

# Structure-Preserving Model Order Reduction for Network Systems

**Dissertation**

zur Erlangung des akademischen Grades

**doctor rerum naturalium**  
**(Dr. rer. nat.)**

von **M. Sc. Petar Mlinarić**

geb. am **25.04.1990** in Zagreb, Kroatien

genehmigt durch die Fakultät für Mathematik  
der Otto-von-Guericke-Universität Magdeburg

Gutachter: **Prof. Dr. Peter Benner**

**Prof. Dr. Harry L. Trentelman**

eingereicht am: **30.10.2019**

Verteidigung am: **23.01.2020**



This thesis considers structure-preserving system-theoretic model order reduction for certain structured input-output systems, particularly network systems. In the first part, our focus lies on the clustering-based approach to reduce network systems. Therein, we begin by considering clustering-based model order reduction for linear multi-agent systems. This approach finds a reduced model whose dynamics evolve over a smaller network. To measure the reduction error, we use the  $\mathcal{H}_2$ -norm and consider the  $\mathcal{H}_2$ -optimal clustering problem. Since clustering is generally a difficult combinatorial problem, we propose a framework based on relaxing the discrete problem to find an  $\mathcal{H}_2$ -suboptimal clustering. Following on this, we directly extend the framework to a class of nonlinear multi-agent systems.

Next, we derive upper bounds for  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  clustering-based reduction errors for linear multi-agent systems based on almost equitable partitions. These results generalize work for multi-agent systems with single-integrator agents. Using a similar approach for power systems, which are a special class of nonlinear multi-agent systems, we find conditions for exact clustering-based reduction.

Additionally, we study subsystem reduction for network systems. We propose a balancing-based approach that guarantees stability preservation under a small-gain condition. Furthermore, we consider the  $\mathcal{H}_2$ -optimal subsystem reduction problem. We derive Gramian-based first-order necessary  $\mathcal{H}_2$ -optimality conditions and use a gradient-based optimization method to fulfill them.

Finally, we apply the structure-preserving  $\mathcal{H}_2$ -optimal model reduction approach for network systems to other structured systems. In particular, we consider  $\mathcal{H}_2$ -optimal model order reduction of second-order systems, port-Hamiltonian systems, time-delay systems and  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal model order reduction of parametric systems. Here, we also derive Gramian-based  $\mathcal{H}_2$ -optimality conditions and use an optimization approach to construct a reduced model. For some structured systems, we also derive interpolatory  $\mathcal{H}_2$ -optimality conditions under additional assumptions on the reduced model.



Diese Arbeit beschäftigt sich mit der strukturerhaltenden systemtheoretischen Modellordnungsreduktion für bestimmte strukturierte Input-Output-Systeme, insbesondere Netzwerksysteme. Im ersten Teil liegt unser Fokus auf dem clusterbasierten Ansatz zur Reduktion von Netzwerksystemen. Zunächst betrachten wir dabei die clusterbasierte Modellreduktion für lineare Multiagentensysteme. Dieser Ansatz findet ein reduziertes Modell, dessen Dynamik sich über ein kleineres Netzwerk entwickelt. Zum Messen des Reduktionsfehlers verwenden wir die  $\mathcal{H}_2$ -Norm und betrachten das  $\mathcal{H}_2$ -optimale Clustering-Problem. Da das Clustering im Allgemeinen ein schwieriges kombinatorisches Problem ist, schlagen wir ein Framework vor, welches darauf basiert, das diskrete Problem zu relaxieren, um ein  $\mathcal{H}_2$ -suboptimales Clustering zu finden. Im Anschluss daran erweitern wir das Framework direkt auf eine Klasse nichtlinearer Multiagentensysteme.

Als nächstes leiten wir obere Schranken für  $\mathcal{H}_2$  und  $\mathcal{H}_\infty$  clusterbasierte Reduktionsfehler für lineare Multiagentensysteme her, basierend auf nahezu gerechten Partitionierungen. Diese Ergebnisse verallgemeinern die Arbeit für Multiagentensysteme mit Single-Integrator-Agenten. Unter Verwendung eines ähnlichen Ansatzes für Energiesysteme, die eine spezielle Klasse nichtlinearer Multiagentensysteme darstellen, finden wir Bedingungen für eine genaue clusterbasierte Reduktion.

Zusätzlich untersuchen wir die Subsystemreduktion für Netzwerksysteme. Wir schlagen einen balancierenden Ansatz vor, der die Erhaltung der Stabilität unter Small-Gain-Bedingungen gewährleistet. Weiterhin betrachten wir das  $\mathcal{H}_2$ -optimale Subsystemreduktionsproblem. Wir leiten grammschenbasierte notwendige  $\mathcal{H}_2$ -Optimalitätsbedingungen erster Ordnung her und verwenden eine gradientenbasierte Optimierungsmethode um diese zu erfüllen.

Schließlich wenden wir die strukturerhaltende  $\mathcal{H}_2$ -optimale Modellreduktion für Netzwerksysteme auf andere strukturierte Systeme an. Insbesondere betrachten wir  $\mathcal{H}_2$ -optimale Modellordnungsreduktion von Systemen zweiter Ordnung, Port-Hamiltonscher Systeme, zeitverzögerter Systeme und  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimale Modellordnungsreduktion von parametrischen Systemen. Auch hier leiten wir grammschenbasierte  $\mathcal{H}_2$ -Optimalitätsbedingungen her und verwenden einen Optimierungsansatz um ein reduziertes Modell zu konstruieren. Für einige strukturierte Systeme leiten wir auch interpolatorische  $\mathcal{H}_2$ -Optimalitätsbedingungen unter zusätzlichen Annahmen an das reduzierte Modell

---

her.

<b>List of Figures</b>	<b>xi</b>
<b>List of Tables</b>	<b>xiii</b>
<b>List of Algorithms</b>	<b>xv</b>
<b>List of Acronyms</b>	<b>xvii</b>
<b>List of Symbols</b>	<b>xix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Outline of the thesis . . . . .	2
<b>2 Mathematical Preliminaries</b>	<b>5</b>
2.1 Linear algebra . . . . .	6
2.1.1 Eigenvalues and eigenvectors of matrices and matrix pairs . . . . .	6
2.1.2 Kronecker product, vectorization, and matrix equations . . . . .	7
2.2 Functional analysis . . . . .	8
2.2.1 Fréchet and Gateaux differentiability . . . . .	8
2.2.2 Implicit function theorem . . . . .	11
2.2.3 Constrained optimization and Lagrange multipliers . . . . .	12
2.3 Linear time-invariant systems . . . . .	13
2.3.1 Solutions and stability . . . . .	14
2.3.2 Controllability, observability, and Gramians . . . . .	15
2.3.3 Hardy spaces and system norms . . . . .	18
2.4 Model order reduction . . . . .	19
2.4.1 Projection-based model reduction . . . . .	19
2.4.2 Balanced truncation . . . . .	20
2.4.3 Rational interpolation . . . . .	20
2.4.4 $\mathcal{H}_2$ -optimal model order reduction . . . . .	22
2.4.4.1 Interpolation-based approach . . . . .	22
2.4.4.2 Gramian-based approach . . . . .	24

2.4.4.3	Comparison of the two approaches . . . . .	28
2.4.5	Model order reduction of unstable systems . . . . .	29
2.5	Graph theory . . . . .	30
2.5.1	Basic concepts . . . . .	31
2.5.2	Graph partitions . . . . .	33
2.6	Linear multi-agent systems . . . . .	34
2.6.1	System description . . . . .	34
2.6.2	Clustering-based model order reduction . . . . .	36
<b>3</b>	<b>Suboptimal Clustering-Based Model Order Reduction</b>	<b>39</b>
3.1	Introduction . . . . .	39
3.2	$\mathcal{H}_2$ -suboptimal clustering . . . . .	40
3.2.1	Single-integrator agents . . . . .	40
3.2.2	QR decomposition-based clustering . . . . .	41
3.2.3	Clustering by k-means algorithm . . . . .	42
3.2.4	Computing the $\mathcal{H}_2$ -error . . . . .	44
3.2.5	Extension to higher-order agents . . . . .	45
3.2.6	Numerical example . . . . .	45
3.3	Clustering for nonlinear multi-agent systems . . . . .	49
3.3.1	Nonlinear multi-agent systems . . . . .	50
3.3.2	Clustering by projection . . . . .	50
3.4	Conclusions . . . . .	51
<b>4</b>	<b>Graph Symmetries and Equitable Partitions in Clustering-Based Model Order Reduction</b>	<b>53</b>
4.1	Introduction . . . . .	54
4.2	Error bounds for clustering-based model order reduction of linear multi-agent systems . . . . .	54
4.2.1	Introduction . . . . .	54
4.2.2	Preliminaries . . . . .	55
4.2.3	Problem formulation . . . . .	58
4.2.4	Graph partitions and reduction by clustering . . . . .	58
4.2.5	$\mathcal{H}_2$ -error bounds . . . . .	60
4.2.6	$\mathcal{H}_\infty$ -error bounds . . . . .	68
4.2.6.1	The single integrator case . . . . .	68
4.2.6.2	The general case with symmetric agent dynamics . . . . .	71
4.2.7	Towards a priori error bounds for general graph partitions . . . . .	74
4.2.7.1	The single integrator case . . . . .	74
4.2.7.2	The general case . . . . .	78
4.2.8	Numerical examples . . . . .	80
4.3	Exact clustering-based model order reduction for nonlinear power systems	84
4.3.1	Introduction . . . . .	84
4.3.2	System description . . . . .	84



4.3.3	Synchronization of generator pair . . . . .	88
4.3.4	Synchronization of generator partition . . . . .	92
4.3.5	Clustering of power systems . . . . .	95
4.3.6	Illustrative example . . . . .	96
4.4	Conclusion . . . . .	97
<b>5</b>	<b>Subsystem Reduction for Interconnected Systems</b>	<b>99</b>
5.1	Introduction . . . . .	99
5.2	Stability-preserving balancing-based model order reduction . . . . .	100
5.2.1	Preliminaries . . . . .	100
5.2.2	Bounded real balanced truncation . . . . .	101
5.2.3	Stability-preserving model order reduction . . . . .	102
5.2.4	Numerical example . . . . .	104
5.3	$\mathcal{H}_2$ -optimal subsystem reduction . . . . .	105
5.3.1	Interconnected systems . . . . .	106
5.3.2	Multi-agent systems . . . . .	109
5.4	Conclusion . . . . .	112
<b>6</b>	<b><math>\mathcal{H}_2</math>-Optimal Model Order Reduction of Further Structured Systems</b>	<b>113</b>
6.1	Introduction . . . . .	113
6.2	Second-order systems . . . . .	114
6.2.1	Wilson-type conditions . . . . .	114
6.2.2	Interpolatory conditions . . . . .	118
6.3	Port-Hamiltonian systems . . . . .	123
6.3.1	Wilson-type conditions . . . . .	125
6.3.2	Interpolatory conditions . . . . .	126
6.4	Linear parametric systems . . . . .	128
6.4.1	$\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal model order reduction . . . . .	128
6.4.2	Interpolatory conditions . . . . .	131
6.5	Linear time-delay systems . . . . .	132
6.5.1	Wilson-type conditions . . . . .	134
6.6	Conclusion . . . . .	138
<b>7</b>	<b>Conclusions and Outlook</b>	<b>139</b>
7.1	Summary . . . . .	139
7.2	Future research perspectives . . . . .	140
	<b>Bibliography</b>	<b>141</b>
	<b>Statement of Scientific Cooperations</b>	<b>151</b>
	<b>Ehrenerklärung</b>	<b>153</b>



## LIST OF FIGURES

2.1	An undirected, weighted, connected graph . . . . .	32
3.1	The undirected, weighted graph from [MTC14] . . . . .	46
3.2	Relative $\mathcal{H}_2$ -errors for all partitions with five clusters . . . . .	48
3.3	Graph with partition . . . . .	48
3.4	Graph partition (3.9) with five clusters. The assignment of vertices to clusters is represented with different patterns. . . . .	49
4.1	A path graph on 5 vertices and its closest graph such that the partition $\{\{1, 2, 3\}, \{4, 5\}\}$ is almost equitable. . . . .	77
4.2	Ratios of $\mathcal{H}_2$ (left) and $\mathcal{H}_\infty$ (right) upper bounds and corresponding true errors, for a fixed almost equitable partition and all possible sets of leaders. In both figures, the sets of leaders are sorted such that the ratio is increasing (in particular, the ordering of the sets of leaders is not the same). . . . .	81
4.3	True $\mathcal{H}_2$ (left) and $\mathcal{H}_\infty$ (right) errors and upper bounds, for a fixed set of leaders and all partitions with five clusters. In each figure, partitions were sorted such that the true errors are increasing. . . . .	81
4.4	First 1000 true errors and upper bounds from Figure 4.3. . . . .	82
4.5	Relative error of $\mathcal{L}$ by $\mathcal{L}_{\text{AEP}}$ in Frobenius norm for all partitions with five clusters. The partitions are ordered such that the errors are increasing. . . . .	82
4.6	Comparison with error bounds from Ishizaki et al. [IKIA14, IKG <sup>+</sup> 15, IKI16a]. The first column shows the $\mathcal{H}_2$ errors and bounds, the second column the $\mathcal{H}_\infty$ errors and bounds. The first row contains values for all partitions with five clusters, the second row only the first 1000 best ones. . . . .	83
4.7	Power system consisting of generators (circles) and buses (vertical bars), where the $i$ th generator is only connected to the $i$ th bus. See Table 4.1 for the notation. . . . .	85
4.8	Partition $\{\{1, 2\}, \{3, 4, 5\}\}$ applied to the original power system in Figure 4.7 with $\chi_i = \chi_{ij} = 1$ for all $i, j$ . . . . .	97
4.9	Reduced power system obtained by clustering the system in Figure 4.8 with $M = D = I_5$ , $f = 0$ , and $E = \mathbf{1}_5$ . . . . .	97

4.10	Initial value response of the original power system from Figure 4.7 and a reduced system obtained by clustering with partition $\{\{1, 2, 3\}, \{4, 5\}\}$ . Original system's parameters are $\chi_i = \chi_{ij} = 1$ for all $i, j$ , $M = D = I_5$ , $f = 0$ , and $E = \mathbb{1}_5$ . The initial value is $\delta(0) = (0, 0.1, 0.2, 0.3, 0.4)$ and $\dot{\delta}(0) = 0$ . . . . .	98
5.1	Interconnected string-beam example from [RS07] . . . . .	105
5.2	Magnitude plot of the full-order and error systems . . . . .	105

## LIST OF TABLES

3.1	Top 20 partitions with 5 clusters for reducing the multi-agent system in Figure 3.1 . . . . .	47
4.1	Notation . . . . .	86



## LIST OF ALGORITHMS

2.1	Balancing-free square root balanced truncation method . . . . .	21
2.2	Iterative rational Krylov algorithm (IRKA) . . . . .	23
2.3	Two-sided iteration algorithm (TSIA) . . . . .	26
3.1	Clustering using QR decomposition with column pivoting [ZHD <sup>+</sup> 01, §3]	42
3.2	QR decomposition with column pivoting for matrices with block-columns	45
5.1	Bounded real balanced truncation [OJ88] . . . . .	101
5.2	Balancing method for network systems preserving stability and structure	104





## LIST OF ACRONYMS

AEP	almost equitable partition
BRBT	bounded real balanced truncation
BT	balanced truncation
FOM	full-order model
IRKA	iterative rational Krylov algorithm
LTI	linear time-invariant
MIMO	multiple-input multiple-output
MOR	model order reduction
ODE	ordinary differential equation
PLTI	parameterized LTI
ROM	reduced-order model
SISO	single-input single-output
TSIA	two-sided iteration algorithm



## LIST OF SYMBOLS

$\mathbb{R}, \mathbb{C}$	fields of real and complex numbers
$\mathbb{C}_+, \mathbb{C}_-$	open right/left complex half plane
$\mathbb{R}_+, \mathbb{R}_-$	strictly positive/negative real line
$\mathbb{R}^n, \mathbb{C}^n$	vector space of real/complex $n$ -tuples
$\mathbb{R}^{m \times n}, \mathbb{C}^{m \times n}$	set of real/complex $m \times n$ matrices
$ S $	cardinality of set $S$
$ \xi $	absolute value of real or complex scalar
$\arg \xi$	argument of complex scalar
$\mathbf{i}$	imaginary unit ( $\mathbf{i}^2 = -1$ )
$\operatorname{Re}(A), \operatorname{Im}(A)$	real and imaginary part of a complex quantity $A = \operatorname{Re}(A) + \mathbf{i} \operatorname{Im}(A) \in \mathbb{C}^{m \times n}$
$\bar{A}$	$:= \operatorname{Re}(A) - \mathbf{i} \operatorname{Im}(A)$ , complex conjugate of $A \in \mathbb{C}^{m \times n}$
$a_{ij}$	the $(i, j)$ -th entry of $A$
$A_{i:j,:}, A_{:,k:\ell}$	rows $i, \dots, j$ of $A$ and columns $k, \dots, \ell$ of $A$
$A_{i:j,k:\ell}$	rows $i, \dots, j$ of columns $k, \dots, \ell$ of $A$
$\operatorname{im}(A)$	$:= \operatorname{span}\{A_{:,1}, A_{:,2}, \dots, A_{:,n}\}$ , subspace generated by the columns of $A \in \mathbb{C}^{m \times n}$
$\ker(A)$	kernel of $A \in \mathbb{C}^{m \times n}$
$A^T$	the transpose of $A$
$A^*$	$:= (\bar{A})^T$ , the complex conjugate transpose
$A^{-1}$	the inverse of invertible $A$
$A^{-T}, A^{-*}$	the inverse of $A^T, A^*$

$A^+$	the Moore-Penrose pseudoinverse of $A$
$I_n, I_{n,r}$	identity matrix of dimension $n$ , first $r$ columns of $I_n$
$e_i$	the $i$ th column of the identity matrix
$\mathbb{1}_r$	vector of ones in $\mathbb{R}^r$
$\text{diag}(\alpha_1, \alpha_2, \dots, \alpha_n)$	$n \times n$ diagonal matrix with the $\alpha_i$ 's on the diagonal
$\text{diag}(A_1, A_2, \dots, A_n)$	block diagonal matrix with the $A_i$ 's as diagonal blocks
$\sigma(A), \sigma(A, E)$	spectrum of matrix $A$ /matrix pair $(A, E)$
$\lambda_i(A), \lambda_i(A, E)$	$i$ -th eigenvalue of $A/(A, E)$
$\rho(A, E)$	$:= \max_i  \lambda_i(A, E) $ , spectral radius of $(A, E)$
$\sigma_{\max}(A)$	the largest singular value of $A$
$\text{tr}(A)$	$:= \sum_{i=1}^n a_{ii}$ , trace of $A$
$\ u\ _2$	$:= \sqrt{\sum_{i=1}^n  u_i ^2}$ for $u \in \mathbb{C}^n$
$\ A\ _2$	$:= \sup\{\ Au\ _2 : \ u\ _2 = 1\}$ , induced matrix 2-norm
$\ A\ _F$	$:= \sqrt{\sum_{i,j}  a_{ij} ^2} = \sqrt{\text{tr}(A^*A)}$ , Frobenius norm of matrix $A \in \mathbb{C}^{m \times n}$
$\ u\ , \ A\ $	Euclidean vector or induced matrix norm $\ \cdot\ _2$
$A \succ (\succcurlyeq, \prec, \preccurlyeq) 0$	$A$ is self-adjoint positive definite (positive semidefinite, negative definite, negative semidefinite)
$A > (\geq) B$	element-wise partial ordering $a_{ij} > (\geq) b_{ij}$ , for all $i, j$
$A \circ B$	Hadamard (element-wise) product of $A$ and $B$
$A \otimes B$	Kronecker product of $A$ and $B$ (see Definition 2.7)
$\text{vec}(A)$	vectorization of $A$ (see Definition 2.7)
$\text{col}(a_1, a_2, \dots, a_n)$	vector in $\mathbb{C}^{nm}$ formed by stacking $a_1, a_2, \dots, a_n \in \mathbb{C}^m$
$\partial_{x_i} f := \frac{\partial}{\partial x_i} f$	partial derivative of $f$ with respect to $x_i$

---

$\partial_{x_i x_j} f := \frac{\partial^2}{\partial x_i \partial x_j} f$	$:= \partial_{x_i} \partial_{x_j} f$ , second order partial derivative of $f$ with respect to $x_i$ and $x_j$
$\partial_{x_i}^2 f := \partial_{x_i x_i} f$	second order partial derivative of $f$ with respect to $x_i$
$\dot{f} := \partial_t f := \frac{\partial}{\partial t} f$	derivative of $f$ with respect to time
$\ddot{f} := \partial_t^2 f := \frac{\partial^2}{\partial t^2} f$	second derivative of $f$ with respect to time
$Df(x)$	Fréchet differential of $f$ at $x$ (see Definition 2.10)
$D^G f(x)$	Gateaux differential of $f$ at $x$ (see Definition 2.14)
$df(x; h)$	directional derivative of $f$ at $x$ in direction $h$ (see Definition 2.14)
$\nabla f(x)$	gradient of $f$ at $x$ (see Definition 2.16)



---

**Contents**

---

1.1 Motivation . . . . .	1
1.2 Outline of the thesis . . . . .	2

---

## 1.1 Motivation

Large-scale dynamical systems, consisting of many ordinary differential equations (ODEs), appear often in applications. Such systems typically arise from detailed discretizations of partial differential equations (PDEs). Another source are interconnected systems, forming large-scale network systems, see, e.g., [New10, EFHO10, BFF<sup>+</sup>14] for applications in complex networks, smart-grids, distributed systems, transportation networks, biological networks, and networked multi-agent systems. Due to the increasing demand for computational resources to analyze, simulate, or control such large-scale systems, there is an interest in using model order reduction (MOR) methods. The idea behind MOR is to find a reduced-order model (ROM) with a much smaller number of ODEs, making it easier to simulate, while reproducing the original model sufficiently accurately.

Direct application of established MOR techniques, such as balanced truncation, Hankel-norm approximation, and Krylov subspace methods, see, e.g., [Ant05, BMS05, BOCW17], to structured models generally leads to a loss of structure. Additionally, for multi-agent systems, properties such as consensus and synchrony are important to preserve in the reduced model (see [Mor05, OSM03, LDCH10, MZ10]). Structure-preserving MOR methods allow preserving the physical interpretation of the model, which can also improve accuracy. Additionally, analysis or optimization methods tailored to specific model structures allows reusing them for the reduced model.

MOR techniques specifically for networked multi-agent systems with first-order agents have been proposed in [IKIA14, IKG<sup>+</sup>15, CKS16]. Extensions to second-order agents have been considered in [II15, CKS17] and to more general higher-order agents in [IKI16a,

[BSJ16, MTC13]. Some of these methods are based on clustering nodes in the network. With clustering, the idea is to partition the set of nodes in the network graph into disjoint sets called clusters, and to associate with each cluster a single, new, node in the reduced network, thus reducing the number of nodes and connections and the complexity of the network topology.

Complementarily, methods for subsystem reduction in more general network systems have been proposed in [RS07, RS08b, VVD08, SM09], with applications in, e.g., multi-body systems. There, the network in the reduced model stays the same and only parts of the system in the nodes are reduced.

