Then it holds

$$\|G\|_{\mathcal{H}_\infty} \geq \sup_{\substack{u \in \mathcal{L}_2([0, \infty), \mathbb{R}^m) \\ u \neq 0}} \frac{\|y\|_{\mathcal{L}_2}}{\|u\|_{\mathcal{L}_2}}.$$

*Proof.* It can be shown that an asymptotically stable system with an input  $u \in \mathcal{L}_2([0, \infty), \mathbb{R}^m)$  results in an output  $y \in \mathcal{L}_2([0, \infty), \mathbb{R}^p)$ . (This can be proven using Young's convolution inequality [Bog07, Theorem 3.9.4].)

With  $U(s) := \mathcal{L}\{u\}(s)$  and  $Y(s) := \mathcal{L}\{y\}(s)$  and using Plancherel's Theorem

we obtain

$$\begin{aligned}
\|y\|_{\mathcal{L}_2}^2 &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \|Y(i\omega)\|_2^2 d\omega \\
&= \frac{1}{2\pi} \int_{-\infty}^{\infty} \|G(i\omega)U(i\omega)\|_2^2 d\omega \\
&\leq \frac{1}{2\pi} \int_{-\infty}^{\infty} \|G(i\omega)\|_2^2 \|U(i\omega)\|_2^2 d\omega \\
&\leq \sup_{\omega \in \mathbb{R}} \|G(i\omega)\|_2^2 \cdot \frac{1}{2\pi} \int_{-\infty}^{\infty} \|U(i\omega)\|_2^2 d\omega \\
&= \sup_{\omega \in \mathbb{R}} \|G(i\omega)\|_2^2 \cdot \|u\|_{\mathcal{L}_2}^2.
\end{aligned}$$

□

It can also be shown that

$$\|G\|_{\mathcal{H}_\infty} = \sup_{\substack{u \in \mathcal{L}_2([0, \infty), \mathbb{R}^m) \\ u \neq 0}} \frac{\|y\|_{\mathcal{L}_2}}{\|u\|_{\mathcal{L}_2}},$$

i. e., the bound is tight. The proof of this is more lengthy, and therefore, it is omitted. There are also several algorithms for computing the  $\mathcal{H}_\infty$ -norm. The most established ones are based on an iteration on structured matrices or pencils. They are too involved to be discussed at this point.

---



## CHAPTER 3

---

### Eigenvalue-Based Approaches

---

Consider a linear system  $[A, B, C, D] \in \Sigma_{n,m,p}$  with transfer function  $G(s) \in \mathbb{R}(s)^{p \times m}$ . Assume that we have a partition of the system as

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad B = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}, \quad C = [C_1 \quad C_2],$$

where  $A_{ij} \in \mathbb{R}^{n_i \times n_j}$ ,  $B_i \in \mathbb{R}^{n_i \times m}$ , and  $C_i \in \mathbb{R}^{p \times n_i}$  for  $i, j = 1, 2$ . Then the system  $[A_{11}, B_1, C_1, D] \in \Sigma_{n_1,m,p}$  is called a *truncation* of the original system  $[A, B, C, D] \in \Sigma_{n,m,p}$ . Assume that it has the transfer function  $G_1(s) \in \mathbb{R}(s)^{p \times m}$ . The goal is to find a good truncation in the following sense:

- a) The state-space dimension  $n_1$  is small compared to  $n$ .
- b) The output  $y_1$  of  $[A_{11}, B_1, C_1, D]$  is similar to output  $y$  of  $[A, B, C, D]$  for the same input  $u$ , i. e.,  $\|y - y_1\|$  is small in some suitable norm. This norm can be often estimated using the norm of  $G(s) - G_1(s)$  such as the  $\mathcal{H}_2$ -norm or  $\mathcal{H}_\infty$ -norm.
- c) If the original model is asymptotically stable, then also the reduced one should be asymptotically stable. In particular, both transfer functions should be in  $\mathcal{RH}_\infty^{p \times m}$ .

It is important to note that without any further assumptions, nothing can be said about asymptotic stability, controllability, or observability of the reduced-order system, even if the original system is asymptotically stable, controllable, or observable.

**Example:** Consider the system  $[A, B, C, D] \in \Sigma_{2,1,1}$  with

$$A = \begin{bmatrix} 1 & \frac{5}{4} \\ -\frac{7}{4} & -2 \end{bmatrix}, \quad B = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad C = [0 \quad 1], \quad D = 5.$$

Then we have  $\Lambda(A) = \{-\frac{1}{4}, -\frac{3}{4}\}$ , i. e., the system is asymptotically stable. Moreover, we have

$$\begin{aligned} \text{rank} [B \quad AB] &= \text{rank} \begin{bmatrix} 0 & \frac{5}{4} \\ 1 & -2 \end{bmatrix} = 2, \\ \text{rank} \begin{bmatrix} C \\ CA \end{bmatrix} &= \text{rank} \begin{bmatrix} 0 & 1 \\ -\frac{7}{4} & -2 \end{bmatrix} = 2, \end{aligned}$$

this means that the system is controllable and observable.

Taking the truncation for  $n_1 = 1$ , we obtain the reduced-order model  $[1, 0, 0, 5] \in \Sigma_{1,1,1}$  which is unstable, uncontrollable, and unobservable.

Most often, good truncations are achieved by performing a state-space transformation (Note that this does not change the transfer function!). Let  $T = [T_1 \quad T_2] \in \mathbb{R}^{n \times n}$  be an invertible matrix with  $T_1 \in \mathbb{R}^{n \times n_1}$  and  $T_2 \in \mathbb{R}^{n \times n_2}$  be given and define  $T^{-1} := [W_1 \quad W_2]^T$  with  $W_1 \in \mathbb{R}^{n \times n_1}$  and  $W_2 \in \mathbb{R}^{n \times n_2}$ . Then we consider the transformed system  $[T^{-1}AT, T^{-1}B, CT, D] \in \Sigma_{n,m,p}$  and obtain the truncation (keeping the first  $n_1$  rows and columns) by setting  $[A_{11}, B_1, C_1, D] = [W_1^T AT_1, W_1^T B, CT_1, D] \in \Sigma_{n_1,m,p}$ .

Note that the above is “*model reduction by projection*”: We assume that the state  $x(\cdot)$  lives approximately in low-dimensional subspace  $\text{im } T_1$ . With  $x(t) \approx T_1 x_1(t)$  we obtain

$$\begin{aligned} T_1 \dot{x}_1(t) &\approx AT_1 x_1(t) + Bu(t), \\ y_1(t) &= CT_1 x_1(t) + Du(t). \end{aligned}$$

Next we “make the state equation square” again by imposing a *Petrov-Galerkin condition*

$$\text{im } W_1 \perp (T_1 \dot{x}_1(t) - (AT_1 x_1(t) + Bu(t))).$$

This results in

$$W_1^T T_1 \dot{x}_1(t) = W_1^T AT_1 x_1(t) + W_1^T Bu(t).$$

By choosing  $T_1$  and  $W_1$  bi-orthogonal, i. e.,  $W_1^T T_1 = I_{n_1}$ , we obtain an ODE as state equation. This bi-orthogonality is automatically fulfilled by the above construction of the truncation. It remains to choose good projection matrices  $T_1$  and  $W_1$ . This principle can also be generalized to non-linear systems.

### 3.1 Modal Truncation

In this chapter we discuss eigenvalue-based methods for model reduction. Assume that we have given a system  $[A, B, C, D] \in \Sigma_{n,m,p}$ . Assume that we have given a state-space transformation  $T \in \mathbb{R}^{n \times n}$  such that

$$T^{-1}AT = \begin{bmatrix} A_{11} & 0 \\ 0 & A_{22} \end{bmatrix}, \quad T^{-1}B = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}, \quad CT = \begin{bmatrix} C_1 \\ C_2 \end{bmatrix}. \quad (3.1)$$

Then for the transfer function we obtain

$$\begin{aligned} G(s) &= C(sI_n - A)^{-1}B + D = (C_1(sI_{n_1} - A_{11})^{-1}B_1 + D) \\ &\quad + (C_2(sI_{n_2} - A_{22})^{-1}B_2) =: G_1(s) + G_2(s). \end{aligned}$$

If we can determine the above decomposition such that  $n_1 \ll n$  and  $\|G_2\|$  is small, we get

$$\|G - G_1\| = \|G_2\|$$

and therefore,  $[A_{11}, B_1, C_1, D] \in \Sigma_{n_1,m,p}$  is a good reduced-order model. This process is called *modal truncation (or modal approximation, modal reduction)*. Here we discuss the computation of such reduced-order models.

**Theorem 3.1:** Assume that the system  $[A, B, C, D] \in \Sigma_{n,m,p}$  is asymptotically stable (controllable, stabilizable, observable, detectable). Then the reduced-order model  $[A_{11}, B_1, C_1, D] \in \Sigma_{n_1,m,p}$  in (3.1) is asymptotically stable (controllable, stabilizable, observable, detectable).

*Proof.* Exercise. □

**Theorem 3.2:** Let the system  $[A, B, C, D] \in \Sigma_{n,m,p}$  be asymptotically stable with transfer function  $G \in \mathcal{RH}_{\infty}^{p \times m}$  and assume that  $A$  is diagonalizable. Assume that there is an invertible matrix  $T \in \mathbb{C}^{n \times n}$  such that

$$T^{-1}AT = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}, \quad T^{-1}B = \begin{bmatrix} \widehat{b}_1^T \\ \vdots \\ \widehat{b}_n^T \end{bmatrix}, \quad CT = [\widehat{c}_1 \quad \dots \quad \widehat{c}_n]$$

and let the reduced-order model be  $[\widetilde{A}, \widetilde{B}, \widetilde{C}, \widetilde{D}] \in \Sigma_{r,m,p}$  with

$$\widetilde{A} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_r \end{bmatrix}, \quad \widetilde{B} = \begin{bmatrix} \widehat{b}_1^T \\ \vdots \\ \widehat{b}_r^T \end{bmatrix}, \quad \widetilde{C} = [\widehat{c}_1 \quad \dots \quad \widehat{c}_r], \quad \widetilde{D} = D,$$

and with the transfer function  $\tilde{G} \in \mathcal{RH}_\infty^{p \times m}$ . Then it holds

$$\|G - \tilde{G}\|_{\mathcal{H}_\infty} \leq \sum_{j=r+1}^n \frac{\|\hat{c}_j\|_2 \cdot \|\hat{b}_j\|_2}{|\operatorname{Re}(\lambda_j)|}.$$

*Proof.* We have

$$G(s) = \sum_{j=1}^n \frac{1}{s - \lambda_j} \hat{c}_j \hat{b}_j^\top + D$$

and

$$\tilde{G}(s) = \sum_{j=1}^r \frac{1}{s - \lambda_j} \hat{c}_j \hat{b}_j^\top + D.$$

Therefore, we have

$$\|G - \tilde{G}\|_{\mathcal{H}_\infty} = \left\| \sum_{j=r+1}^n \frac{1}{\cdot - \lambda_j} \hat{c}_j \hat{b}_j^\top \right\|_{\mathcal{H}_\infty} \leq \sum_{j=r+1}^n \left\| \frac{1}{\cdot - \lambda_j} \hat{c}_j \hat{b}_j^\top \right\|_{\mathcal{H}_\infty}.$$

Moreover, it holds

$$\begin{aligned} \left\| \frac{1}{\cdot - \lambda_j} \hat{c}_j \hat{b}_j^\top \right\|_{\mathcal{H}_\infty} &= \sup_{\omega \in \mathbb{R}} \left\| \frac{1}{i\omega - \lambda_j} \hat{c}_j \hat{b}_j^\top \right\|_2 \\ &= \|\hat{c}_j \hat{b}_j^\top\|_2 \cdot \sup_{\omega \in \mathbb{R}} \left| \frac{1}{i\omega - \lambda_j} \right| \\ &= \|\hat{c}_j\|_2 \cdot \|\hat{b}_j\|_2 \cdot \left| \frac{1}{\operatorname{Re}(\lambda_j)} \right|, \end{aligned}$$

where the latter equality follows from the fact that  $i\omega - \lambda_j$  is minimized for  $\omega = \operatorname{Im}(\lambda_j)$ .  $\square$

**Remark 3.3:** a) In classical modal truncation, the eigenvalues are ordered with respect to distance to the imaginary axis, i. e.,

$$0 > \operatorname{Re}(\lambda_1) \geq \operatorname{Re}(\lambda_2) \geq \dots \geq \operatorname{Re}(\lambda_n).$$

There are good numerical algorithms for approximating eigenvalues closest to the imaginary axis. However, the error bound suggests to order the eigenvalues such that

$$\frac{\|\hat{c}_1\|_2 \cdot \|\hat{b}_1\|_2}{|\operatorname{Re}(\lambda_1)|} \geq \frac{\|\hat{c}_2\|_2 \cdot \|\hat{b}_2\|_2}{|\operatorname{Re}(\lambda_2)|} \geq \dots \geq \frac{\|\hat{c}_n\|_2 \cdot \|\hat{b}_n\|_2}{|\operatorname{Re}(\lambda_n)|}.$$

There are also algorithms that handle this sorting of the eigenvalues (see Section 3.2).

- b) Modal truncation generates good local approximations of the transfer function. This means that the reduced-order model has a good approximation quality near those values on the imaginary axis that are close to some  $\lambda \in \Lambda(A_{11})$  and can have a worse approximation quality near an eigenvalue  $\lambda \in \Lambda(A)$ , if  $\lambda$  is close to the imaginary axis, but  $\lambda \notin \Lambda(A_{11})$ .
- c) There are problems if  $A$  is non-diagonalizable or if  $T$  is ill-conditioned, i. e., the condition number  $\kappa(T) := \|T\|_2 \cdot \|T^{-1}\|_2$  is large. Then  $\|\hat{b}_j\|_2$  and  $\|\hat{c}_j\|_2$  can be large, even if  $B$  and  $C$  are of moderate norm. In this case, the state-space dimension of the reduced-order model often has to be increased to achieve a good approximation error.
- d) The transformation matrix  $T$  can be chosen to be real by treating complex conjugate eigenvalues as pairs. This results in a real reduced-order model.

### 3.2 The Dominant Pole Algorithm

As mentioned above it is desirable to order the eigenvalues such that

$$\frac{\|\hat{c}_1\|_2 \cdot \|\hat{b}_1\|_2}{|\operatorname{Re}(\lambda_1)|} \geq \frac{\|\hat{c}_2\|_2 \cdot \|\hat{b}_2\|_2}{|\operatorname{Re}(\lambda_2)|} \geq \dots \geq \frac{\|\hat{c}_n\|_2 \cdot \|\hat{b}_n\|_2}{|\operatorname{Re}(\lambda_n)|}.$$

The *dominant pole algorithm* that we will discuss now is doing exactly this. First we show that the vectors  $\hat{b}_j$  and  $\hat{c}_j$  have a special structure.

**Lemma 3.4:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  be an asymptotically stable system with transfer function  $G \in \mathcal{RH}_\infty^{p \times m}$ . Assume that  $A$  is diagonalizable. Then it holds

$$G(s) = \sum_{j=1}^n \frac{R_j}{s - \lambda_j} + D$$

with the *residues*  $R_j = (Cx_j)(v_j^H B)$ , where  $x_j, v_j \in \mathbb{C}^n$  denote the right and left eigenvectors of  $A$  associated with the eigenvalue  $\lambda_j$  for  $j = 1, \dots, n$ . Moreover, here we assume the normalization condition  $v_j^H x_j = 1$  for  $j = 1, \dots, n$ .

*Proof.* Let  $T \in \mathbb{C}^{n \times n}$  be such that  $T^{-1}AT = \operatorname{diag}(\lambda_1, \dots, \lambda_n)$ . Then we have that

$$T = [x_1 \quad \dots \quad x_n], \quad T^{-1} = \begin{bmatrix} v_1^H \\ \vdots \\ v_n^H \end{bmatrix}.$$

Using the notation of Theorem 3.2, we obtain

$$\begin{bmatrix} \hat{b}_1 \\ \vdots \\ \hat{b}_n \end{bmatrix} := T^{-1}B = \begin{bmatrix} v_1^H B \\ \vdots \\ v_n^H B \end{bmatrix}, \quad [\hat{c}_1 \ \dots \ \hat{c}_n] := CT = [Cx_1 \ \dots \ Cx_n],$$

which gives the result.  $\square$

We now derive the *dominant pole algorithm* for SISO systems [RM06b]. The case of MIMO systems is conceptually slightly different, see also Remark 3.5. So, assume that we have given an asymptotically stable system  $[A, b, c^T, 0] \in \Sigma_{n,1,1}$  with the transfer function  $G(s) \in \mathbb{R}(s)$ . Here we set the feedthrough term  $D = 0$  to simplify the presentation, but it is no problem to include it as well. Then we have

$$G(s) = \frac{Y(s)}{U(s)} = c^T (sI_n - A)^{-1} b,$$

$$G^H(s) = \frac{\bar{Y}(s)}{\bar{U}(s)} = b^T (sI_n - A)^{-H} c,$$

where  $U(s)$  and  $Y(s)$  are the Laplace transforms of  $u$  and  $y$ , respectively. This can be reformulated

$$\begin{bmatrix} sI_n - A & -b \\ c^T & 0 \end{bmatrix} \begin{bmatrix} X(s) \\ U(s) \end{bmatrix} = \begin{bmatrix} 0 \\ Y(s) \end{bmatrix}, \quad (3.2)$$

$$\begin{bmatrix} (sI_n - A)^H & c \\ -b^T & 0 \end{bmatrix} \begin{bmatrix} V(s) \\ \bar{U}(s) \end{bmatrix} = \begin{bmatrix} 0 \\ \bar{Y}(s) \end{bmatrix},$$

with auxiliary vectors  $X(s)$  and  $V(s)$  (where  $X(s)$  is the Laplace transform of the state of  $[A, b, c^T, 0] \in \Sigma_{n,1,1}$ ). If  $\lambda \in \mathbb{C}$  is a pole of  $G(s)$ , then  $\lim_{s \rightarrow \lambda} |G(s)| = \infty$  and one can choose  $\lim_{s \rightarrow \lambda} U(s) = 0$ , while  $Y(s) \equiv 1$ . This yields that  $\lim_{s \rightarrow \lambda} X(s) = x$  and  $\lim_{s \rightarrow \lambda} V(s) = v$  are right and left eigenvectors of  $A$  associated with the eigenvalue  $\lambda$  and the normalization conditions  $c^T x = 1$  and  $-b^T v = 1$ .

We want to determine the most *dominant poles* of  $G(s)$ , i. e., those  $\lambda_j \in \Lambda(A)$ , where  $|R_j|/|\operatorname{Re}(\lambda_j)|$  is the largest. We do this iteratively in a search in possibly growing subspaces. Assume that we have subspaces spanned by  $\hat{X} \in \mathbb{C}^{n \times k}$  and  $\hat{V} \in \mathbb{C}^{n \times k}$  for some  $k \ll n$ . Then we can project the eigenvalue problem for the matrix pencil  $sI_n - A \in \mathbb{R}[s]^{n \times n}$  to a small eigenvalue problem for the matrix pencil

$$s\hat{V}^H \hat{X} - \hat{V}^H A \hat{X} \in \mathbb{C}[s]^{k \times k}.$$

Assume that this pencil has only semi-simple eigenvalues and that  $\hat{V}^H \hat{X}$  is invertible. Then one could alternatively consider the eigenvalue problem for the projected matrix  $(\hat{V}^H \hat{X})^{-1} \hat{V}^H A \hat{X} \in \mathbb{C}^{k \times k}$ . For this matrix pencil, we can easily

determine all eigenvalues  $\tilde{\lambda}_j \in \mathbb{C}$  and the associated right and left eigenvectors  $\tilde{x}_j \in \mathbb{C}^k$  and  $\tilde{v}_j \in \mathbb{C}^k$  for  $j = 1, \dots, k$ . Then we obtain the eigenvalue and eigenvector approximations for the original problem as

$$\hat{\lambda}_j = \tilde{\lambda}_j, \quad \hat{x}_j = \hat{X}\tilde{x}_j, \quad \hat{v}_j = \hat{V}\tilde{v}_j, \quad j = 1, \dots, k.$$

These approximations can now be sorted according to our dominance measure, i. e., we sort the eigenvalues such that

$$\frac{|c^\top \hat{x}_1 \hat{v}_1^\mathsf{H} b|}{|\operatorname{Re}(\hat{\lambda}_1)|} \geq \frac{|c^\top \hat{x}_2 \hat{v}_2^\mathsf{H} b|}{|\operatorname{Re}(\hat{\lambda}_2)|} \geq \dots \geq \frac{|c^\top \hat{x}_k \hat{v}_k^\mathsf{H} b|}{|\operatorname{Re}(\hat{\lambda}_k)|}.$$

So  $\hat{\lambda}_1$  is our current approximation for the most dominant pole. If  $\|A\hat{x}_1 - \hat{\lambda}_1\hat{x}_1\|_2 < \varepsilon$  (or  $\|\hat{v}_1^\mathsf{H}A - \hat{\lambda}_1\hat{v}_1^\mathsf{H}\|_2 < \varepsilon$ ) for some small tolerance  $\varepsilon > 0$ , then we assume that the eigenvalue  $\hat{\lambda}_1$  and the corresponding eigenvectors have converged.

If this is not the case, we expand the matrices  $\hat{X}$  and  $\hat{V}$  in order to enrich the spaces  $\operatorname{im} \hat{X}$  and  $\operatorname{im} \hat{V}$  in which we search for the eigenvectors. This is done by plugging in our current dominant pole approximation  $\hat{\lambda}_1$  into (3.2) and compute  $\hat{x} := X(\hat{\lambda}_1)$  and  $\hat{v} := V(\hat{\lambda}_1)$  (with  $Y(s) := 1$ ). Then the expanded projection matrices are  $\begin{bmatrix} \hat{X} & \hat{x} \end{bmatrix}$  and  $\begin{bmatrix} \hat{V} & \hat{v} \end{bmatrix}$ . For numerical stability, it is advised to orthogonalize their columns afterwards.

On the other hand, if  $\hat{\lambda}_1$  has converged to an eigenvalue  $\lambda_1$  with right and left eigenvectors  $x_1, v_1 \in \mathbb{C}^n$ , then we want to ensure that we do not find it again in the next iterations. So we want to deflate this eigenvalue and its eigenvectors. This is done by projecting the system  $[A, b, c^\top, 0]$ , namely we replace it by  $[A, \tilde{b}, \tilde{c}^\top, 0]$  with

$$\tilde{b} := \left( I_n - \frac{x_1 v_1^\mathsf{H}}{v_1^\mathsf{H} x_1} \right) b, \quad \tilde{c}^\top := c^\top \left( I_n - \frac{x_1 v_1^\mathsf{H}}{v_1^\mathsf{H} x_1} \right).$$

First of all, note that the matrix  $I_n - \frac{x_1 v_1^\mathsf{H}}{v_1^\mathsf{H} x_1}$  is a projector. Projecting the system like this has the effect that the residue of the deflated eigenvalue is zero, since  $\tilde{c}^\top x_1 = 0$  and  $v_1^\mathsf{H} \tilde{b} = 0$  and the residues of the other eigenvalues remain unchanged (exercise!). Therefore, the already converged eigenvalues are not found again, since their dominance value is set to zero. To summarize this section we formulate the above results as Algorithm 3.1.

**Algorithm 3.1** Dominant pole algorithm

**Input:** Asymptotically stable system  $[A, b, c^\top, 0] \in \Sigma_{n,1,1}$  with transfer function  $G(s) \in \mathbb{R}(s)$ ; an initial pole estimate  $\lambda \in \mathbb{C}$ , tolerance  $\varepsilon > 0$ , number of desired dominant poles  $k$ .

**Output:**  $k$  dominant poles  $\Lambda = \{\lambda_1, \dots, \lambda_k\} \subset \mathbb{C}$  of  $G(s)$  with the associated right and left eigenvectors of  $A$  stored in  $R, L \in \mathbb{C}^{n \times k}$ .

- 1: Initialize  $k_{\text{found}} := 0$ ,  $\Lambda = \{\}$ ,  $R = [ ]$ ,  $L = [ ]$ ,  $\hat{X} = [ ]$ ,  $\hat{V} = [ ]$ .
- 2: **while**  $k > k_{\text{found}}$  **do**
- 3:   Solve the linear system

$$\begin{bmatrix} \lambda I_n - A & -b \\ c^\top & 0 \end{bmatrix} \begin{bmatrix} \hat{x} \\ u \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

for  $\hat{x} \in \mathbb{C}^n$ .

- 4:   Solve the linear system

$$\begin{bmatrix} (\lambda I_n - A)^\text{H} & c \\ -b^\top & 0 \end{bmatrix} \begin{bmatrix} \hat{v} \\ u \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

for  $\hat{v} \in \mathbb{C}^n$ .

- 5:   Expand the search spaces: Set  $\hat{X} := [\hat{X} \ \hat{x}]$  and  $\hat{V} := [\hat{V} \ \hat{v}]$  and orthogonalize.
- 6:   Compute the eigenvalues and eigenvectors of the matrix pencil  $s\hat{V}^\text{H}\hat{X} - \hat{V}^\text{H}A\hat{X} \in \mathbb{C}[s]^{\ell \times \ell}$  and compute eigenvalue and eigenvector approximations, sort them according to the dominance measure and store them as  $\hat{\Lambda} := \{\hat{\lambda}_1, \dots, \hat{\lambda}_\ell\}$ ,  $\hat{X} := [\hat{x}_1 \ \dots \ \hat{x}_\ell]$ ,  $\hat{V} := [\hat{v}_1 \ \dots \ \hat{v}_\ell]$ .
- 7:   **while**  $\|A\hat{x}_1 - \hat{\lambda}_1\hat{x}_1\| < \varepsilon$  **do**
- 8:     Deflate the found eigenvalue: Set

$$\begin{aligned} k_{\text{found}} &:= k_{\text{found}} + 1, & \lambda_{k_{\text{found}}} &:= \hat{\lambda}_1, \\ \Lambda &:= \Lambda \cup \{\lambda_{k_{\text{found}}}\}, & R &:= [R \ \hat{x}_1], & L &:= [L \ \hat{v}_1], \\ b &:= \left( I_n - \frac{\hat{x}_1\hat{v}_1^\text{H}}{\hat{v}_1^\text{H}\hat{x}_1} \right) b, & c^\top &:= c^\top \left( I_n - \frac{\hat{x}_1\hat{v}_1^\text{H}}{\hat{v}_1^\text{H}\hat{x}_1} \right). \end{aligned}$$

- 9:     Set

$$\begin{aligned} \hat{\Lambda} &:= \{\hat{\lambda}_2, \dots, \hat{\lambda}_\ell\} =: \{\hat{\lambda}_1, \dots, \hat{\lambda}_{\ell-1}\}, \\ \hat{X} &:= [\hat{x}_2 \ \dots \ \hat{x}_\ell] =: [\hat{x}_1 \ \dots \ \hat{x}_{\ell-1}], \\ \hat{V} &:= [\hat{v}_2 \ \dots \ \hat{v}_\ell] =: [\hat{v}_1 \ \dots \ \hat{v}_{\ell-1}]. \end{aligned}$$

- 10:   **end while**
- 11:   Set the new pole estimate  $\lambda = \hat{\lambda}_1$ .
- 12: **end while**

**Remark 3.5:** a) The dominant pole algorithm presented here is a subspace accelerated version of an algorithm that was originally designed as a Newton method to find roots of  $G^{-1}(s)$ .

b) It is not guaranteed that the method finds the most dominant poles, but it often works well in practice (in particular, if there are only a few very dominant poles). Convergence of poles can be enhanced by using a Newton scheme or a Rayleigh quotient iteration to update the pole estimates.

c) The projections in (3.2) should not be constructed explicitly. It is rather advised to compute the action of the projection on a vector  $z \in \mathbb{C}^n$  if needed. This means that we compute

$$\left( I_n - \frac{xv^H}{v^Hx} \right) z = z - \frac{v^H z}{v^H x} \cdot x$$

using two inner products and one scaled vector addition.

d) The algorithm can be modified to deal with MIMO systems [RM06a]. The most drastic changes are in lines 3 and 4 of Algorithm 3.1, where we replace the linear systems by

$$(\lambda I_n - A)\hat{x} = Bu, \quad (\lambda I_n - A)^H \hat{v} = Cw,$$

where  $u \in \mathbb{C}^m$  and  $w \in \mathbb{C}^p$  are chosen to be the right and left singular vectors of  $G(\lambda)$  corresponding to its largest singular value.

e) The algorithm can also be modified to output real  $R$  and  $L$  in order to obtain a real reduced-order model. For this, pairs of complex conjugate eigenvalues have to be deflated together.



---

## Balancing-Based Approaches

---

In this chapter we discuss another kind of transformation that simultaneously transforms the controllability and observability Gramians to diagonal form. Then we can sort the transformed states according to their input or output energy and truncate those which are hard to control or hard to observe.

### 4.1 Input and Output Energy

Consider a system  $[A, B, C, D] \in \Sigma_{n,m,p}$ . Here we consider the system for  $t \in \mathbb{R}$  and assume that  $x(-\infty) = 0$ . Assume that we have an input  $u \in \mathcal{L}_2((-\infty, 0], \mathbb{R}^m)$  steering the state to  $x(0) = x_0 \in \mathbb{R}^n$ . Then

$$E_u := \left( \int_{-\infty}^0 \|u(\tau)\|_2^2 d\tau \right)^{1/2} = \|u\|_{\mathcal{L}_2((-\infty, 0], \mathbb{R}^m)}$$

is called the *input energy* and if  $y \in \mathcal{L}_2([0, \infty), \mathbb{R}^p)$ , then

$$E_y := \left( \int_0^{\infty} \|y(\tau)\|_2^2 d\tau \right)^{1/2} = \|y\|_{\mathcal{L}_2([0, \infty), \mathbb{R}^p)}$$

is called the *output energy*. In many applications these can be interpreted as actual physical energies of the system.

For the initial state  $x(0) = x_0 \in \mathbb{R}^n$  we define

$$E_u(x_0) := \inf_{\substack{u \in \mathcal{L}_2((-\infty, 0], \mathbb{R}^m) \\ x(-\infty)=0, x(0)=x_0}} \|u\|_{\mathcal{L}_2((-\infty, 0], \mathbb{R}^m)}, \quad (4.1)$$

which is the minimal energy needed to steer the system from the state zero to the state  $x_0$  in an arbitrary time. If  $E_u(x_0)$  is small, then the state  $x_0$  is easy to reach, otherwise it is hard to reach. Note that  $E_u(x_0) = \infty$  is also possible. Then the state  $x_0$  is unreachable and the system is uncontrollable.