## 1.2 Outline of the thesis

This thesis is structured as follows. In Chapter 2, we review different topics used throughout the thesis. We begin with the necessary topics from linear algebra, particularly properties of eigenvalues and the Kronecker product. Then we give an overview of some topics from functional analysis. In particular, we will need differentiability concepts for functions between arbitrary Hilbert spaces in Chapter 5 and Chapter 6. Since the considered MOR methods are based on systems theory, we review system properties such as stability, controllability, and observability. Additionally, we cover Hardy spaces, particularly the  $\mathcal{H}_2$  space. Next, we provide overview for some projection-based MOR methods, including balanced truncation and interpolatory methods. For multi-agents systems, we revise some basic concepts from graph theory. Then, we cover modeling and some properties of linear multi-agent systems used in Chapter 3 and Chapter 4.

In Chapter 3, we study clustering-based MOR of multi-agent systems. The goal is to find good partitioning of the nodes of the graph over which the dynamics is evolving. Clustering then produces a multi-agent system evolving over a smaller graph. First, we focus on linear multi-agent systems and consider the problem of  $\mathcal{H}_2$ -optimal clustering-based MOR. Since clustering is generally a very difficult combinatorial problem, we relax the discrete problem to a continuous  $\mathcal{H}_2$ -optimal MOR problem. To recover a discrete solution, we interpret it as a problem of approximating a subspace. We propose using a clustering algorithm over the sets of rows of a matrix used to project the original model. Then, we generalize this to a framework of combining a projection-based method with a clustering algorithm and apply it to nonlinear multi-agent systems. The results for linear multi-agent systems with single-integrator agents are published in

P. Mlinarić, S. Grundel, and P. Benner, Efficient Model Order Reduction for Multi-Agent Systems Using QR Decomposition-Based Clustering, *Proceedings of the 54<sup>th</sup> IEEE Conference on Decision and Control (CDC)*, pp. 4794–4799, December 2015.

An extension for linear multi-agent systems with more general agents is published in

P. Mlinarić, S. Grundel, and P. Benner, Clustering-Based Model Order Reduction for Multi-Agent Systems with General Linear Time-Invariant



Agents, *Proceedings of the 22<sup>nd</sup> International Symposium on Mathematical Theory of Networks and Systems (MTNS)*, pp. 230–235, July 2016.

The results for nonlinear multi-agent systems will be published in a forthcoming paper.

In Chapter 4, we turn to more theoretical consideration of error due to clustering. First, we look at linear multi-agent systems and derive  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  error bounds when using an almost equitable partition (AEP), extending known  $\mathcal{H}_2$ -error expressions for multi-agent systems with single-integrator agents. We also propose an extension to arbitrary partitions using the distance to a graph for which the partition becomes almost equitable. These results are published in

H.-J. Jongsma, P. Mlinarić, S. Grundel, P. Benner, and H. L. Trentelman, Model Reduction of Linear Multi-Agent Systems by Clustering with  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  Error Bounds, *Mathematics of Control, Signals, and Systems*, vol. 30, April 2018.

Next, we consider power systems, which are a type of nonlinear multi-agent systems. There, we derive equivalent conditions for clustering with zero error. These conditions involve graph symmetries and equitable partitions. The results are published in

P. Mlinarić, T. Ishizaki, A. Chakraborty, S. Grundel, P. Benner, and J. Imura, Synchronization and Aggregation of Nonlinear Power Systems with Consideration of Bus Network Structures, *Proceedings of the European Control Conference (ECC)*, pp. 2266–2271, June 2018.

In Chapter 5, we study the problem of subsystem reduction for linear network systems. This is a complementary approach to clustering, where the underlying graph is preserved and only the subsystems at the nodes are reduced. In the first part, we extend a balancing-based MOR method which preserves stability for network systems satisfying a certain small-gain condition. Using the known a priori  $\mathcal{H}_\infty$  error bound, it allows automatic choice of the order of the reduced subsystems. The results are published in

P. Benner, S. Grundel, and P. Mlinarić, Stability Preserving Model Reduction for Linearly Coupled Linear Time-Invariant Systems, *Proceedings in Applied Mathematics and Mechanics*, vol. 16, pp. 817–818, October 2016.

The second part is about  $\mathcal{H}_2$ -optimal subsystem reduction. Using the Gramian-based formulation of the  $\mathcal{H}_2$ -error, we derive gradients with respect to matrices defining the ROM. Therefore, we also obtain Wilson-type necessary optimality conditions. These results are subject of a forthcoming article.

In Chapter 6, we use the ideas from the second part of Chapter 5 to other structure-preserving  $\mathcal{H}_2$ -optimal MOR problems. In particular, we consider structure-preserving MOR for second-order systems, port-Hamiltonian systems, and time-delay systems. Additionally, we also consider  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal MOR for parametric systems. We derive Wilson-type necessary optimality conditions and for some systems also the interpolatory optimality conditions. The preliminary results for parametric systems are published in

M. Hund, P. Mlinarić, and J. Saak, An  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -Optimal Model Order Reduction Approach for Parametric Linear Time-Invariant Systems, *Proceedings in Applied Mathematics and Mechanics*, vol. 18, pp. e201800084, December 2018.

Finally, in Chapter 7 we summarize the results of the thesis and discuss possible future directions.

# CHAPTER 2

## MATHEMATICAL PRELIMINARIES

### Contents

---

2.1	Linear algebra . . . . .	6
2.1.1	Eigenvalues and eigenvectors of matrices and matrix pairs . . . . .	6
2.1.2	Kronecker product, vectorization, and matrix equations . . . . .	7
2.2	Functional analysis . . . . .	8
2.2.1	Fréchet and Gateaux differentiability . . . . .	8
2.2.2	Implicit function theorem . . . . .	11
2.2.3	Constrained optimization and Lagrange multipliers . . . . .	12
2.3	Linear time-invariant systems . . . . .	13
2.3.1	Solutions and stability . . . . .	14
2.3.2	Controllability, observability, and Gramians . . . . .	15
2.3.3	Hardy spaces and system norms . . . . .	18
2.4	Model order reduction . . . . .	19
2.4.1	Projection-based model reduction . . . . .	19
2.4.2	Balanced truncation . . . . .	20
2.4.3	Rational interpolation . . . . .	20
2.4.4	$\mathcal{H}_2$ -optimal model order reduction . . . . .	22
2.4.4.1	Interpolation-based approach . . . . .	22
2.4.4.2	Gramian-based approach . . . . .	24
2.4.4.3	Comparison of the two approaches . . . . .	28
2.4.5	Model order reduction of unstable systems . . . . .	29
2.5	Graph theory . . . . .	30
2.5.1	Basic concepts . . . . .	31
2.5.2	Graph partitions . . . . .	33
2.6	Linear multi-agent systems . . . . .	34
2.6.1	System description . . . . .	34
2.6.2	Clustering-based model order reduction . . . . .	36

---

## 2.1 Linear algebra

Here, we recall some definitions and properties related to eigenvalues and matrix equations.

### 2.1.1 Eigenvalues and eigenvectors of matrices and matrix pairs

We begin with the definition of eigenvalues, right and left eigenvectors, spectrum, and spectral radius of a matrix.

**Definition 2.1 ([GV13, Section 7.1.1]):**

A scalar  $\lambda \in \mathbb{C}$  is an *eigenvalue* of a matrix  $A \in \mathbb{C}^{n \times n}$  if there exists a nonzero vector  $x \in \mathbb{C}^n$  such that  $Ax = \lambda x$ . The vector  $x$  is called a (*right*) *eigenvector* of  $A$  corresponding to the eigenvalue  $\lambda$ . A nonzero vector  $y \in \mathbb{C}^n$  is a *left eigenvector* of  $A$  corresponding to the eigenvalue  $\lambda$  if  $y^*A = \lambda y^*$ . The set of eigenvalues of  $A$  is called the *spectrum* and is denoted by  $\sigma(A)$ . The *spectral radius* of the matrix  $A$  is  $\rho(A) := \max_{\lambda \in \sigma(A)} |\lambda|$ .  $\diamond$

Next, we continue with diagonalizability of a matrix and simultaneous diagonalizability of two or more matrices.

**Definition 2.2 ([HJ85, Definition 1.3.6]):**

The matrix  $A \in \mathbb{C}^{n \times n}$  is said to be *diagonalizable* if there exists an invertible  $T \in \mathbb{C}^{n \times n}$  such that  $T^{-1}AT$  is a diagonal matrix.  $\diamond$

**Definition 2.3 ([HJ85, Definition 1.3.11, 1.3.18]):**

Two matrices  $A, B \in \mathbb{C}^{n \times n}$  are said to be *simultaneously diagonalizable* if there exists an invertible  $T \in \mathbb{C}^{n \times n}$  such that  $T^{-1}AT$  and  $T^{-1}BT$  are both diagonal matrices.

A *simultaneously diagonalizable* family  $\mathcal{F} \subseteq \mathbb{C}^{n \times n}$  is a family for which there is a single invertible matrix  $T \in \mathbb{C}^{n \times n}$  such that  $T^{-1}AT$  is diagonal for every  $A \in \mathcal{F}$ .  $\diamond$

The following theorem gives necessary and sufficient conditions for simultaneous diagonalizability. Recall that two matrices  $A$  and  $B$  are said to commute if  $AB = BA$ . A family  $\mathcal{F} \subseteq \mathbb{C}^{n \times n}$  of matrices is a commuting family if every pair of matrices from  $\mathcal{F}$  commute [HJ85, Definition 1.3.16].

**Theorem 2.4 ([HJ85, Theorem 1.3.12, 1.3.19]):**

Let  $A, B \in \mathbb{C}^{n \times n}$  be diagonalizable. Then  $A$  and  $B$  commute if and only if they are simultaneously diagonalizable.

Let  $\mathcal{F} \subseteq \mathbb{C}^{n \times n}$  be a family of diagonalizable matrices. Then  $\mathcal{F}$  is a commuting family if and only if it is a simultaneously diagonalizable family.  $\diamond$

Next definition is for eigenvalues and eigenvectors of a matrix pencil  $A - \lambda B$  (or matrix pair  $(A, B)$ ) with invertible  $B$  matrix.

**Definition 2.5 ([GV13, Section 7.7]):**

A scalar  $\lambda \in \mathbb{C}$  is an *eigenvalue* of a matrix pencil  $A - \lambda B$ ,  $A, B \in \mathbb{C}^{n \times n}$ ,  $B$  invertible, if there exists a nonzero vector  $x \in \mathbb{C}^n$  such that  $Ax = \lambda Bx$ . The vector  $x$  is called a (*right*) *eigenvector* of  $A - \lambda B$  corresponding to the eigenvalue  $\lambda$ . A nonzero vector  $y \in \mathbb{C}^n$  is a *left eigenvector* of  $A - \lambda B$  corresponding to the eigenvalue  $\lambda$  if  $y^* A = \lambda B y^*$ . The set of eigenvalues of  $A - \lambda B$  is denoted by  $\sigma(A, B)$ .  $\diamond$

Notice that  $\sigma(A, B) = \sigma(B^{-1}A)$ .

The following result gives bounds for the eigenvalues of a “projected” symmetric matrix, where we use that the eigenvalues of a symmetric are real ([HJ85, Theorem 2.5.6]) to label them in increasing order.

**Theorem 2.6 (Interlacing property, [MN99, Section 11.10]):**

Let  $A \in \mathbb{R}^{n \times n}$  be a symmetric matrix with eigenvalues  $\lambda_1(A) \leq \lambda_2(A) \leq \dots \leq \lambda_n(A)$ ,  $P \in \mathbb{R}^{n \times r}$  a matrix with orthonormal columns, and  $B = P^T A P$ . Then the eigenvalues  $\lambda_1(B) \leq \lambda_2(B) \leq \dots \leq \lambda_r(B)$  of  $B$  satisfy

$$\lambda_i(A) \leq \lambda_i(B) \leq \lambda_{n-r+i}(A),$$

for  $i = 1, 2, \dots, r$ .  $\diamond$

### 2.1.2 Kronecker product, vectorization, and matrix equations

We start with the definition of the Kronecker product and vectorization operator.

**Definition 2.7 ([GV13, Section 1.3.6]):**

For  $A = [a_{ij}] \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{p \times q}$ , the Kronecker product  $A \otimes B \in \mathbb{R}^{mp \times nq}$  and vectorization  $\text{vec}(A) \in \mathbb{R}^{mn}$  are defined by

$$A \otimes B := \begin{bmatrix} a_{11}B & \cdots & a_{1n}B \\ \vdots & \ddots & \vdots \\ a_{m1}B & \cdots & a_{mn}B \end{bmatrix}, \quad \text{vec}(A) := \begin{bmatrix} A_{:,1} \\ A_{:,2} \\ \vdots \\ A_{:,n} \end{bmatrix}. \quad \diamond$$

Next, we state some properties of the Kronecker product and its relation to the vectorization operator.

**Proposition 2.8 ([GV13, Sections 1.3.6 and 1.3.7]):**

We have for any scalar  $\alpha$  and all matrices  $A, B, C, D$  of compatible dimensions:

1.  $(A \otimes B) \otimes C = A \otimes (B \otimes C)$ ,
2.  $(\alpha A) \otimes B = A \otimes (\alpha B) = \alpha(A \otimes B)$ ,
3.  $(A + B) \otimes C = A \otimes C + B \otimes C$ ,
4.  $A \otimes (B + C) = A \otimes B + A \otimes C$ ,

5.  $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$
6.  $(A \otimes B)^T = A^T \otimes B^T$ ,
7.  $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$  (for invertible  $A$  and  $B$ ),
8.  $\text{vec}(ABC) = (C^T \otimes A) \text{vec}(B)$ ,
9.  $\sigma(A \otimes B) = \{\lambda\mu \mid \lambda \in \sigma(A), \mu \in \sigma(B)\}$  (for square matrices  $A$  and  $B$ ).  $\diamond$

Note that the Kronecker product is not commutative, but the result is equal up to a permutation of rows and columns.

We will be interested in matrix equations of the form

$$AXB^T + CXD^T + E = 0, \tag{2.1}$$

where  $A, C \in \mathbb{R}^{n \times n}$ ,  $B, D \in \mathbb{R}^{m \times m}$ ,  $X, E \in \mathbb{R}^{n \times m}$ , and  $B$  and  $C$  are invertible. Vectorization of (2.1) gives

$$(B \otimes A + D \otimes C) \text{vec}(X) = -\text{vec}(E).$$

Clearly, the above equation has a unique solution if and only if the matrix  $B \otimes A + D \otimes C$  is invertible. The following theorem gives the equivalent condition.

**Theorem 2.9 ([Chu87, Theorem 1]):**

The matrix equation (2.1) has a unique solution if and only if  $\lambda + \mu \neq 0$  for all  $\lambda \in \sigma(A, C)$  and  $\mu \in \sigma(D, B)$ .  $\diamond$

Notice that if  $\sigma(A, C) \subset \mathbb{C}_-$  and  $\sigma(D, B) \subset \mathbb{C}_-$ , then the condition in Theorem 2.9 is satisfied.

## 2.2 Functional analysis

Here, we give necessary basics of calculus in Banach spaces. In particular, we present differentiability, implicit function theorem, and Lagrange multiplier method. We base the presentation on the textbooks [Zei85a, Col12].

### 2.2.1 Fréchet and Gateaux differentiability

We only consider vector spaces over the field of real numbers  $\mathbb{R}$ . For normed vector spaces  $X, Y$ , we use  $B(X, Y)$  to denote the space of bounded linear operators  $A: X \rightarrow Y$ . The dual space of  $X$  is denoted  $X^* := B(X, \mathbb{R})$ .

We begin with the definition of Fréchet differentiability.

**Definition 2.10** ([Col12, Section 2.2], [Zei85a, Definition 4.5]):

Let  $X$  and  $Y$  be normed vector spaces with norms  $\|\cdot\|_X$  and  $\|\cdot\|_Y$ , respectively, and  $f: U \rightarrow Y$  a function, where  $U$  is an open subset of  $X$  and  $x$  an element of  $U$ .

The function  $f$  is *Fréchet differentiable at  $x$*  if there exists an operator  $A \in B(X, Y)$  such that

$$\lim_{h \rightarrow 0} \frac{\|f(x+h) - f(x) - Ah\|_Y}{\|h\|_X} = 0.$$

If it exists, this  $A$  is called the *Fréchet differential of  $f$  at  $x$*  and is denoted by  $Df(x)$ .  $\diamond$

It can be seen that the operator  $A$  in the above definition is unique (see [Col12, Proposition 2.2]), which justifies using the notation  $Df(x)$ .

For a function  $r: U \subseteq X \rightarrow Y$ , defined on some open neighborhood of zero, we write  $r(h) = o(\|h\|^n)$  to mean

$$\lim_{h \rightarrow 0} \frac{r(h)}{\|h\|^n} = 0.$$

Then, for a function  $f$  which is Fréchet differentiable at  $x$ , we can write

$$f(x+h) = f(x) + Df(x)h + o(\|h\|).$$

Similar properties as for finite-dimensional spaces hold here.

**Proposition 2.11** ([Col12, Proposition 2.7], [Zei85a, Proposition 4.9]):

Let  $X$  and  $Y$  be normed vector spaces,  $U$  an open subset of  $X$ , and  $x$  an element of  $U$ . If  $f, g: U \rightarrow Y$  are Fréchet differentiable at  $x$ , then  $f + g$  is differentiable at  $x$ , as is  $\alpha f$ , for any  $\alpha \in \mathbb{R}$ , and

$$D(f+g)(x) = Df(x) + Dg(x), \quad D(\alpha f)(x) = \alpha Df(x). \quad \diamond$$

**Proposition 2.12** ([Col12, Theorem 2.1], [Zei85a, Proposition 4.10]):

Let  $X, Y, Z$  be normed vector spaces,  $U$  an open subset of  $X$ ,  $V$  an open subset of  $Y$ ,  $x$  an element of  $U$ , and  $f: U \rightarrow Y$  and  $g: V \rightarrow Z$  functions such that  $f(U) \subseteq V$ . If  $f$  is Fréchet differentiable at  $x$  and  $g$  is Fréchet differentiable at  $f(x)$ , then  $g \circ f$  is Fréchet differentiable at  $x$  and

$$D(g \circ f)(x) = Dg(f(x)) Df(x). \quad \diamond$$

The following is the extension of continuously differentiability to functions between normed vector spaces.

**Definition 2.13** ([Col12, Section 2.4], [Zei85a, Definition 4.22]):

Let  $X$  and  $Y$  be normed vector spaces and  $f: U \rightarrow Y$  a function, where  $U$  is an open subset of  $X$ . The function  $f$  is of *class  $C^1$*  if it is Fréchet differentiable at every  $x \in U$  and  $Df: U \rightarrow B(X, Y)$  is continuous.  $\diamond$

Next is the extension of directional derivatives.

**Definition 2.14** ([Col12, Section 2.1], [Zei85a, Definition 4.5]):

Let  $X$  and  $Y$  be normed vector spaces and  $f: U \rightarrow Y$  a function, where  $U$  is an open subset of  $X$  containing  $x \in X$ .

1. For  $h \in X$ , if the limit

$$\lim_{t \rightarrow 0} \frac{f(x + th) - f(x)}{t}$$

exists, we call it the *directional derivative of  $f$  at  $x$  in direction  $h$*  and denote it by  $df(x; h)$ .

2. The function  $f$  is *Gateaux differentiable at  $x$*  if  $df(x; h)$  exists for all  $h \in X$  and there exists an operator  $A \in B(X, Y)$  such that  $df(x; h) = Ah$  for all  $h \in X$ . If it exists, the operator  $A$  is called the *Gateaux differential of  $f$  at  $x$*  and is denoted by  $D^G f(x)$ . ◇

Similarly as for Fréchet differential, the Gateaux differential of a function is unique if it exists.

The following result gives the relation between Fréchet and Gateaux differentials.

**Proposition 2.15** ([Zei85a, Proposition 4.8]):

Let  $X$  and  $Y$  be normed vector spaces and  $f: U \rightarrow Y$  a function, where  $U$  is an open subset of  $X$  containing  $x \in X$ . Then,

1. if  $f$  is Fréchet differentiable at  $x$ , then it is also Gateaux differentiable at  $x$  and  $Df(x) = D^G f(x)$ ,
2. if  $f$  is Gateaux differentiable in some neighborhood of  $x$  and  $D^G f$  is continuous at  $x$ , then  $f$  is Fréchet differentiable at  $x$  and  $Df(x) = D^G f(x)$ . ◇

This result provides an alternative way of proving a function is Fréchet differentiable and finding its differential. In particular, assuming  $t \mapsto f(x + th)$  is continuously differentiable, we have that

$$df(x; h) = \lim_{t \rightarrow 0} \frac{f(x + th) - f(x)}{t} = \left. \frac{d}{dt} f(x + th) \right|_{t=0}.$$

Therefore, methods for computing derivatives of functions of a real variable can be used to find a candidate for the Fréchet differential.

Notice that if  $Y = \mathbb{R}$ , then  $Df(x) \in B(X, \mathbb{R}) = X^*$ . Therefore, if  $X$  is a Hilbert space, functional  $Df(x)$  can be identified with an element of  $X$  by the Riesz representation theorem ([Col12, Theorem 6.4]).



**Definition 2.16 ([Col12, Section 6.4]):**

Let  $X$  be a Hilbert space with inner product  $\langle \cdot, \cdot \rangle_X$ ,  $U$  an open subset of  $X$ , and  $f: U \rightarrow \mathbb{R}$  a function. Let  $f$  be Fréchet differentiable at  $x \in U$ . The element  $a \in X$ , such that  $Df(x)h = \langle a, h \rangle_X$  for all  $h \in X$ , is called the *gradient of  $f$  at  $x$*  and is denoted by  $\nabla f(x)$ .  $\diamond$

Next is the extension of partial differentials.

**Definition 2.17 ([Zei85a, Definition 4.13]):**

Let  $X, Y, Z$  be Banach spaces,  $U$  an open subset of  $X \times Y$ ,  $(x_0, y_0)$  an element of  $U$ , and  $f: U \rightarrow Z$  a function.

Define an open subset  $U_{y_0} = \{x \in X : (x, y_0) \in U\} \subseteq X$  and a function  $g: U_{y_0} \rightarrow Z$  with  $g(x) = f(x, y_0)$  for all  $x \in U_{y_0}$ . Let  $g$  be Fréchet differentiable at  $x_0$ . Then  $Dg(x_0)$  is called the *partial Fréchet differential of  $f$  at  $(x_0, y_0)$  with respect to the first variable  $x$*  and is denoted by  $D_x f(x_0, y_0)$ .  $\diamond$

The differential  $D_y f(x_0, y_0)$  is defined similarly.

The following result states the relation between Fréchet and partial Fréchet differentials.

**Proposition 2.18 ([Zei85a, Proposition 4.14]):**

Let  $X, Y, Z$  be Banach spaces,  $U$  an open subset of  $X \times Y$ ,  $(x_0, y_0)$  an element of  $U$ , and  $f: U \rightarrow Z$  a function.

1. If  $f$  is Fréchet differentiable at  $(x_0, y_0)$ , then its partial Fréchet differentials  $D_x f(x_0, y_0)$  and  $D_y f(x_0, y_0)$  also exist and

$$Df(x_0, y_0)(h, k) = D_x f(x_0, y_0)h + D_y f(x_0, y_0)k,$$

for all  $h \in X$  and  $k \in Y$ .

2. If  $f$  has partial Fréchet differentials  $D_x f$  and  $D_y f$  in some neighborhood of  $(x_0, y_0)$  and if these are continuous at  $(x_0, y_0)$ , then  $f$  is Fréchet differentiable at  $(x_0, y_0)$ .
3. The function  $f$  is continuously Fréchet differentiable in a neighborhood of  $(x_0, y_0)$  if and only if all partial Fréchet differentials are continuous in a neighborhood of  $(x_0, y_0)$ .  $\diamond$

Partial differential of a function of more variables  $f(x_1, x_2, \dots, x_n)$  are defined similarly and the above result can be directly extended. Furthermore, partial Gateaux differential, partial directional derivatives, and partial gradients can be defined in a similar way.

### 2.2.2 Implicit function theorem

We can now state the implicit function theorem.

**Theorem 2.19** ([Zei85a, Theorem 4.B], [Col12, Theorem 8.2]):

Let  $X, Y, Z$  be Banach spaces,  $U$  an open subset of  $X \times Y$ ,  $(x_0, y_0)$  an element of  $U$ , and  $f: U \rightarrow Z$  a function. Suppose that

1.  $f(x_0, y_0) = 0$ ,
2.  $D_y f$  exists on  $U$  and  $D_y f(x_0, y_0)$  is bijective,
3.  $f$  and  $D_y f$  are continuous at  $(x_0, y_0)$ .

Then there exist open subsets  $U_X \subseteq X$  and  $U_Y \subseteq Y$  and a bijection  $g: U_X \rightarrow U_Y$  such that  $x_0 \in U_X$ ,  $y_0 \in U_Y$ ,  $g(x_0) = y_0$ , and  $f(x, g(x)) = 0$  for all  $x \in U_X$ .

Additionally, if  $f$  is of class  $C^1$  on a neighborhood of  $(x_0, y_0)$ , then  $g$  is of class  $C^1$  on a neighborhood of  $x_0$  and  $Dg(x_0) = -D_y f(x_0, g(x_0))^{-1} D_x f(x_0, g(x_0))$ .  $\diamond$

### 2.2.3 Constrained optimization and Lagrange multipliers

We conclude with the theorem about Lagrange multipliers. First, we need to define when is a function a submersion [Zei85b, Definition 43.15].

**Definition 2.20:**

Let  $X, Y$  be Banach spaces,  $U$  an open subset of  $X$ ,  $x_0$  an element of  $U$ , and  $g: U \rightarrow Y$  a function. The function  $g$  is a *submersion* at  $x_0$  if

1.  $g$  is of class  $C^1$  in a neighborhood of  $x_0$ ,
2.  $Dg(x_0): X \rightarrow Y$  is surjective,
3. the null space  $\ker(Dg(x_0))$  splits  $X$ , i.e., there exists a continuous projection operator  $P$  of  $X$  on  $\ker(Dg(x_0))$ .  $\diamond$

Note that the third condition is immediately satisfied if  $X$  is a Hilbert space.

Next is the theorem about constrained local minima.

**Theorem 2.21** ([Zei85b, Theorem 43.D]):

Let  $X, Y$  be real Banach spaces,  $U$  an open subset of  $X$ ,  $x_0$  an element of  $U$ ,  $f: U \rightarrow \mathbb{R}$  and  $g: U \rightarrow Y$  functions.

Suppose  $f$  is Fréchet differentiable at  $x_0$  and  $g$  is a submersion at  $x_0$ , and that  $f$  has a constrained local minimum at  $x_0$  with respect to  $\{x \in U : g(x) = 0\}$ . Then there exists  $\lambda \in Y^*$  such that

$$Df(x_0)h - \lambda(Dg(x_0)h) = 0,$$

for all  $h \in X$ .  $\diamond$

For Hilbert spaces, we can use gradients.

**Corollary 2.22:**

Let  $X, Y$  be real Hilbert spaces,  $U$  an open subset  $X$ ,  $x_0$  an element of  $U$ ,  $f: U \rightarrow \mathbb{R}$  and  $g: U \rightarrow Y$  functions.

Suppose  $f$  is Fréchet differentiable at  $x_0$ ,  $g$  is of class  $C^1$  in a neighborhood of  $x_0$ ,  $Dg(x_0)$  is surjective, and that  $f$  has a constrained local minimum at  $x_0$  with respect to  $\{x \in U : g(x) = 0\}$ . Then there exists  $\lambda \in Y$  such that

$$\langle \nabla f(x_0), h \rangle - \langle \lambda, Dg(x_0)h \rangle = 0,$$

for all  $h \in X$ . ◇

## 2.3 Linear time-invariant systems

We consider finite-dimensional, continuous-time, linear time-invariant (LTI) systems of the form

$$\begin{aligned} E\dot{x}(t) &= Ax(t) + Bu(t), & x(0) &= x_0, \\ y(t) &= Cx(t) + Du(t), \end{aligned} \tag{2.2}$$

with *system matrices*  $E, A \in \mathbb{R}^{n \times n}$ , *input matrix*  $B \in \mathbb{R}^{n \times m}$ , *output matrix*  $C \in \mathbb{R}^{p \times n}$ , and *feedthrough matrix*  $D \in \mathbb{R}^{p \times m}$ . Furthermore,  $t \in \mathbb{R}$  is *time*,  $x_0 \in \mathbb{R}^n$  is the initial condition, while  $x(t) \in \mathbb{R}^n$ ,  $u(t) \in \mathbb{R}^m$ , and  $y(t) \in \mathbb{R}^p$  are respectively the *state*, *input*, and *output* of the system. We call dimension  $n$  of the state  $x(t)$  the *order* of the system (2.2). Throughout this thesis, we assume  $E$  to be invertible. Thereby, the state equation in (2.2) is a system of ODEs after multiplying from the left by  $E^{-1}$ .

If there is only one input and one output, i.e., if  $m = p = 1$ , we will refer to such a system as a *single-input single-output (SISO) system*. Otherwise, we will call it a *multiple-input multiple-output (MIMO) system*.

We will use  $(E; A, B, C, D)$  to denote the system (2.2). If  $D$  is a zero matrix, we will also use  $(E; A, B, C)$ . Furthermore, we will use  $(A, B, C, D)$  and  $(A, B, C)$  to mean  $(I_n; A, B, C, D)$  and  $(I_n; A, B, C)$ , respectively.

In this thesis, we focus on continuous-time systems. Results for discrete-time systems

$$\begin{aligned} Ex(k+1) &= Ax(k) + Bu(k), & x(0) &= x_0, \\ y(k) &= Cx(k) + Du(k), \end{aligned}$$

should be directly extendable due to similarities to continuous-time systems (see [Ant05] for more details).

### 2.3.1 Solutions and stability

It can be verified that, for given initial condition  $x_0$  and input  $u$ , the state and output of the system (2.2) satisfy

$$\begin{aligned} x(t) &= e^{tE^{-1}A}x_0 + \int_0^t e^{\tau E^{-1}A}E^{-1}Bu(t-\tau) \, d\tau, \\ y(t) &= Ce^{tE^{-1}A}x_0 + \int_0^t Ce^{\tau E^{-1}A}E^{-1}Bu(t-\tau) \, d\tau + Du(t). \end{aligned}$$

The system (2.2) can also be solved in the frequency-domain. Assuming  $X$ ,  $U$ , and  $Y$  exist as Laplace transforms of  $x$ ,  $u$ , and  $y$ , by applying Laplace transform to the system (2.2), we find that

$$\begin{aligned} sEX(s) - Ex_0 &= AX(s) + BU(s), \\ Y(s) &= CX(s) + DU(s). \end{aligned}$$

After eliminating  $X(s)$ , we obtain

$$Y(s) = C(sE - A)^{-1}Ex_0 + (C(sE - A)^{-1}B + D)U(s).$$

The function  $H(s) = C(sE - A)^{-1}B + D$  is called the *transfer function*, and it characterizes the input-output relationship when  $x_0 = 0$ .

For any pair of invertible matrices  $S, T \in \mathbb{R}^{n \times n}$ , if we define a new state  $\tilde{x}(t) = T^{-1}x(t)$ , we obtain an equivalent system

$$\begin{aligned} SET\tilde{x}(t) &= SAT\tilde{x}(t) + SBu(t), \quad \tilde{x}(0) = T^{-1}x_0, \\ y(t) &= CT\tilde{x}(t) + Du(t), \end{aligned}$$

in the sense that the input and output remain the same. Notice that also the transfer function  $H$  is independent of the choice of  $S$  and  $T$ . It is of interest to see which other system properties are invariant under this transformations.

Important properties of LTI systems are stability and asymptotic stability.

**Definition 2.23** ([Ant05, Section 5.8]):

The autonomous system  $E\dot{x}(t) = Ax(t)$  is called *stable* if all of its solutions  $x$  are bounded for positive time, i.e., the set  $\{x(t) \mid t > 0\}$  is bounded. If additionally  $\lim_{t \rightarrow \infty} x(t) = 0$ , it is called *asymptotically stable*.

The system  $(E; A, B, C, D)$  is (*asymptotically*) *stable* if the corresponding autonomous system  $E\dot{x}(t) = Ax(t)$  is (asymptotically) stable.  $\diamond$

It is known that the system  $(E; A, B, C, D)$  is

- asymptotically stable if and only if  $\sigma(A, E) \subset \mathbb{C}_-$ .
- stable if and only if  $\sigma(A, E) \subset \overline{\mathbb{C}_-}$  and all purely imaginary eigenvalues are semi-simple (their geometric multiplicity and algebraic multiplicity are equal).

We call a matrix  $A$  (a matrix pair  $(A, E)$ ) *asymptotically stable* or *Hurwitz* if  $\sigma(A) \subset \mathbb{C}_-$  ( $\sigma(A, E) \subset \mathbb{C}_-$ ). Additionally, we call a transfer function  $H(s) = C(sE - A)^{-1}B + D$  *asymptotically stable* or *Hurwitz* if its poles are all in  $\mathbb{C}_-$ .

Notice that if the system  $(E; A, B, C, D)$  is asymptotically stable, then its transfer function is also asymptotically stable, but the converse does not hold (e.g., if  $B$  or  $C$  is a zero matrix).

### 2.3.2 Controllability, observability, and Gramians

Further important properties are reachability and observability. We first discuss reachability and controllability (see [DP00, Section 2.2] and [Ant05, Section 4.2.1]).

**Definition 2.24:**

Let  $(E; A, B, C, D)$  be a system with  $A \in \mathbb{R}^{n \times n}$  and  $B \in \mathbb{R}^{n \times m}$ . A state  $x_1 \in \mathbb{R}^n$  is *reachable by time  $t > 0$*  if there exists an input  $u \in \mathcal{L}_2(0, t)$  and a trajectory  $x$  such that  $x(0) = 0$ ,  $x(t) = x_1$ , and  $E\dot{x}(\tau) = Ax(\tau) + Bu(\tau)$  for almost all  $\tau \in (0, t)$ .

The set  $\mathcal{R}(t) \subseteq \mathbb{R}^n$  of all reachable states by time  $t > 0$  is called the *reachable set at time  $t$* .

The  $n \times nm$  matrix

$$\mathcal{C}(E; A, B) = \mathcal{C}(E^{-1}A, E^{-1}B) = [E^{-1}B \quad E^{-1}AE^{-1}B \quad \dots \quad (E^{-1}A)^{n-1}E^{-1}B]$$

is called the *controllability matrix*.

For  $t > 0$ , the *controllability Gramian* is the  $n \times n$  matrix

$$P(t) = \int_0^t e^{E^{-1}At} E^{-1}BB^T E^{-T} e^{A^T E^{-T}t} dt. \quad \diamond$$

We have the following result relating the reachable set, controllability matrix, and controllability Gramian.

**Proposition 2.25 ([DP00, Theorem 2.2]):**

For all  $t > 0$ , equality

$$\mathcal{R}(t) = \text{im}(\mathcal{C}(E; A, B)) = \text{im}(P(t))$$

holds. \(\diamond\)

Since  $\mathcal{C}(E; A, B)$  does not depend on  $t$ , we have that  $\mathcal{R}(t)$  and  $\text{im}(P(t))$  are also constant.

**Definition 2.26:**

The system  $(E; A, B, C, D)$  is *reachable* if  $\mathcal{R}(t) = \mathbb{R}^n$ , for some  $t > 0$ . \(\diamond\)

For continuous-time systems, reachability is equivalent to controllability in the following sense (see [Ant05, Theorem 4.18]).

**Definition 2.27:**

The system  $(E; A, B, C, D)$  is *controllable* if for any pair of states  $x_0, x_1 \in \mathbb{R}^n$  there exist a time  $t > 0$ , an input  $u \in \mathcal{L}_2(0, t)$ , and trajectory  $x$  such that  $x(0) = x_0$ ,  $x(t) = x_1$ , and  $E\dot{x}(\tau) = Ax(\tau) + Bu(\tau)$  for almost all  $\tau \in (0, t)$ .  $\diamond$

Thus, we will use both terms. The following theorem gives equivalent conditions for controllability.

**Theorem 2.28 ([Ant05, Theorem 4.15]):**

The following are equivalent:

1. The system  $(E; A, B, C, D)$  is controllable.
2. The controllability matrix is of full rank:  $\text{rank}(\mathcal{C}(E; A, B)) = n$ .
3. The controllability Gramian  $P(t)$  is positive definite for all  $t > 0$ .
4. If  $v$  is a left eigenvector of  $(A, E)$ , then  $v^*B \neq 0$ .
5.  $\text{rank}([\lambda E - A \quad B]) = n$  for all  $\lambda \in \mathbb{C}$ .  $\diamond$

The following proposition gives an interesting property of the controllability matrix.

**Proposition 2.29 ([DP00, Proposition 2.11]):**

Suppose  $E, A \in \mathbb{R}^{n \times n}$  and  $B \in \mathbb{R}^{n \times m}$ . Then  $\text{im}(\mathcal{C}(E; A, B))$  is invariant under  $E^{-1}A$ .  $\diamond$

In particular, we can see that  $\text{im}(\mathcal{C}(E; A, B))$  is the smallest subspace which contains  $\text{im}(E^{-1}B)$  and is invariant under  $E^{-1}A$ .

Next, we discuss observability (see [DP00, Section 2.4] and [Ant05, Section 4.2.2]).

**Definition 2.30:**

Let  $(E; A, B, C, D)$  be a system with  $A \in \mathbb{R}^{n \times n}$  and  $C \in \mathbb{R}^{p \times n}$ . A state  $x_0 \in \mathbb{R}^n$  is *unobservable* if  $Ce^{tE^{-1}A}x_0 = 0$  for all  $t > 0$ .

The set  $\mathcal{N} \subseteq \mathbb{R}^n$  of all unobservable states is called the *unobservable set*. If  $\mathcal{N} = \{0\}$ , the system is called *observable*.

The  $np \times n$  matrix

$$\mathcal{O}(E; A, C) = \mathcal{C}(E^{-1}A, C) = \begin{bmatrix} C \\ CE^{-1}A \\ \vdots \\ C(E^{-1}A)^{n-1} \end{bmatrix}$$

is called the *observability matrix*.

For  $t > 0$ , the *observability Gramian* is the  $n \times n$  matrix

$$Q(t) = \int_0^t E^{-T} e^{A^T E^{-T} t} C^T C e^{E^{-1} A t} E^{-1} dt. \quad \diamond$$

Here is a result relating the unobservable set and observability matrix.

**Proposition 2.31** ([DP00, Theorem 2.20]):

The unobservable set satisfies

$$\mathcal{N} = \ker(\mathcal{O}(E; A, C)). \quad \diamond$$

The following gives equivalent conditions for observability. Notice the analogy with Theorem 2.28.

**Theorem 2.32** ([Ant05, Theorem 4.26]):

The following are equivalent:

1. The system  $(E; A, B, C, D)$  is observable.
2. The observability matrix is of full rank:  $\text{rank}(\mathcal{O}(E; A, C)) = n$ .
3. The observability Gramian  $Q(t)$  is positive definite for all  $t > 0$ .
4. If  $v$  is a right eigenvector of  $(A, E)$ , then  $Cv \neq 0$ .
5.  $\text{rank}([\lambda E^T - A^T \quad C^T]) = n$  for all  $\lambda \in \mathbb{C}$ . \(\diamond\)

We introduced controllability and observability Gramians which depend on a finite time  $t$ . Important concepts are infinite Gramians [MS05, Section 3.2.4].

**Definition 2.33:**

For an asymptotically stable system  $(E; A, B, C, D)$  the (*infinite*) *controllability Gramian* is

$$P := \int_0^\infty e^{E^{-1}At} E^{-1} B B^T E^{-T} e^{A^T E^{-T}t} dt$$

and the (*infinite*) *observability Gramian* is

$$Q := \int_0^\infty E^{-T} e^{A^T E^{-T}t} C^T C e^{E^{-1}At} E^{-1} dt. \quad \diamond$$

They can also be represented in the frequency domain.

**Proposition 2.34** ([Ant05, Section 4.3], [MS05, Section 3.2.4]):

Under the assumptions of Definition 2.33, the Gramians satisfy

$$\begin{aligned} P &= \frac{1}{2\pi} \int_{-\infty}^\infty (\boldsymbol{\omega}E - A)^{-1} B B^T (-\boldsymbol{\omega}E - A)^{-T} d\boldsymbol{\omega}, \\ Q &= \frac{1}{2\pi} \int_{-\infty}^\infty (-\boldsymbol{\omega}E - A)^{-T} C^T C (\boldsymbol{\omega}E - A)^{-1} d\boldsymbol{\omega}. \end{aligned} \quad \diamond$$

Solving Lyapunov equations can be used to find the Gramians.

**Proposition 2.35** ([Ant05, Proposition 4.27]):

Under the assumptions of Definition 2.33, the Gramians are the unique solutions to the following Lyapunov equations:

$$\begin{aligned} APE^T + EPA^T + BB^T &= 0, \\ A^TQE + E^TQA + C^TC &= 0. \end{aligned} \quad \diamond$$

Compared to Definition 2.33 or Proposition 2.34, Proposition 2.35 provides an approach to efficient computation of Gramians (see surveys [BS13, Sim16]).

### 2.3.3 Hardy spaces and system norms

The transfer function  $H(s) = C(sE - A)^{-1}B + D$  is a rational matrix function. Under some assumptions,  $H$  is an element of Hardy spaces  $\mathcal{H}_2^{p \times m}$  and/or  $\mathcal{H}_\infty^{p \times m}$ .

**Definition 2.36** ([Ant05, Section 5.1.3]):

Hardy space  $\mathcal{H}_2^{p \times m}$ :

$$\mathcal{H}_2^{p \times m} := \left\{ H : \mathbb{C}_+ \rightarrow \mathbb{C}^{p \times m} \mid \begin{array}{l} H \text{ is analytic and} \\ \sup_{\xi > 0} \int_{-\infty}^{\infty} \|H(\xi + \boldsymbol{\nu}\omega)\|_F^2 d\omega < \infty \end{array} \right\},$$

with norm

$$\|H\|_{\mathcal{H}_2} := \sqrt{\sup_{\xi > 0} \int_{-\infty}^{\infty} \|H(\xi + \boldsymbol{\nu}\omega)\|_F^2 d\omega}.$$

Hardy space  $\mathcal{H}_\infty^{p \times m}$ :

$$\mathcal{H}_\infty^{p \times m} := \left\{ H : \mathbb{C}_+ \rightarrow \mathbb{C}^{p \times m} \mid \begin{array}{l} H \text{ is analytic and} \\ \sup_{s \in \mathbb{C}_+} \|H(s)\|_2 < \infty \end{array} \right\}$$

with norm

$$\|H\|_{\mathcal{H}_\infty} := \sup_{s \in \mathbb{C}_+} \|H(s)\|_2. \quad \diamond$$

For simplicity, we will often write  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  instead of  $\mathcal{H}_2^{p \times m}$  and  $\mathcal{H}_\infty^{p \times m}$ .

It can be seen that if the system is asymptotically stable, then  $H$  is an element of  $\mathcal{H}_\infty$ . If additionally  $D = 0$ , then  $H$  is also an element of  $\mathcal{H}_2$ .

**Remark 2.37:**

Functions from Hardy spaces can be extended to the imaginary axis. Furthermore,  $\mathcal{H}_2$  can be shown to be a Hilbert space with inner product

$$\langle H, G \rangle_{\mathcal{H}_2} := \frac{1}{2\pi} \int_{-\infty}^{\infty} \text{tr}(H(\boldsymbol{\nu}\omega)^* G(\boldsymbol{\nu}\omega)) d\omega$$

and  $\mathcal{H}_\infty$  a Banach space with norm

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

For more details, see [Ant05] and [ZDG96]. \(\diamond\)



The following results states that the Gramians can be used to compute the  $\mathcal{H}_2$ -norm.

**Proposition 2.38** ([Ant05, Section 5.5.1]):

For an asymptotically stable system (2.2),

$$\|H\|_{\mathcal{H}_2}^2 = \text{tr}(CPC^T) = \text{tr}(B^TQB). \quad \diamond$$

## 2.4 Model order reduction

### 2.4.1 Projection-based model reduction

We consider the system (2.2). Here, we assume homogeneous initial conditions, i.e.,  $x_0 = 0$ . For some recent work on MOR approaches for inhomogeneous initial conditions, see [HRA11, BGM17].

We can write the system (2.2) in variational form (see [ABG10, Section 2.2]):

Find  $x(t)$  contained in  $\mathbb{R}^n$  such that

$$E\dot{x}(t) - Ax(t) - Bu(t) \perp \mathbb{R}^n.$$

Then the associated output is  $y(t) = Cx(t) + Du(t)$ .

Petrov-Galerkin projection consists of choosing two  $r$ -dimensional subspaces  $\mathcal{V}, \mathcal{W} \subset \mathbb{R}^n$ , where  $r < n$  is the reduced order, and defining the ROM by

Find  $v(t)$  contained in  $\mathcal{V}$  such that

$$E\dot{v}(t) - Av(t) - Bu(t) \perp \mathcal{W}.$$

Then the associated output is  $\hat{y}(t) = Cv(t) + Du(t)$ .

Choosing  $V, W \in \mathbb{R}^{n \times r}$  such that  $\text{im}(V) = \mathcal{V}$  and  $\text{im}(W) = \mathcal{W}$ , we have  $v(t) = V\hat{x}(t)$  for some  $\hat{x}(t) \in \mathbb{R}^r$  and  $W^T(EV\dot{\hat{x}}(t) - AV\hat{x}(t) - Bu(t)) = 0$ . A state-space form of the ROM is

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

with

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

In the following sections, we give an overview of some particular projection-based methods.

### 2.4.2 Balanced truncation

This method is based on the controllability and observability energy functionals

$$E_c(x_0) := \inf \left\{ \|u\|_{\mathcal{L}_2(-\infty,0)}^2 \mid E\dot{x}(t) = Ax(t) + Bu(t), x(-\infty) = 0, x(0) = x_0 \right\},$$

$$E_o(x_0) := \left\| Ce^{tE^{-1}A}x_0 \right\|_{\mathcal{L}_2(0,\infty)}^2.$$

We can interpret  $E_c(x_0)$  as the energy necessary to reach state  $x_0$  and  $E_o(x_0)$  as the observed energy from the system when it starts from state  $x_0$ . The following result relates energy functionals to Gramians.

**Proposition 2.39:**

For a controllable system  $(E; A, B, C, D)$ , we have

$$E_c(x_0) = x_0^T P^{-1}x_0 \quad \text{and} \quad E_o(x_0) = x_0^T E^T Q E x_0. \quad \diamond$$

For a controllable and observable system, it is possible to transform the system such that  $\tilde{P}$  and  $\tilde{E}^T \tilde{Q} \tilde{E}$  become equal to a matrix  $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$  where  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$  (see [Ant05, Lemma 7.3]). If there exists an  $r$  such that  $\sigma_r \gg \sigma_{r+1}$ , then we say that, in the new coordinate system, the states  $e_1, \dots, e_r$  are easy to reach and easy to observe, while the states  $e_{r+1}, \dots, e_n$  are difficult to reach and difficult to observe. The idea is then to truncate the latter states. Balanced truncation (BT), using the balancing-free square root method, is described in Algorithm 2.1. The following theorem gives the a priori error bound which allows automatic choice of the reduced order.