Now assume that  $x(0) = x_0$  and that  $u|_{[0,\infty)} = 0$ . Then we have  $y(t) = Ce^{At}x_0$ . We define

$$E_y(x_0) := \|y\|_{\mathcal{L}_2([0,\infty),\mathbb{R}^p)} = \|Ce^{A\cdot}x_0\|_{\mathcal{L}_2([0,\infty),\mathbb{R}^p)},$$

which is the output energy gained from the state  $x_0$ . If  $E_y(x_0)$  is large, then  $x_0$  is easy to observe, otherwise it is hard to observe. If  $E_y(x_0) = 0$ , the the state  $x_0$  is unobservable, and therefore the system is unobservable.

The next theorem shows that the  $E_u(x_0)$  and  $E_y(x_0)$  can be expressed by the controllability and observability Gramians, respectively, which make them feasible for numerical computations.

**Theorem 4.1:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  be asymptotically stable and controllable. Then the following two statements are satisfied:

a) It holds

$$E_u(x_0)^2 = x_0^\top P^{-1}x_0,$$

where  $P$  is the controllability Gramian of the system. Moreover,  $u_*(t) := B^\top e^{-A^\top t} P^{-1}x_0$  is a trajectory for which the infimum in (4.1) is attained.

b) It holds

$$E_y(x_0)^2 = x_0^\top Qx_0,$$

where  $Q$  is the observability Gramian of the system.

*Proof.* a) Let  $(x, u)$  be an arbitrary solution trajectory with  $x(-\infty) = 0$ ,  $x(0) = x_0$ , and  $E_u < \infty$ . Then we have

$$x(0) = \int_{-\infty}^0 e^{-A\tau} Bu(\tau) d\tau.$$

We show that  $E_u \geq E_{u_*}$  for the above defined  $u_*$ . Define  $v := u - u_*$ . Then we have

$$\begin{aligned} \int_{-\infty}^0 u_*(\tau)^\top v(\tau) d\tau &= x_0^\top P^{-1} \left( \int_{-\infty}^0 e^{-A\tau} Bu(\tau) d\tau \right. \\ &\quad \left. - \underbrace{\int_{-\infty}^0 e^{-A\tau} BB^\top e^{-A^\top \tau} d\tau}_{=P} P^{-1}x_0 \right) \\ &= x_0^\top P^{-1}(x_0 - x_0) = 0. \end{aligned}$$

Hence we obtain

$$\begin{aligned}
E_u^2 &= \int_{-\infty}^0 u(\tau)^\top u(\tau) d\tau \\
&= \int_{-\infty}^0 (v(\tau) + u_*(\tau))^\top (v(\tau) + u_*(\tau)) d\tau \\
&= \underbrace{\int_{-\infty}^0 v(\tau)^\top v(\tau) d\tau}_{\geq 0} + 2 \underbrace{\int_{-\infty}^0 u_*(\tau)^\top v(\tau) d\tau}_{=0} + \underbrace{\int_{-\infty}^0 u_*(\tau)^\top u_*(\tau) d\tau}_{\geq 0} \\
&\geq E_{u_*}^2.
\end{aligned}$$

Moreover, we have

$$\begin{aligned}
E_{u_*}^2 &= x_0^\top P^{-1} \int_{-\infty}^0 e^{-A\tau} B B^\top e^{-A^\top \tau} d\tau P^{-1} x_0 = x_0^\top P^{-1} P P^{-1} x_0 \\
&= x_0^\top P^{-1} x_0.
\end{aligned}$$

b) We have

$$E_y(x_0)^2 = \int_0^\infty y(\tau)^\top y(\tau) d\tau = x_0^\top \int_0^\infty e^{A^\top \tau} C^\top C e^{A\tau} d\tau x_0 = x_0^\top Q x_0.$$

This completes the proof. □

Now consider an eigendecomposition of  $P$ , i. e.,  $P = U\Sigma U^\top$  with orthogonal  $U = [u_1 \ \dots \ u_n]$  and  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ , where  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ . Then the energy needed to reach the state  $x_0 = u_i$  from  $x(-\infty) = 0$  is given by  $E_u(u_i)^2 = u_i^\top P^{-1} u_i = 1/\sigma_i$ . Thus eigenvectors of  $P$  corresponding to large eigenvalues are easy to reach and eigenvectors of  $P$  corresponding to small eigenvalues are hard to reach. The eigenvectors corresponding to zero eigenvalues are unreachable. Analogously, the eigenvectors corresponding to large eigenvalues of  $Q$  are easy to observe, the ones corresponding to small eigenvalues are hard to observe and those corresponding to zero eigenvalues are unobservable.

## 4.2 Balancing Transformations and Balanced Truncation

We motivate the concept of balancing transformations by means of an example. The application of these transformations then leads to the method of balanced truncation that was discussed first in [Moo81].

---

**Example:** Consider the parameter-dependent system  $[A(\alpha), B(\alpha), C(\alpha), D] \in \Sigma_{2,1,1}$  with

$$A(\alpha) = \begin{bmatrix} -1 & -\frac{4}{\alpha} \\ 4\alpha & -2 \end{bmatrix}, \quad B(\alpha) = \begin{bmatrix} 1 \\ 2\alpha \end{bmatrix}, \quad C(\alpha) = \begin{bmatrix} -1 & \frac{2}{\alpha} \end{bmatrix}, \quad D = 0$$

for  $\alpha > 0$ . This system is asymptotically stable, controllable, and observable. Thus the controllability and observability Gramians  $P(\alpha)$  and  $Q(\alpha)$  are symmetric positive definite, where

$$P(\alpha) = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \alpha^2 \end{bmatrix}, \quad Q(\alpha) = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{\alpha^2} \end{bmatrix}.$$

The eigenvectors of  $P$  and  $Q$  are  $e_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$  and  $e_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ .

Assume that the state function  $x(\cdot) = \begin{bmatrix} x_1(\cdot) \\ x_2(\cdot) \end{bmatrix}$  is expressed as a (time-dependent) linear combination of the eigenvectors of  $P(\alpha)$ , in our case we get

$$x(t) = \beta_1(\alpha, t)e_1 + \beta_2(\alpha, t)e_2.$$

Intuitively, if a state  $e_i$  is hard to reach, then its coefficient  $\beta_i$  is negligible, so truncating it should not change the system's dynamics drastically. In our example we have two cases:

- a)  $\alpha \ll 1$ : In this case,  $e_2$  is much harder to reach than  $e_1$ . Thus we truncate  $x_2$  and obtain the reduced-order model  $[-1, 1, -1, 0] \in \Sigma_{1,1,1}$ .
- b)  $\alpha \gg 1$ : In this case,  $e_1$  is much harder to reach than  $e_2$ . Thus we truncate  $x_1$  and obtain the reduced-order model  $[-2, 2\alpha, \frac{2}{\alpha}, 0] \in \Sigma_{1,1,1}$ .

Alternatively, we could express  $x(\cdot)$  as a (time-dependent) linear combination of the eigenvectors of  $Q(\alpha)$ , in our case we get

$$x(t) = \gamma_1(\alpha, t)e_1 + \gamma_2(\alpha, t)e_2.$$

Similarly as above, if a state  $e_i$  is hard to observe, then its coefficient  $\gamma_i$  is negligible, so truncating it should not change the system's dynamics too much. Again we have two cases in our example:

- a)  $\alpha \ll 1$ : In this case,  $e_1$  is much harder to observe than  $e_2$ . Thus we truncate  $x_1$  and obtain the reduced-order model  $[-2, 2\alpha, \frac{2}{\alpha}, 0] \in \Sigma_{1,1,1}$ .
- b)  $\alpha \gg 1$ : In this case,  $e_2$  is much harder to observe than  $e_1$ . Thus we truncate  $x_2$  and obtain the reduced-order model  $[-1, 1, -1, 0] \in \Sigma_{1,1,1}$ .

**Remark 4.2:** a) Both approaches lead to different reduced-order models. In general, this would be OK, but in the above example this leads to contradictory reduced-order models.

b) The behavior in the example can be explained as follows: If  $\alpha \ll 1$ , then  $e_2$  is very hard to reach, but at the same time it is also very easy to observe and thus has a considerable influence on the output function.

c) The transfer function of the system is  $G(s) = \frac{3s+8}{s^2+3s+18}$ , but the reduced-order model depends on  $\alpha$ , which should not be the case.

The solution of the above problems is to truncate states that are simultaneously hard to reach *and* hard to observe. In general, finding these states is difficult, but it is easy if  $P = Q$ .

**Definition 4.3:** An asymptotically stable system  $[A, B, C, D] \in \Sigma_{n,m,p}$  with controllability Gramian  $P$  and observability Gramian  $Q$  is called *balanced*, if  $P = Q = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$ .

If a system is not balanced, then we can find a state-space transformation that balances the system. Before, we have to check how state-space transformations affect the Gramians.

**Lemma 4.4:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  be asymptotically stable. Let  $T \in \mathbb{R}^{n \times n}$  be invertible and define  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] := [T^{-1}AT, T^{-1}B, CT, D]$ . Then

a)  $P$  is the controllability Gramian of  $[A, B, C, D]$ , if and only if  $\tilde{P} := T^{-1}PT^{-T}$  is the controllability Gramian of  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}]$ ;

b)  $Q$  is the observability Gramian of  $[A, B, C, D]$ , if and only if  $\tilde{Q} := T^TQT$  is the observability Gramian of  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}]$ .

*Proof.* Exercise. □

Now we show how to balance a system using so-called *balancing transformations*.

**Theorem 4.5:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  be asymptotically stable, controllable, and observable. Then there exists an invertible matrix  $T \in \mathbb{R}^{n \times n}$  such that  $[T^{-1}AT, T^{-1}B, CT, D]$  is balanced.

*Proof.* By assumption, for the Gramians  $P$  and  $Q$  we have  $P > 0$  and  $Q > 0$ . Thus, there exist Cholesky decompositions  $P = RR^T$  and  $Q = LL^T$ ,

---

where  $R$  and  $L$  are lower triangular and invertible. Now consider the singular value decomposition  $L^T R = U \Sigma V^T$  with orthogonal  $U, V \in \mathbb{R}^{n \times n}$  and  $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$ ,  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$ . Since  $L$  and  $R$  are invertible, so is  $L^T R$  and therefore, we have  $\sigma_n > 0$ .

Set  $T := RV\Sigma^{-\frac{1}{2}}$ . Since  $I_n = \Sigma^{-\frac{1}{2}}U^T L^T R V \Sigma^{-\frac{1}{2}}$  we find  $T^{-1} = \Sigma^{-\frac{1}{2}}U^T L^T$ . For the controllability Gramian  $\tilde{P}$  of the transformed system we have

$$\begin{aligned}\tilde{P} &= T^{-1} P T^{-T} \\ &= \Sigma^{-\frac{1}{2}} U^T L^T R R^T L U \Sigma^{-\frac{1}{2}} \\ &= \Sigma^{-\frac{1}{2}} U^T U \Sigma V^T V \Sigma U^T U \Sigma^{-\frac{1}{2}} = \Sigma.\end{aligned}$$

Analogously, for the transformed observability Gramian  $\tilde{Q}$  we obtain

$$\begin{aligned}\tilde{Q} &= T^T Q T \\ &= \Sigma^{-\frac{1}{2}} V^T R^T L L^T R V \Sigma^{-\frac{1}{2}} \\ &= \Sigma^{-\frac{1}{2}} V^T V \Sigma U^T U \Sigma V^T V \Sigma^{-\frac{1}{2}} = \Sigma = \tilde{P}.\end{aligned}$$

i. e., the transformed system is balanced. □

**Example:** We revisit the above example. We have

$$R = \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 \\ 0 & \alpha \end{bmatrix}, \quad L = \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 \\ 0 & \frac{1}{\alpha} \end{bmatrix}.$$

Moreover, it holds

$$L^T R = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

The balancing transformation is given by

$$T = RV\Sigma^{-\frac{1}{2}} = \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 \\ 0 & \alpha \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \sqrt{2} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ \alpha & 0 \end{bmatrix}.$$

Now the balanced system  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] \in \Sigma_{2,1,1}$  is given by

$$\begin{aligned}\tilde{A} &= \begin{bmatrix} 0 & \frac{1}{\alpha} \\ 1 & 0 \end{bmatrix} \begin{bmatrix} -1 & -\frac{4}{\alpha} \\ 4\alpha & -2 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ \alpha & 0 \end{bmatrix} = \begin{bmatrix} -2 & 4 \\ -4 & -1 \end{bmatrix}, \\ \tilde{B} &= \begin{bmatrix} 0 & \frac{1}{\alpha} \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 2\alpha \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \\ \tilde{C} &= \begin{bmatrix} -1 & \frac{2}{\alpha} \\ \alpha & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ \alpha & 0 \end{bmatrix} = \begin{bmatrix} 2 & -1 \end{bmatrix}, \\ \tilde{D} &= 0.\end{aligned}$$

The transformed system does not depend on  $\alpha$ . The eigenvector  $e_2$  of  $\tilde{P} = \tilde{Q}$  ( $e_1$  in the old coordinates) is harder to reach and harder to observe than  $e_1$ .

This leads to Algorithm 4.1 for model reduction that is called *balanced truncation*.

---

**Algorithm 4.1** Balanced truncation (basic version)

---

**Input:** Asymptotically stable and minimal system  $[A, B, C, D] \in \Sigma_{n,m,p}$ , desired reduced order  $r$ .

**Output:** Reduced-order model  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,p}$ .

- 1: Solve the Lyapunov equations

$$AP + PA^T = -BB^T, \quad A^TQ + QA = -C^TC$$

for  $P > 0$  and  $Q > 0$ .

- 2: Compute Cholesky factorization  $P = RR^T$  and  $Q = LL^T$ .
- 3: Compute the singular value decomposition  $L^TR = U\Sigma V^T$ .
- 4: Set  $T := RV\Sigma^{-\frac{1}{2}}$  (and  $T^{-1} = \Sigma^{-\frac{1}{2}}U^TL^T$ ).
- 5: Do the balancing transformation

$$[T^{-1}AT, T^{-1}B, CT, D] = \left[ \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}, [C_1 \ C_2], D \right]$$

and set the reduced-order model as  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,p}$ .

---

### 4.3 Hankel Operator and Hankel Singular Values

In this section we want to discuss the foundation for the analysis of the balanced truncation algorithm introduced above. For this, we need the Hankel operator and the Hankel singular values [Ant05, Sec. 5.4]. Consider the state equation  $\dot{x}(t) = Ax(t) + Bu(t)$  with  $x(-\infty) = 0$  and an input  $u \in \mathcal{L}_2((-\infty, 0], \mathbb{R}^m)$  that acts on the negative time-horizon leading to  $x(0) = x_0$ . By switching off the input at  $t = 0$ , the output equation  $y(t) = Cx(t) + Du(t)$  gives an output signal  $y \in \mathcal{L}_2([0, \infty), \mathbb{R}^p)$  on the positive time-horizon. This defines an operator

$$\mathcal{H} : \mathcal{L}_2((-\infty, 0], \mathbb{R}^m) \rightarrow \mathcal{L}_2([0, \infty), \mathbb{R}^p), \quad u \mapsto y,$$

which is called the *Hankel operator* of the system  $[A, B, C, D] \in \Sigma_{n,m,p}$ . We have

$$x_0 = \int_{-\infty}^0 e^{-A\tau} Bu(\tau) d\tau, \quad y(t) = Ce^{At}x_0,$$


---

and thus we obtain

$$(\mathcal{H}u)(t) = y(t) = \int_{-\infty}^0 C e^{A(t-\tau)} B u(\tau) d\tau.$$

**Lemma 4.6:** If  $A$  is asymptotically stable, then  $\mathcal{H}$  is a bounded linear operator.

*Proof.* Exercise. □

**Definition 4.7:** Let  $\mathcal{V}$ ,  $\mathcal{W}$  be two linear spaces with the inner products  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$  and  $\langle \cdot, \cdot \rangle_{\mathcal{W}}$ , respectively. Furthermore, let  $\mathcal{L} : \mathcal{V} \rightarrow \mathcal{W}$  be a linear operator. Then  $\mathcal{L}^* : \mathcal{W} \rightarrow \mathcal{V}$  is called the *adjoint* of  $\mathcal{L}$ , if

$$\langle \mathcal{L}v, w \rangle_{\mathcal{W}} = \langle v, \mathcal{L}^*w \rangle_{\mathcal{V}} \quad \text{for all } v \in \mathcal{V}, w \in \mathcal{W}.$$

For  $\mathcal{H}$  as above and  $u \in \mathcal{L}_2((-\infty, 0], \mathbb{R}^m)$ ,  $y \in \mathcal{L}_2([0, \infty), \mathbb{R}^p)$  we obtain

$$\begin{aligned} \langle \mathcal{H}u, y \rangle_{\mathcal{L}_2([0, \infty), \mathbb{R}^p)} &= \int_0^{\infty} ((\mathcal{H}u)(t))^{\top} y(t) dt \\ &= \int_0^{\infty} \int_{-\infty}^0 u(\tau)^{\top} B^{\top} e^{A^{\top}(t-\tau)} C^{\top} y(t) d\tau dt \\ &= \int_{-\infty}^0 u(\tau)^{\top} \int_0^{\infty} B^{\top} e^{A^{\top}(t-\tau)} C^{\top} y(t) dt d\tau \\ &= \langle u, \mathcal{H}^*y \rangle_{\mathcal{L}_2((-\infty, 0], \mathbb{R}^m)}. \end{aligned}$$

Therefore, we have

$$\mathcal{H}^* : \mathcal{L}_2([0, \infty), \mathbb{R}^p) \rightarrow \mathcal{L}_2((-\infty, 0], \mathbb{R}^m), \quad y \mapsto \int_0^{\infty} B^{\top} e^{A^{\top}(t-\cdot)} C^{\top} y(t) dt.$$

**Definition 4.8:** Let  $\mathcal{V}$ ,  $\mathcal{W}$  be two linear spaces with the inner products  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$  and  $\langle \cdot, \cdot \rangle_{\mathcal{W}}$ , respectively. Let  $\mathcal{L} : \mathcal{V} \rightarrow \mathcal{W}$  be a linear operator with adjoint  $\mathcal{L}^* : \mathcal{W} \rightarrow \mathcal{V}$ . Then  $\sigma \in \mathbb{R}^+$  is called a *singular value* of  $\mathcal{L}$ , if  $\sigma^2$  is an eigenvalue of  $\mathcal{L}^* \mathcal{L}$ , i. e., there exists a  $v \in \mathcal{V} \setminus \{0\}$  such that  $\mathcal{L}^* \mathcal{L}v = \sigma^2 v$ .

Note that, if  $\mathcal{L}^* \mathcal{L}v = \lambda v$ , then

$$\lambda \|v\|_{\mathcal{V}}^2 = \lambda \langle v, v \rangle_{\mathcal{V}} = \langle v, \mathcal{L}^* \mathcal{L}v \rangle_{\mathcal{V}} = \langle \mathcal{L}v, \mathcal{L}v \rangle_{\mathcal{W}} = \|\mathcal{L}v\|_{\mathcal{W}}^2,$$

i. e.,  $\lambda$  is real and nonnegative.

---

**Definition 4.9:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  be asymptotically stable and let  $\mathcal{H}$  be its Hankel operator. Then the positive singular values of  $\mathcal{H}$  are called *Hankel singular values*.

We want to compute the Hankel singular values using the state-space matrices  $A, B, C, D$  only. This will be the goal of the following considerations.

**Theorem 4.10:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  be asymptotically stable. Let  $P$  and  $Q$  its controllability and observability Gramians and  $\mathcal{H}$  its Hankel operator. Then the Hankel singular values are exactly the (positive) square-roots of the eigenvalues of  $PQ$ .

*Proof.* We have

$$y(t) := (\mathcal{H}u)(t) = \int_{-\infty}^0 Ce^{A(t-\tau)} Bu(\tau) d\tau = Ce^{At}z$$

for

$$z := \int_{-\infty}^0 e^{-A\tau} Bu(\tau) d\tau. \quad (4.2)$$

Then we get

$$(\mathcal{H}^*y)(t) = \int_0^{\infty} B^T e^{A^T(\tau-t)} C^T y(\tau) d\tau = B^T e^{-A^T t} \int_0^{\infty} e^{A^T \tau} C^T y(\tau) d\tau.$$

This leads to

$$(\mathcal{H}^*\mathcal{H}u)(t) = (\mathcal{H}^*y)(t) = B^T e^{-A^T t} \int_0^{\infty} e^{A^T \tau} C^T Ce^{A\tau} z d\tau = B^T e^{-A^T t} Qz.$$

Assume that  $\sigma > 0$  is a singular value of  $\mathcal{H}$ . Then there exists an eigenfunction  $u \in \mathcal{L}_2((-\infty, 0], \mathbb{R}^m)$  of  $\mathcal{H}^*\mathcal{H}$  corresponding to an eigenvalue  $\sigma^2 > 0$ , i. e.,

$$(\mathcal{H}^*\mathcal{H}u)(t) = B^T e^{-A^T t} Qz = \sigma^2 u(t).$$

This gives

$$u(t) = \frac{1}{\sigma^2} B^T e^{-A^T t} Qz. \quad (4.3)$$

With (4.2) we get

$$z = \int_{-\infty}^0 e^{-A\tau} B \frac{1}{\sigma^2} B^T e^{-A^T \tau} Qz d\tau = \frac{1}{\sigma^2} PQz,$$

i. e.,  $\sigma^2$  is an eigenvalue of  $PQ$ .

Now assume that  $\sigma^2$  is an eigenvalue of  $PQ$  with an eigenvector  $z \in \mathbb{R}^n \setminus \{0\}$ . Define  $u \in \mathcal{L}_2((-\infty, 0], \mathbb{R}^m)$  as in (4.3). Then we have

$$\begin{aligned}
(\mathcal{H}^* \mathcal{H}u)(t) &= B^\top e^{-A^\top t} \int_0^\infty e^{A^\top \tau} C^\top \int_{-\infty}^0 C e^{A(\tau-s)} B u(s) ds d\tau \\
&= B^\top e^{-A^\top t} \int_0^\infty e^{A^\top \tau} C^\top \int_{-\infty}^0 C e^{A(\tau-s)} B \frac{1}{\sigma^2} B^\top e^{-A^\top s} Q z ds d\tau \\
&= B^\top e^{-A^\top t} \int_0^\infty e^{A^\top \tau} C^\top C e^{A\tau} \frac{1}{\sigma^2} \int_{-\infty}^0 e^{-As} B B^\top e^{-A^\top s} Q z ds d\tau \\
&= B^\top e^{-A^\top t} \int_0^\infty e^{A^\top \tau} C^\top C e^{A\tau} \underbrace{\frac{1}{\sigma^2} P Q z}_{=z} d\tau \\
&= B^\top e^{-A^\top t} Q z = \sigma^2 u(t),
\end{aligned}$$

i. e.,  $\sigma$  is a singular value of  $\mathcal{H}$ . □

Note that for a minimal system, there exist the Cholesky factorizations  $P = RR^\top$  and  $Q = LL^\top$ . Thus, if  $\sigma^2$  is an eigenvalue of  $PQ$ , then we have

$$PQz = (RR^\top)(LL^\top)z = \sigma^2 z$$

This is equivalent to

$$(L^\top R)(R^\top L)L^\top z = \sigma^2 L^\top z,$$

which implies that  $\sigma$  is a singular value of  $L^\top R$ . Therefore, we obtain the Hankel singular values as a side product when computing the balancing transformation in Algorithm 4.1.

## 4.4 Properties of Balanced Truncation

In this section we analyze properties of Algorithm 4.1, see also [Ant05, Sec. 7.2]. In particular, we will derive an error bound using the Hankel singular values.

**Theorem 4.11:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  be asymptotically stable and minimal. Apply Algorithm 4.1 to obtain the reduced-order model  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,p}$ . Assume  $\sigma_r > \sigma_{r+1}$  for the Hankel singular values  $\sigma_i$ ,  $i = 1, \dots, n$ . Then the reduced-order model  $[A_{11}, B_1, C_1, D]$  is asymptotically stable, minimal, and balanced with the Gramians  $P_{11} = Q_{11} = \text{diag}(\sigma_1, \dots, \sigma_r) =: \Sigma_1$ .

*Proof.* Since the system  $[A, B, C, D] \in \Sigma_{n,m,p}$  is minimal, the balancing transformation with  $T$  leads to the transformed Gramians

$$\tilde{P} = \tilde{Q} = \text{diag}(\sigma_1, \dots, \sigma_n) =: \text{diag}(\Sigma_1, \Sigma_2) > 0.$$

Then (in balanced coordinates), the Lyapunov equations

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} + \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} \begin{bmatrix} A_{11}^\top & A_{21}^\top \\ A_{12}^\top & A_{22}^\top \end{bmatrix} = - \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} \begin{bmatrix} B_1^\top & B_2^\top \end{bmatrix}, \quad (4.4)$$

$$\begin{bmatrix} A_{11}^\top & A_{21}^\top \\ A_{12}^\top & A_{22}^\top \end{bmatrix} \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} + \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = - \begin{bmatrix} C_1^\top \\ C_2^\top \end{bmatrix} \begin{bmatrix} C_1 & C_2 \end{bmatrix} \quad (4.5)$$

are satisfied. If the reduced-order model is asymptotically stable, i. e.,  $\Lambda(A_{11}) \subset \mathbb{C}^-$ , then  $\Sigma_1 > 0$  is the controllability and observability Gramian of the reduced-order model, i. e., the reduced-order model is minimal and balanced. Now we show that we indeed have  $\Lambda(A_{11}) \subset \mathbb{C}^-$ . Let  $\lambda \in \mathbb{C}$  be an eigenvalue of  $A_{11}^\top$  with eigenvector  $v \in \mathbb{C}^r$ . Then we obtain

$$-\underbrace{\|B_1^\top v\|_2^2}_{\leq 0} = v^H A_{11} \Sigma_1 v + v^H \Sigma_1 A_{11}^\top v = 2 \operatorname{Re}(\lambda) \underbrace{v^H \Sigma_1 v}_{> 0}.$$

This implies  $\operatorname{Re}(\lambda) \leq 0$ . It remains to show that  $A_{11}$  has no eigenvalues on the imaginary axis. Therefore, assume that there exist imaginary eigenvalues. Let  $i\omega \in i\mathbb{R}$  be an imaginary eigenvalue and  $\{v_1, \dots, v_q\} \subset \mathbb{C}^r$  be an orthonormal basis of  $\ker(A_{11} - i\omega I_r)$  and define  $V = [v_1 \ \dots \ v_q]$ . Then we have

$$(A_{11} - i\omega I_r)V = 0, \quad V^H(A_{11}^\top + i\omega I_r) = 0.$$

Moreover, we have

$$(A_{11} - i\omega I_r)\Sigma_1 + \Sigma_1(A_{11}^\top + i\omega I_r) = -B_1 B_1^\top, \quad (4.6)$$

$$(A_{11}^\top + i\omega I_r)\Sigma_1 + \Sigma_1(A_{11} - i\omega I_r) = -C_1^\top C_1. \quad (4.7)$$

Multiplying (4.7) with  $V^H$  from the left and with  $V$  from the right gives

$$\underbrace{V^H(A_{11}^\top + i\omega I_r)\Sigma_1 V}_{=0} + \underbrace{V^H \Sigma_1 (A_{11} - i\omega I_r) V}_{=0} = -V^H C_1^\top C_1 V,$$

resulting in  $C_1 V = 0$ . Multiplying (4.7) with  $V$  from the right yields

$$(A_{11}^\top + i\omega I_r)\Sigma_1 V + \underbrace{\Sigma_1(A_{11} - i\omega I_r)V}_{=0} = -C_1^\top \underbrace{C_1 V}_{=0},$$

and thus  $(A_{11}^\top + i\omega I_r)\Sigma_1 V = 0$ . Now multiplying (4.6) with  $V^H \Sigma_1$  from the left and with  $\Sigma_1 V$  from the right results in

$$\underbrace{V^H \Sigma_1 (A_{11} - i\omega I_r) \Sigma_1^2 V}_{=0} + \underbrace{V^H \Sigma_1^2 (A_{11}^\top + i\omega I_r) \Sigma_1 V}_{=0} = -V^H \Sigma_1 B_1 B_1^\top \Sigma_1 V,$$

giving  $B_1^T \Sigma_1 V = 0$ . By multiplying (4.6) with  $\Sigma_1 V$  from the right, we obtain

$$(A_{11} - i\omega I_r) \Sigma_1^2 V + \underbrace{\Sigma_1 (A_{11}^T + i\omega I_r) \Sigma_1 V}_{=0} = -B_1 \underbrace{B_1^T \Sigma_1 V}_{=0},$$

so we have  $(A_{11} - i\omega I_r) \Sigma_1^2 V = 0$ . Since  $V$  spans  $\ker(A_{11} - i\omega I_r)$ , we have

$$\Sigma_1^2 V = V \Xi \quad \text{for } \Xi \in \mathbb{C}^{q \times q} \text{ with } \Lambda(\Xi) \subseteq \Lambda(\Sigma_1^2). \quad (4.8)$$

Multiplying the (2,1) block of (4.4) by  $\Sigma_1 V$  from the right yields

$$A_{21} \Sigma_1^2 V + \Sigma_2 A_{12}^T \Sigma_1 V = -B_2 B_1^T \Sigma_1 V = 0.$$

On the other hand, multiplying the (2,1) block of (4.5) by  $V$  from the right results in

$$A_{12}^T \Sigma_1 V + \Sigma_2 A_{21} V = -C_2^T C_1 V = 0.$$

Using (4.8) and both of the last two equations we get

$$A_{21} V \Xi = A_{21} \Sigma_1^2 V = -\Sigma_2 A_{12}^T \Sigma_1 V = \Sigma_2^2 A_{21} V,$$

hence

$$(A_{21} V) \Xi - \Sigma_2^2 (A_{21} V) = 0.$$

This is a *Sylvester matrix equation* with the unknown  $A_{21} V$ . Since by (4.8),  $\Lambda(\Xi) \cap \Lambda(\Sigma_2^2) = \emptyset$ , it is uniquely solvable (see exercise!) and thus we have  $A_{21} V = 0$ . Now we have

$$\tilde{A} \begin{bmatrix} V \\ 0 \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} V \\ 0 \end{bmatrix} = \begin{bmatrix} A_{11} V \\ A_{21} V \end{bmatrix} = i\omega \begin{bmatrix} V \\ 0 \end{bmatrix}.$$

Thus,  $i\omega$  is an imaginary eigenvalue of  $A$ , contradicting its asymptotic stability.  $\square$

In the next theorem we will move towards an error bound for balanced truncation.