**Theorem 2.40 ([Ant05, Theorem 7.9]):**

Let  $(E; A, B, C, D)$  be an asymptotically stable, controllable, and observable system of order  $n$  and  $(\hat{E}; \hat{A}, \hat{B}, \hat{C}, \hat{D})$  a ROM obtained by BT of order  $r < n$ . If  $\sigma_r \neq \sigma_{r+1}$ , then the ROM is asymptotically stable, controllable, observable, and

$$\|H - \hat{H}\|_{\mathcal{H}_\infty} \leq 2 \sum_{i=r+1}^n \sigma_i,$$

where  $H$  and  $\hat{H}$  are the transfer functions of the two systems. \(\diamond\)

### 2.4.3 Rational interpolation

MOR using rational interpolation, also called moment matching or Krylov methods, is described here. The main idea is formulated in the following theorem.

**Theorem 2.41 ([ABG10, Theorem 2]):**

Let  $H(s) = C(sE - A)^{-1}B + D$  and  $\hat{H}(s) = \hat{C}(s\hat{E} - \hat{A})^{-1}\hat{B} + \hat{D}$  be two transfer functions with  $\hat{E} = W^T E V$ ,  $\hat{A} = W^T A V$ ,  $\hat{B} = W^T B$ ,  $\hat{C} = C V$ , and  $\hat{D} = D$ . Furthermore, let  $\sigma, \mu \in \mathbb{C}$  be such that  $sE - A$  and  $s\hat{E} - \hat{A}$  are invertible for  $s = \sigma, \mu$ . If  $b \in \mathbb{C}^m$  and  $c \in \mathbb{C}^p$  are fixed nontrivial vectors, then

**Algorithm 2.1:** Balancing-free square root balanced truncation method

**Input:** Asymptotically stable system  $(E; A, B, C, D)$  of order  $n$ , reduced order  $r$ .

**Output:** Reduced-order model  $(\widehat{E}; \widehat{A}, \widehat{B}, \widehat{C}, \widehat{D})$ .

- 1 For Gramians  $P$  and  $Q$  solving Lyapunov equations

$$\begin{aligned} APE^T + EPA^T + BB^T &= 0, \\ A^TQE + E^TQA + C^TC &= 0, \end{aligned}$$

find Cholesky decompositions  $P = Z_P Z_P^T$  and  $Q = Z_Q Z_Q^T$ , with  $Z_P, Z_Q \in \mathbb{R}^{n \times n}$ , or low-rank approximations  $Z_P \in \mathbb{R}^{n \times r_P}$ ,  $Z_Q \in \mathbb{R}^{n \times r_Q}$ , with  $r \leq \min(r_P, r_Q)$ .

- 2 Compute a singular value decomposition

$$Z_Q^T E Z_P = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix}$$

where  $\Sigma_1 \in \mathbb{R}^{r \times r}$ .

- 3 Use thin QR decompositions to find  $V, W \in \mathbb{R}^{n \times r}$  with orthonormal columns such that  $\text{im}(V) = \text{im}(Z_P V_1)$  and  $\text{im}(W) = \text{im}(Z_Q U_1)$ .
- 4 Project to get reduced matrices

$$\widehat{E} = W^T E V, \quad \widehat{A} = W^T A V, \quad \widehat{B} = W^T B, \quad \widehat{C} = C V, \quad \widehat{D} = D.$$

1. if

$$\left( (\sigma E - A)^{-1} E \right)^i (\sigma E - A)^{-1} B b \in \text{im}(V) \text{ for } i = 0, 1, \dots, N - 1, \quad (2.3)$$

then

$$H^{(k)}(\sigma) b = \widehat{H}^{(k)}(\sigma) b \text{ for } k = 0, 1, \dots, N - 1,$$

2. if

$$\left( (\mu E - A)^{-T} E^T \right)^j (\mu E - A)^{-T} C^T c \in \text{im}(W) \text{ for } j = 0, 1, \dots, M - 1, \quad (2.4)$$

then

$$c^T H^{(k)}(\mu) = c^T \widehat{H}^{(k)}(\mu) \text{ for } k = 0, 1, \dots, M - 1,$$

3. if both (2.3) and (2.4) hold and  $\sigma = \mu$ , then

$$c^T H^{(k)}(\sigma) b = c^T \widehat{H}^{(k)}(\sigma) b \text{ for } k = 1, 2, \dots, M + N - 1. \quad \diamond$$

In particular, if

$$\begin{aligned} \text{im}(V) &= \text{im}\left([\!(\sigma_1 E - A)^{-1} B b_1 \quad (\sigma_2 E - A)^{-1} B b_2 \quad \cdots \quad (\sigma_r E - A)^{-1} B b_r\!] \right), \\ \text{im}(W) &= \text{im}\left([\!(\mu_1 E - A)^{-T} C^T c_1 \quad (\mu_2 E - A)^{-T} C^T c_2 \quad \cdots \quad (\mu_r E - A)^{-T} C^T c_r\!] \right), \end{aligned}$$

then we have

$$H(\sigma_i) b_i = \widehat{H}(\sigma_i) b_i \text{ and } c_i^T H(\mu_i) = c_i^T \widehat{H}(\mu_i), \text{ for } i = 1, 2, \dots, r.$$

If additionally  $\sigma_i = \mu_i$ , then also

$$c_i^T H'(\sigma_i) b_i = c_i^T \widehat{H}'(\sigma_i) b_i, \text{ for } i = 1, 2, \dots, r.$$

This case is particularly relevant for  $\mathcal{H}_2$ -optimal MOR, where the ROM necessarily satisfies such interpolatory conditions. More details follow in the next section.

## 2.4.4 $\mathcal{H}_2$ -optimal model order reduction

There are a few approaches for  $\mathcal{H}_2$ -optimal MOR of first-order systems. We will in particular focus on two: interpolation-based and Gramian-based.

Meier and Luenberger [ML67] derived interpolatory necessary optimality conditions for SISO systems. Gugercin, Beattie, and Antoulas [GAB08, ABG10] generalized this to MIMO systems and proposed iterative rational Krylov algorithm (IRKA).

Wilson [Wil70] developed Gramian-based optimality conditions in form of coupled matrix equation. Xu and Zeng [XZ11] used this to propose two-sided iteration algorithm (TSIA). The same algorithm was proposed in parallel by Van Dooren, Gallivan, and Absil [VDGA08, VDGA10]. Benner, Köhler, and Saak [BKS11] investigate some implementation issues in TSIA related to Sylvester equations and orthonormalization.

Let us take the first-order system (2.2) as the full-order model, with  $E$  invertible and  $\lambda E - A$  an asymptotically stable matrix pencil. The  $\mathcal{H}_2$ -optimal MOR problem is finding a ROM  $(\widehat{E}; \widehat{A}, \widehat{B}, \widehat{C}, \widehat{D})$  with transfer function  $\widehat{H}(s) = \widehat{C}(s\widehat{E} - \widehat{A})^{-1} \widehat{B} + \widehat{D}$ , such that the  $\mathcal{H}_2$ -error  $\|H - \widehat{H}\|_{\mathcal{H}_2}$  is locally minimized. Here,  $\widehat{E}, \widehat{A} \in \mathbb{R}^{r \times r}$ ,  $\widehat{B} \in \mathbb{R}^{r \times m}$ ,  $\widehat{C} \in \mathbb{R}^{p \times r}$ , and  $\widehat{D} \in \mathbb{R}^{p \times m}$ , with  $r < n$ . For the  $\mathcal{H}_2$ -error to be defined, it is necessary that  $\widehat{D} = D$ . In particular, we can assume  $D = \widehat{D} = 0$ .

### 2.4.4.1 Interpolation-based approach

The motivation for the interpolation-based approach are the interpolatory necessary optimality conditions, given in the following theorem. A transfer function  $H$  is called real if  $\overline{H(s)} = H(\overline{s})$  for all  $s$  in the domain of  $H$ .

---

**Algorithm 2.2:** Iterative rational Krylov algorithm (IRKA)
 

---

**Input:** System  $(E; A, B, C)$ , initial shifts  $\sigma_i$  and tangential directions  $b_i, c_i$ ,  
 $i = 1, 2, \dots, r$ .

**Output:** Reduced-order model  $(\widehat{E}; \widehat{A}, \widehat{B}, \widehat{C})$ .

1 **while** *not converged* **do**

2     Find  $V, W \in \mathbb{R}^{n \times r}$  with orthonormal columns such that

$$\begin{aligned} \operatorname{im}(V) &= \operatorname{im}\left([\left(\sigma_1 E - A\right)^{-1} B b_1 \quad \dots \quad \left(\sigma_r E - A\right)^{-1} B b_r\right], \\ \operatorname{im}(W) &= \operatorname{im}\left([\left(\sigma_1 E - A\right)^{-\mathrm{T}} C^{\mathrm{T}} c_1 \quad \dots \quad \left(\sigma_r E - A\right)^{-\mathrm{T}} C^{\mathrm{T}} c_r\right]. \end{aligned}$$

3     Project  $\widehat{E} = W^{\mathrm{T}} E V$ ,  $\widehat{A} = W^{\mathrm{T}} A V$ ,  $\widehat{B} = W^{\mathrm{T}} B$ ,  $\widehat{C} = C V$ .

4     Compute  $\operatorname{diag}(\lambda_i)$ ,  $X, Y \in \mathbb{C}^{r \times r}$  such that  $Y^{\mathrm{T}} \widehat{A} X = \operatorname{diag}(\lambda_i)$  and  
 $Y^{\mathrm{T}} \widehat{E} X = I$ .

5     Update interpolation points and tangential directions:

$$\sigma_i = -\lambda_i, \quad b_i = \widehat{B}^{\mathrm{T}} Y e_i, \quad c_i = \widehat{C} X e_i, \quad \text{for } i = 1, 2, \dots, r.$$


---

**Theorem 2.42** ([[ABG10](#), Theorem 5]):

Suppose  $H \in \mathcal{H}_2$  and  $\widehat{H}(s) = \sum_{i=1}^r \frac{c_i b_i^{\mathrm{T}}}{s - \lambda_i}$ , with  $\lambda_i$  pairwise distinct, are real transfer functions. Let  $\widehat{H}$  be a locally  $\mathcal{H}_2$ -optimal ROM of order  $r$ . Then

$$H(-\lambda_i) b_i = \widehat{H}(-\lambda_i) b_i, \quad (2.5a)$$

$$c_i^{\mathrm{T}} H(-\lambda_i) = c_i^{\mathrm{T}} \widehat{H}(-\lambda_i), \quad (2.5b)$$

$$c_i^{\mathrm{T}} H'(-\lambda_i) b_i = c_i^{\mathrm{T}} \widehat{H}'(-\lambda_i) b_i, \quad (2.5c)$$

for  $i = 1, 2, \dots, r$ . ◇

This results states that the transfer function  $\widehat{H}$  of a locally  $\mathcal{H}_2$ -optimal ROM is a bitangential Hermite interpolant, at the reflected poles of  $\widehat{H}$ , of the full-order model's transfer function  $H$ .

For given  $\lambda_i, b_i, c_i$ , the interpolation conditions in (2.5) can be achieved by a Petrov-Galerkin projection (see Theorem 2.41) where  $V$  and  $W$  span tangential rational Krylov subspaces

$$\begin{aligned} \operatorname{im}(V) &= \operatorname{im}\left([\left(-\lambda_1 E - A\right)^{-1} B b_1 \quad \dots \quad \left(-\lambda_r E - A\right)^{-1} B b_r\right], \\ \operatorname{im}(W) &= \operatorname{im}\left([\left(-\lambda_1 E - A\right)^{-\mathrm{T}} C^{\mathrm{T}} c_1 \quad \dots \quad \left(-\lambda_r E - A\right)^{-\mathrm{T}} C^{\mathrm{T}} c_r\right]. \end{aligned}$$

The difficulty is that  $\lambda_i, b_i, c_i$  are not given in advance. Gugercin et al. [[GAB08](#), [ABG10](#)] proposed IRKA (see Algorithm 1 in [[ABG10](#)]). The pseudocode is presented in Algorithm 2.2.

### 2.4.4.2 Gramian-based approach

We use  $\widehat{P}$  and  $\widehat{Q}$  to denote the controllability and observability Gramians of  $(\widehat{E}; \widehat{A}, \widehat{B}, \widehat{C})$ , which solve reduced Lyapunov equations

$$\widehat{A}\widehat{P}\widehat{E}^T + \widehat{E}\widehat{P}\widehat{A}^T + \widehat{B}\widehat{B}^T = 0, \quad (2.6a)$$

$$\widehat{A}^T\widehat{Q}\widehat{E} + \widehat{E}^T\widehat{Q}\widehat{A} + \widehat{C}^T\widehat{C} = 0. \quad (2.6b)$$

Furthermore, we have that the controllability and observability Gramians of the error system

$$\underbrace{\begin{bmatrix} E & 0 \\ 0 & \widehat{E} \end{bmatrix}}_{E_{\text{err}}} \underbrace{\begin{bmatrix} \dot{x}(t) \\ \dot{\widehat{x}}(t) \end{bmatrix}}_{\dot{x}_{\text{err}}} = \underbrace{\begin{bmatrix} A & 0 \\ 0 & \widehat{A} \end{bmatrix}}_{A_{\text{err}}} \underbrace{\begin{bmatrix} x(t) \\ \widehat{x}(t) \end{bmatrix}}_{x_{\text{err}}} + \underbrace{\begin{bmatrix} B \\ \widehat{B} \end{bmatrix}}_{B_{\text{err}}} u(t),$$

$$y(t) - \widehat{y}(t) = \underbrace{\begin{bmatrix} C & -\widehat{C} \end{bmatrix}}_{C_{\text{err}}} \underbrace{\begin{bmatrix} x(t) \\ \widehat{x}(t) \end{bmatrix}}_{x_{\text{err}}},$$

are

$$P_{\text{err}} = \begin{bmatrix} P & \widetilde{P} \\ \widetilde{P}^T & \widehat{P} \end{bmatrix} \quad \text{and} \quad Q_{\text{err}} = \begin{bmatrix} Q & \widetilde{Q} \\ \widetilde{Q}^T & \widehat{Q} \end{bmatrix},$$

where  $\widetilde{P}, \widetilde{Q} \in \mathbb{R}^{n \times r}$  solve Sylvester equations

$$A\widetilde{P}\widehat{E}^T + E\widetilde{P}\widehat{A}^T + B\widehat{B}^T = 0, \quad (2.7a)$$

$$A^T\widetilde{Q}\widehat{E} + E^T\widetilde{Q}\widehat{A} - C^T\widehat{C} = 0. \quad (2.7b)$$

Gramian-based approach uses Wilson conditions [Wil70], given in the following theorem. We include a proof because the result is generalized to include invertible matrices  $E$  and  $\widehat{E}$ . Additionally, the proof uses Lagrange multiplier method, instead of computing the derivatives directly as in [Wil70] or [VDGA08]. We will use the following lemma in the proof.

**Lemma 2.43:**

Let  $A \in \mathbb{R}^{n \times m}$ ,  $B \in \mathbb{R}^{n \times n}$ ,  $C \in \mathbb{R}^{m \times m}$ , and  $f, g: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$  such that  $f(X) = \text{tr}(AX^T)$  and  $g(X) = \text{tr}(BXCX^T)$ . Then  $\nabla f(X_0) = A$  and  $\nabla g(X_0) = BX_0C + B^T X_0 C^T$  for arbitrary  $X_0 \in \mathbb{R}^{n \times m}$ .  $\diamond$

*Proof.* We have

$$f(X_0 + H) = \text{tr}\left(A(X_0 + H)^T\right) = \text{tr}(AX_0^T) + \text{tr}(AH^T) = f(X_0) + \langle A, H \rangle,$$

from which it directly follows that  $f$  is Fréchet differentiable at  $X_0$  and  $\nabla f(X_0) = A$ .

For  $g$  we have

$$\begin{aligned} g(X_0 + H) &= \text{tr}\left(B(X_0 + H)C(X_0 + H)^T\right) \\ &= \text{tr}(BX_0CX_0^T) + \text{tr}(BX_0CH^T) + \text{tr}(BHCX_0^T) + \text{tr}(BHCH^T) \\ &= g(X_0) + \langle BX_0C + B^T X_0C^T, H \rangle + o(\|H\|), \end{aligned}$$

so  $g$  is also Fréchet differentiable at  $X_0$  and  $\nabla g(X_0) = BX_0C + B^T X_0C^T$ .  $\square$

**Theorem 2.44:**

Let  $(\widehat{E}; \widehat{A}, \widehat{B}, \widehat{C})$  be a locally  $\mathcal{H}_2$ -optimal ROM for  $(E; A, B, C)$ . Then

$$\begin{aligned} \widetilde{Q}^T E \widetilde{P} + \widehat{Q} \widehat{E} \widehat{P} &= 0, \\ \widetilde{Q}^T A \widetilde{P} + \widehat{Q} \widehat{A} \widehat{P} &= 0, \\ \widetilde{Q}^T B + \widehat{Q} \widehat{B} &= 0, \\ C \widetilde{P} - \widehat{C} \widehat{P} &= 0. \end{aligned} \tag{2.8} \quad \diamond$$

*Proof.* The proof consists of applying the Lagrange multiplier method to the optimization problem

$$\begin{aligned} &\text{minimize} && \text{tr}(C_{\text{err}} P_{\text{err}} C_{\text{err}}^T), \\ &\text{subject to} && (2.7a), (2.6a). \end{aligned}$$

The Lagrange function is

$$\begin{aligned} \mathcal{L}(\widehat{E}, \widehat{A}, \widehat{B}, \widehat{C}, \widetilde{P}, \widehat{P}, \widetilde{\Lambda}, \widehat{\Lambda}) &= \text{tr}\left(CPC^T - 2C\widetilde{P}\widehat{C}^T + \widehat{C}\widehat{P}\widehat{C}^T\right) \\ &\quad + \text{tr}\left(\widetilde{\Lambda}^T \left(A\widetilde{P}\widehat{E}^T + E\widetilde{P}\widehat{A}^T + B\widehat{B}^T\right)\right) \\ &\quad + \text{tr}\left(\widehat{\Lambda}^T \left(\widehat{A}\widehat{P}\widehat{E}^T + \widehat{E}\widehat{P}\widehat{A}^T + \widehat{B}\widehat{B}^T\right)\right), \end{aligned}$$

where  $\widetilde{\Lambda} \in \mathbb{R}^{r \times r}$  and  $\widehat{\Lambda} \in \mathbb{R}^{r \times r}$  are the Lagrange multipliers. The gradients of  $\mathcal{L}$  with respect to the Gramians are

$$\begin{aligned} \nabla_{\widetilde{P}} \mathcal{L} &= -2C^T \widehat{C} + A^T \widetilde{\Lambda} \widehat{E} + E^T \widetilde{\Lambda} \widehat{A}, \\ \nabla_{\widehat{P}} \mathcal{L} &= \widehat{C}^T \widehat{C} + \widehat{A}^T \widehat{\Lambda} \widehat{E} + \widehat{E}^T \widehat{\Lambda} \widehat{A}. \end{aligned}$$

From  $\nabla_{\widetilde{P}} \mathcal{L} = 0$  and  $\nabla_{\widehat{P}} \mathcal{L} = 0$  it follows that  $\widetilde{\Lambda} = 2\widetilde{Q}$  and  $\widehat{\Lambda} = \widehat{Q}$ . The Lagrange function now simplifies to

$$\begin{aligned} \mathcal{L} &= \text{tr}\left(CPC^T - 2C\widetilde{P}\widehat{C}^T + \widehat{C}\widehat{P}\widehat{C}^T\right) \\ &\quad + \text{tr}\left(2\widetilde{Q}^T A\widetilde{P}\widehat{E}^T + 2\widetilde{Q}^T E\widetilde{P}\widehat{A}^T + 2\widetilde{Q}^T B\widehat{B}^T\right) \\ &\quad + \text{tr}\left(2\widehat{Q}\widehat{A}\widehat{P}\widehat{E}^T + \widehat{Q}\widehat{B}\widehat{B}^T\right). \end{aligned}$$

---

**Algorithm 2.3:** Two-sided iteration algorithm (TSIA)

---

**Input:** System  $(E, A, B, C)$  and initial reduced-order model  $(\widehat{E}, \widehat{A}, \widehat{B}, \widehat{C})$ .

**Output:** Reduced-order model  $(\widehat{E}, \widehat{A}, \widehat{B}, \widehat{C})$  approximately satisfying (2.8).

- 1 **while** *not converged* **do**
  - 2     Solve  $A\widetilde{P}\widehat{E}^T + E\widetilde{P}\widehat{A}^T + B\widehat{B}^T = 0$  and  $A^T\widetilde{Q}\widehat{E} + E^T\widetilde{Q}\widehat{A} - C^T\widehat{C} = 0$ .
  - 3     Find  $V, W \in \mathbb{R}^{n \times r}$  with orthonormal columns such that  $\text{im}(V) = \text{im}(\widetilde{P})$  and  $\text{im}(W) = \text{im}(\widetilde{Q})$ .
  - 4     Project  $\widehat{E} = W^T E V$ ,  $\widehat{A} = W^T A V$ ,  $\widehat{B} = W^T B$ ,  $\widehat{C} = C V$ .
- 

Finally, the gradients of the Lagrange function with respect to the reduced matrices are

$$\begin{aligned}\nabla_{\widehat{E}}\mathcal{L} &= 2\widetilde{Q}^T A\widetilde{P} + 2\widehat{Q}\widehat{A}\widehat{P}, \\ \nabla_{\widehat{A}}\mathcal{L} &= 2\widetilde{Q}^T E\widetilde{P} + 2\widehat{Q}\widehat{E}\widehat{P}, \\ \nabla_{\widehat{B}}\mathcal{L} &= 2\widetilde{Q}^T B + 2\widehat{Q}\widehat{B}, \\ \nabla_{\widehat{C}}\mathcal{L} &= -2C\widetilde{P} + 2\widehat{C}\widehat{P},\end{aligned}$$

which completes the proof. □

Based on these optimality conditions, Xu and Zeng [XZ11] proposed TSIA, given in Algorithm 2.3.

Let  $(\widehat{A}, \widehat{E})$  be diagonalizable and  $\Lambda, X, Y \in \mathbb{C}^{r \times r}$  such that  $Y^T \widehat{E} X = I_r$  and  $Y^T \widehat{A} X = \Lambda$ . Then also  $\widehat{A}x_i = \lambda_i \widehat{E}x_i$  and  $\widehat{A}^T y_i = \lambda_i \widehat{E}^T y_i$ , where  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_r)$ ,  $X = [x_1 \ x_2 \ \dots \ x_r]$ , and  $Y = [y_1 \ y_2 \ \dots \ y_r]$ . Then we have

$$\begin{aligned}A\widetilde{P}\widehat{E}^T y_i + E\widetilde{P}\widehat{A}^T y_i + B\widehat{B}^T y_i &= 0, \\ A^T\widetilde{Q}\widehat{E}x_i + E^T\widetilde{Q}\widehat{A}x_i - C^T\widehat{C}x_i &= 0,\end{aligned}$$

which implies

$$\begin{aligned}(A + \lambda_i E)\widetilde{P}\widehat{E}^T y_i + B\widehat{B}^T y_i &= 0, \\ (A^T + \lambda_i E^T)\widetilde{Q}\widehat{E}x_i - C^T\widehat{C}x_i &= 0.\end{aligned}$$

Therefore,

$$\begin{aligned}\widetilde{P}\widehat{E}^T y_i &= (-\lambda_i E - A)^{-1} B\widehat{B}^T y_i, \\ \widetilde{Q}\widehat{E}x_i &= -(-\lambda_i E - A)^{-T} C^T \widehat{C}x_i,\end{aligned}$$

which shows the connection to IRKA. This can also be used to prove interpolatory conditions in Theorem 2.42 using the Wilson conditions from Theorem 2.44, as given in the following theorem (see [VDGA08, Theorem 4.1]).



**Theorem 2.45:**

Let  $\widehat{H}(s) = \widehat{C}(s\widehat{E} - \widehat{A})^{-1}\widehat{B} = \sum_{i=1}^r \frac{c_i b_i^T}{s - \lambda_i}$ , with  $\lambda_i$  pairwise distinct,  $S^T \widehat{E} T = I_r$ ,  $S^T \widehat{A} T = \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_r)$ ,  $b_i^T = s_i^T \widehat{B}$ ,  $s_i = S e_i$ ,  $c_i = C t_i$ , and  $t_i = T e_i$ . Then

$$\begin{aligned} (C\widetilde{P} - \widehat{C}\widehat{P}) \widehat{E}^T s_i &= (H(-\lambda_i) - \widehat{H}(-\lambda_i)) b_i, \\ -t_i^T \widehat{E}^T (\widetilde{Q}^T B + \widehat{Q}\widehat{B}) &= c_i^T (H(-\lambda_i) - \widehat{H}(-\lambda_i)), \\ t_i^T \widehat{E}^T (\widetilde{Q}^T E\widetilde{P} + \widehat{Q}\widehat{E}\widehat{P}) \widehat{E}^T s_i &= c_i^T (H'(-\lambda_i) - \widehat{H}'(-\lambda_i)) b_i, \\ t_i^T \widehat{E}^T (\widetilde{Q}^T E\widetilde{P} + \widehat{Q}\widehat{E}\widehat{P}) \widehat{E}^T s_j &= c_i^T \left( \frac{H(-\lambda_i) - H(-\lambda_j)}{(-\lambda_i) - (-\lambda_j)} - \frac{\widehat{H}(-\lambda_i) - \widehat{H}(-\lambda_j)}{(-\lambda_i) - (-\lambda_j)} \right) b_j, \\ t_i^T \widehat{E}^T (\widetilde{Q}^T A\widetilde{P} + \widehat{Q}\widehat{A}\widehat{P}) \widehat{E}^T s_i &= c_i^T ([sH(s)]'|_{s=-\lambda_i} - [s\widehat{H}(s)]'|_{s=-\lambda_i}) b_i, \\ t_i^T \widehat{E}^T (\widetilde{Q}^T A\widetilde{P} + \widehat{Q}\widehat{A}\widehat{P}) \widehat{E}^T s_j &= c_i^T \left( \frac{(-\lambda_i)H(-\lambda_i) - (-\lambda_j)H(-\lambda_j)}{(-\lambda_i) - (-\lambda_j)} \right. \\ &\quad \left. - \frac{(-\lambda_i)\widehat{H}(-\lambda_i) - (-\lambda_j)\widehat{H}(-\lambda_j)}{(-\lambda_i) - (-\lambda_j)} \right) b_j, \end{aligned}$$

for  $i, j = 1, 2, \dots, r$ ,  $i \neq j$ . ◇

*Proof.* From  $S^T \widehat{E} T = I_r$  and  $S^T \widehat{A} T = \Lambda$ , we get  $\widehat{A} t_i = \lambda_i \widehat{E} t_i$  and  $s_i^T \widehat{A} = \lambda_i s_i^T \widehat{E}$ . Next, from

$$\begin{aligned} A\widetilde{P}\widehat{E}^T s_i + E\widetilde{P}\widehat{A}^T s_i + B\widehat{B}^T s_i &= 0, \\ \widehat{A}\widehat{P}\widehat{E}^T s_i + \widehat{E}\widehat{P}\widehat{A}^T s_i + \widehat{B}\widehat{B}^T s_i &= 0, \\ A^T \widetilde{Q}\widehat{E} t_i + E^T \widetilde{Q}\widehat{A} t_i - C^T \widehat{C} t_i &= 0, \\ \widehat{A}^T \widehat{Q}\widehat{E} t_i + \widehat{E}^T \widehat{Q}\widehat{A} t_i + \widehat{C}^T \widehat{C} t_i &= 0, \end{aligned}$$

we find

$$\begin{aligned} \widetilde{P}\widehat{E}^T s_i &= (-\lambda_i E - A)^{-1} B b_i, \\ \widehat{P}\widehat{E}^T s_i &= (-\lambda_i \widehat{E} - \widehat{A})^{-1} \widehat{B} b_i, \\ \widetilde{Q}\widehat{E} t_i &= -(-\lambda_i E - A)^{-T} C^T c_i, \\ \widehat{Q}\widehat{E} t_i &= (-\lambda_i \widehat{E} - \widehat{A})^{-T} \widehat{C}^T c_i. \end{aligned}$$

Now it follows that

$$\begin{aligned} (C\widetilde{P} - \widehat{C}\widehat{P}) \widehat{E}^T s_i &= C(-\lambda_i E - A)^{-1} B b_i - \widehat{C} (-\lambda_i \widehat{E} - \widehat{A})^{-1} \widehat{B} b_i \\ &= (H(-\lambda_i) - \widehat{H}(-\lambda_i)) b_i, \end{aligned}$$

and

$$\begin{aligned} -t_i^\top \widehat{E}^\top \left( \widetilde{Q}^\top B + \widehat{Q} \widehat{B} \right) &= c_i^\top C(-\lambda_i E - A)^{-1} B - c_i^\top \widehat{C} \left( -\lambda_i \widehat{E} - \widehat{A} \right)^{-1} \widehat{B} \\ &= c_i^\top \left( H(-\lambda_i) - \widehat{H}(-\lambda_i) \right). \end{aligned}$$

From

$$\begin{aligned} H'(s) &= -C(sE - A)^{-1} E(sE - A)^{-1} B, \\ H(s_1) - H(s_2) &= -(s_1 - s_2) C(s_1 E - A)^{-1} E(s_2 E - A)^{-1} B, \\ [sH(s)]' &= -C(sE - A)^{-1} A(sE - A)^{-1} B, \\ s_1 H(s_1) - s_2 H(s_2) &= -(s_1 - s_2) C(s_1 E - A)^{-1} A(s_2 E - A)^{-1} B, \end{aligned}$$

we find

$$\begin{aligned} &t_i^\top \widehat{E}^\top \left( \widetilde{Q}^\top E \widetilde{P} + \widehat{Q} \widehat{E} \widehat{P} \right) \widehat{E}^\top s_i \\ &= -c_i^\top C(-\lambda_i E - A)^{-1} E(-\lambda_i E - A)^{-1} B b_i + c_i^\top \widehat{C}(-\lambda_i \widehat{E} - \widehat{A})^{-1} \widehat{E}(-\lambda_i \widehat{E} - \widehat{A})^{-1} \widehat{B} b_i \\ &= c_i^\top \left( H'(-\lambda_i) - \widehat{H}'(-\lambda_i) \right) b_i, \\ &t_i^\top \widehat{E}^\top \left( \widetilde{Q}^\top E \widetilde{P} + \widehat{Q} \widehat{E} \widehat{P} \right) \widehat{E}^\top s_j \\ &= -c_i^\top C(-\lambda_i E - A)^{-1} E(-\lambda_j E - A)^{-1} B b_j + c_i^\top \widehat{C}(-\lambda_i \widehat{E} - \widehat{A})^{-1} \widehat{E}(-\lambda_j \widehat{E} - \widehat{A})^{-1} \widehat{B} b_j \\ &= c_i^\top \left( \frac{H(-\lambda_i) - H(-\lambda_j)}{(-\lambda_i) - (-\lambda_j)} - \frac{\widehat{H}(-\lambda_i) - \widehat{H}(-\lambda_j)}{(-\lambda_i) - (-\lambda_j)} \right) b_j, \\ &t_i^\top \widehat{E}^\top \left( \widetilde{Q}^\top A \widetilde{P} + \widehat{Q} \widehat{A} \widehat{P} \right) \widehat{E}^\top s_i \\ &= -c_i^\top C(-\lambda_i E - A)^{-1} A(-\lambda_i E - A)^{-1} B b_i + c_i^\top \widehat{C}(-\lambda_i \widehat{E} - \widehat{A})^{-1} \widehat{A}(-\lambda_i \widehat{E} - \widehat{A})^{-1} \widehat{B} b_i \\ &= c_i^\top \left( [sH(s)]' \Big|_{s=-\lambda_i} - [s\widehat{H}(s)]' \Big|_{s=-\lambda_i} \right) b_i, \\ &t_i^\top \widehat{E}^\top \left( \widetilde{Q}^\top A \widetilde{P} + \widehat{Q} \widehat{A} \widehat{P} \right) \widehat{E}^\top s_j \\ &= -c_i^\top C(-\lambda_i E - A)^{-1} A(-\lambda_j E - A)^{-1} B b_j + c_i^\top \widehat{C}(-\lambda_i \widehat{E} - \widehat{A})^{-1} \widehat{A}(-\lambda_j \widehat{E} - \widehat{A})^{-1} \widehat{B} b_j \\ &= c_i^\top \left( \frac{(-\lambda_i)H(-\lambda_i) - (-\lambda_j)H(-\lambda_j)}{(-\lambda_i) - (-\lambda_j)} - \frac{(-\lambda_i)\widehat{H}(-\lambda_i) - (-\lambda_j)\widehat{H}(-\lambda_j)}{(-\lambda_i) - (-\lambda_j)} \right) b_j, \end{aligned}$$

which concludes the proof.  $\square$

#### 2.4.4.3 Comparison of the two approaches

The interpolation-based and Gramian-based approaches differ in two important ways.

The first is the assumption of diagonalizability. Notice that Theorem 2.42 has the assumption that the ROM has pairwise distinct poles, while Theorem 2.44 does not. Diagonalizability of  $\widehat{E}^{-1}\widehat{A}$  is a generic property, in the sense that the set of non-diagonalizable matrices forms a set of measure zero in the set of all matrices, which justifies this assumption for unstructured first-order systems. But for structured systems, where the assumption of simultaneous diagonalizability of two or more matrices appears, will no longer be a generic property. Additionally, as illustrated by Van Dooren, Gallivan, and Absil [VDGA10], even in the case of first-order systems, if the optimal ROM has a higher-order pole, this can cause numerical issues in IRKA, while TSIA will still converge. The difference in the algorithms is in solving sparse-dense Sylvester equations, where the (generalized) Schur decomposition of  $\widehat{A}$  and  $\widehat{E}$  is sufficient (details in [BKS11]).

The second difference is that the Gramian-based approach gives necessary optimality conditions which are equations that the reduced matrices need to satisfy. In particular, for first-order systems, it is clear from Theorem 2.44 that the ROM is necessarily obtained by Petrov-Galerkin projection, assuming  $\widehat{P}$  and  $\widehat{Q}$  are invertible. This is not immediately clear from Theorem 2.42, where the interpolatory necessary optimality conditions do not explicitly give the equations for the reduced matrices.

One advantage of the interpolation-based approach is its direct application to MOR of infinite-dimensional systems where it is possible to evaluate the transfer function and its derivative. This was used to develop the *transfer function IRKA* in [BG12]. In the Gramian-based approach, access to matrices  $E, A, B, C$  is necessary. One possibility could be to extend these results to infinite-dimensional systems where  $E, A, B, C$  become operators. Then the corresponding Sylvester equations could be solved in a similar way as Lyapunov equations in [ORW13], but this is outside the scope of this thesis.

### 2.4.5 Model order reduction of unstable systems

So far, we assumed the system to be asymptotically stable. Since we will also need to consider systems which are not asymptotically stable, we need to see if and how the previously discussed MOR methods could be applied to such systems. Therefore, here we consider a system  $(E; A, B, C, D)$  with an invertible  $E$ , but which is not asymptotically stable. We will not assume minimality (i.e., both controllability and observability), which means that its transfer function  $H$  may be asymptotically stable.

In the following chapters, we will be interested in finding ROMs for which the  $\mathcal{H}_2$ -error  $\|H - \widehat{H}\|_{\mathcal{H}_2}$  or the  $\mathcal{H}_\infty$ -error  $\|H - \widehat{H}\|_{\mathcal{H}_\infty}$  are small. Particularly, for these quantities to be defined, it is necessary that the transfer function  $H - \widehat{H}$  is asymptotically stable. This implies that the ROM's transfer function  $\widehat{H}$  needs to have the same unstable part as the original  $H$ . Therefore, we will consider additive decomposition into the asymptotically stable and unstable part as in [Enn85].

There are works ([YCDAGX93, Zil91, BNBG10]) about MOR for unstable systems

using shifting, but this approach does not guarantee preservation of unstable poles. The work in [ZSW99] gives an extension for BT, assuming there are no poles on the imaginary axis. There is also work extending  $\mathcal{H}_2$ -optimal MOR to unstable systems ([MBG10, BBG19]), which also assumes there are no poles on the imaginary axis or cannot preserve unstable poles.

Let  $H = H_- + H_+$ , where  $H_-$  has poles in  $\mathbb{C}_-$  and  $H_+$  in  $\overline{\mathbb{C}_+}$ . Furthermore, let  $S$  and  $T$  be such that

$$S^T E T = \begin{bmatrix} E_- & 0 \\ 0 & E_+ \end{bmatrix}, \quad S^T A T = \begin{bmatrix} A_- & 0 \\ 0 & A_+ \end{bmatrix}, \quad S^T B = \begin{bmatrix} B_- \\ B_+ \end{bmatrix}, \quad \text{and } C T = [C_- \quad C_+],$$

where

$$T = [T_- \quad T_+] \in \mathbb{C}^{n \times n} \text{ and } S = [S_- \quad S_+] \in \mathbb{C}^{n \times n},$$

with  $\sigma(A_-, E_-) \subset \mathbb{C}_-$  and  $\sigma(A_+, E_+) \subset \overline{\mathbb{C}_+}$ . Notice that  $H_-(s) = C_-(sE_- - A_-)^{-1}B_-$  and  $H_+(s) = C_+(sE_+ - A_+)^{-1}B_+$ . To find a ROM  $\widehat{H} = \widehat{H}_- + H_+$ , we can apply Petrov-Galerkin projection to the asymptotically stable part  $(E_-; A_-, B_-, C_-)$  to get  $(\widehat{E}_-; \widehat{A}_-, \widehat{B}_-, \widehat{C}_-)$  with

$$\widehat{E}_- = W_-^T E_- V_-, \quad \widehat{A}_- = W_-^T A_- V_-, \quad \widehat{B}_- = W_-^T B_-, \quad \widehat{C}_- = C_- V_-.$$

Then we have

$$\begin{aligned} \widehat{E} &= \begin{bmatrix} \widehat{E}_- & 0 \\ 0 & E_+ \end{bmatrix} \\ &= \begin{bmatrix} W_-^T & 0 \\ 0 & I \end{bmatrix} [S_- \quad S_+]^T E [T_- \quad T_+] \begin{bmatrix} V_- & 0 \\ 0 & I \end{bmatrix} \\ &= [S_- W_- \quad S_+]^T E [T_- V_- \quad T_+]. \end{aligned}$$

In a similar way, we find  $\widehat{A} = W^T A V$ ,  $\widehat{B} = W^T B$ ,  $\widehat{C} = C V$  for  $V = [T_- V_- \quad T_+]$  and  $W = [S_- W_- \quad S_+]$ . Therefore, Petrov-Galerkin projection of  $(E_-; A_-, B_-, C_-)$  can be represented as a Petrov-Galerkin projection of  $(E; A, B, C)$ . BT can also be applied to the asymptotically stable part and the  $\mathcal{H}_\infty$ -error bound will hold since  $\|H - \widehat{H}\|_{\mathcal{H}_\infty} = \|H_- - \widehat{H}_-\|_{\mathcal{H}_\infty}$ . Alternatively, IRKA can be applied to  $H_-$  to find  $\widehat{H}_-$  which is a local minimum for  $\|H_- - \widehat{H}_-\|_{\mathcal{H}_2}$ . This implies that  $\widehat{H} = \widehat{H}_- + H_+$  is a local minimum for  $\|H - \widehat{H}\|_{\mathcal{H}_2}$ . Notice that  $H - \widehat{H} = H_- - \widehat{H}_-$ , which can be used to compute the  $\mathcal{H}_2$  or  $\mathcal{H}_\infty$  error.

## 2.5 Graph theory

Parts of this thesis are about systems defined over networks. Thus, we present some basic concepts from graph theory here. Notation is based on [ME10] and [GR01].

### 2.5.1 Basic concepts

A graph  $\mathfrak{G}$  consists of a *vertex set*  $\mathfrak{V}$  and an *edge set*  $\mathfrak{E}$  encoding the relation between vertices. *Undirected* graphs are those for which the edge set is a subset of the set of all unordered pairs of vertices, i.e.,  $\mathfrak{E} \subseteq \{\{i, j\} : i, j \in \mathfrak{V}, i \neq j\}$ . On the other hand, a graph is *directed* if  $\mathfrak{E} \subseteq \{(i, j) : i, j \in \mathfrak{V}, i \neq j\}$ . We think of an edge  $(i, j)$  as an arrow starting from vertex  $i$  and ending at  $j$ .

**Remark 2.46:**

Notice that we exclude graphs with multiple copies of the same edge, i.e., *multigraphs*. Furthermore, we exclude graphs containing *self-loops*, i.e., edges of the form  $\{i\}$  or  $(i, i)$ . Therefore, we will only consider *simple* graphs.  $\diamond$

We will only consider *finite graphs*, i.e., graphs with a finite number of vertices  $\mathfrak{n} := |\mathfrak{V}|$ . Without loss of generality, let  $\mathfrak{V} = \{1, 2, \dots, \mathfrak{n}\}$ .

For an undirected graph, a *path* of length  $\ell$  is a sequence of distinct vertices  $i_0, i_1, \dots, i_\ell$  such that  $\{i_k, i_{k+1}\} \in \mathfrak{E}$  for  $k = 0, 1, \dots, \ell - 1$ . For a directed graph, a *directed path* of length  $\ell$  is a sequence of distinct vertices  $i_0, i_1, \dots, i_\ell$  such that  $(i_k, i_{k+1}) \in \mathfrak{E}$  for  $k = 0, 1, \dots, \ell - 1$ . An undirected graph is *connected* if there is a path between any two distinct vertices  $i, j \in \mathfrak{V}$ . A directed graph is *strongly connected* if there is a directed path between any two distinct vertices  $i, j \in \mathfrak{V}$ .

We can associate weights to edges of a graph by a *weight function*  $\mathfrak{w} : \mathfrak{E} \rightarrow \mathbb{R}$ . If  $\mathfrak{w}(\epsilon) > 0$  for all  $\epsilon \in \mathfrak{E}$ , the tuple  $\mathfrak{G} = (\mathfrak{V}, \mathfrak{E}, \mathfrak{w})$  is called a *weighted graph*. In the following, we will focus on weighted graphs. In particular, we will directly generalize concepts for unweighted graphs from [ME10, GR01], as was done in [MTC14].

The *adjacency matrix*  $\mathfrak{A} = [\mathfrak{a}_{ij}]_{i, j \in \mathfrak{V}} \in \mathbb{R}^{\mathfrak{n} \times \mathfrak{n}}$  of an undirected graph is defined component-wise by

$$\mathfrak{a}_{ij} := \begin{cases} \mathfrak{w}(\{i, j\}), & \text{if } \{i, j\} \in \mathfrak{E}, \\ 0, & \text{otherwise,} \end{cases}$$

and for a directed graph as

$$\mathfrak{a}_{ij} := \begin{cases} \mathfrak{w}((j, i)), & \text{if } (j, i) \in \mathfrak{E}, \\ 0, & \text{otherwise.} \end{cases}$$

For every vertex  $i \in \mathfrak{V}$ , its *in-degree* is  $\mathfrak{d}_i := \sum_{j=1}^{\mathfrak{n}} \mathfrak{a}_{ij}$ . The diagonal matrix  $\mathfrak{D} := \text{diag}(\mathfrak{d}_1, \mathfrak{d}_2, \dots, \mathfrak{d}_\mathfrak{n})$  is called the *in-degree matrix*. Notice that  $\mathfrak{D} = \text{diag}(\mathfrak{A}\mathbf{1})$ .

Let  $\epsilon_1, \epsilon_2, \dots, \epsilon_{|\mathfrak{E}|}$  be all the edges of  $\mathfrak{G}$  in some order. The *incidence matrix*  $\mathfrak{R} \in \mathbb{R}^{\mathfrak{n} \times |\mathfrak{E}|}$  of a directed graph  $\mathfrak{G}$  is defined component-wise

$$[\mathfrak{R}]_{ik} := \begin{cases} -1, & \text{if } \epsilon_k = (i, j) \text{ for some } j \in \mathfrak{V}, \\ 1, & \text{if } \epsilon_k = (j, i) \text{ for some } j \in \mathfrak{V}, \\ 0, & \text{otherwise.} \end{cases}$$

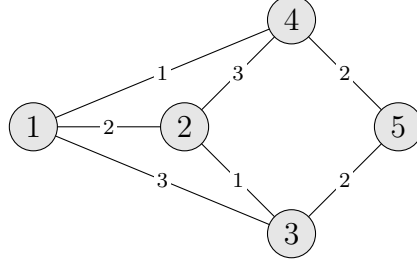


Figure 2.1: An undirected, weighted, connected graph

If  $\mathfrak{G}$  is undirected, we assign some orientation to every edge to define a directed graph  $\mathfrak{G}^o$ , and define the incidence matrix of  $\mathfrak{G}$  to be the incidence matrix of  $\mathfrak{G}^o$ . The *weight matrix* is defined as  $\mathfrak{W} := \text{diag}(\mathfrak{w}(\mathfrak{e}_1), \mathfrak{w}(\mathfrak{e}_2), \dots, \mathfrak{w}(\mathfrak{e}_{|\mathfrak{E}|}))$ .

The (*in-degree*) *Laplacian matrix*  $\mathfrak{L}$  is defined by  $\mathfrak{L} := \mathfrak{D} - \mathfrak{A}$ . For undirected graphs, it can be checked that  $\mathfrak{L} = \mathfrak{R}\mathfrak{W}\mathfrak{R}^T$ , using

$$\mathfrak{R}\mathfrak{W}\mathfrak{R}^T = \sum_{\{i,j\} \in \mathfrak{E}} \mathfrak{a}_{ij} (e_i - e_j)(e_i - e_j)^T,$$

which is independent of the order of edges defining  $\mathfrak{R}$  and  $\mathfrak{W}$  or the orientation of edges in  $\mathfrak{G}^o$ . From the definition of  $\mathfrak{L}$ , it directly follows that the sum of each row in  $\mathfrak{L}$  is zero, i.e.,  $\mathfrak{L}\mathbf{1} = 0$ . From  $\mathfrak{L} = \mathfrak{R}\mathfrak{W}\mathfrak{R}^T$ , we immediately see that, for undirected weighted graphs, the Laplacian matrix  $\mathfrak{L}$  is symmetric positive semidefinite.