**Theorem 4.12:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  with transfer function  $G \in \mathcal{RH}_\infty^{p \times m}$  be asymptotically stable and balanced with the controllability and observability Gramians  $P = Q = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$ . Let  $\sigma_r > \sigma_{r+1} = \dots = \sigma_n$ . Let  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,p}$  be the reduced-order model of order  $r$  obtained by Algorithm 4.1 with transfer function  $\tilde{G} \in \mathcal{RH}_\infty^{p \times m}$ . Then it holds that

$$\|G - \tilde{G}\|_{\mathcal{H}_\infty} \leq 2\sigma_{r+1}$$

(independently of the multiplicity of  $\sigma_{r+1}$ ).

*Proof.* Define

$$\Sigma_1 := \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r), \quad \Sigma_2 := \text{diag}(\sigma_{r+1}, \sigma_{r+2}, \dots, \sigma_n),$$

and the error transfer function

$$E(s) = G(s) - \tilde{G}(s) = C(sI_n - A)^{-1}B + D - (C_1(sI_r - A_{11})^{-1}B_1 + D),$$

and consider its realization  $[\hat{A}, \hat{B}, \hat{C}, \hat{D}] \in \Sigma_{n+r, m, p}$  with

$$\hat{A} = \begin{bmatrix} A_{11} & 0 & 0 \\ 0 & A_{11} & A_{12} \\ 0 & A_{21} & A_{22} \end{bmatrix}, \quad \hat{B} = \begin{bmatrix} B_1 \\ B_1 \\ B_2 \end{bmatrix}, \quad \hat{C} = [-C_1 \quad C_1 \quad C_2], \quad \hat{D} = 0.$$

Using the state-space transformation

$$T := \begin{bmatrix} I_r & I_r & 0 \\ I_r & -I_r & 0 \\ 0 & 0 & I_{n-r} \end{bmatrix} \quad \text{with} \quad T^{-1} = \frac{1}{2} \begin{bmatrix} I_r & I_r & 0 \\ I_r & -I_r & 0 \\ 0 & 0 & 2I_{n-r} \end{bmatrix},$$

we obtain the alternative realization

$$[\hat{A}_1, \hat{B}_1, \hat{C}_1, \hat{D}_1] := [T^{-1}\hat{A}T, T^{-1}\hat{B}, \hat{C}T, \hat{D}] \in \Sigma_{n+r, m, p}$$

with

$$\hat{A}_1 = \begin{bmatrix} A_{11} & 0 & \frac{1}{2}A_{12} \\ 0 & A_{11} & -\frac{1}{2}A_{12} \\ A_{21} & -A_{21} & A_{22} \end{bmatrix}, \quad \hat{B}_1 = \begin{bmatrix} B_1 \\ 0 \\ B_2 \end{bmatrix}, \\ \hat{C}_1 = [0 \quad -2C_1 \quad C_2], \quad \hat{D}_1 = 0.$$

Now define

$$A_0 := \hat{A}_1, \quad B_0 := \begin{bmatrix} \hat{B}_1 & \hat{B}_2 \end{bmatrix} := \begin{bmatrix} B_1 & 0 \\ 0 & \sigma_{r+1}\Sigma_1^{-1}C_1^\top \\ B_2 & -C_2^\top \end{bmatrix}, \\ C_0 := \begin{bmatrix} \hat{C}_1 \end{bmatrix} := \begin{bmatrix} 0 & -2C_1 & C_2 \\ -2\sigma_{r+1}B_1^\top\Sigma_1^{-1} & 0 & -B_2^\top \end{bmatrix}, \\ D_0 := \begin{bmatrix} 0 & 2\sigma_{r+1}I_p \\ 2\sigma_{r+1}I_m & 0 \end{bmatrix}.$$

Then the transfer function of  $[A_0, B_0, C_0, D_0] \in \Sigma_{r+n, m+p, m+p}$  is given by

$$E_0(s) := C_0(sI_{r+n} - A_0)^{-1}B_0 + D_0 \\ = \begin{bmatrix} \hat{C}_1(sI_{r+n} - \hat{A}_1)^{-1}\hat{B}_1 & \hat{C}_1(sI_{r+n} - \hat{A}_1)^{-1}\hat{B}_2 + 2\sigma_{r+1}I_p \\ \hat{C}_2(sI_{r+n} - \hat{A}_1)^{-1}\hat{B}_1 + 2\sigma_{r+1}I_m & \hat{C}_2(sI_{r+n} - \hat{A}_1)^{-1}\hat{B}_2 \end{bmatrix}$$


---

Since  $A$  and  $A_{11}$  are both asymptotically stable, then also the matrices  $\hat{A}$ ,  $\hat{A}_1$ , and  $A_0$  are asymptotically stable by construction. Moreover,  $[A_0, B_0, C_0, D_0]$  has the controllability Gramian  $P_0 = \text{diag}(\Sigma_1, \sigma_{r+1}^2 \Sigma_1^{-1}, 2\Sigma_2)$ , because

$$\begin{aligned} A_0 P_0 + P_0 A_0^\top + B_0 B_0^\top &= \begin{bmatrix} A_{11} \Sigma_1 & 0 & A_{12} \Sigma_2 \\ 0 & \sigma_{r+1}^2 A_{11} \Sigma_1^{-1} & -A_{12} \Sigma_2 \\ A_{21} \Sigma_1 & -\sigma_{r+1}^2 A_{21} \Sigma_1^{-1} & 2A_{22} \Sigma_2 \end{bmatrix} \\ &+ \begin{bmatrix} \Sigma_1 A_{11}^\top & 0 & \Sigma_1 A_{21}^\top \\ 0 & \sigma_{r+1}^2 \Sigma_1^{-1} A_{11}^\top & -\sigma_{r+1}^2 \Sigma_1^{-1} A_{21}^\top \\ \Sigma_2 A_{12}^\top & -\Sigma_2 A_{12}^\top & 2\Sigma_2 A_{22}^\top \end{bmatrix} \\ &+ \begin{bmatrix} B_1 B_1^\top & 0 & B_1 B_2^\top \\ 0 & \sigma_{r+1}^2 \Sigma_1^{-1} C_1^\top C_1 \Sigma_1^{-1} & -\sigma_{r+1}^2 \Sigma_1^{-1} C_1^\top C_2 \\ B_2 B_1^\top & -\sigma_{r+1}^2 C_2^\top C_1 \Sigma_1^{-1} & B_2 B_2^\top + C_2^\top C_2 \end{bmatrix} = 0, \quad (4.9) \end{aligned}$$

where the latter equality follows from combining (4.4) and (4.5). Moreover, we have

$$B_0 D_0^\top = \begin{bmatrix} 0 & 2\sigma_{r+1} B_1 \\ 2\sigma_{r+1}^2 \Sigma_1^{-1} C_1^\top & 0 \\ -2\sigma_{r+1} C_2^\top & 2\sigma_{r+1} B_2 \end{bmatrix} = -P_0 C_0^\top. \quad (4.10)$$

We have

$$\|E\|_{\mathcal{H}_\infty} = \sup_{\omega \in \mathbb{R}} \|E(i\omega)\|_2 \leq \sup_{\omega \in \mathbb{R}} \|E_0(i\omega)\|_2 = \sup_{\omega \in \mathbb{R}} \left( \lambda_{\max}(E_0(i\omega)E_0(i\omega)^H) \right)^{1/2}.$$

Define the conjugated transfer function  $E_0^\sim(s) := E_0(-\bar{s})^H$  which is realized by  $[-A_0^\top, C_0^\top, -B_0^\top, D_0^\top] \in \Sigma_{r+n, m+p, m+p}$ . Then a realization of  $E_0(s)E_0^\sim(s)$  is given by  $[\tilde{A}_0, \tilde{B}_0, \tilde{C}_0, \tilde{D}_0] \in \Sigma_{2(r+n), m+p, m+p}$  (see homework!) with

$$\begin{aligned} \tilde{A}_0 &= \begin{bmatrix} A_0 & -B_0 B_0^\top \\ 0 & -A_0^\top \end{bmatrix}, \quad \tilde{B}_0 = \begin{bmatrix} B_0 D_0^\top \\ C_0^\top \end{bmatrix}, \quad \tilde{C}_0 = [C_0 \quad -D_0 B_0^\top], \\ \tilde{D}_0 &= D_0 D_0^\top. \end{aligned}$$

Using the state-space transformation

$$\tilde{T} := \begin{bmatrix} I_{r+n} & -P_0 \\ 0 & I_{r+n} \end{bmatrix} \quad \text{with} \quad \tilde{T}^{-1} := \begin{bmatrix} I_{r+n} & P_0 \\ 0 & I_{r+n} \end{bmatrix},$$

we obtain an equivalent realization by  $[\hat{A}_0, \hat{B}_0, \hat{C}_0, \hat{D}_0] \in \Sigma_{2(r+n), m+p, m+p}$

$$\begin{aligned} \hat{A}_0 &:= \tilde{T}^{-1} \tilde{A}_0 \tilde{T} = \begin{bmatrix} A_0 & -A_0 P_0 - P_0 A_0^\top - B_0 B_0^\top \\ 0 & -A_0^\top \end{bmatrix} \stackrel{(4.9)}{=} \begin{bmatrix} A_0 & 0 \\ 0 & -A_0^\top \end{bmatrix}, \\ \hat{B}_0 &:= \tilde{T}^{-1} \tilde{B}_0 = \begin{bmatrix} B_0 D_0^\top + P_0 C_0^\top \\ C_0^\top \end{bmatrix} \stackrel{(4.10)}{=} \begin{bmatrix} 0 \\ C_0^\top \end{bmatrix}, \\ \hat{C}_0 &:= \tilde{C}_0 \tilde{T} = [C_0 \quad -C_0 P_0 - D_0 B_0^\top] \stackrel{(4.10)}{=} [C_0 \quad 0], \\ \hat{D}_0 &:= \tilde{D}_0. \end{aligned}$$

Finally we obtain

$$\begin{aligned}
E_0(s)E_0^\sim(s) &= \widehat{C}_0(sI_{2(r+n)} - \widehat{A}_0)^{-1}\widehat{B}_0 + \widehat{D}_0 \\
&= [C_0 \ 0] \begin{bmatrix} sI_{r+n} - A_0 & 0 \\ 0 & sI_{r+n} + A_0^\top \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ C_0^\top \end{bmatrix} + D_0 D_0^\top \\
&= \begin{bmatrix} 0 & 2\sigma_{r+1}I_p \\ 2\sigma_{r+1}I_m & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 2\sigma_{r+1}I_m \\ 2\sigma_{r+1}I_p & 0 \end{bmatrix} = 4\sigma_{r+1}^2 I_{m+p}.
\end{aligned}$$

This implies

$$\begin{aligned}
\|E\|_{\mathcal{H}_\infty} &\leq \sup_{\omega \in \mathbb{R}} \left( \lambda_{\max}(E_0(i\omega)E_0(i\omega)^\mathsf{H}) \right)^{1/2} \\
&= \sup_{\omega \in \mathbb{R}} \left( \lambda_{\max}(E_0(i\omega)E_0^\sim(i\omega)) \right)^{1/2} = \sqrt{4\sigma_{r+1}^2} = 2\sigma_{r+1}.
\end{aligned}$$

□

Now we can conclude an  $\mathcal{H}_\infty$  error bound for the general case.

**Corollary 4.13:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  with transfer function  $G \in \mathcal{RH}_\infty^{p \times m}$  be asymptotically stable and balanced with the controllability and observability Gramians  $P = Q = \text{diag}(\sigma_1 I_{s_1}, \sigma_2 I_{s_2}, \dots, \sigma_k I_{s_k})$ , where  $\sigma_1 > \sigma_2 > \dots > \sigma_k \geq 0$ . Let  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,p}$  be the reduced-order model of order  $r$  obtained by Algorithm 4.1 with  $r = s_1 + s_2 + \dots + s_\ell$  for some  $\ell \leq k$  and with transfer function  $\tilde{G} \in \mathcal{RH}_\infty^{p \times m}$ . Then it holds that

$$\|G - \tilde{G}\|_{\mathcal{H}_\infty} \leq \sum_{j=\ell+1}^k 2\sigma_j.$$

*Proof.* Denote by  $G_j(s)$  the transfer function of the reduced-order model obtained by Algorithm 4.1 by truncating only the Hankel singular values  $\sigma_{j+1}, \dots, \sigma_k$ . So we have  $G(s) = G_k(s)$  and  $\tilde{G}(s) = G_\ell(s)$ . Now it holds that

$$\begin{aligned}
G(s) - G_\ell(s) &= (G_k(s) - G_{k-1}(s)) + (G_{k-1}(s) - G_{k-2}(s)) + \dots \\
&\quad + (G_{\ell+1}(s) - G_\ell(s)),
\end{aligned}$$

which implies

$$\|G - \tilde{G}\|_{\mathcal{H}_\infty} \leq \sum_{j=\ell+1}^k \|G_j - G_{j-1}\|_{\mathcal{H}_\infty} \leq 2 \sum_{j=\ell+1}^k \sigma_j.$$

□

There are also other error bounds for balanced truncation. An *a-posteriori* error bound is given as follows: Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  with transfer function  $G \in \mathcal{RH}_\infty^{p \times m}$  be asymptotically stable and balanced. Let  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,p}$  be the reduced-order model of order  $r$  obtained by Algorithm 4.1 with transfer function  $\tilde{G} \in \mathcal{RH}_\infty^{p \times m}$ . Moreover, assume that  $Y_1 \in \mathbb{R}^{r \times r}$  and  $Y_2 \in \mathbb{R}^{(n-r) \times r}$  solve the Sylvester equation

$$\begin{bmatrix} A_{11}^\top & A_{21}^\top \\ A_{12}^\top & A_{22}^\top \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} + \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} A_{11} = -C^\top C_1.$$

Then we have

$$\|G - \tilde{G}\|_{\mathcal{H}_2}^2 \leq \text{tr} \left( (B_2 B_2^\top + 2Y_2 A_{12}) \Sigma_2 \right).$$

## 4.5 Numerical Solution of Large-Scale Lyapunov Equations

In this section, we discuss the numerical solution of large-scale Lyapunov equations. Since a Lyapunov equation is a special Sylvester equation, the same conditions for unique solvability apply. This means, that a Lyapunov equation

$$AX + XA^\top = -W$$

has a unique solution, if and only if  $\Lambda(A) \cap \Lambda(-A) = \emptyset$ . Since for balanced truncation,  $A$  is assumed to be asymptotically stable, this condition is fulfilled a-priori. In the following we will derive the alternating directions implicit (ADI) iteration, that was introduced to solve partial differential equations in [PR55]. We will see soon that this method is also suitable for large-scale Lyapunov equations that appear in model reduction. There are many other methods, in particular Krylov subspace methods [Sim07], that are often equally good. For sake of brevity, we will not discuss these here in detail.

### 4.5.1 Derivation of the ADI Iteration

Consider the discrete-time Lyapunov equation

$$X = AXA^\top + W, \quad A \in \mathbb{R}^{n \times n}, \quad W = W^\top \in \mathbb{R}^{n \times n}. \quad (4.11)$$

The existence of a unique solution is ensured if  $|\lambda| < 1$  for all  $\lambda \in \Lambda(A)$  (see exercise). This motivates the basic iteration

$$X_k = AX_{k-1}A^\top + W, \quad k \geq 1, \quad X_0 \in \mathbb{R}^{n \times n}. \quad (4.12)$$

Let  $A$  be diagonalizable, i.e., there exists a nonsingular matrix  $V \in \mathbb{C}^{n \times n}$  such that  $A = V\Lambda V^{-1}$ . Let  $\rho(A) := \max_{\lambda \in \Lambda(A)} |\lambda|$  denote the spectral radius of  $A$ . Since

$$\begin{aligned} \|X_k - X\|_2 &= \|A(X_{k-1} - X)A^T\|_2 = \dots = \|A^k(X_0 - X)(A^T)^k\|_2 \\ &\leq \|A^k\|_2^2 \|X_0 - X\|_2 \leq \|V\|_2^2 \|V^{-1}\|_2^2 \rho(A)^{2k} \|X_0 - X\|_2, \end{aligned} \quad (4.13)$$

this iteration converges because  $\rho(A) < 1$  (fixed point argumentation).

Continuous-time Lyapunov equations can be treated similarly, however, we must first transform the data as pointed out in the next lemma.

**Lemma 4.14:** The continuous-times Lyapunov equation

$$AX + XA^T = -W, \quad \Lambda(A) \subset \mathbb{C}^-$$

is equivalent to the discrete-time Lyapunov equation

$$\begin{aligned} X &= C(p)XC(p)^H + \widetilde{W}(p), \quad C(p) := (A - \bar{p}I_n)(A + pI_n)^{-1}, \\ \widetilde{W}(p) &:= -2 \operatorname{Re}(p)(A + pI_n)^{-1}W(A + pI_n)^{-H} \end{aligned} \quad (4.14)$$

for  $p \in \mathbb{C}^-$ .

*Proof.* Exercise. □

Moreover,  $\rho(C(p)) < 1$  (see exercise). Applying (4.12) to (4.14) gives the *Smith iteration*

$$X_k = C(p)X_{k-1}C(p)^H + \widetilde{W}(p), \quad k \geq 1, \quad X_0 \in \mathbb{R}^{n \times n}. \quad (4.15)$$

Similarly as in (4.13), we have

$$\|X_k - X\|_2 \leq \|V\|_2^2 \|V^{-1}\|_2^2 \rho(C(p))^{2k} \|X_0 - X\|_2.$$

This means that we obtain fast convergence by choosing  $p$  such that  $\rho(C(p)) < 1$  is as small as possible. We will discuss this later in more detail.

By varying the shifts  $p$  in (4.15) in every step, we obtain the *ADI iteration for Lyapunov equations*

$$X_k = C(p_k)X_{k-1}C(p_k)^H + \widetilde{W}(p_k), \quad k \geq 1, \quad X_0 \in \mathbb{R}^{n \times n}, \quad p_k \in \mathbb{C}^-. \quad (4.16)$$

### 4.5.2 The ADI Shift Parameter Problem

One can show, similarly to (4.13), that

$$\|X_k - X\|_2 \leq \|V\|_2^2 \|V^{-1}\|_2^2 \rho(M_k)^2 \|X_0 - X\|_2, \quad M_k := \prod_{i=1}^k C(p_i), \quad (4.17)$$

where  $V$  is a transformation matrix diagonalizing  $A$  (assuming it is diagonalizable). The eigenvalues of the product of the Cayley transformations  $M_k$  are

$$\Lambda(M_k) = \left\{ \prod_{i=1}^k \frac{\lambda - \bar{p}_i}{\lambda + p_i} \mid \lambda \in \Lambda(A) \right\}.$$

Good shifts  $p_1^*, \dots, p_k^*$  should make  $\rho(M_k) < 1$  as small as possible. This motivates the ADI shift parameter problem

$$[p_1^*, \dots, p_k^*] = \operatorname{argmin}_{[p_1, \dots, p_k] \in (\mathbb{C}^-)^k} \max_{\lambda \in \Lambda(A)} \left| \prod_{i=1}^k \frac{\lambda - \bar{p}_i}{\lambda + p_i} \right|. \quad (4.18)$$

In general, this is very hard to solve. For instance, in general,  $\rho(C(p))$  is not differentiable and the problem is very expensive, if  $A$  is a large matrix. However, there are some procedures that work well in practice:

- **Wachspress shifts [Wac13]:** Embed  $\Lambda(A)$  in an elliptic function region that depends on the parameters  $\max_{\lambda \in \Lambda(A)} \operatorname{Re}(\lambda)$ ,  $\min_{\lambda \in \Lambda(A)} \operatorname{Re}(\lambda)$ , and  $\arctan \max_{\lambda \in \Lambda(A)} \left| \frac{\operatorname{Im}(\lambda)}{\operatorname{Re}(\lambda)} \right|$  (or approximations thereof). Then, (4.18) can be solved by employing elliptic integral.
- **Heuristic Penzl shifts [Pen00]:** If  $A$  is a large and sparse matrix,  $\Lambda(A)$  is replaced by a small number of approximate eigenvalues (e.g., Ritz values). Then (4.18) is solved heuristically.
- **Self-generating shifts [BKS15]:** If  $A$  is large and sparse, these shifts are based on projections of  $A$  with the data obtained by previous iterations. These shifts also make use of the right-hand side  $W$ .

### 4.5.3 The Low-Rank Phenomenon

Now we consider

$$AX + XA^T = -BB^T, \quad (4.19)$$

where  $A \in \mathbb{R}^{n \times n}$  and  $n$  is 'large', but  $A$  is sparse, i. e., only a few entries in  $A$  are non-zero. Therefore, multiplication with  $A$  can be performed in  $\mathcal{O}(n)$  rather than  $\mathcal{O}(n^2)$  FLOPS. Also solves with  $A$  or  $A + pI$  can be performed efficiently.

---

However,  $X \in \mathbb{R}^{n \times n}$  is usually dense and thus  $X$  cannot be stored for large  $n$  since we would need  $\mathcal{O}(n^2)$  memory. Thus the question arises whether it is possible to store the solution  $X$  more efficiently. In practice we often have  $B \in \mathbb{R}^{n \times m}$ , where  $m \ll n$ , i. e., the right-hand side  $BB^T$  has a low rank. Recall that if  $(A, B)$  is controllable then  $X = X^T \geq 0$  and  $\text{rank}(X) = n$ .

It is a very common observation in practice that the eigenvalues of  $X$  solving (4.19) decay very rapidly towards zero, and fall early below the machine precision. The following theorem explains this eigenvalue decay [SZ02].

**Theorem 4.15:** Let  $A$  be diagonalizable, i. e., there exists an invertible matrix  $V \in \mathbb{C}^{n \times n}$  such that  $A = V\Lambda V^{-1}$ . Then the eigenvalues of  $X$  solving (4.19) satisfy

$$\frac{\lambda_{km+1}(X)}{\lambda_1(X)} \leq \|V\|_2^2 \|V^{-1}\|_2^2 \rho(M_k)^2$$

for any choice of shift parameters  $p$  used to construct  $M_k$  (in particular, the optimal ones).

**Remark 4.16:** • If the eigenvalues of  $A$  cluster in the complex plane, only a few  $p_k$  in the clusters suffice to get a small  $\rho(M_k)$  and thus  $\lambda_i(X)$  decay fast.

- If  $A$  is normal, then  $\|V\|_2 \|V^{-1}\|_2 = 1$  and the bound gives a good explanation for the decay. The nonnormal case is much harder to understand.
- This bound (and most others) does not precisely incorporate the eigenvectors of  $A$  as well as the precise influence of  $B$ .

**Consequence:** If there is a fast decay of  $\lambda_i(X)$ , then  $X$  can be well approximated as  $X = X^T \approx ZZ^H$ , where  $Z \in \mathbb{C}^{n \times r}$  with  $r \ll n$  is a *low-rank solution factor*. Hence, only  $nr$  memory is required. Thus, in the next subsection we consider algorithms for computing the factor  $Z$  without explicitly forming  $X$ .

#### 4.5.4 The Low-Rank Cholesky Factor ADI Iteration

The idea [Pen00] consists of considering one step of the dense ADI iteration (4.16) and inserting  $X_j = Z_j Z_j^H$ . This leads to

$$\begin{aligned} X_j &= C(p_j)X_{j-1}C(p_j)^H + \widetilde{W}(p_j) \\ &= (A - \bar{p}_j I_n)(A + p_j I_n)^{-1} Z_{j-1} Z_{j-1}^H (A + p_j I_n)^{-H} (A - \bar{p}_j I_n)^H \\ &\quad - 2 \text{Re}(p_j)(A + p_j I_n)^{-1} B B^T (A + p_j I_n)^{-H}. \end{aligned}$$

Note that if  $p_j \in \mathbb{R}_-$ , then  $X_j \in \mathbb{R}^{n \times n}$ . Furthermore, if  $p_j \in \mathbb{C}^-$  and  $p_{j+1} = \bar{p}_j$ , then  $X_{j+1} \in \mathbb{R}^{n \times n}$ . Obviously, we have  $X_j = Z_j Z_j^H$  with

$$Z_j = \left[ \sqrt{-2 \operatorname{Re}(p_j)}(A + p_j I_n)^{-1} B \quad (A - \bar{p}_j I_n)(A + p_j I_n)^{-1} Z_{j-1} \right].$$

With  $Z_0 = 0$  we find a low rank variant the ADI iteration (4.16) forming  $Z_j$  successively (grows by  $m$  columns in each step).

The drawback is that all columns are processed in every step which leads to quickly growing costs (in total  $jm$  linear systems have to be solved to get  $Z_j$ ).

However, there is a remedy to this problem. We have observed that

$$S_i = (A + p_i I_n)^{-1} \text{ and } T_j = (A - \bar{p}_j I_n)$$

commute for all  $i, j$  with each other and themselves.

Now consider  $Z_j$  being the iterate after iteration step  $j$

$$Z_j = \left[ \alpha_j S_j B \quad (T_j S_j) \alpha_{j-1} S_{j-1} B \quad \dots \quad (T_j S_j) \cdots (T_2 S_2) \alpha_1 S_1 B \right]$$

with  $\alpha_i = \sqrt{-2 \operatorname{Re}(p_i)}$ . Due to the commutativity, the order of application of the shifts is not important, and we reverse their application to obtain the following alternative iterate

$$\begin{aligned} \tilde{Z}_j &= \left[ \alpha_1 S_1 B \quad \alpha_2 (T_1 S_1) S_2 B \quad \dots \quad \alpha_j (T_1 S_1) \cdots (T_{j-1} S_{j-1}) S_j B \right] \\ &= \left[ \alpha_1 S_1 B \quad \alpha_2 (T_1 S_2) S_1 B \quad \dots \quad \alpha_j (T_{j-1} S_j) (T_{j-2} S_{j-1}) \cdots (T_1 S_2) S_1 B \right] \\ &= \left[ \alpha_1 V_1 \quad \alpha_2 V_2 \quad \dots \quad \alpha_j V_j \right], \\ V_1 &= S_1 B, \quad V_i = T_{i-1} S_i V_{i-1}, \quad i = 1, \dots, j. \end{aligned}$$

We have  $X_j = \tilde{Z}_j \tilde{Z}_j^H$ , but in this formulation only the new columns are processed. Even more structure is revealed by the Lyapunov residual.

**Theorem 4.17:** The residual at step  $j$  of (4.16), started with  $X_0 = 0$ , is of rank at most  $m$  and given by

$$\begin{aligned} R_j &:= A \tilde{Z}_j \tilde{Z}_j^H + \tilde{Z}_j \tilde{Z}_j^H A^T + B B^T = W_j W_j^H, \\ W_j &= M_j B = C(p_j) W_{j-1} = W_{j-1} - 2 \operatorname{Re}(p_j) V_j, \quad W_0 := B, \end{aligned}$$

where  $M_j := \prod_{i=1}^j C(p_i)$ . Moreover, it holds  $V_j = (A + p_j I_n)^{-1} W_{j-1}$ .

*Proof.* We have

$$\begin{aligned} R_j &= A X_j + X_j A^T + B B^T = A(X_j - X) + (X_j - X) A^T \quad (\text{by (4.19)}) \\ &= A M_j (X_0 - X) M_j^H + M_j (X_0 - X) M_j^H A^T \\ &= -M_j A X M_j^H - M_j X A^T M_j^H \\ &= -M_j (A X + X A^T) M_j^H = M_j B B^T M_j^H. \end{aligned}$$

**Algorithm 4.2** Low-rank ADI (LRCF-ADI) iteration for Lyapunov equations

**Input:**  $A, B$  from (4.19), shifts  $P = \{p_1, \dots, p_{\maxiter}\} \subset \mathbb{C}^-$ , residual tolerance  $\text{tol}$ .

**Output:**  $Z_k$  such that  $X = Z_k Z_k^H$  (approx.) solves (4.19).

- 1: Initialize  $j = 1, W_0 := B, Z_0 := []$ .
- 2: **while**  $\|W_{j-1}\|_2 \geq \text{tol}$  **do**
- 3:   Set  $V_j := (A + p_j I_n)^{-1} W_{j-1}$ .
- 4:   Set  $W_j := W_{j-1} - 2 \operatorname{Re}(p_j) V_j$ .
- 5:   Set  $Z_j := [Z_{j-1} \quad \sqrt{-\operatorname{Re}(p_j)} V_j]$ .
- 6:   Set  $j := j + 1$ .
- 7: **end while**

Moreover, it holds

$$\begin{aligned} V_j &= T_{j-1} S_j V_{j-1} = T_{j-1} S_j T_{j-2} S_{j-1} V_{j-2} = \dots = \\ &= S_j \left( \prod_{k=1}^{j-1} T_k S_k \right) B = S_j M_{j-1} B = (A + p_j I_n)^{-1} W_{j-1}, \end{aligned} \quad (4.20)$$

and

$$W_j = M_j B = S_j T_j W_{j-1} = W_{j-1} - 2 \operatorname{Re}(p_j) S_j W_{j-1} = W_{j-1} - 2 \operatorname{Re}(p_j) V_j.$$

□

Thanks to the above theorem, the norm of the Lyapunov residual can be cheaply computed via  $\|R_j\|_2 = \|W_j W_j^H\|_2 = \|W_j\|_2^2$ . All this leads to Algorithm 4.2 which is also often referred to as low-rank Cholesky-factor ADI (LRCF-ADI) iteration. Algorithm 4.2 produces complex low-rank factors, if some of the shifts are complex, which might be required for problems with nonsymmetric  $A$ .

However, it is still possible to ensure that  $Z_j \in \mathbb{R}^{n \times n_j}$ , see [BKS13].