**Example 2.1:**

For the graph in Figure 2.1, the adjacency matrix is

$$\mathfrak{A} = \begin{bmatrix} 0 & 2 & 3 & 1 & 0 \\ 2 & 0 & 1 & 3 & 0 \\ 3 & 1 & 0 & 0 & 2 \\ 1 & 3 & 0 & 0 & 2 \\ 0 & 0 & 2 & 2 & 0 \end{bmatrix},$$

degree matrix  $\mathfrak{D} = \text{diag}(6, 6, 6, 6, 4)$ , Laplacian and incidence matrix

$$\mathfrak{L} = \begin{bmatrix} 6 & -2 & -3 & -1 & 0 \\ -2 & 6 & -1 & -3 & 0 \\ -3 & -1 & 6 & 0 & -2 \\ -1 & -3 & 0 & 6 & -2 \\ 0 & 0 & -2 & -2 & 4 \end{bmatrix}, \quad \mathfrak{R} = \begin{bmatrix} -1 & -1 & -1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 & -1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & -1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix},$$

and weight matrix  $\mathfrak{W} = \text{diag}(2, 3, 1, 1, 3, 2, 2)$ . ◇

The following theorem states that connectedness of a graph is related to the spectral properties of  $\mathfrak{L}$ .

**Theorem 2.47 ([ME10, Theorem 2.8]):**

Let  $\mathfrak{G} = (\mathfrak{V}, \mathfrak{E}, \mathfrak{w})$  be an undirected weighted graph,  $\mathfrak{L}$  its Laplacian matrix, and  $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  the eigenvalues of  $\mathfrak{L}$ . Then the following statements are equivalent:

1.  $\mathfrak{G}$  is connected,
2.  $\lambda_2 > 0$ ,
3.  $\ker(\mathfrak{L}) = \text{im}(\mathbb{1})$ . ◇

### 2.5.2 Graph partitions

A nonempty subset  $\mathfrak{C} \subseteq \mathfrak{V}$  is called a *cluster* or *cell* of  $\mathfrak{V}$ . A *graph partition*  $\pi$  is a partition of the vertex set  $\mathfrak{V}$ . Vertices  $i$  and  $j$  are called *cellmates* in  $\pi$  if they belong to the same cell of  $\pi$ . The *characteristic vector* of a cluster  $\mathfrak{C} \subseteq \mathfrak{V}$  is the  $n$ -dimensional column vector  $\mathfrak{p}(\mathfrak{C})$  defined as

$$[\mathfrak{p}(\mathfrak{C})]_i = \begin{cases} 1 & \text{if } i \in \mathfrak{C}, \\ 0 & \text{otherwise.} \end{cases}$$

The *characteristic matrix* of a partition  $\pi = \{\mathfrak{C}_1, \mathfrak{C}_2, \dots, \mathfrak{C}_r\}$  of the graph  $\mathfrak{G}$  is the matrix  $\mathfrak{P} \in \mathbb{R}^{n \times r}$  defined by

$$\mathfrak{P} = [\mathfrak{p}(\mathfrak{C}_1) \quad \mathfrak{p}(\mathfrak{C}_2) \quad \dots \quad \mathfrak{p}(\mathfrak{C}_r)].$$

For a given vertex  $i \in \mathfrak{V}$  and a cluster  $\mathfrak{C}_q \in \pi$ , we define the *in-degree of vertex  $i$  with respect to cluster  $\mathfrak{C}_q$*  by

$$\mathfrak{d}(i, \mathfrak{C}_q) = \sum_{k \in \mathfrak{C}_q} \mathfrak{a}_{ik}.$$

The partition  $\pi$  is an *equitable partition* if  $\mathfrak{d}(i, \mathfrak{C}_q) = \mathfrak{d}(j, \mathfrak{C}_q)$  for all vertices  $i, j \in \mathfrak{C}_p$  and any two clusters  $\mathfrak{C}_p, \mathfrak{C}_q \in \pi$  (they can be the same cluster). The partition  $\pi$  is an *AEP* if  $\mathfrak{d}(i, \mathfrak{C}_q) = \mathfrak{d}(j, \mathfrak{C}_q)$  all vertices  $i, j \in \mathfrak{C}_p$  and for any two different clusters  $\mathfrak{C}_p, \mathfrak{C}_q \in \pi$ .

**Example 2.2:**

For the partition  $\pi = \{\{1, 2\}, \{3, 4\}, \{5\}\}$ , the characteristic matrix is

$$\mathfrak{P} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Notice that  $\pi$  is an equitable partition for the graph in Figure 2.1. Furthermore, notice that  $\{\{1, 2, 5\}, \{3, 4\}\}$  is an AEP, but not an equitable partition. ◇

The following theorem gives necessary and sufficient conditions for equitableness and almost-equitableness (see [ME10, Lemma 2.24], [MTC14, Lemma 5], and [CDR07, Proposition 1]).

**Lemma 2.48:**

Let  $\mathfrak{G}$  be an undirected weighted graph with adjacency matrix  $\mathfrak{A}$  and Laplacian matrix  $\mathfrak{L}$ . Furthermore, let  $\pi$  be a partition of the graph  $\mathfrak{G}$ . Then,

- $\pi$  is equitable if and only if  $\text{im}(\mathfrak{P})$  is  $\mathfrak{A}$ -invariant,
- $\pi$  is almost equitable if and only if  $\text{im}(\mathfrak{P})$  is  $\mathfrak{L}$ -invariant. ◇

Notice that  $\{\mathfrak{V}\}$  and  $\{\{1\}, \{2\}, \dots, \{\mathfrak{n}\}\}$  are always AEPs. We will refer to them as *trivial AEPs*.

**Remark 2.49:**

To the best of our knowledge, there is no known polynomial-time algorithm for finding nontrivial AEPs of a given graph. There is a polynomial-time algorithm for finding the coarsest AEP which is finer than a given partition (see [ZCC14]), but there is no guarantee that it will find a nontrivial AEP.

Furthermore, it is not clear whether a given graph has any nontrivial AEPs at all. On the other hand, a graph can have many AEPs, e.g., every partition of a complete unweighted graph is an AEP. ◇

## 2.6 Linear multi-agent systems

### 2.6.1 System description

Our goal here is to unify definitions from Besselink et al. [BSJ16], Cheng et al. [CKS16, CKS18], Ishizaki et al. [IKIA14, IKG<sup>+</sup>15, IKI16a], and Monshizadeh et al. [MTC13, MTC14]. In particular, we focus on LTI multi-agent systems. Additionally, we restrict to multi-agent system on a undirected, weighted, and connected graph  $\mathfrak{G} = (\mathfrak{V}, \mathfrak{E}, \mathfrak{w})$ .

We base the derivation on [BSJ16] and [CKS18]. The dynamics of the  $i$ th agent, for  $i \in \mathfrak{V} = \{1, 2, \dots, \mathfrak{n}\}$ , is

$$\begin{aligned} E\dot{x}_i(t) &= Ax_i(t) + Bv_i(t), \\ z_i(t) &= Cx_i(t), \end{aligned}$$

with system matrices  $E, A \in \mathbb{R}^{n \times n}$ , input matrix  $B \in \mathbb{R}^{n \times m}$ , output matrix  $C \in \mathbb{R}^{p \times n}$ , state  $x_i(t) \in \mathbb{R}^n$ , input  $v_i(t) \in \mathbb{R}^m$ , and output  $z_i(t) \in \mathbb{R}^p$ . We assume the matrix  $E$  to be invertible. The interconnections are

$$\mathfrak{m}_i v_i(t) = K \sum_{j=1}^{\mathfrak{n}} \mathfrak{a}_{ij} (z_j(t) - z_i(t)) + \sum_{k=1}^{\mathfrak{m}} \mathfrak{b}_{ik} u_k(t),$$



for  $i = 1, 2, \dots, \mathbf{n}$ , with inertia  $\mathbf{m}_i > 0$ , coupling matrix  $K \in \mathbb{R}^{m \times p}$ , external inputs  $u_k(t) \in \mathbb{R}^m$ ,  $k = 1, 2, \dots, \mathbf{m}$ , where  $\mathfrak{A} = [\mathbf{a}_{ij}]$  is the adjacency matrix of the graph  $\mathfrak{G}$ . Here, we allow  $m$  and  $p$  to be different, and for this reason also need a (possibly non-square) coupling matrix  $K$ . The outputs are

$$y_\ell(t) = \sum_{j=1}^{\mathbf{n}} \mathbf{c}_{\ell j} z_j(t)$$

for  $\ell = 1, 2, \dots, \mathbf{p}$ . Define

$$\begin{aligned} \mathfrak{M} &:= \text{diag}(\mathbf{m}_i) \in \mathbb{R}^{\mathbf{n} \times \mathbf{n}}, \quad \mathfrak{B} := [\mathbf{b}_{ik}] \in \mathbb{R}^{\mathbf{n} \times \mathbf{m}}, \quad \mathfrak{C} := [\mathbf{c}_{\ell j}] \in \mathbb{R}^{\mathbf{p} \times \mathbf{n}}, \\ x(t) &:= \text{col}(x_i(t)) \in \mathbb{R}^{\mathbf{n}}, \quad v(t) := \text{col}(v_i(t)) \in \mathbb{R}^{\mathbf{m}}, \quad z(t) := \text{col}(z_i(t)) \in \mathbb{R}^{\mathbf{n}}, \\ u(t) &:= \text{col}(u_k(t)) \in \mathbb{R}^{\mathbf{m}}, \quad \text{and } y(t) := \text{col}(y_\ell(t)) \in \mathbb{R}^{\mathbf{p}}. \end{aligned}$$

Then the agent dynamics can be rewritten as

$$\begin{aligned} (I_{\mathbf{n}} \otimes E)\dot{x}(t) &= (I_{\mathbf{n}} \otimes A)x(t) + (I_{\mathbf{n}} \otimes B)v(t), \\ z(t) &= (I_{\mathbf{n}} \otimes C)x(t), \end{aligned}$$

interconnection as

$$(\mathfrak{M} \otimes I_{\mathbf{n}})v(t) = (-\mathfrak{L} \otimes K)z(t) + (\mathfrak{B} \otimes I_{\mathbf{m}})u(t),$$

and output as

$$y(t) = (\mathfrak{C} \otimes I_{\mathbf{p}})z(t).$$

Therefore, we have

$$\begin{aligned} (\mathfrak{M} \otimes E)\dot{x}(t) &= (\mathfrak{M} \otimes A - \mathfrak{L} \otimes BKC)x(t) + (\mathfrak{B} \otimes B)u(t), \\ y(t) &= (\mathfrak{C} \otimes C)x(t). \end{aligned}$$

This can be generalized further by replacing  $B$  in the input matrix  $\mathfrak{B} \otimes B$  or  $C$  in the output matrix  $\mathfrak{C} \otimes C$ . For example, output  $y(t) = (\mathfrak{C} \otimes I_{\mathbf{n}})x(t)$  is used in [IKI16a], with  $\mathfrak{C} = I_{\mathbf{n}}$ . We can consider

$$\begin{aligned} (\mathfrak{M} \otimes E)\dot{x}(t) &= (\mathfrak{M} \otimes A - \mathfrak{L} \otimes BKC)x(t) + (\mathfrak{B} \otimes F)u(t), \\ y(t) &= (\mathfrak{C} \otimes G)x(t), \end{aligned} \tag{2.9}$$

where  $F$  and  $G$  can be different from  $B$  and  $C$  respectively (including the dimensions).

Of particular interest are *leader-follower multi-agent systems* where only some agents (*leaders*) receive external input, while other agents (*followers*) receive no inputs. Let  $\mathbf{m} \in \{1, 2, \dots, \mathbf{n}\}$  be the number of leaders,  $\mathfrak{V}_L = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_\mathbf{m}\} \subseteq \mathfrak{V}$  the set of leaders, and  $\mathfrak{V}_F = \mathfrak{V} \setminus \mathfrak{V}_L$  the set of followers. Then, with  $\mathfrak{B}$  defined by

$$\mathbf{b}_{ij} := \begin{cases} 1, & \text{if } i = \mathbf{v}_j, \\ 0, & \text{otherwise,} \end{cases}$$

the system (2.9) becomes a *leader-follower multi-agent system*. One important class are multi-agent systems with *single-integrator agents*, i.e., with  $n = 1$ ,  $A = 0$ ,  $B = C = 1$ , and  $K = 1$ . Thus, system (2.9) becomes

$$\begin{aligned}\mathfrak{M}\dot{x}(t) &= -\mathfrak{L}x(t) + \mathfrak{B}u(t), \\ y(t) &= \mathfrak{C}x(t).\end{aligned}\tag{2.10}$$

In particular, [MTC14] uses  $\mathfrak{M} = I_n$  and  $\mathfrak{C} = \mathfrak{W}^{\frac{1}{2}}\mathfrak{R}^T$ , while [IKG<sup>+</sup>15] and [CKS16] use  $\mathfrak{C} = I_n$ .

The property of interest for multi-agent systems is *synchronization*.

**Definition 2.50:**

The system  $(\mathfrak{M} \otimes E)\dot{x}(t) = (\mathfrak{M} \otimes A - \mathfrak{L} \otimes B)x(t)$  is *synchronized* if

$$\lim_{t \rightarrow \infty} (x_i(t) - x_j(t)) = 0,$$

for all  $i, j \in \mathfrak{V}$  and all initial conditions  $x(0) = x_0$ .  $\diamond$

For system (2.10), this is equivalent to  $(\mathfrak{L} \otimes I_n)x(t) \rightarrow 0$ , because the multi-agent system is defined on a connected graph. The following results gives another equivalent condition.

**Proposition 2.51 ([LDCH10, Theorem 1], [MTC13, Lemma 4.2]):**

The system  $(\mathfrak{M} \otimes E)\dot{x}(t) = (\mathfrak{M} \otimes A - \mathfrak{L} \otimes B)x(t)$  is synchronized if and only if  $(A - \lambda B, E)$  is Hurwitz for all nonzero eigenvalues  $\lambda$  of  $(\mathfrak{L}, \mathfrak{M})$ .  $\diamond$

## 2.6.2 Clustering-based model order reduction

By choosing some matrices  $V, W \in \mathbb{R}^{n \times r}$ , we get the ROM for (2.10)

$$\begin{aligned}W^T \mathfrak{M} V \dot{\hat{x}}(t) &= -W^T \mathfrak{L} V \hat{x}(t) + W^T \mathfrak{B} u(t), \\ \hat{y}(t) &= \mathfrak{C} V \hat{x}(t),\end{aligned}\tag{2.11}$$

or, for (2.9),

$$\begin{aligned}(W^T \mathfrak{M} V \otimes E) \dot{\hat{x}}(t) &= (W^T \mathfrak{M} V \otimes A - W^T \mathfrak{L} V \otimes B K C) \hat{x}(t) + (W^T \mathfrak{B} \otimes F) u(t), \\ \hat{y}(t) &= (\mathfrak{C} V \otimes G) \hat{x}(t),\end{aligned}\tag{2.12}$$

which is not necessarily a multi-agent system. Different projection methods are suggested in literature:

1.  $V = \mathfrak{P}, W = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1}$  ([MTC14]),
2.  $V = W = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}}$  ([IKIA14]),
3.  $V = W = \mathfrak{P}$  ([CKS16]),

where  $\mathfrak{P}$  is a characteristic matrix of a partition  $\pi$  of the vertex set  $\mathfrak{V}$ . Since  $\text{im}(V) = \text{im}(W) = \text{im}(\mathfrak{P})$  in all cases, all methods give equivalent ROMs. They each have their advantages. The first preserves structure when  $\mathfrak{M} = I_n$ , in the sense that the ROM represents a multi-agent system defined on a directed, symmetric graph [MTC14]. The second gives the realization for which it is easy to apply the interlacing property (Theorem 2.6), from which it follows that all these methods preserve synchronization. The third preserves the structure in the general case.



# CHAPTER 3

## SUBOPTIMAL CLUSTERING-BASED MODEL ORDER REDUCTION

### Contents

3.1	Introduction . . . . .	39
3.2	$\mathcal{H}_2$ -suboptimal clustering . . . . .	40
3.2.1	Single-integrator agents . . . . .	40
3.2.2	QR decomposition-based clustering . . . . .	41
3.2.3	Clustering by k-means algorithm . . . . .	42
3.2.4	Computing the $\mathcal{H}_2$ -error . . . . .	44
3.2.5	Extension to higher-order agents . . . . .	45
3.2.6	Numerical example . . . . .	45
3.3	Clustering for nonlinear multi-agent systems . . . . .	49
3.3.1	Nonlinear multi-agent systems . . . . .	50
3.3.2	Clustering by projection . . . . .	50
3.4	Conclusions . . . . .	51

## 3.1 Introduction

In this chapter, we study methods for clustering-based MOR for multi-agent systems (see Section 2.6 for an overview). Clustering was proposed in the literature as a way to preserve the multi-agent structure. Here, we will be interested in finding optimal (in a specific error measure) partitions. Since clustering is generally a difficult combinatorial problem (see, e.g., [Sch07]), we will propose a heuristic approach for finding suboptimal partitions.

References [MEB08, RJME09, MEB10] introduce *leader-invariant equitable partitions* and corresponding *quotient graphs*, which can be used for eliminating some uncontrollable states, thus performing no model reduction error. Reference [CM11] proposes *leader partitions*, which introduce “only small model errors”. The authors of [IKIA14] develop a clustering-based  $\mathcal{H}_\infty$  MOR method based on *positive tridiagonalization* and

$\theta$ -reducible clusters, applicable to LTI systems with asymptotically stable and symmetric dynamics matrices.

In [MTC14], the authors focus on *leader-follower linearly diffusively coupled multi-agent systems* with agents having single-integrator dynamics. In particular, these systems have Laplacian-based dynamics, which means they are not asymptotically stable. The authors demonstrate how using partitions for MOR can be transformed to using Petrov-Galerkin projections, while preserving network structure and consensus in the ROM. Further, they derive a simple expression for the relative  $\mathcal{H}_2$ -error when using an AEP and establish a lower bound based on AEPs when using any partition.

The authors of [BSJ16] study *networks of identical passive systems over weighted and directed graphs*, but confine to interconnections with tree structures. Otherwise, these systems would be more general than those in [MTC14]. They present a clustering method relying on the analysis of the corresponding *edge system* to find adjacent subsystems to cluster. Furthermore, they prove that this method preserves the consensus property. Nonetheless, clustering only adjacent subsystems might be restrictive.

The paper [IKG<sup>+</sup>12] presents an efficient clustering-based method for  $\mathcal{H}_2$  MOR of *positive networks*, which include systems with Laplacian-based dynamics. Similarly to [IKIA14], the method is based on  $\theta$ -reducible clusters and an  $\mathcal{H}_2$ -error bound. However, it is not clear how the ROMs resulting from this method compare to the  $\mathcal{H}_2$ -optimal ones.

To the best of our knowledge, there is no efficient method for finding an  $\mathcal{H}_2$ -optimal ROM using graph partitions. In Section 3.2, we propose an efficient  $\mathcal{H}_2$ -suboptimal method, originating from the  $\mathcal{H}_2$ -optimal MOR problem for LTI systems, and illustrate it on an example. Results indicate that the method finds a partition close to the optimal. In Section 3.3, we extend the approach to a class of nonlinear multi-agent systems.

## 3.2 $\mathcal{H}_2$ -suboptimal clustering

The outline of the section is as follows. In Section 3.2.2, we motivate and describe our MOR method, together with the issue of computing the  $\mathcal{H}_2$ -error for systems with Laplacian-based dynamics. We illustrate the method on an example in Section 3.2.6 and conclude with Section 3.4.

### 3.2.1 Single-integrator agents

We first consider multi-agent systems with single-integrator agents as in (2.10). Let

$$\begin{aligned} H(s) &= \mathfrak{C}(s\mathfrak{M} + \mathfrak{L})^{-1}\mathfrak{B}, \\ \hat{H}(s) &= \mathfrak{C}V (sW^T\mathfrak{M}V + W^T\mathfrak{L}V)^{-1}W^T\mathfrak{B} \end{aligned}$$

be the transfer functions of systems (2.10) and (2.11), respectively. We consider the following  $\mathcal{H}_2$ -optimal MOR problem:

$$\underset{V, W \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad \|H - \widehat{H}\|_{\mathcal{H}_2}, \quad (3.1a)$$

$$\text{subject to} \quad V = \mathfrak{P}, \quad (3.1b)$$

$$W = \mathfrak{P}, \quad (3.1c)$$

$$\pi \in \Pi, \quad |\pi| = r, \quad (3.1d)$$

where  $\Pi$  is a set of all partitions of the vertex set  $\mathfrak{V}$ .

To the best of our knowledge, there is no efficient method to exactly solve the optimization problem (3.1). Additionally, since many graph clustering problems are NP-hard (see [Sch07]), we assume the same is true for the above problem.

The problem (3.1) is actually a discrete optimization problem over the set of partitions of the set  $\{1, 2, \dots, n\}$  with  $r$  clusters. Thus, our idea is to solve a relaxed, continuous optimization problem, and use that solution to find a feasible solution for (3.1). We hope that this feasible solution is then close to the optimal solution as it is close to the optimal solution of the relaxed problem.

We relax the problem (3.1) by dropping all constraints (3.1b), (3.1c), and (3.1d). Thus, we obtain the  $\mathcal{H}_2$ -optimal MOR problem for an LTI system. Note that the system (2.10) is not asymptotically stable. See also Section 3.2.4 below for some elaboration on the technicalities associated to this problem.

Applying an  $\mathcal{H}_2$ -optimal method, such as IRKA, to (2.10) will not in general solve the original problem (3.1). Therefore, IRKA will return as a result matrices  $V$  and  $W$  which will not (in general) satisfy the constraints. However, note that in Petrov-Galerkin projection, the subspaces  $\mathcal{V}$  and  $\mathcal{W}$  are enough to determine the transfer function of the ROM. Therefore, the constraints (3.1b) and (3.1c) can be replaced by

$$\text{im}(V) = \text{im}(\mathfrak{P}), \quad (3.2a)$$

$$\text{im}(W) = \text{im}(\mathfrak{P}), \quad (3.2b)$$

without changing  $\widehat{H}$  and the cost (3.1a). The expressions (3.2) motivate us to look for a partition  $\pi$  such that  $\text{im}(\mathfrak{P})$  is “close” to  $\text{im}(V)$  and/or  $\text{im}(W)$ .

### 3.2.2 QR decomposition-based clustering

We know that the condition (3.2a) is equivalent to the existence of a invertible matrix  $Z$  such that  $V = \mathfrak{P}Z$ . In general, condition (3.2a) will not be satisfied, so there will be an error  $\Delta$  such that  $V = \mathfrak{P}Z + \Delta$ . A very similar problem of finding a partition  $\pi$  was encountered in [ZHD<sup>+</sup>01, §3], where a proposed solution is a clustering algorithm based on the *QR decomposition with column pivoting*. Algorithm 3.1 outlines the procedure. The motivation behind it is the result of the following lemma, which says that if  $V = \mathfrak{P}Z$ , the procedure will return the correct result. Therefore, the idea is that if the error  $\Delta$  is small, then Algorithm 3.1 should still return the correct partition.

---

**Algorithm 3.1:** Clustering using QR decomposition with column pivoting [ZHD<sup>+</sup>01, §3]

---

**Input:** Matrix  $V \in \mathbb{R}^{n \times \tau}$  of rank  $\tau$ .

**Output:** Partition  $\pi$  such that  $\text{im}(\mathfrak{P}) \approx \text{im}(V)$ .

- 1 Compute QR decomposition with column pivoting for the matrix  $V^T$ , i.e., find orthogonal  $Q \in \mathbb{R}^{\tau \times \tau}$ , upper-trapezoidal  $R \in \mathbb{R}^{\tau \times n}$ , and a permutation matrix  $P \in \mathbb{R}^{n \times n}$  such that  $V^T P = QR$ .
  - 2 Partition  $R$  as  $\begin{bmatrix} R_{11} & R_{12} \end{bmatrix}$ , with  $R_{11} \in \mathbb{R}^{\tau \times \tau}$  upper-triangular and  $R_{12} \in \mathbb{R}^{\tau \times (n-\tau)}$ .
  - 3 Solve the triangular system  $R_{11}X = R_{12}$ .
  - 4 Compute  $Y = P \begin{bmatrix} I_\tau & X \end{bmatrix}^T = [y_{ij}] \in \mathbb{R}^{n \times \tau}$ .
  - 5 Find a partition  $\pi = \{\mathfrak{C}_1, \mathfrak{C}_2, \dots, \mathfrak{C}_\tau\}$  such that  $i \in \mathfrak{C}_j$  if and only if  $j = \text{argmax}_k |y_{ik}|$ .
- 

**Lemma 3.1:**

Algorithm 3.1 returns  $\pi$  with  $\mathfrak{P}Z$  as input, for an arbitrary partition  $\pi$  and a invertible matrix  $Z$ .  $\diamond$

*Proof.* Let the number of clusters in  $\pi$  be  $\tau$ , i.e.,  $|\pi| = \tau$ . Let us denote the rows of  $Z$  with  $z_1^T, z_2^T, \dots, z_\tau^T$ . Furthermore, without loss of generality we can assume that  $\mathfrak{P} = \text{diag}(\mathbf{1}_{|C_1|}, \mathbf{1}_{|C_2|}, \dots, \mathbf{1}_{|C_\tau|})$ , where  $\mathbf{1}_k \in \mathbb{R}^k$  is a vector of all ones (this structure can be achieved by relabeling the vertices of the graph). Then we have

$$Z^T \mathfrak{P}^T = \begin{bmatrix} z_1 \mathbf{1}_{|C_1|}^T & z_2 \mathbf{1}_{|C_2|}^T & \cdots & z_\tau \mathbf{1}_{|C_\tau|}^T \end{bmatrix}. \quad (3.3)$$

Therefore,  $Z^T \mathfrak{P}^T$  has repeating columns in blocks. If we perform the QR decomposition with column pivoting on (3.3) (ignoring the orthogonal matrix) and then undo the permutation of columns, the result is

$$\begin{bmatrix} \tilde{z}_1 \mathbf{1}_{|C_1|}^T & \tilde{z}_2 \mathbf{1}_{|C_2|}^T & \cdots & \tilde{z}_\tau \mathbf{1}_{|C_\tau|}^T \end{bmatrix}, \quad (3.4)$$

where a permutation of the columns of  $\begin{bmatrix} \tilde{z}_1 & \tilde{z}_2 & \cdots & \tilde{z}_\tau \end{bmatrix}$  gives the upper-triangular matrix  $R_{11}$  from Algorithm 3.1. We see that multiplying (3.4) on the left by  $R_{11}^{-1}$  ( $\mathfrak{P}Z$  having full rank implies that  $R_{11}$  is invertible) and transposing produces  $\mathfrak{P}$ , possibly with permuted columns. Thus we conclude that Algorithm 3.1 returns the partition  $\pi$ .  $\square$

Furthermore, note that the number of floating point operations in Algorithm 3.1 is  $\mathcal{O}(n\tau^2)$  and the size of additional storage is  $\mathcal{O}(n\tau)$ . Therefore, it is efficient in the large-scale setting if  $\tau$  is small compared to  $n$ .

### 3.2.3 Clustering by k-means algorithm

The motivation for using the k-means algorithm ([HW79]) is the following result about how a change in the Petrov-Galerkin subspaces affects the ROM.



**Theorem 3.2 ([BGW12, Theorem 3.3]):**

Let  $V_1, V_2, W_1, W_2 \in \mathbb{R}^{n \times r}$ ,

$$\mathcal{V}_i = \text{im}(V_i), \quad \mathcal{W}_i = \text{im}(W_i), \quad \hat{H}_i(s) = CV_i (sW_i^T EV_i - W_i^T AV_i)^{-1} W_i^T B,$$

for  $i = 1, 2$ . Then

$$\frac{\|\hat{H}_1 - \hat{H}_2\|_{\mathcal{H}_\infty}}{\frac{1}{2}(\|\hat{H}_1\|_{\mathcal{H}_\infty} + \|\hat{H}_2\|_{\mathcal{H}_\infty})} \leq M \max(\sin \Theta(\mathcal{V}_1, \mathcal{V}_2), \sin \Theta(\mathcal{W}_1, \mathcal{W}_2)),$$

where

$$\begin{aligned} M &= 2 \max(M_1, M_2), \\ M_1 &= \frac{\max_{\omega \in \mathbb{R}} \|C\|_2 \left\| V_1 (\omega W_1^T EV_1 - W_1^T AV_1)^{-1} W_1^T B \right\|_2 \|\hat{H}_1(\omega)\|_2^{-1}}{\min_{\omega \in \mathbb{R}} \cos \Theta(\ker(W_2^T (\omega E - A)^{-1}), \mathcal{V}_2)}, \\ M_2 &= \frac{\max_{\omega \in \mathbb{R}} \|CV_2 (\omega W_2^T EV_2 - W_2^T AV_2)^{-1} W_2^T\|_2 \|B\|_2 \|\hat{H}_2(\omega)\|_2^{-1}}{\min_{\omega \in \mathbb{R}} \cos \Theta(\text{im}((\omega E - A)^{-1} V_1), \mathcal{W}_1)}, \end{aligned}$$

and  $\Theta(\mathcal{M}, \mathcal{N})$  is the largest principal angle between subspaces  $\mathcal{M}, \mathcal{N} \subseteq \mathbb{R}^n$ .  $\diamond$

The angle between two subspaces  $\mathcal{V}_1, \mathcal{V}_2 \subseteq \mathbb{R}^n$  is defined by (see [BGW12, Section 3.1])

$$\sin \Theta(\mathcal{V}_1, \mathcal{V}_2) := \sup_{v_1 \in \mathcal{V}_1} \inf_{v_2 \in \mathcal{V}_2} \frac{\|v_2 - v_1\|_2}{\|v_1\|_2}.$$

If  $\dim \mathcal{V}_1 = \dim \mathcal{V}_2$ , then we have

$$\sin \Theta(\mathcal{V}_1, \mathcal{V}_2) = \sin \Theta(\mathcal{V}_2, \mathcal{V}_1) = \|(I - V_1 V_1^T) V_2\|_2,$$

where  $\mathcal{V}_1 = \text{im}(V_1)$ ,  $\mathcal{V}_2 = \text{im}(V_2)$ , and both  $V_1$  and  $V_2$  have orthonormal columns. If additionally  $V_1 = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}}$  and  $V_2 = V$  from IRKA, then

$$\begin{aligned} (\sin \Theta(\mathcal{V}_1, \mathcal{V}_2))^2 &\leq \left\| \left( I - \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \right) V \right\|_F^2 \\ &= \left\| \left( I - [\mathfrak{p}(\mathfrak{C}_1) \ \cdots \ \mathfrak{p}(\mathfrak{C}_r)] \begin{bmatrix} |\mathfrak{C}_1|^{-1} & & \\ & \ddots & \\ & & |\mathfrak{C}_r|^{-1} \end{bmatrix} \begin{bmatrix} \mathfrak{p}(\mathfrak{C}_1)^T \\ \vdots \\ \mathfrak{p}(\mathfrak{C}_r)^T \end{bmatrix} \right) V \right\|_F^2 \\ &= \left\| \left( I - \sum_{i=1}^r \frac{1}{|\mathfrak{C}_i|} \mathfrak{p}(\mathfrak{C}_i) \mathfrak{p}(\mathfrak{C}_i)^T \right) V \right\|_F^2 \\ &= \sum_{i=1}^r \left\| V_{\mathfrak{C}_i, :} - \mathfrak{p}(\mathfrak{C}_i) \frac{1}{|\mathfrak{C}_i|} \mathbb{1}_{|\mathfrak{C}_i|}^T V_{\mathfrak{C}_i, :} \right\|_F^2 \\ &= \sum_{i=1}^r \sum_{p \in \mathfrak{C}_i} \left\| V_{p, :} - \frac{1}{|\mathfrak{C}_i|} \sum_{q \in \mathfrak{C}_i} V_{q, :} \right\|_F^2, \end{aligned}$$

which is equal to the k-means cost functional. Therefore, applying the k-means algorithm to the rows of  $V$  will minimize an upper bound on the largest principal angle between  $\text{im}(V)$  and  $\text{im}(\mathfrak{P})$ .

The advantage of using k-means compared to QR decomposition-based in that the latter can only, given  $V \in \mathbb{R}^{n \times \mathfrak{r}}$ , return a partition with  $\mathfrak{r}$  clusters. On the other hand, k-means clustering can return a partition with any number of clusters, which makes it more efficient if the number of clusters is relatively large.

### 3.2.4 Computing the $\mathcal{H}_2$ -error

When using an AEP, the relative  $\mathcal{H}_2$ -error can be computed directly [MTC14, Theorem 6]. In other cases, we need to solve Lyapunov equations to compute the  $\mathcal{H}_2$ -norms.

As discussed in Section 2.4.5, we can use the eigenspaces to find the asymptotically stable part of the system and use it to compute the  $\mathcal{H}_2$ -norm. We are looking for an invertible matrix  $T$  such that

$$T^T \mathfrak{M} T = \begin{bmatrix} \mathfrak{M}_- & 0 \\ 0 & \mathfrak{M}_+ \end{bmatrix} \text{ and } T^T \mathfrak{L} T = \begin{bmatrix} \mathfrak{L}_- & 0 \\ 0 & \mathfrak{L}_+ \end{bmatrix},$$

where  $\sigma(-\mathfrak{L}_-, \mathfrak{M}_-) \subset \mathbb{C}_-$  and  $\sigma(-\mathfrak{L}_+, \mathfrak{M}_+) \subset \overline{\mathbb{C}_+}$ . We see that if

$$T = \begin{bmatrix} T_- & \mathbf{1}_n \end{bmatrix},$$

then

$$T^T \mathfrak{M} T = \begin{bmatrix} T_-^T \mathfrak{M} T_- & T_-^T \mathfrak{M} \mathbf{1}_n \\ \mathbf{1}_n^T \mathfrak{M} T_- & \mathbf{1}_n^T \mathfrak{M} \mathbf{1}_n \end{bmatrix} \text{ and } T^T \mathfrak{L} T = \begin{bmatrix} T_-^T \mathfrak{L} T_- & 0 \\ 0 & 0 \end{bmatrix}.$$

To have  $T_-^T \mathfrak{M} \mathbf{1}_n = 0$ , we need for the columns of  $T_-$  to be orthogonal to  $\mathfrak{M} \mathbf{1}_n$ . Additionally,  $T_-$  should be such that both  $T_-^T \mathfrak{M} T_-$  and  $T_-^T \mathfrak{L} T_-$  are sparse. We choose

$$T_- = \begin{bmatrix} \alpha_1 & & & & \\ -\beta_1 & \alpha_2 & & & \\ & -\beta_2 & \ddots & & \\ & & \ddots & \alpha_{n-1} & \\ & & & -\beta_{n-1} & \end{bmatrix}$$

for some  $\alpha_i, \beta_i > 0$ ,  $i = 1, 2, \dots, n-1$ , which we determine next. From  $e_i^T T_-^T \mathfrak{M} \mathbf{1}_n = 0$ , we find  $\alpha_i \mathbf{m}_i = \beta_i \mathbf{m}_{i+1}$ . If we additionally set  $\alpha_i^2 + \beta_i^2 = 1$ , we get

$$\alpha_i = \frac{\mathbf{m}_{i+1}}{\sqrt{\mathbf{m}_i^2 + \mathbf{m}_{i+1}^2}} \text{ and } \beta_i = \frac{\mathbf{m}_i}{\sqrt{\mathbf{m}_i^2 + \mathbf{m}_{i+1}^2}}.$$

---

**Algorithm 3.2:** QR decomposition with column pivoting for matrices with block-columns

---

**Input:** Matrix  $X \in \mathbb{R}^{n \times n}$  of full rank and  $\mathfrak{r} < n$ .

**Output:** Orthogonal matrix  $Q$ , upper-triangular matrix  $R$ , and permutation matrix  $P$  such that  $XP = QR$

- 1 Denote  $X = [X_1 \ X_2 \ \cdots \ X_n]$ , where  $X_i \in \mathbb{R}^{n \times n}$ .
  - 2 Find a block-column  $X_i$  with the largest Frobenius norm and swap it with  $X_1$ .
  - 3 Perform QR decomposition with column pivoting on  $X_1$ , i.e. find an orthogonal  $Q_1 \in \mathbb{R}^{n \times n}$ , an upper-triangular  $R_1 \in \mathbb{R}^{n \times n}$ , and a permutation matrix  $P_1 \in \mathbb{R}^{n \times n}$  such that  $X_1 P_1 = Q_1 R_1$ .
  - 4 Multiply all block-columns in  $X$  from the right by  $P_1$ .
  - 5 Multiply  $X$  from the left by  $Q_1^T$ .
  - 6 Repeat the procedure for  $X_{n+1:n\mathfrak{r}, n+1:n}$ , which computes the matrices  $Q_i, R_i$ , and  $P_i$ , for  $i \in \{2, 3, \dots, \mathfrak{r}\}$ .
  - 7 Return  $Q = Q_1 Q_2 \cdots Q_{\mathfrak{r}}$ ,  $R = X$ , and  $P$  with all of the column permutations recorded.
- 

### 3.2.5 Extension to higher-order agents

For multi-agent systems (2.9) with agents of order  $n$ , the constraints (3.2) become

$$\text{im}(V) = \text{im}(\mathfrak{B} \otimes I_n), \quad (3.5a)$$

$$\text{im}(W) = \text{im}(\mathfrak{B} \otimes I_n). \quad (3.5b)$$

QR decomposition-based clustering can then be extended as in Algorithm 3.2 by clustering the block-columns of  $V^T$  (or  $W^T$ ). For the k-means algorithm, we can show in a similar way as in the single-integrator case that clustering the block-rows leads to minimizing an upper bound of the largest principal angle. Therefore, k-means can be directly applied to the vectorized block-rows of  $V$  or  $W$ .

Additionally, we need to consider synchronization preservation. In the single-integrator case, clustering using any partition preserves synchronization. In the general case, using Proposition 2.51, we need that  $(A - \hat{\lambda}, E)$  is Hurwitz for all nonzero eigenvalues  $\hat{\lambda}$  of  $(\hat{\mathcal{L}}, \hat{\mathcal{M}})$ . This is true for AEP because  $\sigma(\hat{\mathcal{L}}, \hat{\mathcal{M}})$  is a subset of  $\sigma(\mathcal{L}, \mathcal{M})$ . Note that if  $(A - \lambda, E)$  is Hurwitz for all  $\lambda \in [\lambda_2, \dots, \lambda_n]$ , where  $0 = \lambda_1 < \lambda_2 \leq \dots \leq \lambda_n$  are the eigenvalues of  $(\mathcal{L}, \mathcal{M})$ , then from Theorem 2.6 we get that every partition preserves synchronization.

### 3.2.6 Numerical example

We use the example from [MTC14]. It is a leader-follower multi-agent system, defined on a undirected, weighted, connected graph with 10 vertices, shown in Figure 3.1. Agents 6 and 7 are the leaders of the multi-agent system. The Laplacian and input

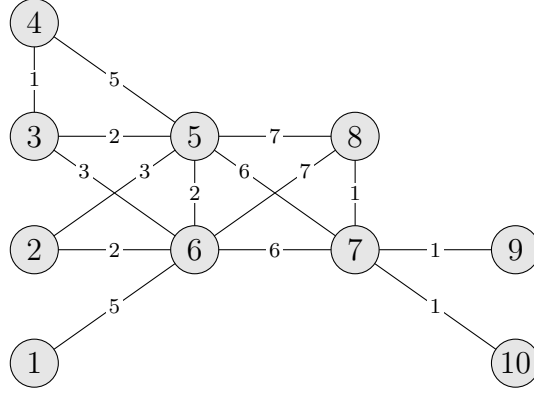


Figure 3.1: The undirected, weighted graph from [MTC14]

matrices are

$$\mathfrak{L} = \begin{bmatrix} 5 & 0 & 0 & 0 & 0 & -5 & 0 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 & -3 & -2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 6 & -1 & -2 & -3 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 6 & -5 & 0 & 0 & 0 & 0 & 0 \\ 0 & -3 & -2 & -5 & 25 & -2 & -6 & -7 & 0 & 0 \\ -5 & -2 & -3 & 0 & -2 & 25 & -6 & -7 & 0 & 0 \\ 0 & 0 & 0 & 0 & -6 & -6 & 15 & -1 & -1 & -1 \\ 0 & 0 & 0 & 0 & -7 & -7 & -1 & 15 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -1 & 0 & 0 & 1 \end{bmatrix}, \quad \mathfrak{B} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

We chose an edge ordering and orientation such that the incidence and edge-weights matrices are

$$\mathfrak{R} = \begin{bmatrix} -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & -1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & -1 & -1 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

and  $\mathfrak{W} = \text{diag}(5, 3, 2, 1, 2, 3, 5, 2, 6, 7, 6, 7, 1, 1, 1)$ .

For this example, we focus on partitions with five clusters. There are in total 42,525 such partitions. Table 3.1 shows the best 20 partitions of size 5 and their relative  $\mathcal{H}_2$ -errors. Figure 3.2 shows all relative  $\mathcal{H}_2$ -errors for all partitions with five clusters.

Table 3.1: Top 20 partitions with 5 clusters for reducing the multi-agent system in Figure 3.1

Rank	Relative $\mathcal{H}_2$ -error	Partition
1	0.128053	$\{\{1, 8\}, \{2, 3, 4, 9, 10\}, \{5\}, \{6\}, \{7\}\}$
2	0.131311	$\{\{1, 2, 3, 4\}, \{5, 8\}, \{6\}, \{7\}, \{9, 10\}\}$
3	0.137466	$\{\{1, 2, 3, 4, 9, 10\}, \{5\}, \{6\}, \{7\}, \{8\}\}$
4	0.137473	$\{\{1, 3, 8\}, \{2, 4, 9, 10\}, \{5\}, \{6\}, \{7\}\}$
5	0.143700	$\{\{1, 5, 8\}, \{2, 3, 4\}, \{6\}, \{7\}, \{9, 10\}\}$
6	0.145900	$\{\{1, 2, 3\}, \{4, 9, 10\}, \{5, 8\}, \{6\}, \{7\}\}$
7	0.146196	$\{\{1, 8\}, \{2, 3, 4, 9\}, \{5, 10\}, \{6\}, \{7\}\}$
8	0.146196	$\{\{1, 8\}, \{2, 3, 4, 10\}, \{5, 9\}, \{6\}, \{7\}\}$
9	0.147022	$\{\{1, 2, 3, 8\}, \{4, 9, 10\}, \{5\}, \{6\}, \{7\}\}$
10	0.149240	$\{\{1, 8, 9\}, \{2, 3, 4, 10\}, \{5\}, \{6\}, \{7\}\}$
11	0.149240	$\{\{1, 8, 10\}, \{2, 3, 4, 9\}, \{5\}, \{6\}, \{7\}\}$
12	0.149654	$\{\{1, 8\}, \{2, 4, 9, 10\}, \{3, 5\}, \{6\}, \{7\}\}$
13	0.150440	$\{\{1, 5\}, \{2, 3, 4, 9, 10\}, \{6\}, \{7\}, \{8\}\}$
14	0.150654	$\{\{1, 3\}, \{2, 4, 9, 10\}, \{5, 8\}, \{6\}, \{7\}\}$
15	0.151684	$\{\{1, 2, 8\}, \{3, 4, 9, 10\}, \{5\}, \{6\}, \{7\}\}$
16	0.153100	$\{\{1, 2, 3, 4, 9\}, \{5, 8\}, \{6\}, \{7\}, \{10\}\}$
17	0.153100	$\{\{1, 2, 3, 4, 10\}, \{5, 8\}, \{6\}, \{7\}, \{9\}\}$
18	0.153819	$\{\{1\}, \{2, 3, 4, 9, 10\}, \{5, 8\}, \{6\}, \{7\}\}$
19	0.154374	$\{\{1, 3, 8, 9\}, \{2, 4, 10\}, \{5\}, \{6\}, \{7\}\}$
20	0.154374	$\{\{1, 3, 8, 10\}, \{2, 4, 9\}, \{5\}, \{6\}, \{7\}\}$

The partition

$$\{\{1, 2, 3, 4\}, \{5, 6\}, \{7\}, \{8\}, \{9, 10\}\} \quad (3.6)$$

is an AEP of the graph in Figure 3.1 (see Figure 3.3) used in [MTC14]. It has five clusters, so we compare the relative  $\mathcal{H}_2$ -errors of different ROMs of order  $\mathfrak{r} = 5$ . The partition (3.6) is actually the only AEP with five clusters, among the total of five AEPs of the graph in Figure 3.1.

From [MTC14, Theorem 6], it follows that for the AEP in (3.6), the relative  $\mathcal{H}_2$ -error is

$$\frac{\|H - \widehat{H}_{\text{AEP}}\|_{\mathcal{H}_2}}{\|H\|_{\mathcal{H}_2}} = \sqrt{\frac{(1 - \frac{1}{2}) + (1 - \frac{1}{1})}{2(1 - \frac{1}{10})}} \approx 0.527046, \quad (3.7)$$

where  $\widehat{H}_{\text{AEP}}$  is the transfer function of the ROM using the graph partition (3.6).

Next, we used IRKA to find a ROM of order  $\mathfrak{r} = 5$ . It found a ROM with relative

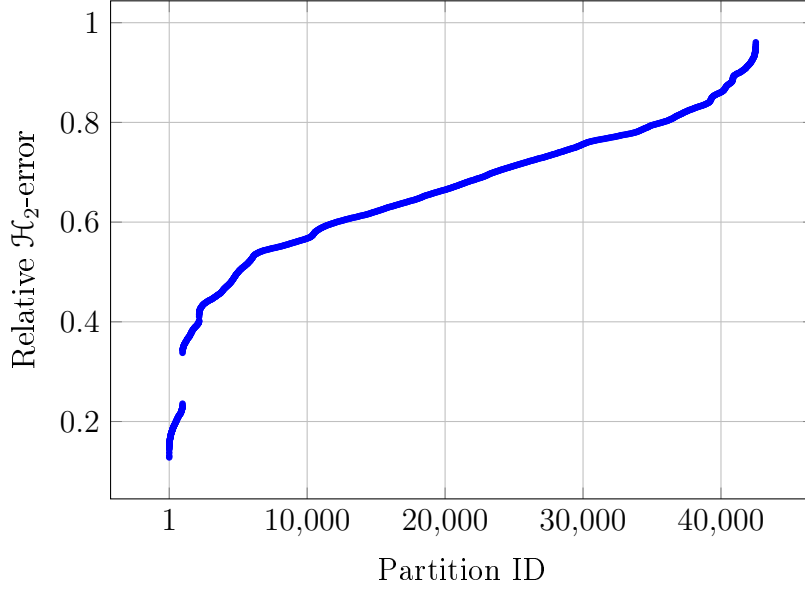


Figure 3.2: Relative  $\mathcal{H}_2$ -errors for all partitions with five clusters

$\mathcal{H}_2$ -error of

$$\frac{\|H - \hat{H}_{\text{IRKA}}\|_{\mathcal{H}_2}}{\|H\|_{\mathcal{H}_2}} \approx 0.0330412, \quad (3.8)$$

which is almost 16 times better than (3.7).