**Definition 4.18:** A set of shift parameters  $P$  is called proper if for all  $p \in P$ , also  $\bar{p} \in P$ .

**Theorem 4.19:** Assume  $P = \{p_1, \dots, p_k\}$  to be a set of proper shifts and assume w. l. o. g. that  $p_{j+1} = \bar{p}_j \notin \mathbb{R}$ . Then for  $V_j, V_{j+1}$  it holds

$$\begin{aligned} V_{j+1} &= \bar{V}_j + 2\beta_j \operatorname{Im}(V_j), \\ W_{j+1} &= W_{j-1} - 4 \operatorname{Re}(p_j) (\operatorname{Re}(V_j) + \beta_j \operatorname{Im}(V_j)), \end{aligned}$$

with  $\beta_j = \frac{\operatorname{Re}(p_j)}{\operatorname{Im}(p_j)}$ .

**Question:** Why does that help?  $V_j, V_{j+1}$  are still complex. Consider

$$Z_{j+1} = [Z_{j-1} \quad \alpha_j V_j \quad \alpha_j V_{j+1}].$$

This gives

$$X_{j+1} = Z_{j+1} Z_{j+1}^H = Z_{j-1} Z_{j-1}^H + \alpha_j^2 \hat{Z} \hat{Z}^H,$$

with

$$\begin{aligned} \hat{Z} &= [V_j \quad V_{j+1}] = [\operatorname{Re}(V_j) + i \operatorname{Im}(V_j) \quad \operatorname{Re}(V_j) + 2\beta_j \operatorname{Im}(V_j) - i \operatorname{Im}(V_j)] \\ &= [\operatorname{Re}(V_j) \quad \operatorname{Im}(V_j)] \underbrace{\begin{bmatrix} I_m & I_m \\ iI_m & (2\beta_j - i)I_m \end{bmatrix}}_{=:N}. \end{aligned}$$

This yields

$$\hat{Z} \hat{Z}^H = [\operatorname{Re}(V_j) \quad \operatorname{Im}(V_j)] N N^H [\operatorname{Re}(V_j) \quad \operatorname{Im}(V_j)]^H$$

and

$$\begin{aligned} 0 < N N^H &= \begin{bmatrix} 2I_m & 2\beta_j I_m \\ 2\beta_j I_m & (4\beta_j^2 + 1)I_m \end{bmatrix} \\ &= \underbrace{\begin{bmatrix} I_m & 0 \\ \beta_j I_m & I_m \end{bmatrix}}_{=:L} \underbrace{\begin{bmatrix} 2I_m & 0 \\ 0 & 2(\beta_j^2 + 1)I_m \end{bmatrix}}_{=: \Gamma > 0} \begin{bmatrix} I_m & \beta_j I_m \\ 0 & I_m \end{bmatrix}. \end{aligned}$$

Therefore, we can alternatively choose the following  $\check{Z}$  instead of  $\hat{Z}$  to obtain the same  $Z_{j+1}$ , namely

$$\begin{aligned} \check{Z} &:= [\operatorname{Re}(V_j) \quad \operatorname{Im}(V_j)] L \Gamma^{\frac{1}{2}} \\ &= \sqrt{2} \begin{bmatrix} \operatorname{Re}(V_j) + \beta_j \operatorname{Im}(V_j) & \sqrt{(\beta_j^2 + 1)} \operatorname{Im}(V_j) \end{bmatrix} \in \mathbb{R}^{n \times 2m}, \end{aligned}$$

in other words,  $Z_{j+1}$  is constructed to be real. This leads to the real version of the LRCF-ADI iteration.

#### 4.5.5 Balanced Truncation Using the LRCF-ADI Method

Algorithm 4.1 can now be modified by including low-rank methods for solving the Lyapunov equations. The result is Algorithm 4.3.

**Remark 4.20:** The  $\mathcal{H}_\infty$  error bound for balanced truncation does not necessarily hold anymore. First of all, we do not compute all Hankel singular values. Hence, the ones which have not been computed, can only be estimated using

**Algorithm 4.3** Balanced truncation (pro version)

**Input:** Asymptotically stable system  $[A, B, C, D] \in \Sigma_{n,m,p}$ , desired maximum reduced order  $r_{\max}$ , ADI residual tolerance  $\text{tol}$ .

**Output:** Reduced-order model  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] \in \Sigma_{r,m,p}$  with  $r \leq r_{\max}$ .

1: Solve the Lyapunov equations

$$AP + PA^T = -BB^T, \quad A^TQ + QA = -C^TC$$

using (the real version of) Algorithm 4.2 to determine two low-rank factors  $R \in \mathbb{R}^{n \times r_P}$  and  $L \in \mathbb{R}^{n \times r_Q}$  such that  $P \approx RR^T$  and  $Q \approx LL^T$  and with residuals less than  $\text{tol}$ .

2: Set  $r := \min\{r_{\max}, r_P, r_Q\}$ .

3: Compute the SVD of  $L(:, 1:r)^T R(:, 1:r) = U\Sigma V^T$ .

4: Set  $T := RV\Sigma^{-\frac{1}{2}}$  and  $W := LU\Sigma^{-\frac{1}{2}}$ .

5: Balance and truncate to obtain the reduced-order model

$$[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] := [W^TAT, W^TB, CT, D].$$

the smallest singular value of  $\Sigma$  in Step 3 of Algorithm 4.3. Moreover, the singular values contained in  $\Sigma$  may be corrupted by the approximation errors done when computing  $R$  and  $L$ . Therefore, the reduction error can only be estimated in practice.



## CHAPTER 5

---

### Passivity-Preserving Balancing-Based Model Reduction

---

In this chapter we will focus on some aspects of structure-preservation. In this chapter we consider *passive* systems which often appear in the modeling of electrical circuits, power network, mechanical systems, and many more. See Chapter 1 for some examples. Thus, when doing model reduction, one would like to obtain a passive reduced-order model in order to preserve the physical properties in the model. In this chapter, we will first define passivity and show that each passive system admits a positive real transfer function. Thereafter we will discuss a passivity-preserving model reduction scheme using alternative energy functionals. This will lead to the method of positive real balanced truncation which we will analyze afterwards. Since positive real balanced truncation relies on algebraic Riccati equations rather than Lyapunov equations, we will also treat the numerical solution of large-scale algebraic Riccati equations. Many of the results presented here can be found in the famous works by Jan C. Willems [Wil71, Wil72a, Wil72b].

#### 5.1 Passivity and Positive Real Transfer Functions

First we define passivity for a LTI systems.

**Definition 5.1:** Let  $[A, B, C, D] \in \Sigma_{n,m,m}$  be given. Then the system is called

passive, if

$$\int_0^T y(\tau)^\top u(\tau) d\tau \geq 0 \quad (5.1)$$

holds for all  $T \geq 0$  and all solution trajectories  $(x, u, y) \in \mathcal{L}_2([0, T], \mathbb{R}^{n+2m})$  of the system with  $x(0) = 0$ .

The expression on the left-hand side of (5.1) can be understood as the energy that is supplied to the system in the time interval  $[0, T]$ . Therefore, passivity of a dynamical system is the property that for each solution, more energy has to be supplied than energy that can be extracted from the system. So the system cannot internally produce energy. Passivity is connected to two energy functionals:

a) the (virtual) available storage  $V^- : \mathbb{R}^n \rightarrow \mathbb{R}$  with

$$V^-(x_0) := \sup \left\{ - \int_0^\infty 2y(\tau)^\top u(\tau) d\tau \mid (x, u, y) \in \mathcal{L}_2([0, \infty), \mathbb{R}^{n+2m}) \right. \\ \left. \text{is a solution of } [A, B, C, D] \text{ with } x(0) = x_0 \right\};$$

b) the required supply  $V^+ : \mathbb{R}^n \rightarrow \mathbb{R}$  with

$$V^+(x_0) := \inf \left\{ \int_{-\infty}^0 2y(\tau)^\top u(\tau) d\tau \mid (x, u, y) \in \mathcal{L}_2((-\infty, 0], \mathbb{R}^{n+2m}) \right. \\ \left. \text{is a solution of } [A, B, C, D] \text{ with } x(0) = x_0 \right\}.$$

The value of  $V^+(x_0)$  is the least amount of energy that has to be supplied to the system to reach the state  $x_0$ . On the other hand, the value of  $V^-(x_0)$  is the maximum amount of energy that can be extracted from the system by stabilizing solution trajectories.

We will later see that under some conditions, the functionals  $V^+$ ,  $V^-$  are so-called *storage functions*. A storage function is a function  $V : \mathbb{R}^n \rightarrow \mathbb{R}^+$  with  $V(0) = 0$  that fulfills the *dissipation inequality*

$$V(x_1) - V(x_0) \leq \int_{t_0}^{t_1} 2y(\tau)^\top u(\tau) d\tau, \quad (5.2)$$

where  $(x, u, y) \in \mathcal{L}_2([t_0, t_1], \mathbb{R}^{n+2m})$  is a solution trajectory of the system with  $x(t_0) = x_0$  and  $x(t_1) = x_1$ . If  $V$  is differentiable, then the dissipation inequality can be formulated in its differential form

$$V'(x(t)) \cdot \dot{x}(t) \leq 2y(t)^\top u(t), \quad (5.3)$$

where  $V' : \mathbb{R}^n \rightarrow \mathbb{R}^{1 \times n}$  is the Jacobian of  $V$ .

Now we consider the special case of *quadratic* storage functions, i. e.,  $V(x) = x^\top P x$  for some symmetric positive semi-definite matrix  $P \in \mathbb{R}^{n \times n}$ . For such,  $V'(x) = 2x^\top P$  and (5.3) gives

$$\begin{aligned} V'(x(t)) \cdot \dot{x}(t) &= 2x(t)^\top P \dot{x}(t) \\ &= x(t)^\top P (Ax(t) + Bu(t)) + (Ax(t) + Bu(t))^\top P x(t) \\ &= \begin{bmatrix} x(t) \\ u(t) \end{bmatrix}^\top \begin{bmatrix} A^\top P + PA & PB \\ B^\top P & 0 \end{bmatrix} \begin{bmatrix} x(t) \\ u(t) \end{bmatrix} \\ &\leq 2y(t)^\top u(t) \\ &= (Cx(t) + Du(t))^\top u(t) + u(t)^\top (Cx(t) + Du(t)) \\ &= \begin{bmatrix} x(t) \\ u(t) \end{bmatrix}^\top \begin{bmatrix} 0 & C^\top \\ C & D + D^\top \end{bmatrix} \begin{bmatrix} x(t) \\ u(t) \end{bmatrix}. \end{aligned}$$

Therefore, each quadratic storage function can be expressed by a solution  $P \geq 0$  of the *linear matrix inequality (LMI)*

$$\begin{bmatrix} A^\top P + PA & PB - C^\top \\ B^\top P - C & -D - D^\top \end{bmatrix} \leq 0, \quad P = P^\top. \quad (5.4)$$

It can be shown that that if the system is controllable, then the LMI (5.4) has two extremal solutions  $P^+ \in \mathbb{R}^{n \times n}$  and  $P^- \in \mathbb{R}^{n \times n}$  such that  $P^- \leq P \leq P^+$  for each solution  $P \in \mathbb{R}^{n \times n}$  of the LMI and with the properties

$$V^-(x_0) = x_0^\top P^- x_0, \quad V^+(x_0) = x_0^\top P^+ x_0.$$

We have the following theorem connecting all these concepts.

**Theorem 5.2:** Let  $[A, B, C, D] \in \Sigma_{n,m,m}$  be controllable. Then the following statements are equivalent:

- The system  $[A, B, C, D]$  is passive.
- It holds that  $V^+(x_0) \geq 0$  for all  $x_0 \in \mathbb{R}^n$ .
- There exists a function  $V : \mathbb{R}^n \rightarrow \mathbb{R}^+$  that satisfies the dissipation inequality (5.2).

Moreover, whenever one of the above conditions is fulfilled, then we have

$$-\infty < V^-(x_0) \leq V^+(x_0) < \infty.$$

*Proof.* We show “a)  $\Rightarrow$  b)”: Assume that b) is not satisfied, i. e., there exist an  $\varepsilon > 0$  and a solution trajectory  $(x, u, y) \in \mathcal{L}_2((-\infty, 0], \mathbb{R}^{n+2m})$  with  $x(0) = x_0$

(and  $\lim_{t \rightarrow -\infty} x(t) = 0$ ) such that

$$\int_{-\infty}^0 2y(\tau)^\top u(\tau) d\tau < -\varepsilon.$$

It can be shown that for each  $\varepsilon > 0$  there exist a  $T > 0$  and a solution trajectory  $(\tilde{x}, \tilde{u}, \tilde{y}) \in \mathcal{L}_2((-\infty, 0], \mathbb{R}^{n+2m})$  with compact support in the interval  $[-T, 0]$  and with  $\tilde{x}(0) = x_0$  such that

$$\left| \int_{-\infty}^0 2y(\tau)^\top u(\tau) d\tau - \int_{-\infty}^0 2\tilde{y}(\tau)^\top \tilde{u}(\tau) d\tau \right| < \frac{\varepsilon}{2}.$$

This gives

$$\begin{aligned} \int_{-\infty}^0 2\tilde{y}(\tau)^\top \tilde{u}(\tau) d\tau &= \int_{-T}^0 2\tilde{y}(\tau)^\top \tilde{u}(\tau) d\tau \\ &= \int_0^T 2\tilde{y}(\tau - T)^\top \tilde{u}(\tau - T) d\tau < -\frac{\varepsilon}{2}. \end{aligned}$$

Therefore, condition a) is violated for the solution trajectory  $(\tilde{x}(\cdot - T), \tilde{u}(\cdot - T), \tilde{y}(\cdot - T)) \in \mathcal{L}_2([0, T], \mathbb{R}^{n+2m})$ .

The statement “b)  $\Rightarrow$  c)” follows from the fact that  $V^+$  is a storage function.

Now we show “c)  $\Rightarrow$  a)”: From the dissipation inequality with  $t_0 = 0$ ,  $x(t_0) = 0$ ,  $t_1 = T$ , and the condition  $V(0) = 0$  we obtain

$$0 \leq V(x_1) \leq \int_0^T 2y(\tau)^\top u(\tau) d\tau,$$

which gives the result.

The last inequality follows from b) since for all  $T > 0$  and all solution trajectories  $(x, u, y) \in \mathcal{L}_2((-\infty, T], \mathbb{R}^{n+2m})$  with  $x(0) = x_0$  we obtain

$$-\int_0^T 2y(\tau)^\top u(\tau) d\tau \leq \int_{-\infty}^0 2y(\tau)^\top u(\tau) d\tau.$$

With  $T \rightarrow \infty$ , taking the supremum on the left-hand side and the infimum on the right-hand side gives  $V^-(x_0) \leq V^+(x_0)$  for all  $x_0 \in \mathbb{R}^n$ . The finiteness of both functionals then follows from controllability, since every point  $x_0$  can be reached by a solution trajectory.  $\square$

The passivity property is equivalent to a structural property of its transfer function, namely, they are positive real.

**Definition 5.3:** Let  $[A, B, C, D] \in \Sigma_{n,m,m}$  be given with the transfer function  $G(s) \in \mathbb{R}(s)^{m \times m}$ . Then  $G(s)$  is called *positive real*, if  $G(s)$  has no poles in  $\mathbb{C}^+$  and

$$\Psi(\lambda) := G(\lambda) + G(\lambda)^H \geq 0 \quad \forall \lambda \in \mathbb{C}^+. \quad (5.5)$$

The following famous theorem (called the positive real lemma, sometimes also the Kalman-Yakubovich-Popov(-Anderson) lemma) makes a connection between solvability of the LMI (5.4) and positive realness.

**Theorem 5.4:** Let  $[A, B, C, D] \in \Sigma_{n,m,m}$  be given with the transfer function  $G(s) \in \mathbb{R}(s)^{m \times m}$  and let  $\Psi$  be as in (5.5). Then the following statements are satisfied:

- If the LMI (5.4) has a solution  $P > 0$ , then  $G(s)$  is positive real.
- If the system  $[A, B, C, D]$  is minimal and  $G(s)$  is positive real, then there exists a solution  $P > 0$  of the LMI (5.4).

*Proof.* We prove statement a): Let  $P > 0$  be a solution of the LMI (5.4). Let  $v \in \mathbb{C}^n$  be an eigenvector corresponding to an eigenvalue  $\lambda \in \mathbb{C}$  of  $A$ . Then we have

$$v^H A^T P v + v^H P A v = \bar{\lambda} v^H P v + \lambda v^H P v = 2 \operatorname{Re}(\lambda) \underbrace{v^H P v}_{>0} \leq 0.$$

Therefore, we have  $\operatorname{Re}(\lambda) \leq 0$  and thus,  $G(s)$  has no poles in  $\mathbb{C}^+$ .

Moreover, we have that

$$\begin{aligned} \begin{bmatrix} A & B \end{bmatrix} \begin{bmatrix} (\lambda I_n - A)^{-1} B \\ I_m \end{bmatrix} &= A(\lambda I_n - A)^{-1} B + B \\ &= (A + \lambda I_n - A)(\lambda I_n - A)^{-1} B \\ &= \lambda(\lambda I_n - A)^{-1} B. \end{aligned}$$

With  $\lambda \in \mathbb{C}^+$  and using the above identity we obtain

$$\begin{aligned} & \begin{bmatrix} (\lambda I_n - A)^{-1} B \\ I_m \end{bmatrix}^H \begin{bmatrix} A^T P + P A & P B \\ B^T P & 0 \end{bmatrix} \begin{bmatrix} (\lambda I_n - A)^{-1} B \\ I_m \end{bmatrix} \\ &= \begin{bmatrix} (\lambda I_n - A)^{-1} B \\ I_m \end{bmatrix}^H \left( \begin{bmatrix} P A & P B \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} A^T P & 0 \\ B^T P & 0 \end{bmatrix} \right) \begin{bmatrix} (\lambda I_n - A)^{-1} B \\ I_m \end{bmatrix} \\ &= (\lambda + \bar{\lambda}) B^T (\lambda I_n - A)^{-H} P (\lambda I_n - A)^{-1} B \geq 0. \end{aligned}$$

Now we get

$$\begin{aligned}\Psi(\lambda) &= \begin{bmatrix} (\lambda I_n - A)^{-1} B \\ I_m \end{bmatrix}^H \begin{bmatrix} 0 & C^T \\ C & D + D^T \end{bmatrix} \begin{bmatrix} (\lambda I_n - A)^{-1} B \\ I_m \end{bmatrix} \\ &\geq \begin{bmatrix} (\lambda I_n - A)^{-1} B \\ I_m \end{bmatrix}^H \begin{bmatrix} -A^T P - PA & -PB + C^T \\ -B^T P + C & D + D^T \end{bmatrix} \begin{bmatrix} (\lambda I_n - A)^{-1} B \\ I_m \end{bmatrix} \\ &\geq 0,\end{aligned}$$

where the latter inequality follows from the fact, that  $P$  solves (5.4).

The proof of statement b) is quite complicated and technical and therefore, we omit it here. The proof can be found in [Ran96].  $\square$

**Remark 5.5:** There are many relaxations of the assumptions of the positive real lemma. For instance, it can be shown that  $G(s)$  is already positive real, if there exists a solution  $P \geq 0$  of the LMI (5.4). However, the techniques for the proof get more involved, see, e. g., [AV73].

From the above theorem, the following corollary is immediate.

**Corollary 5.6:** The system  $[A, B, C, D] \in \Sigma_{n,m,p}$  is passive if and only if its transfer function  $G(s)$  is positive real.

*Proof.* Let  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] \in \Sigma_{\tilde{n},m,p}$  be a minimal realization of  $G(s)$ . Then,  $[A, B, C, D]$  is passive if and only if  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}]$  is passive, since both generate the same input/output pairs. From Theorem 5.2 and controllability, this is equivalent to the existence of a matrix  $\tilde{P} \geq 0$  such that the LMI

$$\begin{bmatrix} \tilde{A}^T \tilde{P} + \tilde{P} \tilde{A} & \tilde{P} \tilde{B} - \tilde{C}^T \\ \tilde{B}^T \tilde{P} - \tilde{C} & -\tilde{D} - \tilde{D}^T \end{bmatrix} \leq 0, \quad \tilde{P} = \tilde{P}^T$$

is satisfied. From Theorem 5.4 and Remark 5.5, this is equivalent to positive realness of  $G(s)$ .  $\square$

## 5.2 Positive Real Balanced Truncation

Now we want to derive a balancing-type algorithm for passivity-preserving model reduction. Assume that  $P \geq 0$  is a solution of the LMI (5.4). Then there exist matrices  $K \in \mathbb{R}^{q \times n}$  and  $L \in \mathbb{R}^{q \times m}$  such that

$$\begin{bmatrix} A^T P + PA & PB - C^T \\ B^T P - C & -D - D^T \end{bmatrix} = - \begin{bmatrix} K^T \\ L^T \end{bmatrix} \begin{bmatrix} K & L \end{bmatrix}.$$

Under the assumption that  $D + D^\top$  is invertible (then  $q \geq m$ ), we can apply the Schur complement on both sides and obtain

$$\begin{aligned} A^\top P + PA + (PB - C^\top)(D + D^\top)^{-1}(B^\top P - C) \\ = -K^\top K + K^\top L(L^\top L)^{-1}L^\top K. \end{aligned}$$

Let  $L = U \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} V^\top$  with orthogonal matrices  $U \in \mathbb{R}^{q \times q}$ ,  $V \in \mathbb{R}^{m \times m}$  and an invertible diagonal matrix  $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_m)$  with  $\sigma_1 \geq \dots \geq \sigma_m > 0$  be given. Then we have

$$\begin{aligned} & -K^\top K + K^\top L(L^\top L)^{-1}L^\top K \\ & = -K^\top K + K^\top U \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} V^\top \left( V \begin{bmatrix} \Sigma & 0 \end{bmatrix} U^\top U \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} V^\top \right)^{-1} V \begin{bmatrix} \Sigma & 0 \end{bmatrix} U^\top K \\ & \quad - K^\top K + K^\top U \begin{bmatrix} I_m & 0 \\ 0 & 0 \end{bmatrix} U^\top K \leq 0. \end{aligned}$$

Therefore, each solution  $P$  of the LMI (5.4) satisfies the *algebraic Riccati inequality*

$$A^\top P + PA + (PB - C^\top)(D + D^\top)^{-1}(B^\top P - C) \leq 0, \quad P = P^\top.$$

It can be seen that the extremal elements of its solution set,  $P^+$  and  $P^-$ , even satisfy the *algebraic Riccati equation (ARE)*

$$A^\top P + PA + (PB - C^\top)(D + D^\top)^{-1}(B^\top P - C) = 0, \quad P = P^\top. \quad (5.6)$$

Recall that if the system  $[A, B, C, D]$  is controllable and passive, then the extremal solutions  $P^+$  and  $P^-$  exist and  $P^+ \geq 0$ . It can be shown (see [Obe91, Sect. 6]) that if the system is also observable, then we even have

$$0 < P^- \leq P^+. \quad (5.7)$$

Furthermore, it can be shown (exercise!) that if  $P > 0$  is a solution of the ARE, then  $Q = P^{-1}$  is a solution of the *dual ARE*

$$AQ + QA^\top + (QC^\top - B)(D + D^\top)^{-1}(CQ - B^\top) = 0, \quad Q = Q^\top. \quad (5.8)$$

Therefore, if we have (5.7), then there exist a minimal solution  $Q^- \in \mathbb{R}^{n \times n}$  and a maximal solution  $Q^+ \in \mathbb{R}^{n \times n}$  with

$$0 < Q^- \leq Q \leq Q^+ \quad \text{for all solutions } Q \text{ of (5.8),}$$

and with  $Q^- = (P^+)^{-1}$  and  $Q^+ = (P^-)^{-1}$ . The minimal solutions  $P^-$  and  $Q^-$  are called the *positive real (observability and controllability) Gramians* and they are now subject to our balancing procedure. These will attain the role of the controllability and observability Gramian from the previous chapter.

**Definition 5.7:** Let  $[A, B, C, D] \in \Sigma_{n,m,m}$  be minimal and passive with positive real Gramians  $P^- > 0$  and  $Q^- > 0$ . Then the system is called *positive real balanced*, if  $P^- = Q^- = \Sigma$ . In this case, the eigenvalues of the matrix  $\Sigma$  are called the *positive real characteristic values*.

Now we discuss the associated balancing transformations.

**Theorem 5.8:** Let  $[A, B, C, D] \in \Sigma_{n,m,m}$  be minimal and passive with positive real Gramians  $P^- > 0$  and  $Q^- > 0$ . Then there exists an invertible matrix  $T \in \mathbb{R}^{n \times n}$  such that  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] := [T^{-1}AT, T^{-1}B, CT, D]$  with the positive real Gramians  $\tilde{P}^- > 0$  and  $\tilde{Q}^- > 0$  is positive real balanced.

*Proof.* Since  $P^- > 0$  and  $Q^- > 0$ , there exist Cholesky decompositions  $P^- = RR^T$  and  $Q^- = LL^T$ , where  $R$  and  $L$  are lower triangular and invertible. Now consider the singular value decomposition  $L^T R = U\Sigma V^T$  with orthogonal  $U, V \in \mathbb{R}^{n \times n}$  and  $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$ ,  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$ . Since  $L$  and  $R$  are invertible, so is  $L^T R$  and therefore, we have  $\sigma_n > 0$ .

It is easy to check that as for the case of standard balancing, we have that

$$\tilde{P}^- = T^T P^- T, \quad \tilde{Q}^- = T^{-1} Q^- T^{-T}.$$

Now the rest of the proof is similar to the proof of Theorem 4.5. With  $T := LU\Sigma^{-\frac{1}{2}}$  we find  $T^{-1} = \Sigma^{-\frac{1}{2}}V^T R^T$  which make the system positive real balanced (exercise!).  $\square$

This leads to Algorithm 5.1 for model reduction that is called *positive real balanced truncation* and which has first been considered in [Obe91, Sect. 6].

**Remark 5.9:** a) If the passive system  $[A, B, C, D] \in \Sigma_{n,m,m}$  is not minimal, then the algebraic Riccati equations (5.6) and (5.8) may not have (minimal) solutions. However, for many problems the existence of minimal solutions can be derived from the structure of the models, such as for electrical circuit models.

b) If the matrix  $D + D^T$  is not invertible, then the algebraic Riccati equations cannot be formed. In this case, one has to resort to *Lur'e equations* such as

$$\begin{aligned} A^T P + PA &= -K^T K, & P &= P^T, \\ PB - C^T &= -K^T L, \\ D + D^T &= L^T L, \end{aligned} \tag{5.9}$$

**Algorithm 5.1** Positive real balanced truncation (basic version)

**Input:** Minimal and passive system  $[A, B, C, D] \in \Sigma_{n,m,m}$  with invertible  $D + D^\top$ , desired reduced order  $r$ .

**Output:** Passive reduced-order model  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,m}$ .

- 1: Compute the minimal (and positive definite) solutions  $P^-$  and  $Q^-$  of the algebraic Riccati equations

$$\begin{aligned} A^\top P + PA + (PB - C^\top)(D + D^\top)^{-1}(B^\top P - C) &= 0, & P &= P^\top, \\ AQ + QA^\top + (QC^\top - B)(D + D^\top)^{-1}(CQ - B^\top) &= 0, & Q &= Q^\top. \end{aligned}$$

- 2: Compute Cholesky factorization  $P^- = RR^\top$  and  $Q^- = LL^\top$ .
- 3: Compute the singular value decomposition  $L^\top R = U\Sigma V^\top$ .
- 4: Set  $T := LU\Sigma^{-\frac{1}{2}}$  (and  $T^{-1} = \Sigma^{-\frac{1}{2}}V^\top R^\top$ ).
- 5: Do the balancing transformation

$$[T^{-1}AT, T^{-1}B, CT, D] = \left[ \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}, [C_1 \ C_2], D \right]$$

and set the reduced-order model as  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,m}$ .

which has to be solved for the triple  $(P, K, L) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{q \times n} \times \mathbb{R}^{q \times m}$ , where  $q$  is as small as possible among all such triples solving (5.9). This minimal rank property is motivated by the fact, that the algebraic Riccati inequality turns to an equation, if and only if  $L \in \mathbb{R}^{q \times m}$  is invertible, i. e.,  $q = m$ . This is the smallest possible rank, since  $D + D^\top$  was assumed to be invertible in this case. In the general case, the solutions of the LMI which have this minimal rank property are called *rank-minimizing solutions*. It can be shown that this minimal rank is

$$q = \text{rank}_{\mathbb{R}(s)}(G(s) + G^\sim(s)).$$

Extremal solutions of (5.9) are always rank-minimizing.

### 5.3 Analysis of the Method

In this section, we give a brief analysis of the properties of positive real balanced truncation (mainly without the proofs). First of all, we see that the reduced-order model is again passive.

**Theorem 5.10:** Let  $[A, B, C, D] \in \Sigma_{n,m,m}$  with invertible  $D + D^\top$  be asymptotically stable, minimal, and passive. Apply Algorithm 5.1 to obtain the reduced-order model  $[A_{11}, B_1, C_1, D] \in \Sigma_{r,m,m}$ . Assume further, that for the positive real characteristic values  $\sigma_1, \sigma_2, \dots, \sigma_n > 0$  sorted in decreasing order, it holds that  $\sigma_r > \sigma_{r+1}$ . Then the reduced-order model  $[A_{11}, B_1, C_1, D]$  is asymptotically stable, minimal, passive, and positive real balanced with the positive real Gramians  $P_{11} = Q_{11} = \Sigma_1 := \text{diag}(\sigma_1, \dots, \sigma_r)$ .

*Proof.* The preservation of passivity is easy to see, since it holds that  $\Sigma_1 > 0$  solves the reduced ARE

$$A_{11}^\top \tilde{P} + \tilde{P} A_{11} + (\tilde{P} B_1 - C_1^\top)(D + D^\top)^{-1}(B_1^\top \tilde{P} - C_1) = 0, \quad \tilde{P} = \tilde{P}^\top.$$

Therefore, the corresponding LMI (5.4) has a positive definite solution and thus by Theorem 5.4, the reduced transfer function  $G(s) = C_1(sI_r - A_{11})^{-1}B_1 + D$  is positive real. By Corollary 5.6, the reduced-order model is passive. The proof of asymptotic stability and minimality is quite involved and therefore, it is omitted here.  $\square$

Next we want to address error bounds. In contrast to standard balancing, there are no a priori error bounds in the  $\mathcal{H}_\infty$ -norm, even if the original and the reduced-order model both have transfer functions in  $\mathcal{RH}_\infty^{m \times m}$ . Instead, one has to resort to the so-called *gap metric*. The following has been taken from [GO13].

**Definition 5.11:** Let  $\mathcal{V}_1$  and  $\mathcal{V}_2$  be two closed subspaces (“closed” means that any sequence of elements of the subspace has its limit in this subspace) of a Hilbert space  $\mathcal{H}$  with induced norm  $\|\cdot\|_{\mathcal{H}}$ . Then the *gap* between  $\mathcal{V}_1$  and  $\mathcal{V}_2$  is defined by

$$g(\mathcal{V}_1, \mathcal{V}_2) := \|\Pi_{\mathcal{V}_1} - \Pi_{\mathcal{V}_2}\|_{\mathcal{L}(\mathcal{H}, \mathcal{H})},$$

where  $\Pi_{\mathcal{V}_1}, \Pi_{\mathcal{V}_2} : \mathcal{H} \rightarrow \mathcal{H}$  denote the orthogonal projections onto the spaces  $\mathcal{V}_1, \mathcal{V}_2$  (which exist by the closedness assumption).

It can be shown that  $g$  makes the set of all closed subspaces of  $\mathcal{H}$  to a (complete) metric space. Moreover, it can be shown that

$$g(\mathcal{V}_1, \mathcal{V}_2) = \max \{ \vec{g}(\mathcal{V}_1, \mathcal{V}_2), \vec{g}(\mathcal{V}_2, \mathcal{V}_1) \},$$

where

$$\vec{g}(\mathcal{V}_1, \mathcal{V}_2) = \|(I - \Pi_{\mathcal{V}_2})\Pi_{\mathcal{V}_1}\|_{\mathcal{L}(\mathcal{H}, \mathcal{H})} = \sup_{v \in \mathcal{V}_1, \|v\|_{\mathcal{H}}=1} \text{dist}(v, \mathcal{V}_2).$$

is the *directed gap*. Now we apply these concepts to spaces related to linear systems.

**Definition 5.12:** Let  $\Sigma := [A, B, C, D] \in \Sigma_{n,m,p}$  be a given asymptotically stable system with transfer function  $G \in \mathcal{RH}_\infty^{p \times m}$ . Let  $u \in \mathcal{L}_2([0, \infty), \mathbb{R}^m)$  be given and define  $U \in \mathcal{H}_2^m$  by  $U(s) := \mathcal{L}\{u\}(s)$ . Furthermore, define the *multiplication operator*

$$M_G : \mathcal{H}_2^m \rightarrow \mathcal{H}_2^p, \quad (M_G U)(s) = G(s)U(s) \quad \forall s \in \mathbb{C}^+.$$

The *graph* of the dynamical system  $\Sigma$  is then defined by

$$\mathcal{G}(\Sigma) := \left\{ \begin{bmatrix} I_m \\ M_G \end{bmatrix} U \mid U \in \mathcal{H}_2^m \right\} = \text{im} \begin{bmatrix} I_m \\ M_G \end{bmatrix},$$

which is a closed subspace of the Hilbert space  $\mathcal{H}_2^{p+m}$ .

The *gap metric* between two asymptotically stable systems  $\Sigma_1, \Sigma_2$  is defined by

$$\delta(\Sigma_1, \Sigma_2) := g(\mathcal{G}(\Sigma_1), \mathcal{G}(\Sigma_2)).$$

The gap metric can be interpreted as the distance of the subspaces of input/output trajectories generated by two dynamical systems in frequency domain. With the gap metric, we can now obtain the following error bounds.

**Theorem 5.13:** Let  $\Sigma := [A, B, C, D] \in \Sigma_{n,m,m}$  be an asymptotically stable, minimal, and passive system with transfer function  $G \in \mathcal{RH}_\infty^{m \times m}$  and let  $\tilde{\Sigma} := [A_{11}, B_1, C_1, D] \in \Sigma_{r,m,m}$  be the reduced-order model obtained by positive real balanced truncation with transfer function  $\tilde{G} \in \mathcal{RH}_\infty^{m \times m}$ . Assume further, that for the positive real characteristic values  $\sigma_1, \sigma_2, \dots, \sigma_n > 0$  sorted in decreasing order, it holds that  $\sigma_r > \sigma_{r+1}$ . Then we have the (a priori) gap metric error bound

$$\delta(\Sigma, \tilde{\Sigma}) \leq \sum_{j=r+1}^n \sigma_j.$$

Moreover, there is the (a posteriori)  $\mathcal{H}_\infty$  error bound

$$\|G - \tilde{G}\|_{\mathcal{H}_\infty} \leq 2 \min \left\{ \left(1 + \|G\|_{\mathcal{H}_\infty}^2\right) \left(1 + \|\tilde{G}\|_{\mathcal{H}_\infty}\right), \right. \\ \left. \left(1 + \|G\|_{\mathcal{H}_\infty}\right) \left(1 + \|\tilde{G}\|_{\mathcal{H}_\infty}^2\right) \right\} \sum_{j=r+1}^n \sigma_j.$$

*Proof.* Omitted. □

**Remark 5.14:** a) The  $\mathcal{H}_\infty$  error bound is only an a posteriori error bound, since it requires the knowledge of the reduced-order model. Therefore, its use in practice is limited.

b) The gap metric can be expressed by the *normalized coprime factorizations* of  $G(s)$  and  $\tilde{G}(s)$ , see [Geo88]. A (right) normalized coprime factorization is given by  $G(s) = N(s)M(s)^{-1}$ , where  $\begin{bmatrix} M \\ N \end{bmatrix} \in \mathcal{RH}_\infty^{(p+m) \times m}$  and there exist  $Y \in \mathcal{RH}_\infty^{m \times m}$ ,  $Z \in \mathcal{RH}_\infty^{m \times p}$  such that the *Bézout identity*

$$Y(s)M(s) + Z(s)N(s) = I_m$$

with the normalization condition

$$M^\sim(s)M(s) + N^\sim(s)N(s) = I_m$$

is satisfied. If  $G(s) = N(s)M(s)^{-1}$  and  $\tilde{G}(s) = \tilde{N}(s)\tilde{M}(s)^{-1}$  are normalized coprime factorizations of  $G(s)$  and  $\tilde{G}(s)$ , respectively, then for the directed gap we have

$$\vec{g}(\mathcal{G}(\Sigma), \mathcal{G}(\tilde{\Sigma})) = \inf_{H \in \mathcal{H}_\infty^{m \times m}} \left\| \begin{bmatrix} M \\ N \end{bmatrix} - \begin{bmatrix} \tilde{M} \\ \tilde{N} \end{bmatrix} H \right\|_{\mathcal{H}_\infty}.$$

From this, some bounds for the gap metric can be derived. Efficient methods for its computation however, seem to be widely unexplored, except for [Geo88].

c) The gap metric error bound can also be expressed by the gap of subspaces of  $\mathcal{L}_2([0, \infty), \mathbb{R}^{p+m})$  in the time domain. This analysis makes use of the *behavior approach of systems theory* which was developed by Jan C. Willems in the early 90s.

## 5.4 Numerical Solution of Large-Scale Algebraic Riccati Equations

### 5.4.1 Newton's Method for Solving Algebraic Riccati Equations

In this section we discuss the numerical solution of algebraic Riccati equations of the form

$$A^\top P + PA + (PB - C^\top)(D + D^\top)^{-1}(B^\top P - C) = 0, \quad P = P^\top.$$

This ARE can be rewritten as

$$\begin{aligned}\mathcal{R}(P) &:= F + \hat{A}^\top P + P\hat{A} + PGP = 0, \quad P = P^\top, \quad \text{where} \\ \hat{A} &:= A - B(D + D^\top)^{-1}C, \\ F &:= C^\top(D + D^\top)^{-1}C \geq 0, \\ G &:= B(D + D^\top)^{-1}B^\top \geq 0.\end{aligned}\tag{5.10}$$

Remember that we want to compute the minimal solution  $P^-$  of (5.10). Under the assumption that the pair  $(A, B)$  is stabilizable (equivalently, the pair  $(\hat{A}, B)$  is stabilizable), it can even be shown that  $P^-$  is the unique *stabilizing solution* of (5.10), that is

$$\Lambda(\hat{A} + GP^-) = \Lambda\left(A - B(D + D^\top)^{-1}(C - B^\top P^-)\right) \subset \mathbb{C}^-.$$

We consider (5.10) as a nonlinear system of equations and apply Newton's method which has first been considered in [Kle68]. For this, we need to evaluate the (Fréchet) derivative of  $\mathcal{R}(P)$  with respect to  $P$ .

**Definition 5.15** (Fréchet differentiability, Fréchet derivative): Let  $(\mathcal{X}, \|\cdot\|_{\mathcal{X}})$  and  $(\mathcal{Y}, \|\cdot\|_{\mathcal{Y}})$  be two normed linear spaces and let  $\mathcal{U} \subset \mathcal{X}$  be an open subset. An operator  $\mathcal{F} : \mathcal{U} \rightarrow \mathcal{Y}$  is called *Fréchet differentiable* at  $X \in \mathcal{U}$  if there exists a bounded linear operator  $\mathcal{F}'(X) : \mathcal{X} \rightarrow \mathcal{Y}$  such that

$$\lim_{\|N\|_{\mathcal{X}} \rightarrow 0} \frac{1}{\|N\|_{\mathcal{X}}} \|\mathcal{F}(X + N) - \mathcal{F}(X) - (\mathcal{F}'(X))(N)\|_{\mathcal{Y}} = 0.$$

The operator  $\mathcal{F}'(X)$  is called *Fréchet derivative* of  $\mathcal{F}$  at  $X$ . The map  $\mathcal{F}' : \mathcal{U} \rightarrow \mathcal{L}(\mathcal{X}, \mathcal{Y})$  with  $X \mapsto \mathcal{F}'(X)$  is called *Fréchet derivative* of  $\mathcal{F}$  on  $\mathcal{U}$ .

Let us see whether  $\mathcal{R}(\cdot)$  is Fréchet differentiable and (if yes) determine its Fréchet derivative. If the Fréchet derivative exists it is given by

$$\begin{aligned}(\mathcal{R}'(P))(N) &= \lim_{h \rightarrow 0} \frac{1}{h} (\mathcal{R}(P + hN) - \mathcal{R}(P)) \\ &= \lim_{h \rightarrow 0} \frac{1}{h} \left( F + \hat{A}^\top(P + hN) + (P + hN)\hat{A} \right. \\ &\quad \left. + (P + hN)G(P + hN) - (F + \hat{A}^\top P + P\hat{A} + PGP) \right) \\ &= \lim_{h \rightarrow 0} \frac{1}{h} (h\hat{A}^\top N + hN\hat{A} + hPGN + hNGP + h^2NGN) \\ &= \lim_{h \rightarrow 0} (\hat{A}^\top N + N\hat{A} + PGN + NGP + hNGN) \\ &= \hat{A}^\top N + N\hat{A} + PGN + NGP \\ &= (\hat{A} + GP)^\top N + N(\hat{A} + GP).\end{aligned}$$

**Algorithm 5.2** Newton's method for the algebraic Riccati equation**Input:**  $\hat{A}$ ,  $F$ ,  $G$  as in (5.10) and initial value  $P_0$  such that  $\Lambda(\hat{A} + GP_0) \subset \mathbb{C}^-$ .**Output:** Stabilizing (and minimal) solution  $P^-$  solving (5.10).

- 1: **for**  $j = 1, 2, \dots$  **do**
- 2:   Set  $\hat{A}_j := \hat{A} + GP_{j-1}$ .
- 3:   Solve  $\hat{A}_j^\top N_{j-1} + N_{j-1} \hat{A}_j = -\mathcal{R}(P_{j-1})$  for  $N_{j-1}$ .
- 4:   Set  $P_j := P_{j-1} + N_{j-1}$ .
- 5: **end for**
- 6: Set  $P^- := P_j$ .

In other words, the Fréchet derivative of a Riccati operator is a Lyapunov operator. Now the Newton iteration is given by

$$(\mathcal{R}'(P_{j-1}))(N_{j-1}) = -\mathcal{R}(P_{j-1}), \quad P_j = P_{j-1} + N_{j-1}, \quad j = 1, 2, \dots$$

and the iteration is summarized in Algorithm 5.2. This formulation of the algorithm has the disadvantage that  $\mathcal{R}(P_{j-1})$  is evaluated in every iteration. Therefore, let us revisit the computation of the update  $N_{j-1}$ . We know that

$$\begin{aligned} (\hat{A} + GP_{j-1})^\top N_{j-1} + N_{j-1}(\hat{A} + GP_{j-1}) \\ = -F - \hat{A}^\top P_{j-1} - P_{j-1} \hat{A} - P_{j-1} G P_{j-1}. \end{aligned} \quad (5.11)$$

Plugging in  $N_{j-1} = P_j - P_{j-1}$  then gives

$$\begin{aligned} (\hat{A} + GP_{j-1})^\top (P_j - P_{j-1}) + (P_j - P_{j-1})(\hat{A} + GP_{j-1}) \\ = -F - \hat{A}^\top P_{j-1} - P_{j-1} \hat{A} - P_{j-1} G P_{j-1}. \end{aligned}$$

Some manipulations and rearrangements of the terms finally lead to

$$(\hat{A} + GP_{j-1})^\top P_j + P_j (\hat{A} + GP_{j-1}) = -F + P_{j-1} G P_{j-1}. \quad (5.12)$$

This leads to Kleinman's formulation of the Newton iteration which is given in Algorithm 5.3. The question arises whether Algorithm 5.3 converges to the right solution. The following theorem makes this clear, see also [LR95] for a proof.

**Theorem 5.16:** Consider the ARE (5.10) with stabilizable  $(A, B)$ . Let  $P^-$  be its unique stabilizing solution. Let further  $P_0 \in \mathbb{R}^{n \times n}$  be stabilizing, i. e., it holds that  $\Lambda(\hat{A} + GP_0) \subset \mathbb{C}^-$ . Then the iterates  $P_j$ ,  $j = 1, 2, \dots$  fulfill the following statements:

- a) The matrix  $P_j$  is stabilizing.
- b) It holds that  $P_1 \leq \dots \leq P_j \leq P_{j+1} \leq \dots \leq P^-$ .

---

**Algorithm 5.3** Newton-Kleinman iteration for the algebraic Riccati equation

---

**Input:**  $\hat{A}$ ,  $F$ ,  $G$  as in (5.10) and initial value  $P_0$  such that  $\Lambda(\hat{A} + GP_0) \subset \mathbb{C}^-$ .

**Output:** Stabilizing (and minimal) solution  $P^-$  solving (5.10).

- 1: **for**  $j = 1, 2, \dots$  **do**
  - 2:   Set  $\hat{A}_j := \hat{A} + GP_{j-1}$  and  $F_j := -F + P_{j-1}GP_{j-1}$ .
  - 3:   Solve  $\hat{A}_j^\top P_j + P_j \hat{A}_j = F_j$ .
  - 4: **end for**
  - 5: Set  $P^- := P_j$ .
- 

c) It holds that  $\lim_{j \rightarrow \infty} P_j = P^-$ .

d) The convergence is globally quadratic, i. e., there exists a constant  $\gamma > 0$  such that

$$\|P^- - P_j\| \leq \gamma \|P^- - P_{j-1}\|^2, \quad j = 1, 2, \dots$$

*Proof.* The proof of this theorem needs a few technical results from the solution theory of Lyapunov and Riccati equations. Therefore, we omit it here.  $\square$

**Remark 5.17:** a) If  $\hat{A}$  is not asymptotically stable (otherwise  $P_0 = 0$  is stabilizing), then the computation of a stabilizing  $P_0$  usually costs as much as another iteration step since this requires the solution of one additional Lyapunov equation (e. g., in Bass' algorithm).

b) It can be proven that if  $\Lambda(\hat{A} + GP_{j-1}) \subset \mathbb{C}^-$ , then it holds that  $\Lambda(\hat{A} + G(P_{j-1} + tN_{j-1})) \subset \mathbb{C}^-$  for all  $t \in [0, 2]$ . This motivates line search algorithms to optimize the step length after computing the direction  $N_{j-1}$  in Algorithm 5.2. This means that we set  $P_j := P_{j-1} + tN_{j-1}$  where  $t$  is chosen as

$$t = \operatorname{argmin}_{\tau \in [0, 2]} \|\mathcal{R}(P_{j-1} + \tau N_{j-1})\|_F.$$

The computation of  $t$  is usually much cheaper than the actual Newton step which can drastically accelerate the iteration [BB98].

---

### 5.4.2 The Low-Rank Newton-Kleinman Method

We now aim at applying the Newton-Kleinman iteration for large-scale AREs and derive a low-rank formulation. Recall that the ARE attains the form

$$\begin{aligned} \mathcal{R}(P) := & C^\top (D + D^\top)^{-1} C + \left( A - B(D + D^\top)^{-1} C \right)^\top P \\ & + P \left( A - B(D + D^\top)^{-1} C \right) + PB(D + D^\top)^{-1} B^\top P = 0, \end{aligned}$$

where  $A \in \mathbb{R}^{n \times n}$  is large and sparse,  $B, C^\top \in \mathbb{R}^{n \times m}$ ,  $D + D^\top > 0$ , and  $m \ll n$ . Thus, the constant and the quadratic term admit low-rank factorizations  $C^\top (D + D^\top)^{-1} C = \tilde{C}^\top \tilde{C}$  and  $B(D + D^\top)^{-1} B^\top = \tilde{B} \tilde{B}^\top$ . Inserting this into (5.12) gives the iteration scheme

$$\begin{aligned} & \left( A - B(D + D^\top)^{-1} (C - B^\top P_{j-1}) \right)^\top P_j + \\ & P_j \underbrace{\left( A - B(D + D^\top)^{-1} (C - B^\top P_{j-1}) \right)}_{=: \hat{A}_j} = -\tilde{C}^\top \tilde{C} + P_{j-1} \tilde{B} \tilde{B}^\top P_{j-1} \quad (5.13) \end{aligned}$$

Unfortunately, we cannot solve the Lyapunov equation (5.13) directly with the low-rank ADI method, since its right-hand side may be indefinite. However, we can split it into two Lyapunov equations

$$\begin{aligned} \hat{A}_j^\top P_{1,j} + P_{1,j} \hat{A}_j &= -\tilde{C}^\top \tilde{C}, \\ \hat{A}_j^\top P_{2,j} + P_{2,j} \hat{A}_j &= -P_{j-1} \tilde{B} \tilde{B}^\top P_{j-1}, \end{aligned}$$

where we start with the initial values  $P_{1,0} = P_0$  and  $P_{2,0} = 0$  [RS10]. Then we obtain the iterate  $P_j = P_{1,j} - P_{2,j}$  by linearity of the Lyapunov operator. This results in the *low-rank Newton-Kleinman method for algebraic Riccati equations*, see also [BS13].

One problem remains: even if  $A$  is sparse and  $B, C^\top$  are thin, the feedback

$$\hat{A}_j := A - B(D + D^\top)^{-1} \underbrace{(C - B^\top P_{j-1})}_{=: K_j} \quad (5.14)$$

is usually dense. This means that we should never explicitly form (5.14).

There are several ways to solve linear systems with the system matrix  $\hat{A}_j - p_j I_n$  efficiently in the low-rank ADI method:

- a) **Application of an iterative solver:** This option only requires multiplications with  $\hat{A}_j$ . Since  $B$  and  $K_j$  have only a few columns and rows, respectively, these can be carried out efficiently. On the other hand, the convergence of iterative solvers is often slow, as long as no good preconditioner is available.

b) **Application of the Sherman-Morrison-Woodbury identity:** It holds that

$$\left(A - B(D + D^T)^{-1}K_j\right)^{-1} = A^{-1} + A^{-1}B(D + D^T - K_jA^{-1}B)^{-1}K_jA^{-1}.$$

Then a linear system solve with  $\hat{A}_j$  only requires two sparse (block) solves with  $A$  and one small dense solve with the matrix  $D + D^T - K_jA^{-1}B$ .

---



## CHAPTER 6

---

### Interpolatory Model Reduction

---

In this section we will discuss interpolatory methods for model reduction. Here we take a more transfer function point of view and try to find reduced representations of it. This reduction is then normally done by evaluating the original transfer function (and its derivatives) in a number of points in the complex plane and then to construct a *rational interpolant* that match this information. Mostly, this rational interpolant has a realization of low order which will be our reduced-order model.

#### 6.1 Moment Matching

##### 6.1.1 Moments

Consider the LTI system  $[A, B, C, D] \in \Sigma_{n,m,p}$  with transfer function  $G(s) \in \mathbb{R}(s)^{p \times m}$ . Since  $G$  is rational and proper, its poles are contained in  $\Lambda(A)$ . Therefore,  $G$  is analytic in a neighborhood of all  $s_0 \in \mathbb{C} \setminus \Lambda(A)$ . Hence it can be locally expanded into a Taylor series at the expansion point  $s_0$ . Thus, for finite  $s_0$  we obtain

$$G(s) = \sum_{k=0}^{\infty} (s - s_0)^k M_k(s_0)$$

for some neighborhood of  $s_0$  and some matrices  $M_0(s_0), M_1(s_0), M_2(s_0), \dots$ . On the other hand, for  $s_0 = \infty$  we obtain the Laurent series

$$G(s) = \sum_{k=0}^{\infty} s^{-k} M_k(\infty).$$

The matrices  $M_k(s_0)$  are called the ( $k$ -th) *moments* of  $G$  at  $s_0$ . For the case  $s_0 = \infty$  they are also called the *Markov parameters* of the transfer function. Now we want to determine the moments. For this we need the following lemma.

**Lemma 6.1** (Neumann series): Let  $A \in \mathbb{C}^{n \times n}$  with spectral radius  $\rho(A) < 1$  be given. Then  $I_n - A$  is invertible and it holds that

$$(I_n - A)^{-1} = \sum_{k=0}^{\infty} A^k.$$

For finite  $s_0$  we have

$$\begin{aligned} G(s) &= C((s - s_0)I_n - A + s_0I_n)^{-1}B + D \\ &= C((s - s_0)(s_0I_n - A)^{-1} + I_n)^{-1}(s_0I_n - A)^{-1}B + D \\ &= C(I_n - (s_0 - s)(s_0I_n - A)^{-1})^{-1}(s_0I_n - A)^{-1}B + D \\ &= \sum_{k=0}^{\infty} C(s_0I_n - A)^{-k-1}B(s_0 - s)^k + D, \end{aligned}$$

where we have used, that by Lemma 6.1 and  $s$  sufficiently close to  $s_0$  it holds that

$$(I_n - (s_0 - s)(s_0I_n - A)^{-1})^{-1} = \sum_{k=0}^{\infty} (s_0 - s)^k (s_0I_n - A)^{-k}.$$

Therefore, we have

$$M_k(s_0) = \begin{cases} C(s_0I_n - A)^{-1}B + D, & \text{if } k = 0, \\ (-1)^k C(s_0I_n - A)^{-k-1}B, & \text{if } k \geq 1. \end{cases}$$

Moreover, for  $s_0 = \infty$  and sufficiently large  $s$ , we have

$$\begin{aligned} G(s) &= C(sI_n - A)^{-1}B + D \\ &= \frac{1}{s} C \underbrace{\left(I_n - \frac{1}{s}A\right)^{-1}}_{= \sum_{k=0}^{\infty} \frac{1}{s^k} A^k} B + D \\ &= \sum_{k=1}^{\infty} C A^{k-1} B \frac{1}{s^k} + D. \end{aligned}$$

Therefore, we have

$$M_k(\infty) = \begin{cases} D, & \text{if } k = 0, \\ CA^{k-1}B, & \text{if } k \geq 1. \end{cases}$$

Model reduction by moment matching consists of finding a reduced-order model  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] \in \Sigma_{r,m,p}$  with  $r \ll n$  such that for a given  $s_0 \in (\mathbb{C} \cup \{\infty\}) \setminus \Lambda(A)$  the corresponding moments  $\tilde{M}_k(s_0)$ ,  $k = 1, 2, \dots$  fulfill

$$\tilde{M}_k(s_0) = M_k(s_0), \quad k = 0, 1, \dots, \ell$$

for  $\ell$  as large as possible. The moment matching problem for  $s_0 = 0$  is also called *Padé approximation problem*, for  $s_0 = \infty$  it is also called the *partial realization problem*. Since for finite  $s_0$  we have

$$\begin{aligned} G(s) &= \sum_{k=0}^{\infty} (s - s_0)^k M_k(s_0) = \sum_{k=0}^{\ell} (s - s_0)^k M_k(s_0) + \mathcal{O}((s - s_0)^{\ell+1}), \\ \tilde{G}(s) &= \sum_{k=0}^{\infty} (s - s_0)^k \tilde{M}_k(s_0) = \sum_{k=0}^{\ell} (s - s_0)^k M_k(s_0) + \mathcal{O}((s - s_0)^{\ell+1}), \end{aligned}$$

we obtain

$$G(s) - \tilde{G}(s) = \mathcal{O}((s - s_0)^{\ell+1}).$$

Similarly, for  $s_0 = \infty$  we get

$$G(s) - \tilde{G}(s) = \mathcal{O}(s^{-(\ell+1)}).$$

### 6.1.2 One-Sided Moment Matching

For ease of notation we will now have a look at moment matching for SISO systems. In particular we discuss efficient ways to generate the reduced-order model without explicitly forming the moments.

**Definition 6.2:** Let  $[A, B, C, D] \in \Sigma_{n,m,p}$  be given. The *generalized controllability matrices* of  $[A, B, C, D]$  for  $s_0 \in (\mathbb{C} \cup \{\infty\}) \setminus \Lambda(A)$  are given by

$$\mathcal{C}_k(s_0) = \begin{bmatrix} (s_0 I_n - A)^{-1} B & (s_0 I_n - A)^{-2} B & \dots & (s_0 I_n - A)^{-k} B \end{bmatrix}, \\ k = 1, \dots, n,$$

if  $s_0 \in \mathbb{C} \setminus \Lambda(A)$ , and by

$$\mathcal{C}_k(\infty) = [B \quad AB \quad \dots \quad A^{k-1}B], \quad k = 1, \dots, n,$$

if  $s_0 = \infty$ .

Note that  $\mathcal{C}_n(\infty)$  is the Kalman controllability matrix that can be used for checking controllability of a linear system. We have the following lemma.

**Lemma 6.3:** Let  $s_0 \in (\mathbb{C} \cup \{\infty\}) \setminus \Lambda(A)$  be given. Then the pair  $(A, B) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times 1}$  is controllable if and only if it holds that  $\text{rank } \mathcal{C}_k(s_0) = k$  for all  $k = 1, 2, \dots, n$ .

*Proof.* Consider the case that  $s_0 = \infty$ . Then controllability of  $(A, B)$  is equivalent to

$$\text{rank } \mathcal{C}_n(\infty) = \text{rank} \begin{bmatrix} B & AB & \dots & A^{n-1}B \end{bmatrix} = n.$$

This is furthermore equivalent to

$$\text{rank } \mathcal{C}_k(\infty) = k, \quad k = 1, \dots, n.$$

On the other hand, for  $s_0 \in \mathbb{C} \setminus \Lambda(A)$ , the result follows as above by noticing that  $(A, B)$  is controllable, if and only if

$$\begin{aligned} \mathbb{C}^n &= \text{span} \{B, AB, \dots, A^{n-1}B\} \\ &= \text{span} \{B, (s_0 I_n - A)B, \dots, (s_0 I_n - A)^{n-1}B\} \\ &= (s_0 I_n - A)^{-n} \text{span} \{B, (s_0 I_n - A)B, \dots, (s_0 I_n - A)^{n-1}B\} \\ &= \text{span} \{(s_0 I_n - A)^{-1}B, (s_0 I_n - A)^{-2}B, \dots, (s_0 I_n - A)^{-n}B\}. \end{aligned}$$

□

These matrices determine projection matrices that result in reduced-order models with matched moments. In particular, we have the following theorem.

**Theorem 6.4:** Let  $[A, B, C, D] \in \Sigma_{n,1,1}$  be controllable with the moments  $M_k(s_0)$ ,  $k = 0, 1, \dots$  for some given  $s_0 \in (\mathbb{C} \cup \{\infty\}) \setminus \Lambda(A)$ . Assume that  $\ell \in \{1, 2, \dots, n\}$  and that  $\tilde{T} \in \mathbb{C}^{\ell \times \ell}$  is invertible. Set  $T := \mathcal{C}_\ell(s_0) \tilde{T} \in \mathbb{C}^{n \times \ell}$  and choose  $W \in \mathbb{C}^{n \times \ell}$  such that  $W^H T = I_\ell$ . Define  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] := [W^H A T, W^H B, C T, D] \in \Sigma_{\ell,1,1}$  with the moments  $\tilde{M}_k(s_0)$ ,  $k = 0, 1, \dots$ . Then the first  $\ell$  moments are matched, i. e., it holds that

$$M_k(s_0) = \tilde{M}_k(s_0), \quad k = 0, 1, \dots, \ell - 1.$$

*Proof.* We can assume w. l. o. g. that  $\tilde{T} = I_\ell$ , because the moments  $\tilde{M}_k(s_0)$ ,  $k = 1, 2, \dots$  are invariant under state-space transformations. Since the system is controllable, we have that  $\text{rank } T = \ell$  and there exists a matrix  $W \in \mathbb{C}^{n \times \ell}$

such that  $W^H T = I_\ell$ . Here we will only do the proof for finite  $s_0$ . The proof for  $s_0 = \infty$  is similar, and therefore, we omit it here. We have that

$$\begin{aligned} s_0 I_\ell - \tilde{A} &= W^H (s_0 I_n - A) T \\ &= W^H (s_0 I_n - A) \left[ (s_0 I_n - A)^{-1} B \quad \dots \quad (s_0 I_n - A)^{-\ell} B \right] \\ &= \left[ W^H B \quad W^H (s_0 I_n - A)^{-1} B \quad \dots \quad W^H (s_0 I_n - A)^{-\ell+1} B \right] \\ &=: \left[ \tilde{B} \quad e_1 \quad \dots \quad e_{\ell-1} \right]. \end{aligned}$$

Therefore, for  $k \in \{0, 1, \dots, \ell - 1\}$  we have

$$(s_0 I_\ell - \tilde{A})^{k+1} = \left[ * \quad \dots \quad * \quad \tilde{B} \quad e_1 \quad \dots \quad e_{\ell-1-k} \right],$$

which gives  $(s_0 I_\ell - \tilde{A})^{-(k+1)} \tilde{B} = e_{k+1}$ .

So, for the moments we have,

$$\begin{aligned} \tilde{M}_0(s_0) &= \tilde{C} (s_0 I_\ell - \tilde{A})^{-1} \tilde{B} + D \\ &= C T e_1 + D \\ &= C \mathcal{C}_\ell(s_0) e_1 + D \\ &= C (s_0 I_n - A)^{-1} B + D = M_0(s_0), \end{aligned}$$

and for  $k = 1, \dots, \ell - 1$  we obtain

$$\begin{aligned} \tilde{M}_k(s_0) &= (-1)^k \tilde{C} (s_0 I_\ell - \tilde{A})^{-(k+1)} \tilde{B} \\ &= (-1)^k C T e_{k+1} \\ &= (-1)^k C \mathcal{C}_\ell(s_0) e_{k+1} \\ &= (-1)^k C (s_0 I_n - A)^{-(k+1)} B = M_k(s_0). \end{aligned}$$

□

**Remark 6.5:** a) The assumption that  $[A, B, C, D]$  is controllable can be weakened in the sense that it is only required that  $\ell$  is small enough such that  $\text{rank } \mathcal{C}_\ell(s_0) = \ell$ .

b) On the other hand, if  $\ell = \text{rank } \mathcal{C}_n(s_0)$ , then the reduced system matches all moments.

We are free to choose the matrix  $\tilde{T}$  in Theorem 6.4. The simplest choice  $\tilde{T} = I_\ell$  unfortunately usually leads to very ill-conditioned projection matrices  $T$ , implying that the numerical computation of  $T$  is very sensitive to round-off errors. It is better to choose  $T$  such that it has orthonormal columns. Since  $\text{im } \mathcal{C}(s_0)$  is a Krylov subspace, i. e., a space of the general form

$$\mathcal{K}_\ell(F, v) := \text{span} \left\{ v, Fv, \dots, F^{\ell-1}v \right\},$$

$T$  can be efficiently computed by the (*shift-and-invert*) *Arnoldi method* (see the course on numerical linear algebra).

Since  $T$  has orthonormal columns, we can choose  $W = T$ . In the case  $s_0 = \infty$ , there are further simplifications. The Arnoldi iteration computes  $T \in \mathbb{C}^{n \times \ell}$  with orthonormal columns and a matrix  $H \in \mathbb{C}^{\ell \times \ell}$  in Hessenberg form such that

$$AT = TH + f_{\ell+1}e_\ell^\top$$

for some vector  $f_{\ell+1} \in \mathbb{C}^n$  with  $T^H f_{\ell+1} = 0$ . Thus we obtain

$$\tilde{A} = W^H AT = T^H (TH + f_{\ell+1}e_\ell^\top) = H.$$

Moreover, since  $B$  is the initial vector of the Krylov space  $\mathcal{K}_\ell(A, B)$ , we have

$$\tilde{B} = W^H B = T^H B = \|B\|_2 e_1.$$

Thus,  $\tilde{A}$  and  $\tilde{B}$  are obtained at no extra cost, only the matrix  $\tilde{C}$  has to be computed.

Note that stability and observability may be lost by moment matching as the following example shows.

**Example:** Consider the system  $[A, B, C, D] \in \Sigma_{2,1,1}$  with

$$A = \begin{bmatrix} -1 & 5 \\ 0 & -2 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad C = [1 \quad -1]. \quad (6.1)$$

It is easily checked that this system is stable and observable. For  $s_0 = \infty$  and  $\ell = 1$ , we obtain  $T = \frac{1}{\sqrt{2}}B = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = W$ . This gives

$$\begin{aligned} \tilde{A} &= \frac{1}{2} [1 \quad 1] \begin{bmatrix} -1 & 5 \\ 0 & -2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1, \\ \tilde{C} &= \frac{1}{\sqrt{2}} [1 \quad -1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0. \end{aligned}$$

Therefore, the reduced model is neither stable nor observable.

Note that stability and observability can be shown for symmetric systems with  $A = A^\top$  and  $B = C^\top$ .

By using the generalized observability matrices

$$\mathcal{O}_k(s_0) = \begin{bmatrix} C(s_0 I_n - A)^{-1} \\ C(s_0 I_n - A)^{-2} \\ \vdots \\ C(s_0 I_n - A)^{-k} \end{bmatrix}, \quad k = 1, \dots, n,$$

if  $s_0 \in \mathbb{C} \setminus \Lambda(A)$ , and

$$\mathcal{O}_k(\infty) = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{k-1} \end{bmatrix}, \quad k = 1, \dots, n,$$

if  $s_0 = \infty$ , we obtain the following result which is completely analogous to Theorem 6.4.

**Theorem 6.6:** Let  $[A, B, C, D] \in \Sigma_{n,1,1}$  be observable with the moments  $M_k(s_0)$ ,  $k = 0, 1, \dots$  for some given  $s_0 \in (\mathbb{C} \cup \{\infty\}) \setminus \Lambda(A)$ . Assume that  $\ell \in \{1, 2, \dots, n\}$  and that  $\widetilde{W} \in \mathbb{C}^{\ell \times \ell}$  is invertible. Set  $W := \mathcal{O}_\ell(s_0)^H \widetilde{W} \in \mathbb{C}^{n \times \ell}$  and choose  $T \in \mathbb{C}^{n \times \ell}$  such that  $W^H T = I_\ell$ . Define  $[\widetilde{A}, \widetilde{B}, \widetilde{C}, \widetilde{D}] := [W^H A T, W^H B, C T, D] \in \Sigma_{\ell,1,1}$  with the moments  $\widetilde{M}_k(s_0)$ ,  $k = 0, 1, \dots$ . Then the first  $\ell$  moments are matched, i. e., it holds that

$$M_k(s_0) = \widetilde{M}_k(s_0), \quad k = 0, 1, \dots, \ell - 1.$$

The matrix  $W$  can be computed by the Arnoldi algorithm applied to the matrix pair  $(A^T, C^T)$ . There are also variants of moment matching for MIMO systems which make use of the block Arnoldi method and for matching moments at several interpolation points  $s_0, s_1, \dots$

### 6.1.3 Two-Sided Moment Matching

We have just seen that choosing the projection matrices  $T$  and  $W$  such that  $W^H T = I_\ell$  and

- i)  $\text{im } T = \text{im } \mathcal{C}_\ell(s_0)$  or
- ii)  $\text{im } W = \text{im } \mathcal{O}_\ell(s_0)^H$

guarantees that  $\ell$  moment at  $s_0$  are matched. So what happens if we do both? We will show soon that in this case we can match  $2\ell$  moments. However, this is not always possible since a certain regularity condition is needed which we will prove now.

**Lemma 6.7:** Let  $[A, B, C, D] \in \Sigma_{n,1,1}$  be given with the generalized controllability and observability matrices  $\mathcal{C}_\ell(s_0)$  and  $\mathcal{O}_\ell(s_0)$  for some  $\ell \in \{1, 2, \dots, n\}$ . Then there exist matrices  $T, W \in \mathbb{C}^{n \times \ell}$  with  $W^H T = I_\ell$ ,  $\text{im } T = \text{im } \mathcal{C}_\ell(s_0)$ , and  $\text{im } W = \text{im } \mathcal{O}_\ell(s_0)^H$ , if and only if the matrix  $\mathcal{H}_\ell(s_0) := \mathcal{O}_\ell(s_0) \mathcal{C}_\ell(s_0)$  is invertible.

*Proof.* First we show “ $\Leftarrow$ ”: So assume that  $\mathcal{H}_\ell(s_0)$  is invertible. By choosing  $T := \mathcal{C}_\ell(s_0)$  and  $W := \mathcal{O}_\ell(s_0)^H \mathcal{H}_\ell(s_0)^{-H}$ , the conditions  $\text{im } T = \text{im } \mathcal{C}_\ell(s_0)$  and  $\text{im } W = \text{im } \mathcal{O}_\ell(s_0)^H$  are obviously satisfied. Moreover, we have

$$W^H T = \mathcal{H}_\ell(s_0)^{-1} \underbrace{\mathcal{O}_\ell(s_0) \mathcal{C}_\ell(s_0)}_{=\mathcal{H}_\ell(s_0)} = I_\ell.$$

Now we show “ $\Rightarrow$ ”: So let  $T, W \in \mathbb{C}^{n \times \ell}$  such that  $W^H T = I_\ell$ ,  $\text{im } T = \text{im } \mathcal{C}_\ell(s_0)$ , and  $\text{im } W = \text{im } \mathcal{O}_\ell(s_0)^H$ . Then there exist two invertible matrices  $K_C, K_O \in \mathbb{C}^{\ell \times \ell}$  such that  $T = \mathcal{C}_\ell(s_0) K_C$  and  $W = \mathcal{O}_\ell(s_0)^H K_O^H$ . Therefore, we get

$$I_\ell = W^H T = K_O \mathcal{O}_\ell(s_0) \mathcal{C}_\ell(s_0) K_C = K_O \mathcal{H}_\ell(s_0) K_C.$$

So  $\mathcal{H}_\ell(s_0) = K_O^{-1} K_C^{-1}$  is invertible.  $\square$

**Theorem 6.8:** Let  $[A, B, C, D] \in \Sigma_{n,1,1}$  be given with the moments  $M_k(s_0)$ ,  $k = 0, 1, \dots$  for some given  $s_0 \in (\mathbb{C} \cup \{\infty\}) \setminus \Lambda(A)$ . Let  $T, W \in \mathbb{C}^{n \times \ell}$  with  $W^H T = I_\ell$  and define the reduced-order model by  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] := [W^H A T, W^H B, C T, D] \in \Sigma_{\ell,1,1}$  with the moments  $\tilde{M}_k(s_0)$ ,  $k = 0, 1, \dots$ . If  $\text{im } \mathcal{C}_{\ell_1}(s_0) \subseteq \text{im } T$  and  $\text{im } \mathcal{O}_{\ell_2}(s_0)^H \subseteq \text{im } W$  for some  $\ell_1, \ell_2 \in \mathbb{N}_0$ , then it holds that

$$M_k(s_0) = \tilde{M}_k(s_0), \quad k = 0, 1, \dots, \begin{cases} \ell_1 + \ell_2 - 1, & \text{if } s_0 \in \mathbb{C} \setminus \Lambda(A), \\ \ell_1 + \ell_2, & \text{if } s_0 = \infty. \end{cases}$$

*Proof.* Here we prove the theorem only for the case  $s_0 \in \mathbb{C} \setminus \Lambda(A)$ , the case  $s_0 = \infty$  is analogous. We show the statement in several steps:

**Step 1:** We show that for any  $F \in \mathbb{C}^{\ell \times n}$  with  $FT = I_\ell$  we have  $TFv = v$  for all  $v \in \text{im } T$ . Indeed, if  $v \in \text{im } T$ , there exists a  $z \in \mathbb{C}^\ell$  such that  $v = Tz$ . This gives  $TFv = TFTz = Tz = v$ .

**Step 2:** We show that  $(s_0 I_n - A)^{-k} B = T(s_0 I_\ell - \tilde{A})^{-k} \tilde{B}$  for  $k = 1, \dots, \ell_1$ . First note that with  $F := (s_0 I_\ell - \tilde{A})^{-1} W^H (s_0 I_n - A)$  we have  $FT = (s_0 I_\ell - \tilde{A})^{-1} W^H (s_0 I_n - A) T = I_\ell$ , so we can use Step 1 here.

Now we prove the statement via induction. First we show the case  $k = 1$ : It

holds that

$$\begin{aligned}
T(s_0 I_\ell - \tilde{A})^{-1} \tilde{B} &= T(s_0 I_\ell - \tilde{A})^{-1} W^H B \\
&= T(s_0 I_\ell - \tilde{A})^{-1} W^H (s_0 I_n - A) (s_0 I_n - A)^{-1} B \\
&= TF \underbrace{(s_0 I_n - A)^{-1} B}_{\in \text{im } \mathcal{C}_{\ell_1}(s_0) \subseteq \text{im } T} \\
&= (s_0 I_n - A)^{-1} B.
\end{aligned}$$

Now assume that  $(s_0 I_n - A)^{-k} B = T(s_0 I_\ell - \tilde{A})^{-k} \tilde{B}$ . Then we have

$$\begin{aligned}
T(s_0 I_\ell - \tilde{A})^{-(k+1)} \tilde{B} &= T(s_0 I_\ell - \tilde{A})^{-1} W^H \underbrace{T(s_0 I_\ell - \tilde{A})^{-k} \tilde{B}}_{=(s_0 I_n - A)^{-k} B} \\
&= T(s_0 I_\ell - \tilde{A})^{-1} W^H (s_0 I_n - A) (s_0 I_n - A)^{-(k+1)} B \\
&= TF \underbrace{(s_0 I_n - A)^{-(k+1)} B}_{\in \text{im } \mathcal{C}_{\ell_1}(s_0) \subseteq \text{im } T, \text{ if } k+1 \leq \ell_1} = (s_0 I_n - A)^{-(k+1)} B.
\end{aligned}$$

**Step 3:** We show that  $C(s_0 I_n - A)^{-k} = \tilde{C}(s_0 I_\ell - \tilde{A})^{-k} W^H$  for  $k = 1, 2, \dots, \ell_2$ . This statement can be analogously proven as in Step 2 by replacing  $A$  by  $A^\top$ ,  $\tilde{A}$  by  $\tilde{A}^H$ ,  $B$  by  $C^\top$ ,  $\tilde{B}$  by  $\tilde{C}^H$ ,  $T$  by  $W$ , and  $W$  by  $T$ .

**Step 4:** We show that  $M_k(s_0) = \tilde{M}_k(s_0)$  for  $k = 0, 1, \dots, \ell_1 + \ell_2 - 1$ . For  $k = 0$  we obtain

$$\begin{aligned}
\tilde{M}_0(s_0) &= \tilde{C}(s_0 I_\ell - \tilde{A})^{-1} \tilde{B} + D = C \underbrace{T(s_0 I_\ell - \tilde{A})^{-1} \tilde{B}}_{=(s_0 I_n - A)^{-1} B} + D \\
&= C(s_0 I_n - A)^{-1} B + D = M_0(s_0),
\end{aligned}$$

where we have used Step 2. Now for  $k \in \{1, 2, \dots, \ell_1 + \ell_2 - 1\}$  let  $k_1 \in \{1, 2, \dots, \ell_1\}$  and  $k_2 \in \{1, 2, \dots, \ell_2\}$  be such that  $k_1 + k_2 = k + 1$ . Then we have

$$\begin{aligned}
\tilde{M}_k(s_0) &= (-1)^k \tilde{C}(s_0 I_\ell - \tilde{A})^{-(k+1)} \tilde{B} \\
&= (-1)^k \tilde{C}(s_0 I_\ell - \tilde{A})^{-k_2} W^H T(s_0 I_\ell - \tilde{A})^{-k_1} \tilde{B} \\
&= (-1)^k C(s_0 I_n - A)^{-k_2} (s_0 I_n - A)^{-k_1} B = M_k(s_0).
\end{aligned}$$

This completes the proof.  $\square$

**Remark 6.9:** a) Most often, the reduced-order model is computed for  $\ell_1 = \ell_2 = \ell$ . In case the matrix  $\mathcal{H}_\ell(s_0)$  is invertible, then the corresponding projection matrices  $T$  and  $W$  can be computed by the nonsymmetric Lanczos process.

b) If the system  $[A, B, C, D] \in \Sigma_{n,m,p}$  is a MIMO system, then the nonsymmetric block-Lanczos method can be used. Then we also have that

$$\begin{aligned} \text{span} \left\{ (s_0 I_n - A)^{-1} B, \dots, (s_0 I_n - A)^{-\ell} B \right\} &\subseteq \text{im } T, \\ \text{span} \left\{ (s_0 I_n - A)^{-\text{H}} C^{\text{H}}, \dots, \left( (s_0 I_n - A)^{\ell} \right)^{-\text{H}} C^{\text{H}} \right\} &\subseteq \text{im } W \end{aligned}$$

implies

$$M_k(s_0) = \widetilde{M}_k(s_0), \quad k = 0, 1, \dots, 2\ell - 1,$$

but the projection spaces must have a larger dimension to match the same order of moments. Alternatively, *tangential interpolation* techniques can be employed which yield moment matching properties in certain tangential directions, but normally also give smaller projection matrices. This will be discussed in more detail later.

Algorithm 6.1 employs the nonsymmetric Lanczos algorithm [Saa82] for moment matching. In the literature it is known as the *Padé-via-Lanczos (PVL) algorithm*.

The quantities computed in Algorithm 6.1 have the following properties (exercise!):

- a) The matrices  $W$  and  $T$  are biorthogonal, i. e.,  $W^{\text{H}}T = I_{\ell}$ .
- b) It holds that

$$\begin{aligned} (s_0 I_n - A)^{-1} T &= TH + [0 \ \dots \ 0 \ t_{\ell+1}] \beta_{\ell+1}, \\ (s_0 I_n - A)^{-\text{H}} W &= WH^{\text{H}} + [0 \ \dots \ 0 \ w_{\ell+1}] \overline{\gamma_{\ell+1}}, \end{aligned}$$

where

$$H = \begin{bmatrix} \alpha_1 & \gamma_2 & & & \\ \beta_2 & \alpha_2 & \ddots & & \\ & \ddots & \ddots & \gamma_{\ell} & \\ & & \beta_{\ell} & \alpha_{\ell} & \end{bmatrix}.$$

- c) We have that

$$\text{im } T = \text{im } \mathcal{C}_{\ell}(s_0), \quad \text{im } W = \text{im } \mathcal{O}_{\ell}(s_0)^{\text{H}}.$$

**Remark 6.10:** In Algorithm 6.1 breakdowns may occur. It can be shown that this is the case if and only if the matrix  $\mathcal{H}_k(s_0)$  is singular for some  $k \leq \ell$ . With regard to Lemma 6.7 this means that we cannot construct matrices  $W, T \in \mathbb{C}^{n \times k}$  with  $W^{\text{H}}T = I_k$  such that  $\text{im } T = \text{im } \mathcal{C}_{\ell}(s_0)$  and  $\text{im } W = \text{im } \mathcal{O}_{\ell}(s_0)^{\text{H}}$ .

**Algorithm 6.1** Padé-via-Lanczos for moment matching**Input:**  $[A, B, C, D] \in \Sigma_{n,1,1}$ ,  $s_0 \in \mathbb{C} \setminus \Lambda(A)$ , reduced order  $\ell$ .**Output:** Reduced-order model  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] = [W^H A T, W^H B, C T, D] \in \Sigma_{\ell,1,1}$  with  $M_k(s_0) = \tilde{M}_k(s_0)$  for  $k = 0, 1, \dots, 2\ell - 1$ .

- 1: **if**  $C(s_0 I_n - A)^{-2} B = 0$  **then**
- 2:   Breakdown.
- 3: **else**
- 4:   Set  $t_1 := \frac{1}{\|(s_0 I_n - A)^{-1} B\|_2} (s_0 I_n - A)^{-1} B$ ,  $w_1 := \frac{\|(s_0 I_n - A)^{-1} B\|_2}{C(s_0 I_n - A)^{-2} B} (s_0 I_n - A)^{-H} C^H$ .
- 5:   Set  $\alpha_1 := w_1^H (s_0 I_n - A)^{-1} t_1$ ,  $t_0 := 0$ ,  $w_0 := 0$ ,  $\beta_1 := 0$ ,  $\gamma_1 := 0$ .
- 6: **end if**
- 7: **for**  $k = 1, 2, \dots, \ell$  **do**
- 8:   Set  $r_k := (s_0 I_n - A)^{-1} t_k - \alpha_k t_k - \gamma_k t_{k-1}$ .
- 9:   Set  $s_k := (s_0 I_n - A)^{-H} w_k - \bar{\alpha}_k w_k - \beta_k w_{k-1}$ .
- 10:   **if**  $s_k^H r_k = 0$  **then**
- 11:     Breakdown.
- 12:   **else**
- 13:     Set  $\beta_{k+1} := \|r_k\|_2$ .
- 14:     Set  $\gamma_{k+1} := \frac{1}{\beta_{k+1}} s_k^H r_k$ .
- 15:     Set  $t_{k+1} = \frac{1}{\beta_{k+1}} r_k$ .
- 16:     Set  $w_{k+1} = \frac{1}{\gamma_{k+1}} s_k$ .
- 17:     Set  $\alpha_{k+1} := w_{k+1}^H (s_0 I_n - A)^{-1} t_{k+1}$ .
- 18:   **end if**
- 19: **end for**
- 20: Construct the reduced-order model as

$$\tilde{A} := W^H A T, \quad \tilde{B} := W^H B, \quad \tilde{C} := C T, \quad \tilde{D} := D,$$

where

$$W := [w_1 \ \dots \ w_\ell], \quad T := [t_1 \ \dots \ t_\ell].$$

We conclude this subsection with some general remarks on the difficulties of moment matching model reduction.

**Remark 6.11:** a) Moment matching is a comparably cheap method for model reduction, since for each  $s_0$ , only one (sparse) LU factorization or one preconditioner has to be computed and stored.

b) Computable error estimates or bounds are often very pessimistic or expensive to evaluate. Moreover, there is only a good approximation quality close to the expansion point  $s_0$ .

- c) The expansion point  $s_0$  can often only be chosen heuristically by using a priori knowledge from the system under consideration.
- d) Preservation of physical properties such as stability or passivity can only be guaranteed in special cases. Usually, a post-processing to restore these properties would be required, but this may destroy the moment matching properties.

In the following section we will discuss optimal choices for interpolation points which will cure many of the above mentioned problems.

## 6.2 $\mathcal{H}_2$ -Optimal Interpolation: The Iterative Rational Krylov Algorithm

This section is based on [GAB08]. Consider an LTI system  $[A, B, C, D] \in \Sigma_{n,1,1}$ . From the results of the previous section it is known that choosing  $W, T \in \mathbb{C}^{n \times \ell}$  such that  $W^H T = I_\ell$  with

$$\begin{aligned} \text{im } T &= \text{im} \left[ (s_1 I_n - A)^{-1} B \quad \dots \quad (s_\ell I_n - A)^{-1} B \right], \\ \text{im } W &= \text{im} \left[ (s_1 I_n - A)^{-H} C^T \quad \dots \quad (s_\ell I_n - A)^{-H} C^T \right] \end{aligned}$$

will lead to a reduced-order model  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] := [W^H A T, W^H B, C T, D] \in \Sigma_{\ell,1,1}$  with

$$G(s_k) = \tilde{G}(s_k), \quad G'(s_k) = \tilde{G}'(s_k), \quad k = 1, 2, \dots, \ell.$$

The question arises whether for fixed  $\ell$  we can choose the interpolation points  $s_1, \dots, s_\ell$  such that  $\|G - \tilde{G}\|$  can be minimized in a suitable system norm. It turns out that we can determine reduced-order models which fulfill such an optimality criterion locally in the  $\mathcal{H}_2$ -norm. Recall from Chapter 2 that for SISO systems we have that

$$\|G - \tilde{G}\|_{\mathcal{H}_2} = \sup_{\substack{u \in \mathcal{L}_2([0, \infty), \mathbb{R}) \\ u \neq 0}} \frac{\|y - \tilde{y}\|_{\mathcal{L}_\infty}}{\|u\|_{\mathcal{L}_2}},$$

where  $y, \tilde{y} \in \mathcal{L}_\infty([0, \infty), \mathbb{R})$  are the outputs of the full and reduced-order models obtained by feeding in the input  $u \in \mathcal{L}_2([0, \infty), \mathbb{R})$ . Further recall that

$$\|G - \tilde{G}\|_{\mathcal{H}_2} = \left( \frac{1}{2\pi} \int_{-\infty}^{\infty} |G(i\omega) - \tilde{G}(i\omega)|^2 d\omega \right)^{1/2}$$

for  $G - \tilde{G} \in \mathcal{RH}_2$  which is usually the case, since we have  $D = \tilde{D}$  in projection-based model reduction.

Finding a global optimum is hard, but we are satisfied with a local optimum which can be obtained by checking first-order (necessary) optimality conditions. For this we will need the  $\mathcal{H}_2$  inner product, which for SISO systems  $G, H \in \mathcal{RH}_2$  reduces to

$$\langle G, H \rangle_{\mathcal{H}_2} = \frac{1}{2\pi} \int_{-\infty}^{\infty} G(-i\omega)H(i\omega)d\omega.$$

Now we look for a different way of evaluating this inner product. For this we need the residues of a function  $G(s) \in \mathbb{R}(s)$  with possibly multiple poles. The residue of  $G$  at pole  $\lambda \in \mathbb{C}$  of order  $m$  is given by

$$\text{res}(G, \lambda) := \frac{1}{(m-1)!} \lim_{s \rightarrow \lambda} \frac{d^{m-1}}{ds^{m-1}} ((s-\lambda)^m G(s)).$$

Note that if  $G$  is proper and the pole  $\lambda$  is simple, then the residues can be computed as in Lemma 3.4. Residues are important tools in complex analysis, for instance for the evaluation of integrals. One of the key results in complex analysis is the *residue theorem* which we state next (in a simplified version).

**Theorem 6.12** (Residue theorem): Let  $f : \mathcal{D} \rightarrow \mathbb{C}$  be a meromorphic function (i. e., a function holomorphic in almost all points in  $\mathcal{D}$ ),  $\Gamma \subset \mathcal{D}$  be a simple closed contour in counterclockwise orientation that does not contain a pole of  $f$ , and let  $\mu_1, \dots, \mu_k$  be the poles of  $f$  inside the contour  $\Gamma$ . Then we have

$$\int_{\Gamma} f(s)ds = 2\pi i \sum_{j=1}^k \text{res}(f, \mu_j).$$

**Lemma 6.13:** Let  $G, H \in \mathcal{RH}_2$  and let  $\lambda_1, \dots, \lambda_p$  be the poles of  $G$  and  $\mu_1, \dots, \mu_q$  be the poles of  $H$ . Then it holds that

$$\langle G, H \rangle_{\mathcal{H}_2} = \sum_{j=1}^q \text{res}(G(-\cdot)H(\cdot), \mu_j) = \sum_{j=1}^p \text{res}(H(-\cdot)G(\cdot), \lambda_j).$$

If  $\mu_j$  is a simple pole of  $H$ , then it holds that

$$\text{res}(G(-\cdot)H(\cdot), \mu_j) = G(-\mu_j) \text{res}(H, \mu_j).$$

Moreover, if  $\mu_j$  is a double pole of  $H$ , then we have

$$\text{res}(G(-\cdot)H(\cdot), \mu_j) = G(-\mu_j) \text{res}(H, \mu_j) - G'(-\mu_j)h_0(\mu_j),$$

where  $h_0(\mu_j) = \lim_{s \rightarrow \mu_j} ((s - \mu_j)^2 H(s))$ .

*Proof.* First note that the function  $G(-s)H(s)$  may have poles only in  $\mu_1, \dots, \mu_q$  and  $-\lambda_1, \dots, -\lambda_p$ . For any  $R > 0$  define the semicircular contour

$$\Gamma_R := \underbrace{\{z \mid z = i\omega \text{ with } \omega \in [-R, R]\}}_{=: \Gamma_{1,R}} \cup \underbrace{\{z \mid z = Re^{i\theta} \text{ with } \theta \in [\frac{\pi}{2}, \frac{3\pi}{2}]\}}_{=: \Gamma_{2,R}}.$$

For sufficiently large  $R$ , all poles of  $H$  are inside  $\Gamma_R$ . Since  $G, H \in \mathcal{RH}_2$ , it holds that  $\lim_{|s| \rightarrow \infty} G(s) = \lim_{|s| \rightarrow \infty} H(s) = 0$ , and thus we have  $G(-s)H(s) = \mathcal{O}(|s|^{-2})$  for  $|s| \rightarrow \infty$ . Therefore,  $G(-s)H(s)$  decays faster than the length of the contour  $\Gamma_{2,R}$  grows for  $R \rightarrow \infty$  and thus we have

$$\lim_{R \rightarrow \infty} \int_{\Gamma_{2,R}} G(-s)H(s)ds = 0.$$

By using the residue theorem, we obtain

$$\begin{aligned} \langle G, H \rangle_{\mathcal{H}_2} &= \frac{1}{2\pi} \int_{-\infty}^{\infty} G(-i\omega)H(i\omega)d\omega \\ &= \frac{1}{2\pi i} \lim_{R \rightarrow \infty} \int_{\Gamma_R} G(-s)H(s)ds \\ &= \sum_{j=1}^q \text{res}(G(-\cdot)H(\cdot), \mu_j). \end{aligned}$$

(Note that if a pole  $\mu_j$  of  $H$  is canceled by a zero of  $G(-\cdot)$ , then the above formula is still correct, since then  $\text{res}(G(-\cdot)H(\cdot), \mu_j) = 0$ .) If  $\mu_j$  is a simple pole of  $H$ , then if  $G(-\mu_j) \neq 0$ , it is also a simple pole of  $G(-\cdot)H(\cdot)$  and we get

$$\begin{aligned} \text{res}(G(-\cdot)H(\cdot), \mu_j) &= \lim_{s \rightarrow \mu_j} (s - \mu_j)G(-s)H(s) \\ &= G(-\mu_j) \lim_{s \rightarrow \mu_j} (s - \mu_j)H(s) = G(-\mu_j) \text{res}(H, \mu_j). \end{aligned}$$

On the other hand, if  $G(-\mu_j) = 0$ , then  $G(-\cdot)H(\cdot)$  has no pole at  $\mu_j$  and

$$0 = \text{res}(G(-\cdot)H(\cdot), \mu_j) = G(-\mu_j) \text{res}(H, \mu_j).$$

If  $\mu_j$  is a double pole of  $H$ , then  $G(-\cdot)H(\cdot)$  has no pole, a simple pole, or a double pole at  $\mu_j$ . With

$$\begin{aligned} H(s) &= \frac{H_{-2}}{(s - \mu_j)^2} + \frac{H_{-1}}{s - \mu_j} + \dots, \\ G(-s) &= G_0 + G_1(s - \mu_j) + \dots, \end{aligned}$$

where  $G_0 = G(-\mu_j)$  and  $G_1 = -G'(-\mu_j)$ , it holds that

$$\begin{aligned}
 & \text{res}(G(\cdot)H(\cdot), \mu_j) \\
 &= G_0 \cdot H_{-1} + G_1 \cdot H_{-2} \\
 &= G(-\mu_j) \lim_{s \rightarrow \mu_j} \frac{d}{ds} ((s - \mu_j)^2 H(s)) - G'(-\mu_j) \lim_{s \rightarrow \mu_j} (s - \mu_j)^2 H(s) \\
 &= G(-\mu_j) \text{res}(H, \mu_j) - G'(-\mu_j) h_0(\mu_j).
 \end{aligned}$$

□

Now we get back to the optimization problem. Assume for simplicity that  $G \in \mathcal{RH}_2$  (The case of non-strictly proper  $G$  can be handled similarly.) The question arises how we can check whether  $\tilde{G} \in \mathcal{RH}_2$  with McMillan degree  $\ell$  locally minimizes  $\|G - \tilde{G}\|_{\mathcal{H}_2}$ . First we attempt to answer an easier question: Is there a  $\hat{G} \in \mathcal{RH}_2$  of McMillan degree  $\ell$  with the same poles as  $\tilde{G}$  that yields a smaller  $\mathcal{H}_2$ -error norm?

**Theorem 6.14:** Let  $\boldsymbol{\mu} := \{\mu_1, \dots, \mu_\ell\} \subset \mathbb{C}^-$  be closed under complex conjugation. We define  $\mathcal{M}(\boldsymbol{\mu})$  to be the set of all functions in  $\mathcal{RH}_2$  of McMillan degree at most  $\ell$  whose poles are simple and contained in the set  $\boldsymbol{\mu}$ . Then there are the following facts:

- If  $H \in \mathcal{M}(\boldsymbol{\mu})$ , then  $H$  is the transfer function of a stable and minimal LTI system with state-space dimension at most  $\ell$ .
- The set  $\mathcal{M}(\boldsymbol{\mu})$  is an  $\ell$ -dimensional closed subspace of  $\mathcal{H}_2$ .
- The function  $\hat{G}$  solves the minimization problem

$$\|G - \hat{G}\|_{\mathcal{H}_2} = \min_{H \in \mathcal{M}(\boldsymbol{\mu})} \|G - H\|_{\mathcal{H}_2}, \quad (6.2)$$

if and only if

$$\langle G - \hat{G}, H \rangle_{\mathcal{H}_2} = 0 \quad \text{for all } H \in \mathcal{M}(\boldsymbol{\mu}).$$

Furthermore, the solution  $\hat{G}$  of the minimization problem (6.2) exists and is unique.

*Proof.* Statements a) and b) are easy to show (exercise!).

We only discuss statement c). Let  $\hat{G} \in \mathcal{M}(\boldsymbol{\mu})$  be such that  $\langle G - \hat{G}, H \rangle_{\mathcal{H}_2} = 0$

for all  $H \in \mathcal{M}(\boldsymbol{\mu})$ . Then for  $H \in \mathcal{M}(\boldsymbol{\mu})$  we get

$$\begin{aligned} \|G - H\|_{\mathcal{H}_2}^2 &= \|(G - \hat{G}) + (\hat{G} - H)\|_{\mathcal{H}_2}^2 \\ &= \|G - \hat{G}\|_{\mathcal{H}_2}^2 + \|\hat{G} - H\|_{\mathcal{H}_2}^2 + 2\langle G - \hat{G}, \hat{G} - H \rangle_{\mathcal{H}_2} \\ &= \|G - \hat{G}\|_{\mathcal{H}_2}^2 + \|\hat{G} - H\|_{\mathcal{H}_2}^2, \end{aligned}$$

where the latter equality follows from  $\hat{G} - H \in \mathcal{M}(\boldsymbol{\mu})$ . This proves that  $\hat{G}$  is the unique minimizer.

Let  $\hat{G}$  be a minimizer of (6.2), which is in  $\mathcal{M}(\boldsymbol{\mu})$  by the closedness property. Then we get

$$\|(\hat{G} + \varepsilon H) - G\|_{\mathcal{H}_2}^2 - \|G - \hat{G}\|_{\mathcal{H}_2}^2 = 2\varepsilon \langle \hat{G} - G, H \rangle_{\mathcal{H}_2} + \varepsilon^2 \|H\|_{\mathcal{H}_2}^2 \geq 0 \quad (6.3)$$

If  $\langle \hat{G} - G, H \rangle_{\mathcal{H}_2} \neq 0$  for some  $H \in \mathcal{M}(\boldsymbol{\mu})$ , then for sufficiently small  $\varepsilon > 0$  we have  $|2\varepsilon \langle \hat{G} - G, H \rangle_{\mathcal{H}_2}| > \varepsilon^2 \|H\|_{\mathcal{H}_2}^2$ . However, this is a contradiction to the nonnegativity of (6.3). Therefore, we have  $\langle \hat{G} - G, H \rangle_{\mathcal{H}_2} = 0$  for all  $H \in \mathcal{M}(\boldsymbol{\mu})$ .  $\square$

By Lemma 6.13 we obtain

$$0 = \langle G - \hat{G}, H \rangle_{\mathcal{H}_2} = \sum_{j=1}^{\ell} (G(-\mu_j) - \hat{G}(-\mu_j)) \text{res}(H, \mu_j)$$

for all  $H \in \mathcal{M}(\boldsymbol{\mu})$ , which is equivalent to

$$G(-\mu_j) = \hat{G}(-\mu_j), \quad j = 1, 2, \dots, \ell.$$

So in other words,  $\hat{G}$  has to interpolate  $G$  at the mirror images of its own poles. The next question is whether the poles  $\mu_1, \mu_2, \dots, \mu_\ell$  are optimal.

**Definition 6.15:** Let  $G \in \mathcal{RH}_2$  be given and let  $\delta(\cdot)$  denote the McMillan degree of a rational function. Then the function  $\tilde{G} \in \mathcal{RH}_2$  is a *local minimizer* of the minimization problem

$$\min_{H \in \mathcal{RH}_2, \delta(H) = \ell} \|G - H\|_{\mathcal{H}_2}, \quad (6.4)$$

if

$$\|G - \tilde{G}\|_{\mathcal{H}_2} \leq \|G - H\|_{\mathcal{H}_2}$$

is satisfied for all  $H \in \mathcal{RH}_2$  with  $\delta(H) = \ell$  and  $\|\tilde{G} - H\|_{\mathcal{H}_2} \leq \varepsilon$  for some  $\varepsilon > 0$ .

We will now derive necessary optimality conditions for  $\tilde{G}$  being a local minimizer of the minimization problem (6.4).

**Theorem 6.16:** If  $\tilde{G} \in \mathcal{RH}_2$  is a local minimizer of (6.4) and  $\tilde{G}$  has only simple poles  $\boldsymbol{\mu} := \{\mu_1, \dots, \mu_\ell\}$  (closed under complex conjugation), then it holds that

$$\langle G - \tilde{G}, \tilde{G} \cdot H_1 + H_2 \rangle_{\mathcal{H}_2} = 0$$

for all  $H_1, H_2 \in \mathcal{M}(\boldsymbol{\mu})$ .

*Proof.* Suppose that  $\tilde{G}_\varepsilon \in \mathcal{RH}_2$  with  $\delta(\tilde{G}_\varepsilon) = \ell$  and  $\|\tilde{G} - \tilde{G}_\varepsilon\|_{\mathcal{H}_2} \leq c \cdot \varepsilon$  for a sufficiently small  $\varepsilon > 0$  and some constant  $c > 0$ . Then it holds that

$$\begin{aligned} \|G - \tilde{G}\|_{\mathcal{H}_2}^2 &\leq \|G - \tilde{G}_\varepsilon\|_{\mathcal{H}_2}^2 \\ &= \|(G - \tilde{G}) + (\tilde{G} - \tilde{G}_\varepsilon)\|_{\mathcal{H}_2}^2 \\ &= \|G - \tilde{G}\|_{\mathcal{H}_2}^2 + 2\langle G - \tilde{G}, \tilde{G} - \tilde{G}_\varepsilon \rangle_{\mathcal{H}_2} + \|\tilde{G} - \tilde{G}_\varepsilon\|_{\mathcal{H}_2}^2. \end{aligned}$$

This implies that for all sufficiently small  $\varepsilon > 0$  we have that

$$0 \leq 2\langle G - \tilde{G}, \tilde{G} - \tilde{G}_\varepsilon \rangle_{\mathcal{H}_2} + \|\tilde{G} - \tilde{G}_\varepsilon\|_{\mathcal{H}_2}^2. \quad (6.5)$$

Let  $\mu_1, \dots, \mu_\ell$  be ordered such that  $\mu_1, \mu_2, \dots, \mu_{\ell_R}$  are the real poles and  $\mu_j = \alpha_j \pm i\beta_j$  with  $\beta_j > 0$  for  $j = \ell_R + 1, \dots, \ell_R + \ell_C$  are the nonreal poles. Then any  $H_1 \in \mathcal{M}(\boldsymbol{\mu})$  can be written as a partial fraction expansion as

$$H_1(s) = \sum_{j=1}^{\ell_R} \frac{\gamma_j}{s - \mu_j} + \sum_{j=\ell_R+1}^{\ell_R+\ell_C} \frac{\rho_j(s - \alpha_j) + \tau_j}{(s - \alpha_j)^2 + \beta_j^2}$$

for some  $\gamma_j, \rho_j, \tau_j \in \mathbb{R}$ . Thus we get

$$\begin{aligned} \langle G - \tilde{G}, \tilde{G} \cdot H_1 + H_2 \rangle_{\mathcal{H}_2} &= \sum_{j=1}^{\ell_R} \gamma_j \left\langle G - \tilde{G}, \frac{\tilde{G}}{\cdot - \mu_j} \right\rangle_{\mathcal{H}_2} \\ &\quad + \sum_{j=\ell_R+1}^{\ell_R+\ell_C} \rho_j \left\langle G - \tilde{G}, \frac{(\cdot - \alpha_j)\tilde{G}}{(\cdot - \alpha_j)^2 + \beta_j^2} \right\rangle_{\mathcal{H}_2} \\ &\quad + \sum_{j=\ell_R+1}^{\ell_R+\ell_C} \tau_j \left\langle G - \tilde{G}, \frac{\tilde{G}}{(\cdot - \alpha_j)^2 + \beta_j^2} \right\rangle_{\mathcal{H}_2} + \langle G - \tilde{G}, H_2 \rangle_{\mathcal{H}_2}. \end{aligned}$$

By Theorem 6.14, we have  $\langle G - \tilde{G}, H_2 \rangle_{\mathcal{H}_2} = 0$ . Now we show that also the other summands vanish. Assume that for some  $j \in \{1, 2, \dots, \ell_R\}$  we have that

$$\left\langle G - \tilde{G}, \frac{\tilde{G}}{\cdot - \mu_j} \right\rangle_{\mathcal{H}_2} \neq 0.$$

We write  $\tilde{G}(s) = \frac{p(s)}{(s-\mu_j)q(s)}$  for some polynomials  $p(s), q(s) \in \mathbb{R}[s]$  and define

$$\tilde{G}_\varepsilon(s) := \frac{p(s)}{(s-\mu_j - (\pm\varepsilon))q(s)},$$

where the sign of  $\pm\varepsilon$  is chosen to match the one of  $\left\langle G - \tilde{G}, \frac{\tilde{G}}{\cdot - \mu_j} \right\rangle_{\mathcal{H}_2}$ . Then we have

$$\tilde{G}_\varepsilon(s) = \tilde{G}(s) \pm \varepsilon \frac{p(s)}{(s-\mu_j)^2 q(s)} + \mathcal{O}(\varepsilon^2),$$

which leads to

$$\tilde{G}(s) - \tilde{G}_\varepsilon(s) = \mp \varepsilon \frac{\tilde{G}(s)}{s-\mu_j} + \mathcal{O}(\varepsilon^2)$$

and

$$\langle G - \tilde{G}, \tilde{G} - \tilde{G}_\varepsilon \rangle_{\mathcal{H}_2} = -\varepsilon \underbrace{\left\langle G - \tilde{G}, \frac{\tilde{G}}{\cdot - \mu_j} \right\rangle_{\mathcal{H}_2}}_{=: \tilde{c} > 0} + \mathcal{O}(\varepsilon^2).$$

Together with (6.5) we get

$$0 \leq -2\tilde{c}\varepsilon + c^2\varepsilon^2 + \mathcal{O}(\varepsilon^2) = -2\tilde{c}\varepsilon + \mathcal{O}(\varepsilon^2).$$

But this is a contradiction for sufficiently small  $\varepsilon > 0$ .

Now write  $\tilde{G}(s) = \frac{p(s)}{((s-\alpha_j)^2 + \beta_j^2)q(s)}$  for some  $j \in \{\ell_R + 1, \ell_R + 2, \dots, \ell_R + \ell_C\}$  and  $p(s), q(s) \in \mathbb{R}[s]$ . In an analogous fashion, it can be shown that

$$\begin{aligned} \left\langle G - \tilde{G}, \frac{(\cdot - \alpha_j)\tilde{G}}{(\cdot - \alpha_j)^2 + \beta_j^2} \right\rangle_{\mathcal{H}_2} &= 0, \\ \left\langle G - \tilde{G}, \frac{\tilde{G}}{(\cdot - \alpha_j)^2 + \beta_j^2} \right\rangle_{\mathcal{H}_2} &= 0, \end{aligned}$$

by using

$$\begin{aligned} \tilde{G}_\varepsilon(s) &= \frac{p(s)}{((s-\alpha_j - (\pm\varepsilon))^2 + \beta_j^2)q(s)} \quad \text{and} \\ \tilde{G}_\varepsilon(s) &= \frac{p(s)}{((s-\alpha_j)^2 + \beta_j^2 - (\pm\varepsilon))q(s)}, \end{aligned}$$

respectively. From this we can conclude the claim of the theorem.  $\square$

The question arises whether we can interpret the above theorem as interpolation conditions. With Lemma 6.13 we obtain

$$0 = \langle G - \tilde{G}, \tilde{G} \cdot H_1 \rangle_{\mathcal{H}_2} = \sum_{j=1}^{\ell} \text{res} \left( (G(\cdot) - \tilde{G}(\cdot)) \tilde{G}(\cdot) H_1(\cdot), \mu_j \right).$$

If  $H_1$  has exactly the same poles as  $\tilde{G}$ , then  $\mu_j$  is a double pole of  $\tilde{G} \cdot H_1$  (recall that  $\tilde{G}, H_1 \in \mathcal{M}(\mu)$ ) and  $G(-\mu_j) = \tilde{G}(-\mu_j)$  for all  $j = 1, 2, \dots, \ell$ . This gives

$$\begin{aligned} 0 &= - \sum_{j=1}^{\ell} (G'(-\mu_j) - \tilde{G}'(-\mu_j)) \left( \lim_{s \rightarrow \mu_j} (s - \mu_j)^2 \tilde{G}(s) H_1(s) \right) \\ &= - \sum_{j=1}^{\ell} (G'(-\mu_j) - \tilde{G}'(-\mu_j)) \operatorname{res}(\tilde{G}, \mu_j) \operatorname{res}(H_1, \mu_j). \end{aligned}$$

Since this is true for all  $H_1 \in \mathcal{M}(\mu)$ , we obtain that

$$G'(-\mu_j) = \tilde{G}'(-\mu_j), \quad j = 1, 2, \dots, \ell.$$

The analysis carried out above can also be generalized to deal with transfer functions  $G(s) \in \mathbb{R}(s)^{p \times m}$  of MIMO systems. This will lead to *tangential interpolation conditions*, meaning that the original and the reduced transfer function moments only match in certain tangential directions. We will not discuss this in detail here, but summarize the result in the following theorem.

**Theorem 6.17:** Let an asymptotically stable system  $[A, B, C, D] \in \Sigma_{n,m,p}$  with transfer function  $G \in \mathcal{RH}_{\infty}^{p \times m}$  be given. Let

$$\tilde{G}(s) = \sum_{j=1}^{\ell} \frac{1}{s - \mu_j} \hat{c}_j \hat{b}_j^H + \tilde{D}$$

with  $\hat{b}_j \in \mathbb{C}^m$  and  $\hat{c}_j \in \mathbb{C}^p$ ,  $j = 1, 2, \dots, \ell$  be a local minimizer of the minimization problem

$$\min_{H \in \mathcal{RH}_{\infty}^{p \times m}, \delta(H) = \ell} \|G - H\|_{\mathcal{H}_2}.$$

Then the following statements are satisfied:

- a) It holds that  $D = \tilde{D}$ .
- b) It holds that

$$G(-\mu_j) \hat{b}_j = \tilde{G}(-\mu_j) \hat{b}_j, \quad \hat{c}_j^H G(-\mu_j) = \hat{c}_j^H \tilde{G}(-\mu_j), \quad j = 1, \dots, \ell.$$

- c) It holds that

$$\hat{c}_j^H G'(-\mu_j) \hat{b}_j = \hat{c}_j^H \tilde{G}'(-\mu_j) \hat{b}_j, \quad j = 1, \dots, \ell.$$

Now we want to put this into an algorithm. We want to compute a state-space realization of the reduced transfer function  $\tilde{G}$  without explicitly computing it. Given the interpolation points  $\hat{s}_j$  and the tangential directions  $\hat{b}_j$  and  $\hat{c}_j$  for  $j =$

1, 2, ...,  $\ell$ , we can construct a reduced-order model such that the tangential interpolation conditions

$$G(\hat{s}_j)\hat{b}_j = \tilde{G}(\hat{s}_j)\hat{b}_j, \quad \hat{c}_j^H G(\hat{s}_j) = \hat{c}_j^H \tilde{G}(\hat{s}_j), \quad j = 1, \dots, \ell$$

and

$$\hat{c}_j^H G'(\hat{s}_j)\hat{b}_j = \hat{c}_j^H \tilde{G}'(\hat{s}_j)\hat{b}_j, \quad j = 1, \dots, \ell.$$

are satisfied. We do this by two-sided interpolation with projection matrices  $W, T \in \mathbb{C}^{n \times \ell}$  with  $W^H T = I_\ell$  and

$$\begin{aligned} \text{im } T &= \text{im} \left[ (\hat{s}_1 I_n - A)^{-1} B \hat{b}_1 \quad \dots \quad (\hat{s}_\ell I_n - A)^{-1} B \hat{b}_\ell \right], \\ \text{im } W &= \text{im} \left[ (\hat{s}_1 I_n - A)^{-H} C^H \hat{c}_1 \quad \dots \quad (\hat{s}_\ell I_n - A)^{-H} C^H \hat{c}_\ell \right]. \end{aligned}$$

On the other hand, if a reduced-order system  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] \in \Sigma_{\ell, m, p}$  is given, we compute new interpolation data as follows: Assume that  $\tilde{A}$  is diagonalizable, i. e., there exists an invertible matrix  $X \in \mathbb{C}^{\ell \times \ell}$  such that

$$X^{-1} \tilde{A} X = M = \text{diag}(\mu_1, \dots, \mu_\ell).$$

Then we have

$$\tilde{G}(s) = \tilde{C}(sI_\ell - \tilde{A})^{-1} \tilde{B} + \tilde{D} = \tilde{C} X (sI_\ell - M)^{-1} X^{-1} \tilde{B} + \tilde{D}.$$

With

$$\begin{bmatrix} \hat{b}_1^H \\ \vdots \\ \hat{b}_\ell^H \end{bmatrix} := X^{-1} \tilde{B}, \quad [\hat{c}_1 \quad \dots \quad \hat{c}_\ell] := \tilde{C} X$$

we then obtain

$$\tilde{G}(s) = \sum_{j=1}^{\ell} \frac{1}{s - \mu_j} \hat{c}_j \hat{b}_j^H + \tilde{D}.$$

From this we can construct a new model as above that fulfills the tangential interpolation conditions at the data  $(-\mu_j, \hat{b}_j, \hat{c}_j)$ . This leads to an algorithm that alternates between the construction of a reduced-order model from interpolation data and the generation of new interpolation data from a given reduced-order model. If this iteration converges, then the tangential interpolation conditions are satisfied. This can be interpreted as a fixed-point iteration and results in Algorithm 6.2. The algorithm is usually terminated if the change in the pole locations between two consecutive iterations is below a given tolerance, or if the distance (e. g., measured in the  $\mathcal{H}_2$ -norm) of the reduced transfer functions between two consecutive iterations is sufficiently small.

---

**Algorithm 6.2** Iterative rational Krylov algorithm (IRKA)

**Input:** Asymptotically stable  $[A, B, C, D] \in \Sigma_{n,m,p}$ , desired reduced order  $\ell$ .

**Output:** Reduced-order model  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] \in \Sigma_{\ell,m,p}$ .

- 1: Choose initial interpolation data  $(\mu_j, \hat{b}_j, \hat{c}_j), j = 1, \dots, \ell$ .
- 2: **while** not converged **do**
- 3: Construct the model  $[\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}] := [W^H A T, W^H B, C T, D] \in \Sigma_{\ell,m,p}$  where  $W, T \in \mathbb{C}^{n \times \ell}$  with  $W^H T = I_\ell$  and

$$\begin{aligned} \text{im } T &= \text{im} \begin{bmatrix} (-\mu_1 I_n - A)^{-1} B \hat{b}_1 & \dots & (-\mu_\ell I_n - A)^{-1} B \hat{b}_\ell \end{bmatrix}, \\ \text{im } W &= \text{im} \begin{bmatrix} (-\mu_1 I_n - A)^{-H} C^H \hat{c}_1 & \dots & (-\mu_\ell I_n - A)^{-H} C^H \hat{c}_\ell \end{bmatrix}. \end{aligned}$$

- 4: Compute new interpolation data  $(\mu_j, \hat{b}_j, \hat{c}_j), j = 1, \dots, \ell$ : Compute an invertible matrix  $X \in \mathbb{C}^{\ell \times \ell}$  such that  $X^{-1} \tilde{A} X = \text{diag}(\mu_1, \dots, \mu_\ell)$  and set

$$\begin{bmatrix} \hat{b}_1^H \\ \vdots \\ \hat{b}_\ell^H \end{bmatrix} := X^{-1} \tilde{B}, \quad [\hat{c}_1 \quad \dots \quad \hat{c}_\ell] := \tilde{C} X, \quad \tilde{D} := D.$$

- 5: **end while**
- 

**Remark 6.18:** a) The initial interpolation data is either chosen randomly or by a cheaply computable reduced-order model such as a partial realization.

b) Stability of intermediate reduced-order models is not guaranteed. Unstable intermediate systems may occur, in particular, if the initial interpolation data is far away of the optimal one.

c) If the full-order model contains only real matrices, then also a real reduced-order model can be obtained by choosing the interpolation data such that it is closed under complex conjugation and by a realification of the projection matrices  $W$  and  $T$  by noting that

$$\begin{aligned} \text{im} \begin{bmatrix} (-\mu I_n - A)^{-1} B \hat{b} & (-\bar{\mu} I_n - A)^{-1} B \bar{\hat{b}} \end{bmatrix} \\ = \text{im} \left[ \text{Re} \left( (-\mu I_n - A)^{-1} B \hat{b} \right) \quad \text{Im} \left( (-\mu I_n - A)^{-1} B \hat{b} \right) \right]. \end{aligned}$$

d) Convergence of IRKA cannot be guaranteed, counter-examples are known.

e) In case of convergence, the resulting system may only be locally optimal, but in many cases, it even converges to a very good local or even global optimizer in only a few iterations. Unfortunately, a complete convergence

---

analysis of the method it still an open problem. Partial results are only known for symmetric systems.

- f) The reduced-order models obtained by IRKA are usually competitive to those obtained by balanced truncation.

### 6.3 Interpolation from Data: The Loewner Framework

In this section we discuss another framework for interpolation that is solely based on data obtained from the system under consideration, see [MA07]. In many situations it may be the case that a model is not available, but only some input/output data may given. Then we seek for a model of low order that best reproduces this data. Note that this process is rather called *reduced-order modeling* instead of model reduction, since no model is given. The methods discussed here are related to the field of system identification, where it is the goal to determine a model or its parameters from measurements, to validate it and to assess its robustness with respect to uncertainties in the data via statistical considerations.

Assume that we have given *right or column data*  $(\lambda_i, r_i, w_i) \in \mathbb{C} \times \mathbb{C}^m \times \mathbb{C}^p$ ,  $i = 1, 2, \dots, k$ , and *left or row data*  $(\mu_j, \ell_j, v_j) \in \mathbb{C} \times \mathbb{C}^p \times \mathbb{C}^m$ ,  $j = 1, 2, \dots, q$ . We seek a function  $\tilde{G}(s) \in \mathbb{C}(s)^{p \times m}$  such that

$$\begin{aligned}\tilde{G}(\lambda_i)r_i &= w_i, & i = 1, 2, \dots, k, \\ \ell_j^H \tilde{G}(\mu_j) &= v_j^H, & j = 1, 2, \dots, q.\end{aligned}$$

For simplicity we assume that  $\{\lambda_1, \dots, \lambda_k\} \cap \{\mu_1, \dots, \mu_q\} = \emptyset$ . We reorganize the right data as

$$\begin{aligned}\Lambda &:= \text{diag}(\lambda_1, \dots, \lambda_k) \in \mathbb{C}^{k \times k}, \\ R &:= [r_1 \ r_2 \ \dots \ r_k] \in \mathbb{C}^{m \times k}, \\ W &:= [w_1 \ w_2 \ \dots \ w_k] \in \mathbb{C}^{p \times k},\end{aligned}\tag{6.6}$$

and the left data as

$$\begin{aligned}M &:= \text{diag}(\mu_1, \dots, \mu_q) \in \mathbb{C}^{q \times q}, \\ L^H &:= [\ell_1 \ \ell_2 \ \dots \ \ell_q] \in \mathbb{C}^{p \times q}, \\ V^H &:= [v_1 \ v_2 \ \dots \ v_q] \in \mathbb{C}^{m \times q}.\end{aligned}\tag{6.7}$$

From the data we construct the *Loewner matrix*

$$\mathbb{L} := \begin{bmatrix} \frac{v_1^H r_1 - \ell_1^H w_1}{\mu_1 - \lambda_1} & \cdots & \frac{v_1^H r_k - \ell_1^H w_k}{\mu_1 - \lambda_k} \\ \vdots & \ddots & \vdots \\ \frac{v_q^H r_1 - \ell_q^H w_1}{\mu_q - \lambda_1} & \cdots & \frac{v_q^H r_k - \ell_q^H w_k}{\mu_q - \lambda_k} \end{bmatrix} \in \mathbb{C}^{q \times k}, \quad (6.8)$$

and the *shifted Loewner matrix*

$$\mathbb{L}_\sigma := \begin{bmatrix} \frac{\mu_1 v_1^H r_1 - \ell_1^H w_1 \lambda_1}{\mu_1 - \lambda_1} & \cdots & \frac{\mu_1 v_1^H r_k - \ell_1^H w_k \lambda_k}{\mu_1 - \lambda_k} \\ \vdots & \ddots & \vdots \\ \frac{\mu_q v_q^H r_1 - \ell_q^H w_1 \lambda_1}{\mu_q - \lambda_1} & \cdots & \frac{\mu_q v_q^H r_k - \ell_q^H w_k \lambda_k}{\mu_q - \lambda_k} \end{bmatrix} \in \mathbb{C}^{q \times k}. \quad (6.9)$$

The matrix pencil  $s\mathbb{L} - \mathbb{L}_\sigma \in \mathbb{C}[s]^{q \times k}$  is also called the *Loewner pencil*.

**Lemma 6.19:** With the notation introduced in (6.6), (6.7), (6.8), and (6.9), the two Sylvester equations

$$\mathbb{L}\Lambda - M\mathbb{L} = LW - VR \quad (6.10)$$

and

$$\mathbb{L}_\sigma\Lambda - M\mathbb{L}_\sigma = LW\Lambda - MVR \quad (6.11)$$

are satisfied.

*Proof.* By denoting with  $[\cdot]_{ij}$ , the  $(i, j)$ -th entry of the matrix in the brackets, we obtain that

$$\begin{aligned} [\mathbb{L}\Lambda - M\mathbb{L}]_{ij} &= \frac{v_i^H r_j - \ell_i^H w_j}{\mu_i - \lambda_j} \lambda_j - \mu_i \frac{v_i^H r_j - \ell_i^H w_j}{\mu_i - \lambda_j} \\ &= \frac{v_i^H r_j \lambda_j - \ell_i^H w_j \lambda_j - \mu_i v_i^H r_j + \mu_i \ell_i^H w_j}{\mu_i - \lambda_j} \\ &= \ell_i^H w_j - v_i^H r_j = [LW - VR]_{ij}. \end{aligned}$$

Moreover, we have

$$\begin{aligned} [\mathbb{L}_\sigma\Lambda - M\mathbb{L}_\sigma]_{ij} &= \frac{\mu_i v_i^H r_j - \ell_i^H w_j \lambda_j}{\mu_i - \lambda_j} \lambda_j - \mu_i \frac{\mu_i v_i^H r_j - \ell_i^H w_j \lambda_j}{\mu_i - \lambda_j} \\ &= \frac{\mu_i v_i^H r_j \lambda_j - \ell_i^H w_j \lambda_j^2 - \mu_i^2 v_i^H r_j + \mu_i \ell_i^H w_j \lambda_j}{\mu_i - \lambda_j} \\ &= \ell_i^H w_j \lambda_j - \mu_i v_i^H r_j = [LW\Lambda - MVR]_{ij}. \end{aligned}$$

□

**Remark 6.20:** Let the data be sampled from a (regular) linear descriptor system

$$\begin{aligned}\frac{d}{dt}Ex(t) &= Ax(t) + Bu(t), \\ y(t) &= Cx(t)\end{aligned}$$

with transfer function  $G(s) = C(sE - A)^{-1}B$  and define the matrices

$$\begin{aligned}C_k &:= [(\lambda_1 E - A)^{-1}Br_1 \quad \dots \quad (\lambda_k E - A)^{-1}Br_k] \in \mathbb{C}^{n \times k}, \\ O_q &:= \begin{bmatrix} \ell_1^H C(\mu_1 E - A)^{-1} \\ \vdots \\ \ell_q^H C(\mu_q E - A)^{-1} \end{bmatrix} \in \mathbb{C}^{q \times n},\end{aligned}$$

which are called *generalized tangential controllability matrices* and *generalized tangential observability matrices*, respectively. With these, it follows that

$$\begin{aligned}[\mathbb{L}]_{ij} &= \frac{v_i^H r_j - \ell_i^H w_j}{\mu_i - \lambda_j} \\ &= \frac{\ell_i^H G(\mu_i) r_j - \ell_i^H G(\lambda_j) r_j}{\mu_i - \lambda_j} \\ &= \frac{\ell_i^H C((\mu_i E - A)^{-1} - (\lambda_j E - A)^{-1}) Br_j}{\mu_i - \lambda_j} \\ &= \frac{\ell_i^H C(\mu_i E - A)^{-1}((\lambda_j E - A) - (\mu_i E - A))(\lambda_j E - A)^{-1} Br_j}{\mu_i - \lambda_j} \\ &= -\ell_i^H C(\mu_i E - A)^{-1} E(\lambda_j E - A)^{-1} Br_j,\end{aligned}$$

and similarly

$$\begin{aligned}[\mathbb{L}_\sigma]_{ij} &= \frac{\mu_i v_i^H r_j - \lambda_j \ell_i^H w_j}{\mu_i - \lambda_j} \\ &= \frac{\mu_i \ell_i^H G(\mu_i) r_j - \lambda_j \ell_i^H G(\lambda_j) r_j}{\mu_i - \lambda_j} \\ &= -\ell_i^H C(\mu_i E - A)^{-1} A(\lambda_j E - A)^{-1} Br_j.\end{aligned}$$

This gives

$$\mathbb{L} = -O_q E C_k, \quad \mathbb{L}_\sigma = -O_q A C_k.$$

Now we state and proof a result that gives us the structure of the dynamical system that interpolates the data.

**Theorem 6.21:** Assume that there is given data as in (6.6), (6.7) with  $k = q$ . Let the associated Loewner pencil  $s\mathbb{L} - \mathbb{L}_\sigma$  be regular (i. e.,  $\det(s\mathbb{L} - \mathbb{L}_\sigma)$  is not the zero polynomial) and assume that no  $\lambda_i$ ,  $i = 1, 2, \dots, k$  and  $\mu_j$ ,  $j = 1, 2, \dots, q$  is an eigenvalue of the pencil  $s\mathbb{L} - \mathbb{L}_\sigma$ . Then

$$\tilde{E} = -\mathbb{L}, \quad \tilde{A} = -\mathbb{L}_\sigma, \quad \tilde{B} = V, \quad \tilde{C} = W$$

is an interpolating descriptor system realization, i. e., the function  $\tilde{G}(s) := \tilde{C}(s\tilde{E} - \tilde{A})^{-1}\tilde{B}$  interpolates the given data.

*Proof.* For the proof we make use of the Sylvester equations (6.10) and (6.11). By multiplying the first equation by  $s$  and subtracting it from the second one, we obtain

$$(\mathbb{L}_\sigma - s\mathbb{L})\Lambda - M(\mathbb{L}_\sigma - s\mathbb{L}) = LW(\Lambda - sI_k) - (M - sI_q)VR. \quad (6.12)$$

Multiplying this equation by  $e_i$  from the right and setting  $s = \lambda_i$ , we obtain

$$(\lambda_i I_q - M)(\mathbb{L}_\sigma - \lambda_i \mathbb{L})e_i = (\lambda_i I_q - M)Vr_i,$$

which is equivalent to

$$(\mathbb{L}_\sigma - \lambda_i \mathbb{L})e_i = Vr_i$$

and yields

$$w_i = We_i = W(\mathbb{L}_\sigma - \lambda_i \mathbb{L})^{-1}Vr_i.$$

Therefore, we obtain  $w_i = \tilde{G}(\lambda_i)r_i$ ,  $i = 1, \dots, k$ , i. e., the right data is interpolated. To prove that also the left data is interpolated we multiply (6.12) by  $e_j^T$  from the left and take  $s = \mu_j$ . This gives

$$e_j^T(\mathbb{L}_\sigma - \mu_j \mathbb{L})(\Lambda - \mu_j I_k) = e_j^T LW(\Lambda - \mu_j I_k),$$

which is equivalent to

$$e_j^T(\mathbb{L}_\sigma - \mu_j \mathbb{L}) = \ell_j^H W.$$

Therefore, we obtain

$$v_j^H = e_j^T V = \ell_j^H W(\mathbb{L}_\sigma - \mu_j \mathbb{L})^{-1}V,$$

which yields  $v_j^H = \ell_j^H \tilde{G}(\mu_j)$ ,  $j = 1, \dots, q$ . This completes the proof.  $\square$

A common situation that arises in practice is redundant data, i. e., the case where we have too much data. In this case, the Loewner pencil  $s\mathbb{L} - \mathbb{L}_\sigma$  is singular and thus the transfer function  $\tilde{G}$  does not exist. The question that arises is how we can treat this situation in the Loewner framework.

**Theorem 6.22:** Let data as in (6.6), (6.7) and the associated Loewner matrices (6.8) and (6.9) be given. Assume that

$$\begin{aligned} \text{rank}(\xi\mathbb{L} - \mathbb{L}_\sigma) = \text{rank} \begin{bmatrix} \mathbb{L} & \mathbb{L}_\sigma \end{bmatrix} = \text{rank} \begin{bmatrix} \mathbb{L} \\ \mathbb{L}_\sigma \end{bmatrix} = r \\ \forall \xi \in \{\lambda_1, \dots, \lambda_k\} \cup \{\mu_1, \dots, \mu_q\}. \end{aligned}$$

Consider the economic SVDs

$$\begin{bmatrix} \mathbb{L} & \mathbb{L}_\sigma \end{bmatrix} = Y\Sigma_1\tilde{X}^H, \quad \begin{bmatrix} \mathbb{L} \\ \mathbb{L}_\sigma \end{bmatrix} = \tilde{Y}\Sigma_r X^H \quad (6.13)$$

with  $\Sigma_1, \Sigma_r \in \mathbb{R}^{r \times r}$ , and  $X \in \mathbb{C}^{k \times r}$  and  $Y \in \mathbb{C}^{q \times r}$ . If  $R$  and  $L^H$  both have full column rank, then

$$\tilde{E} = -Y^H\mathbb{L}X, \quad \tilde{A} = -Y^H\mathbb{L}_\sigma X, \quad \tilde{B} = Y^H V, \quad \tilde{C} = WX$$

is an interpolating descriptor system realization, i.e., the function  $\tilde{G}(s) := \tilde{C}(s\tilde{E} - \tilde{A})^{-1}\tilde{B}$  interpolates the given data.

*Proof.* Here we show only show that the right interpolation conditions are satisfied, the proof for the left interpolation conditions is analogous. From (6.13) and since  $X, Y$  have orthonormal columns, we have that  $\mathbb{L} = -\tilde{Y}_1\Sigma_r X^H$  and  $\mathbb{L}_\sigma = -\tilde{Y}_2\Sigma_r X^H$ . Therefore, we obtain

$$\begin{aligned} \mathbb{L}X X^H &= -\tilde{Y}_1\Sigma_r X^H X X^H = -\tilde{Y}_1\Sigma_r X^H = \mathbb{L}, \\ \mathbb{L}_\sigma X X^H &= -\tilde{Y}_2\Sigma_r X^H X X^H = -\tilde{Y}_2\Sigma_r X^H = \mathbb{L}_\sigma. \end{aligned}$$

Moreover, we have

$$\mathbb{L}_\sigma - \mathbb{L}\Lambda = VR,$$

which follows from

$$\begin{aligned} [\mathbb{L}_\sigma - \mathbb{L}\Lambda]_{ij} &= \frac{\mu_i v_i^H r_j - \lambda_j \ell_i^H w_j}{\mu_i - \lambda_j} - \frac{v_i^H r_j - \ell_i^H w_j}{\mu_i - \lambda_j} \lambda_j = v_i^H r_j = [VR]_{ij}, \\ & \quad i = 1, 2, \dots, k, \quad j = 1, 2, \dots, q. \end{aligned}$$

Similarly we have

$$\mathbb{L}_\sigma - M\mathbb{L} = LW.$$

This gives

$$LW X X^H = (\mathbb{L}_\sigma - M\mathbb{L}) X X^H = \mathbb{L}_\sigma - M\mathbb{L} = LW.$$

Since  $L^H$  has full column rank, this implies  $WXX^H = W$ . With the above identities we get

$$\begin{aligned} -\tilde{A}X^H + \tilde{E}X^H\Lambda &= Y^H\mathbb{L}_\sigma XX^H - Y^H\mathbb{L}XX^H\Lambda \\ &= Y^H(\mathbb{L}_\sigma - \mathbb{L}\Lambda) \\ &= Y^HV R. \end{aligned}$$

Then for  $i = 1, 2, \dots, k$  we get

$$\begin{aligned} \tilde{G}(\lambda_i)r_i &= \tilde{C}(\lambda_i\tilde{E} - \tilde{A})^{-1}\tilde{B}r_i \\ &= \tilde{C}(\lambda_i\tilde{E} - \tilde{A})^{-1}Y^HV r_i \\ &= \tilde{C}(\lambda_i\tilde{E} - \tilde{A})^{-1}(-\tilde{A}X^H + \tilde{E}X^H\Lambda)e_i \\ &= \tilde{C}(\lambda_i\tilde{E} - \tilde{A})^{-1}(\lambda_i\tilde{E} - \tilde{A})X^He_i \\ &= \tilde{C}X^He_i = WXX^He_i = We_i = w_i. \end{aligned}$$

□

**Remark 6.23:** In general there are many projections leading to the same transfer function. To make this precise, assume that  $\Phi \in \mathbb{C}^{k \times r}$ ,  $\Psi \in \mathbb{C}^{q \times r}$  are given such that  $X^H\Phi$  and  $\Psi^HY$  are both nonsingular. Then the model given by

$$\hat{E} = -\Phi^H\mathbb{L}\Psi, \quad \hat{A} = -\Phi^H\mathbb{L}_\sigma\Psi, \quad \hat{B} = \Phi^HV, \quad \hat{C} = W\Psi$$

has the transfer function  $\tilde{G}(s) := \tilde{C}(s\tilde{E} - \tilde{A})^{-1}\tilde{B}$  with the notation as in Theorem 6.22.

To illustrate the method we give a simple example.

**Example:** Consider the function  $G(s) = \frac{1}{s^2+1}$ . We sample this function at the points

$$\begin{aligned} \lambda_1 &= 1, & \lambda_2 &= 2, & \lambda_3 &= 3, \\ \mu_1 &= -1, & \mu_2 &= -2, & \mu_3 &= -3, \end{aligned}$$

and obtain the data

$$\begin{aligned} \Lambda &= \text{diag}(1, 2, 3), & R &= [1 \ 1 \ 1], & W &= \left[\frac{1}{2} \ \frac{1}{5} \ \frac{1}{10}\right], \\ M &= \text{diag}(-1, -2, -3), & L^H &= [1 \ 1 \ 1], & V^H &= \left[\frac{1}{2} \ \frac{1}{5} \ \frac{1}{10}\right]. \end{aligned}$$

Then the Loewner pencil  $s\mathbb{L} - \mathbb{L}_\sigma$  is given by

$$\mathbb{L} = \begin{bmatrix} 0 & -\frac{1}{10} & -\frac{1}{10} \\ \frac{1}{10} & 0 & -\frac{1}{50} \\ \frac{1}{10} & \frac{1}{50} & 0 \end{bmatrix}, \quad \mathbb{L}_\sigma = \begin{bmatrix} \frac{1}{2} & \frac{3}{10} & \frac{1}{5} \\ \frac{3}{10} & \frac{1}{2} & \frac{7}{50} \\ \frac{1}{5} & \frac{7}{50} & \frac{1}{10} \end{bmatrix}.$$

It is easily checked by some SVDs that

$$\text{rank}(\xi\mathbb{L} - \mathbb{L}_\sigma) = \text{rank} \begin{bmatrix} \mathbb{L} & \mathbb{L}_\sigma \end{bmatrix} = \text{rank} \begin{bmatrix} \mathbb{L} \\ \mathbb{L}_\sigma \end{bmatrix} = 2$$

$$\forall \xi \in \{-3, -2, -1, 1, 2, 3\}.$$

Therefore, we choose  $\Phi, \Psi \in \mathbb{C}^{3 \times 2}$ . We take

$$\Phi = \Psi = \begin{bmatrix} 5 & -5 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

This gives a reduced-order

$$\hat{E} = -\Phi^H \mathbb{L} \Psi = \begin{bmatrix} 0 & \frac{51}{50} \\ -\frac{51}{50} & 0 \end{bmatrix}, \quad \hat{A} = -\Phi^H \mathbb{L}_\sigma \Psi = \begin{bmatrix} -\frac{157}{10} & \frac{643}{50} \\ \frac{643}{50} & -\frac{53}{5} \end{bmatrix},$$

$$\hat{B} = \Phi^H V = \begin{bmatrix} \frac{27}{10} \\ -\frac{12}{5} \end{bmatrix}, \quad \hat{C} = W \Psi = \begin{bmatrix} \frac{27}{10} & -\frac{12}{5} \end{bmatrix}.$$

Now a simple calculation gives

$$\hat{C}(s\hat{E} - \hat{A})^{-1}\hat{B} = \frac{1}{s^2 + 1} = G(s),$$

so we have reconstructed our original model (which can be realized by a system of state-space dimension two) from the data.

# CHAPTER 7

---

## Outlook

---

In this course on model reduction we have mainly considered system-theoretic methods, i. e., methods that mainly try to approximate the input/output behavior of a dynamical system. In this chapter we will have a brief look onto some further aspects of model reduction that could not be covered in this course in full detail.

### 7.1 Parametric Model Reduction

In industrial applications one often considers dynamical systems that depend on parameters. Then one is often interested in optimizing these parameters which often requires a lot of simulations of the dynamical system for many parameters. In this context one is particularly interested in reduced representations of the model to make these simulations feasible. The great challenge consists of finding a reduced-order model that is a good approximation of the original one for all parameters. A good survey on such methods is [BGW15]. To make this precise, consider a *parametric LTI system*

$$\begin{aligned}\dot{x}(t; p) &= A(p)x(t; p) + B(p)u(t), \\ y(t; p) &= C(p)x(t; p) + D(p)u(t),\end{aligned}$$

where  $A(p) \in \mathbb{R}^{n \times n}$ ,  $B(p) \in \mathbb{R}^{n \times m}$ ,  $C(p) \in \mathbb{R}^{q \times n}$ , and  $D(p) \in \mathbb{R}^{q \times m}$  for all  $p \in \Omega \subset \mathbb{R}^d$ . In this setting it is desirable to have a parameter-affine representation

of the model as

$$\begin{aligned} A(p) &= A_0 + a_1(p)A_1 + \dots + a_{\kappa_A}(p)A_{\kappa_A}, \\ B(p) &= B_0 + b_1(p)B_1 + \dots + b_{\kappa_B}(p)B_{\kappa_B}, \\ C(p) &= C_0 + c_1(p)C_1 + \dots + c_{\kappa_C}(p)C_{\kappa_C} \end{aligned} \quad (7.1)$$

for fixed matrices  $A_0, \dots, A_{\kappa_A} \in \mathbb{R}^{n \times n}$ ,  $B_0, \dots, B_{\kappa_B} \in \mathbb{R}^{n \times m}$ ,  $C_0, \dots, C_{\kappa_C} \in \mathbb{R}^{q \times n}$  and functions  $a_1, \dots, a_{\kappa_A}, b_1, \dots, b_{\kappa_B}, c_1, \dots, c_{\kappa_C} : \Omega \rightarrow \mathbb{C}$ . Note that any system can be written in this form, but for computational efficiency it is desirable to have  $\kappa_A, \kappa_B, \kappa_C \ll n$ . This representation easily allows the construction of reduced-order models via projection, i. e., projection matrices  $W, T \in \mathbb{R}^{n \times r}$  are constructed such that  $W^T T = I_r$  and such that the reduced-order model is given by

$$\begin{aligned} \dot{\tilde{x}}(t; p) &= \tilde{A}(p)\tilde{x}(t; p) + \tilde{B}(p)u(t), \\ \tilde{y}(t; p) &= \tilde{C}(p)\tilde{x}(t; p) + \tilde{D}(p)u(t), \end{aligned}$$

where

$$\begin{aligned} \tilde{A}(p) &= W^T A(p)T = W^T A_0 T + a_1(p)W^T A_1 T + \dots + a_{\kappa_A}(p)W^T A_{\kappa_A} T, \\ \tilde{B}(p) &= W^T B(p) = W^T B_0 + b_1(p)W^T B_1 + \dots + b_{\kappa_B}(p)W^T B_{\kappa_B}, \\ \tilde{C}(p) &= C(p)T = C_0 T + c_1(p)C_1 T + \dots + c_{\kappa_C}(p)C_{\kappa_C} T, \\ \tilde{D}(p) &= D(p). \end{aligned}$$

The main questions that have to be faced here, are the following:

- a) Do we want to determine *global projection matrices*  $W$  and  $T$  that are good for all  $p \in \Omega$  or do we rather determine *local projection matrices*  $W_1, \dots, W_k, T_1, \dots, T_k$  from models for particular parameters  $p_1, \dots, p_k$ ?
- b) If we determine local projection matrices, how do we choose good parameters  $p_1, \dots, p_k$  and how can we use the information to get a reduced-order model for  $p \notin \{p_1, \dots, p_k\}$ ?

There are several ways to attack these questions. If the local projection matrices  $W_1, \dots, W_k, T_1, \dots, T_k$  from models for particular parameters  $p_1, \dots, p_k$  are known, then one can obtain a reduced-order model for some  $p \notin \{p_1, \dots, p_k\}$  by

- a) *transfer function interpolation*: If  $\tilde{G}_j(s) = \tilde{C}(p_j)(sI_n - \tilde{A}(p_j))^{-1}\tilde{B}(p_j) + \tilde{D}(p_j)$ , then one can construct

$$\tilde{G}(s, p) = \sum_{j=1}^k g_j(p)\tilde{G}_j(s)$$

for an arbitrary  $p \in \Omega$ , where  $g_j(\cdot)$  are some interpolation functions (such as Lagrange polynomials). Further, rational interpolation techniques have been developed to find good global projection matrices. These lead to (tangential) interpolation conditions that do not only lead to moment matching at interpolation points  $s_1, \dots, s_\ell$ , but also at the parameters  $p_1, p_2, \dots, p_k$ .

- b) *matrix interpolation*: Instead of interpolating the transfer functions, the reduced functions  $\tilde{A}(\cdot)$ ,  $\tilde{B}(\cdot)$ ,  $\tilde{C}(\cdot)$ ,  $\tilde{D}(\cdot)$  can be obtained via interpolation at the parameters  $p_1, \dots, p_k$ .

Note that the above techniques do not need the affine representation (7.1). On the other hand, the affine representation is explicitly used in determining good global projection matrices by parametric balanced truncation. This method computes global low-rank factorizations of parametric Lyapunov equations, but it still has limitations since it can only be applied to special classes of systems such as systems with pointwise positive definite  $A(p)$  that allow the development of error estimators. These error estimators are important to find good sampling parameters  $p$  in the algorithm.

## 7.2 Sampling-Based Methods

A further class of methods that has not been discussed in this course are sampling-based methods. These methods mainly consist of sampling the solution of the system under consideration for several initial conditions or parameter values and approximating the space in which the solutions live. If a basis for this space for some initial conditions or parameters is known, then the hope is that also the solutions for the other values can be well represented in this basis. Sampling-based methods play a great role in model reduction of PDEs, see, e. g., [Vol13]. Consider for example 1D Burgers' equation

$$\frac{\partial}{\partial t} y(x, t) - \nu \frac{\partial^2}{\partial x^2} y(x, t) + \frac{1}{2} \frac{\partial}{\partial x} y(x, t)^2 = 0, \quad x \in [0, 1], \quad t \geq 0$$

with the boundary and initial conditions

$$\begin{aligned} y(0, t) &= y(1, t) = 0, \quad t \geq 0, \\ y(x, 0) &= y_0(x), \quad x \in [0, 1]. \end{aligned}$$

A finite element discretization with the finite element basis  $\{\varphi_1, \dots, \varphi_N\}$  leads to the spatially discretized model

$$\widehat{M} \frac{d}{dt} \widehat{y}(t) - \nu \widehat{K} \widehat{y}(t) + \widehat{L}(\widehat{y}(t)) = 0.$$


---

Then an approximate solution of the original problem is of the form

$$\hat{y}_h(x, t) = \sum_{i=1}^N \varphi_i(x) \hat{y}_i(t),$$

where  $\hat{y}_i(t)$  is the  $i$ -th component of  $\hat{y}(t)$ . The goal is to replace this model by a reduced-order model

$$\tilde{M} \frac{d}{dt} \tilde{y}(t) - \nu \tilde{K} \tilde{y}(t) + \tilde{L}(\tilde{y}(t)) = 0$$

that can be described by an orthonormal basis  $\{\psi_1, \dots, \psi_k\}$  with  $k \ll N$ . Then an approximate solution of the original problem is of the form

$$\tilde{y}_h(x, t) = \sum_{i=1}^k \psi_i(x) \tilde{y}_i(t),$$

where  $\tilde{y}_i(t)$  is the  $i$ -th component of  $\tilde{y}(t)$ . The question is how to determine the basis  $\{\psi_1, \dots, \psi_k\}$ . Assume that we have given “snapshots” of the full-order solution  $\hat{y}_h(x, t)$ , i. e., we know  $\{\hat{y}_1(\cdot), \dots, \hat{y}_n(\cdot)\} := \{\hat{y}_h(\cdot, t_1), \dots, \hat{y}_h(\cdot, t_n)\} \subset X$ , where  $X$  is assumed to be a Hilbert space with inner product  $\langle \cdot, \cdot \rangle_X$  and induced norm  $\|\cdot\|_X := \langle \cdot, \cdot \rangle_X^{1/2}$ . For given  $k \leq n$ , the method of *proper orthogonal decomposition* (POD) computes an optimal solution of the optimization problem

$$\min_{\{\psi_1, \dots, \psi_k\}} \sum_{i=1}^n \|\tilde{y}_i - \hat{y}_i\|_X^2, \quad \langle \psi_j, \psi_\ell \rangle_X = \delta_{j\ell} \quad \forall j, \ell \in \{1, \dots, k\},$$

where  $\tilde{y}_i(x) = \sum_{j=1}^k \langle \hat{y}_i, \psi_j \rangle_X \psi_j(x)$  is an approximation of  $\hat{y}_i$  using the basis  $\{\psi_1, \dots, \psi_k\}$ . Assume w. l. o. g. that  $X = \mathbb{R}^N$ . (Note that if  $X$  is any finite dimensional Hilbert space, then there exists an isometric isomorphism  $\Phi : (X, \langle \cdot, \cdot \rangle_X) \rightarrow (\mathbb{R}^N, \langle \cdot, \cdot \rangle_{\mathbb{R}^N})$ , where  $\langle \cdot, \cdot \rangle_{\mathbb{R}^N}$  is the standard inner product in  $\mathbb{R}^N$  – so we can exploit the notation and simply denote  $\Phi(y)$  by  $y$ .) So we get an equivalent optimization problem

$$\min_{\{\psi_1, \dots, \psi_k\}} \sum_{i=1}^n \|\tilde{y}_i - \hat{y}_i\|_2^2, \quad \langle \psi_j, \psi_\ell \rangle_{\mathbb{R}^N} = \delta_{j\ell} \quad \forall j, \ell \in \{1, \dots, k\},$$

where  $\tilde{y}_i = \sum_{j=1}^k \langle \hat{y}_i, \psi_j \rangle_{\mathbb{R}^N} \psi_j$ . Let  $Y = [\hat{y}_1 \ \dots \ \hat{y}_n] \in \mathbb{R}^{N \times n}$  and

$$YY^T \tilde{\psi}_j = \lambda_j \tilde{\psi}_j \quad \text{for } j = 1, 2, \dots, n \quad \text{with } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n.$$

Then the POD basis is given by  $\{\psi_1, \dots, \psi_k\} = \{\tilde{\psi}_1, \dots, \tilde{\psi}_k\}$ . In other words, the POD basis consists of the  $k$  dominant left singular vectors of the snapshot matrix  $Y$ . Furthermore, the error is given by

$$\sum_{i=1}^n \left\| \sum_{j=1}^k \langle \hat{y}_i, \psi_j \rangle_{\mathbb{R}^N} \psi_j - \hat{y}_i \right\|_2^2 = \sum_{j=k+1}^n \lambda_j.$$

In the literature several further aspects are discussed such as optimal choices of snapshots for parameter-dependent problems (by developing error estimators) and the application to optimal control problems.



---

## Bibliography

---

- [Ant05] A. C. Antoulas. *Approximation of Large-Scale Dynamical Systems*. SIAM Publications, Philadelphia, PA, USA, 2005.
- [AV73] B. D. O. Anderson and S. Vongpanitlerd. *Network Analysis and Synthesis – A Modern Systems Theory Approach*. Prentice-Hall, Englewood Cliffs, NJ, USA, 1973.
- [BB98] P. Benner and R. Byers. An exact line search method for solving generalized continuous-time algebraic Riccati equations. *IEEE Trans. Automat. Control*, 43(1):101–107, 1998.
- [BGW15] P. Benner, S. Gugercin, and K. Willcox. A survey on projection-based model reduction methods for parametric dynamical systems. *SIAM Rev.*, 57(4):483–531, 2015.
- [BKS13] P. Benner, P. Kürschner, and J. Saak. Efficient handling of complex shift parameters in the low-rank cholesky factor ADI method. *Numer. Algorithms*, 62(2):225–251, 2013.
- [BKS15] P. Benner, P. Kürschner, and J. Saak. Self-generating and efficient shift parameters in ADI methods for large Lyapunov and Sylvester equations. *Electron. Trans. Numer. Anal.*, 43:142–162, 2014–2015.
- [Bog07] V. I. Bogachev. *Measure Theory*. Springer-Verlag, Berlin, Heidelberg, 1st edition, 2007.
- [BS13] P. Benner and J. Saak. Numerical solution of large and sparse continuous time algebraic matrix Riccati and Lyapunov equations: A state of the art survey. *GAMM-Mitt.*, 6(1):32–52, 2013.

- [GAB08] S. Gugercin, A. C. Antoulas, and C. A. Beattie.  $\mathcal{H}_2$  model reduction for large-scale linear dynamical systems. *SIAM J. Matrix Anal. Appl.*, 30(2):609–638, 2008.
- [Geo88] T. T. Georgiou. On the computation of the gap metric. *Systems Control Lett.*, 11(4):253–257, 1988.
- [GO13] C. Guiver and M. R. Opmeer. Error bounds in the gap metric for dissipative balanced approximations. *Linear Algebra Appl.*, 439(12):3659–3698, 2013.
- [Kle68] D. Kleinman. On an iterative technique for Riccati equation computations. *IEEE Trans. Automat. Control*, 13(1):114–115, 1968.
- [LR95] P. Lancaster and L. Rodman. *Algebraic Riccati Equations*. Oxford University Press, New York, 1995.
- [MA07] A. J. Mayo and A. C. Antoulas. A framework for the solution of the generalized realization problem. *Linear Algebra Appl.*, 425(2–3):634–662, 2007.
- [Moo81] B. C. Moore. Principal component analysis in linear systems: Controllability, observability, and model reduction. *IEEE Trans. Automat. Control*, AC-26(1):17–32, 1981.
- [Obe91] R. Ober. Balanced parametrization of classes of linear systems. *SIAM J. Control Optim.*, 29(6):1251–1287, 1991.
- [Pen00] T. Penzl. A cyclic low rank Smith method for large sparse Lyapunov equations. *SIAM J. Sci. Comput.*, 21(4):1401–1418, 2000.
- [PR55] D. Peaceman and H. Rachford. The numerical solution of elliptic and parabolic differential equations. *J. Soc. Indust. Appl. Math.*, 3(1):28–41, 1955.
- [Ran96] A. Rantzer. On the Kalman-Yakubovich-Popov lemma. *Systems Control Lett.*, 28(1):7–10, 1996.
- [RM06a] J. Rommes and N. Martins. Efficient computation of multivariate transfer function dominant poles using subspace acceleration. *IEEE Trans. Power Syst.*, 21(4):1471–1483, 2006.
- [RM06b] J. Rommes and N. Martins. Efficient computation of transfer function dominant poles using subspace acceleration. *IEEE Trans. Power Syst.*, 21(3):1218–1226, 2006.
- [RS10] T. Reis and T. Stykel. Positive real and bounded real balancing for model reduction of descriptor systems. *Internat. J. Control*, 83(1):74–88, 2010.
-

- [Saa82] Y. Saad. The Lanczos biorthogonalization algorithm and other oblique projection methods for solving large unsymmetric systems. *SIAM J. Numer. Anal.*, 19(3):485–506, 1982.
- [Sim07] V. Simoncini. A new iterative method for solving large-scale Lyapunov equations. *SIAM J. Sci. Comput.*, 29(3):1268–1288, 2007.
- [SZ02] D. C. Sorenson and Y. Zhou. Bounds on eigenvalue decay rates and sensitivity of solutions to Lyapunov equations. CAAM Technical Reports, Rice University, 2002.
- [Vol13] S. Volkwein. Proper orthogonal decomposition: Theory and reduced-order modeling, August 2013. Available from <http://www.math.uni-konstanz.de/numerik/personen/volkwein/teaching/POD-Book.pdf>.
- [Wac13] E. Wachspress. *The ADI Model Problem*. Springer-Verlag, New York, NY, USA, 1st edition, 2013.
- [Wil71] J. C. Willems. Least squares stationary optimal control and the algebraic Riccati equation. *IEEE Trans. Automat. Control*, AC-16(6):621–634, 1971.
- [Wil72a] J. C. Willems. Dissipative dynamical systems part I: General theory. *Arch. Ration. Mech. Anal.*, 45(5):321–351, 1972.
- [Wil72b] J. C. Willems. Dissipative dynamical systems part II: Linear systems with quadratic supply rates. *Arch. Ration. Mech. Anal.*, 45(5):352–393, 1972.
- [ZDG96] K. Zhou, J. C. Doyle, and K. Glover. *Robust and Optimal Control*. Prentice-Hall, Englewood Cliffs, NJ, 1996.
-