The partition resulting from Algorithm 3.1 applied to IRKA's  $V$  matrix is

$$\{\{1, 3\}, \{2, 4, 9, 10\}, \{5, 8\}, \{6\}, \{7\}\}, \quad (3.9)$$

and is shown in Figure 3.4. The relative  $\mathcal{H}_2$ -error of a ROM using partition (3.9) is

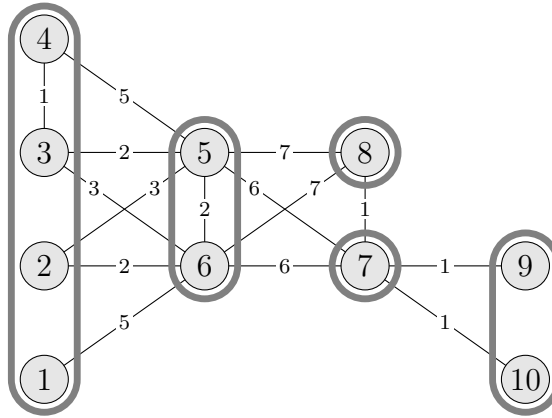


Figure 3.3: Graph with partition

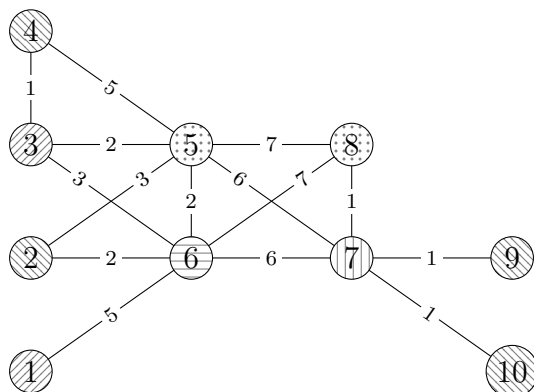


Figure 3.4: Graph partition (3.9) with five clusters. The assignment of vertices to clusters is represented with different patterns.

$$\frac{\|H - \hat{H}_V\|_{\mathcal{H}_2}}{\|H\|_{\mathcal{H}_2}} \approx 0.150654, \quad (3.10)$$

which is more than 4 times worse than (3.8), but 3 times better than (3.7).

We notice by (3.2b) that  $W$  can also be used to find a good partition. In this example, Algorithm 3.1 returns the partition

$$\{\{1, 2, 3, 9, 10\}, \{4, 8\}, \{5\}, \{6\}, \{7\}\}. \quad (3.11)$$

The relative  $\mathcal{H}_2$ -error when using the partition (3.11) is

$$\frac{\|H - \hat{H}_W\|_{\mathcal{H}_2}}{\|H\|_{\mathcal{H}_2}} \approx 0.179746,$$

which is worse than (3.10).

We notice that partition (3.9) is the 14th partition and that the best partition produces about 1.18 times better error. Table 3.1 does not show that partition (3.11) is 192nd and partition (3.6) is 5996th. Thus, the new method gets a lot closer to the optimal partition than the AEP in this example. Using the k-means algorithm to cluster the rows of  $V$  and  $W$  results in the 6th and 27th partition, respectively, which is in this case an improvement over QR decomposition-based clustering.

Furthermore, from Table 3.1 we see that, in all of the top 20 partitions, leaders 6 and 7 appear in singletons. This makes sense, because this way no input is diffused over more agents. However, further research is needed to see if, in general, the best partition has leaders appearing in singletons.

### 3.3 Clustering for nonlinear multi-agent systems

In this section, we extend the approach from the previous section to certain nonlinear multi-agent systems. We describe the class of multi-agent systems in Section 3.3.1. Next, in Section 3.3.2, we extend clustering by projection to this class of systems.

### 3.3.1 Nonlinear multi-agent systems

Here, we consider a class of nonlinear multi-agent systems. In particular, let the dynamics of the  $i$ th agent, for  $i = 1, 2, \dots, \mathbf{n}$ , be

$$\begin{aligned}\dot{x}_i(t) &= f(x_i(t)) + g(x_i(t))v_i(t), \\ z_i(t) &= h(x_i(t)),\end{aligned}$$

with functions  $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $g: \mathbb{R}^n \rightarrow \mathbb{R}^{n \times m}$ ,  $h: \mathbb{R}^n \rightarrow \mathbb{R}^{p \times n}$ , state  $x_i(t) \in \mathbb{R}^n$ , input  $v_i(t) \in \mathbb{R}^m$ , and output  $z_i(t) \in \mathbb{R}^p$ . The interconnections are

$$v_i(t) = \sum_{j=1}^{\mathbf{n}} \mathbf{a}_{ij} k(z_j(t) - z_i(t)) + \sum_{\ell=1}^{\mathbf{m}} \mathbf{b}_{i\ell} u_\ell(t),$$

for  $i = 1, 2, \dots, \mathbf{n}$ , with coupling  $k: \mathbb{R}^p \rightarrow \mathbb{R}^m$ , external input  $u_\ell(t) \in \mathbb{R}^m$ ,  $\ell = 1, 2, \dots, \mathbf{m}$ , where  $\mathfrak{A} = [\mathbf{a}_{ij}]$  is the adjacency matrix of the graph  $\mathfrak{G}$ , and  $\mathfrak{B} = [\mathbf{b}_{i\ell}]$ . Additionally, let the external output be

$$y_i(t) = \sum_{j=1}^{\mathbf{n}} \mathbf{c}_{ij} z_j(t),$$

with  $\mathfrak{C} = [\mathbf{c}_{ij}]$ . We assume functions  $f, g, h, k$  are continuous and that there is a unique global solution  $x(t) = \text{col}(x_1(t), x_2(t), \dots, x_n(t))$  for any admissible  $u(t)$ .

### 3.3.2 Clustering by projection

We want to find the form of the reduced order model obtained from Galerkin projection with  $V = \mathfrak{P} \otimes I_n$ . We can write

$$\begin{aligned}\dot{x}(t) &= F(x(t), u(t)), \\ y(t) &= G(x(t)),\end{aligned}$$

for some functions  $F$  and  $G$ . The reduced model is

$$\begin{aligned}(\mathfrak{P}^T \mathfrak{P} \otimes I_n) \hat{\dot{x}}(t) &= (\mathfrak{P}^T \otimes I_n) F((\mathfrak{P} \otimes I_n) \hat{x}(t), u(t)), \\ \hat{y}(t) &= G((\mathfrak{P} \otimes I_n) \hat{x}(t)).\end{aligned}$$



Premultiplying the first equation with  $e_p^T \otimes I_n$ , we find

$$\begin{aligned}
 |\mathcal{C}_p| \dot{\hat{x}}_p(t) &= \sum_{i \in \mathcal{C}_p} \left( f(\hat{x}_p(t)) + g(\hat{x}_p(t)) \left( \sum_{j=1}^n \mathbf{a}_{ij} k(h(\hat{x}_{C(j)}(t)) - h(\hat{x}_p(t))) + \sum_{\ell=1}^m \mathbf{b}_{i\ell} u_\ell(t) \right) \right) \\
 &= |\mathcal{C}_p| f(\hat{x}_p(t)) \\
 &\quad + g(\hat{x}_p(t)) \left( \sum_{i \in \mathcal{C}_p} \sum_{j=1}^n \mathbf{a}_{ij} k(h(\hat{x}_{C(j)}(t)) - h(\hat{x}_p(t))) + \sum_{i \in \mathcal{C}_p} \sum_{\ell=1}^m \mathbf{b}_{i\ell} u_\ell(t) \right) \\
 &= |\mathcal{C}_p| f(\hat{x}_p(t)) \\
 &\quad + g(\hat{x}_p(t)) \left( \sum_{q=1}^r \sum_{i \in \mathcal{C}_p} \sum_{j \in \mathcal{C}_q} \mathbf{a}_{ij} k(h(\hat{x}_q(t)) - h(\hat{x}_p(t))) + \sum_{\ell=1}^m \sum_{i \in \mathcal{C}_p} \mathbf{b}_{i\ell} u_\ell(t) \right) \\
 &= |\mathcal{C}_p| f(\hat{x}_p(t)) \\
 &\quad + g(\hat{x}_p(t)) \left( \sum_{q=1}^r \hat{\mathbf{a}}_{pq} k(h(\hat{x}_q(t)) - h(\hat{x}_p(t))) + \sum_{\ell=1}^m \hat{\mathbf{b}}_{p\ell} u_\ell(t) \right),
 \end{aligned}$$

for

$$\hat{\mathbf{a}}_{pq} = \sum_{i \in \mathcal{C}_p} \sum_{j \in \mathcal{C}_q} \mathbf{a}_{ij}, \quad \hat{\mathbf{b}}_{p\ell} = \sum_{i \in \mathcal{C}_p} \mathbf{b}_{i\ell}.$$

Defining  $\hat{\mathfrak{A}} := [\hat{\mathbf{a}}_{pq}]$  and  $\hat{\mathfrak{B}} := [\hat{\mathbf{b}}_{p\ell}]$ , we see that  $\hat{\mathfrak{A}} = \mathfrak{P}^T \mathfrak{A} \mathfrak{P}$  and  $\hat{\mathfrak{B}} = \mathfrak{P}^T \mathfrak{B}$ . For the output we have

$$\hat{y}_i(t) = \sum_{j=1}^n \mathbf{c}_{ij} h(\hat{x}_{C(j)}(t)) = \sum_{j=1}^r \sum_{j \in \mathcal{C}_q} \mathbf{c}_{ij} h(\hat{x}_q(t)) = \sum_{j=1}^r \hat{\mathbf{c}}_{iq} h(\hat{x}_q(t)),$$

where

$$\hat{\mathbf{c}}_{iq} = \sum_{j \in \mathcal{C}_q} \mathbf{c}_{ij}.$$

Thus, for  $\hat{\mathcal{C}} := [\hat{\mathbf{c}}_{iq}]$ , we have  $\hat{\mathcal{C}} = \mathcal{C} \mathfrak{P}$ . Therefore, we showed how to construct a ROM of the same structure as the original multi-agent system. Based on this, to find a good partition, we can apply any projection-based MOR method for nonlinear systems (e.g., proper orthogonal decomposition [HV05]) and cluster the (block-)rows of the matrix used to project the system.

## 3.4 Conclusions

In Section 3.2, we presented our method, combining IRKA and a clustering algorithm, for MOR of multi-agent systems using graph partitions. It seems heuristically that

this method is able to create a partition whose  $\mathcal{H}_2$ -error is small. We furthermore elaborated that this method is scalable to large-scale systems. Our numerical test for a small network shows that, among 42 525 partitions, our algorithm found the 14th best approximation whose error is of the same order of magnitude as the error of the best partition. A theoretical foundation that the algorithm always finds a partition in some sense close to optimal remains an open problem for future work.

In Section 3.3, we showed how to extend the approach to nonlinear multi-agent systems by combining a projection-based method and a clustering algorithm.

# CHAPTER 4

## GRAPH SYMMETRIES AND EQUITABLE PARTITIONS IN CLUSTERING-BASED MODEL ORDER REDUCTION

### Contents

4.1	Introduction . . . . .	54
4.2	Error bounds for clustering-based model order reduction of linear multi-agent systems . . . . .	54
4.2.1	Introduction . . . . .	54
4.2.2	Preliminaries . . . . .	55
4.2.3	Problem formulation . . . . .	58
4.2.4	Graph partitions and reduction by clustering . . . . .	58
4.2.5	$\mathcal{H}_2$ -error bounds . . . . .	60
4.2.6	$\mathcal{H}_\infty$ -error bounds . . . . .	68
4.2.6.1	The single integrator case . . . . .	68
4.2.6.2	The general case with symmetric agent dynamics . . . . .	71
4.2.7	Towards a priori error bounds for general graph partitions . . . . .	74
4.2.7.1	The single integrator case . . . . .	74
4.2.7.2	The general case . . . . .	78
4.2.8	Numerical examples . . . . .	80
4.3	Exact clustering-based model order reduction for nonlinear power systems . . . . .	84
4.3.1	Introduction . . . . .	84
4.3.2	System description . . . . .	84
4.3.3	Synchronization of generator pair . . . . .	88
4.3.4	Synchronization of generator partition . . . . .	92
4.3.5	Clustering of power systems . . . . .	95
4.3.6	Illustrative example . . . . .	96
4.4	Conclusion . . . . .	97

## 4.1 Introduction

Here, we consider theoretical aspects of clustering-based MOR, using graph symmetries and equitable partitions.

The work in [MTC14] showed upper bounds for  $\mathcal{H}_2$  model reduction error using AEP for leader-follower multi-agent systems with single-integrator agents. In Section 4.2, we extend these results to more general agents, including also upper bounds for  $\mathcal{H}_\infty$  error and arbitrary partitions.

The results of [RJME09] on clustering using graph symmetries and equitable partitions for linear multi-agent systems are extended in Section 4.3 nonlinear power systems.

## 4.2 Error bounds for clustering-based model order reduction of linear multi-agent systems

### 4.2.1 Introduction

In [MTC14], MOR by clustering was put in the context of MOR by Petrov-Galerkin projection. The results in [MTC14] provide explicit expressions for the  $\mathcal{H}_2$  model reduction error if a leader-follower network with *single integrator agent dynamics* is clustered using an almost equitable partition of the graph. Here, our aim is to generalize and extend the results in [MTC14] to networks where the agent dynamics is given by an *arbitrary multivariable input-state-output system*. We also aim at finding explicit formulas and a priori upper bounds for the model reduction error measured in the  $\mathcal{H}_\infty$ -norm. Finally, we will consider the problem of clustering a network according to arbitrary, not necessarily almost equitable, graph partitions. The main contributions of this section are the following:

1. We derive an a priori upper bound for the  $\mathcal{H}_2$  model reduction error for the case that the agents are represented by an arbitrary input-state-output system.
2. We extend the results in [MTC14] for single integrator dynamics by giving an explicit expression for the  $\mathcal{H}_\infty$  model reduction error in terms of properties of the given graph partition.
3. We establish an a priori upper bound for the  $\mathcal{H}_\infty$  model reduction error for the case that the agents are represented by an arbitrary but *symmetric* input-state-output system.
4. We establish some preliminary results on the model reduction error in case of clustering using an arbitrary, possibly non almost equitable, partition.

The outline of this section is as follows. In Section 4.2.3, we formulate our problem of MOR of leader-follower multi-agent networks. Section 4.2.4 reviews some theory on graph partitions and MOR by clustering, and relates this method to Petrov-Galerkin

projection of the original network. Also preservation of synchronization is discussed here. In Section 4.2.5, we provide a priori error bounds on the  $\mathcal{H}_2$  model reduction error for networks with arbitrary agent dynamics, clustered using almost equitable partitions. In Section 4.2.6, we complement these results by providing upper bounds on the  $\mathcal{H}_\infty$  model reduction error. In Section 4.2.7, the problem of clustering networks according to general partitions is considered and the first steps towards a priori error bounds on both the  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  model reduction errors are made. Numerical examples for which we compare the actual errors with the a priori bounds established here are presented in Section 4.2.8.

## 4.2.2 Preliminaries

In this section, we briefly discuss some basic facts on finite-dimensional linear systems and how they extend to unstable systems.

Consider the input-state-output system

$$\begin{aligned}\dot{x}(t) &= Ax(t) + Bu(t), \\ y(t) &= Cx(t),\end{aligned}\tag{4.1}$$

with  $x(t) \in \mathbb{R}^n$ ,  $u(t) \in \mathbb{R}^m$ ,  $y(t) \in \mathbb{R}^p$ , and transfer function  $H(s) = C(sI - A)^{-1}B$ . If  $A$  is Hurwitz, then the  $\mathcal{H}_2$ -norm can be computed as

$$\|H\|_{\mathcal{H}_2}^2 = \text{tr}(B^T X B),$$

where  $X$  is the unique positive semi-definite solution of the Lyapunov equation

$$A^T X + X A + C^T C = 0.\tag{4.2}$$

For the purposes of the work in this part, we also need to deal with the situation when  $A$  is not Hurwitz. Let  $\mathcal{X}_+(A)$  denote the unstable subspace of  $A$ , i.e., the direct sum of the generalized eigenspaces of  $A$  corresponding to its eigenvalues in the closed right half plane. We state the following proposition.

**Proposition 4.1:**

Assume that  $\mathcal{X}_+(A) \subseteq \ker(C)$ . Then the Lyapunov equation (4.2) has at least one positive semi-definite solution. Among all positive semi-definite solutions, there is exactly one solution, say  $X$ , with the property  $\mathcal{X}_+(A) \subseteq \ker(X)$ . For this particular solution  $X$ , we have  $\|H\|_{\mathcal{H}_2}^2 = \text{tr}(B^T X B)$ .  $\diamond$

*Proof.* Without loss of generality, assume that

$$A = \begin{bmatrix} A_- & 0 \\ 0 & A_+ \end{bmatrix}, \quad B = \begin{bmatrix} B_- \\ B_+ \end{bmatrix}, \quad C = [C_- \quad 0],$$

where  $A_-$  is Hurwitz, and  $A_+$  has all its eigenvalues in the closed right half plane. Let  $X_-$  be the unique solution to the reduced Lyapunov equation

$$A_-^T X_- + X_- A_- + C_-^T C_- = 0.\tag{4.3}$$

Then  $X_- = \int_0^\infty e^{A^T t} C_-^T C_- e^{A-t} dt \succcurlyeq 0$ . Obviously then,  $X = \text{diag}(X_-, 0)$  is a positive semi-definite solution of (4.2). Now, let  $X$  be a positive semi-definite solution to (4.2) with the property that  $\mathcal{X}_+(A) \subseteq \ker(X)$ . Then  $X$  must be of the form  $X = \text{diag}(X_1, 0)$ , and  $X_1$  must satisfy the reduced Lyapunov equation (4.3). Thus  $X = \text{diag}(X_-, 0)$ . Finally,  $H$  is asymptotically stable since  $\mathcal{X}_+(A) \subseteq \ker(C)$ . Moreover,

$$\begin{aligned} \|H\|_{\mathcal{H}_2}^2 &= \text{tr} \left( B^T \int_0^\infty e^{A^T t} C^T C e^{At} dt B \right) \\ &= \text{tr} \left( B_-^T \int_0^\infty e^{A^T t} C_-^T C_- e^{A-t} dt B_- \right) \\ &= \text{tr}(B_-^T X_- B_-) \\ &= \text{tr}(B^T X B). \end{aligned} \quad \square$$

We will now deal with computing the  $\mathcal{H}_\infty$ -norm. The result is a generalization of Lemma 4 in [IKIA14].

**Lemma 4.2:**

Consider the system (4.1). Assume that its transfer function  $H$  has all its poles in the open left half plane. If there exists  $X \in \mathbb{R}^{p \times p}$  such that  $X = X^T$  and  $CA = XC$ , then  $\|H\|_{\mathcal{H}_\infty} = \|H(0)\|_2$ .  $\diamond$

*Proof.* For the first part of the proof, let us assume that  $(A, B, C)$  is minimal. Then, in particular,  $A$  is Hurwitz and  $(A, B)$  is controllable.

Clearly, the inequality  $\|S\|_{\mathcal{H}_\infty} \geq \|S(0)\|_2$  is always satisfied. We will prove that  $\|S\|_{\mathcal{H}_\infty} \leq \|S(0)\|_2$  using the bounded real lemma [Ran96], which states that  $\|S\|_{\mathcal{H}_\infty} \leq \gamma$  if and only if there exists  $P \in \mathbb{R}^{n \times n}$  such that  $P = P^T$  and

$$A^T P + PA + C^T C + \frac{1}{\gamma^2} P B B^T P \preccurlyeq 0.$$

Let us take  $\gamma = \|S(0)\|_2 = \|CA^{-1}B\|_2$ . This implies that

$$CA^{-1}BB^T A^{-T} C^T \preccurlyeq \gamma^2 I_p. \quad (4.4)$$

Defining  $P := -A^{-T} C^T X C A^{-1}$  and using (4.4) yields

$$\begin{aligned} &A^T P + PA + C^T C + \frac{1}{\gamma^2} P B B^T P \\ &= -C^T X C A^{-1} - A^{-T} C^T X C + C^T C \\ &\quad + \frac{1}{\gamma^2} A^{-T} C^T X C A^{-1} B B^T A^{-T} C^T X C A^{-1} \\ &\preccurlyeq -C^T X C A^{-1} - A^{-T} C^T X C + C^T C + A^{-T} C^T X X C A^{-1} \\ &= (X C A^{-1} - C)^T (X C A^{-1} - C) \\ &= 0. \end{aligned}$$

From the bounded real lemma, we conclude that  $\|S\|_{\mathcal{H}_\infty} \leq \|S(0)\|_2$ .

For a non-minimal representation  $(A, B, C)$ , applying the Kalman decomposition, let  $T$  be a invertible matrix such that

$$T^{-1}AT = \begin{bmatrix} A_1 & 0 & A_6 & 0 \\ A_2 & A_3 & A_4 & A_5 \\ 0 & 0 & A_7 & 0 \\ 0 & 0 & A_8 & A_9 \end{bmatrix}, \quad T^{-1}B = \begin{bmatrix} B_1 \\ B_2 \\ 0 \\ 0 \end{bmatrix}, \quad CT = [C_1 \ 0 \ C_2 \ 0],$$

where  $(A_1, B_1, C_1)$  is a minimal representation of  $(A, B, C)$  with  $A_1$  Hurwitz. From,

$$\begin{aligned} (CT)(T^{-1}AT) &= CAT = XCT = X(CT), \\ (CT)(T^{-1}AT) &= [C_1 \ 0 \ C_2 \ 0] \begin{bmatrix} A_1 & 0 & A_6 & 0 \\ A_2 & A_3 & A_4 & A_5 \\ 0 & 0 & A_7 & 0 \\ 0 & 0 & A_8 & A_9 \end{bmatrix} \\ &= [C_1A_1 \ 0 \ C_1A_6 + C_2A_7 \ 0], \\ X(CT) &= X [C_1 \ 0 \ C_2 \ 0] = [XC_1 \ 0 \ XC_2 \ 0], \end{aligned}$$

we find that  $C_1A_1 = XC_1$ . Therefore, the minimal representation satisfies the sufficient condition and using the result obtained above the proof is completed.  $\square$

Continuing our effort to compute the  $\mathcal{H}_\infty$ -norm, we now formulate a lemma that will be instrumental in evaluating a transfer function at the origin. Recall that for a given matrix  $A$ , its Moore-Penrose inverse is denoted by  $A^+$ .

**Lemma 4.3:**

Consider the system (4.1). If  $A$  is symmetric and  $\ker(A) \subseteq \ker(C)$ , then 0 is not a pole of the transfer function  $H$  and we have  $H(0) = -CA^+B$ .  $\diamond$

*Proof.* If  $A$  is invertible, then the conclusion follows immediately. Otherwise, let  $A = U\Lambda U^T$  be an eigenvalue decomposition with orthogonal  $U$  and  $\Lambda = \text{diag}(0, \Lambda_2)$ , where  $\Lambda_2 \in \mathbb{R}^{r \times r}$  and  $r$  is the rank of  $A$ . We denote  $U = [U_1 \ U_2]$ , with  $U_2 \in \mathbb{R}^{n \times r}$ . Then

$$A^+ = U\Lambda^+U^T = [U_1 \ U_2] \begin{bmatrix} 0 & 0 \\ 0 & \Lambda_2^{-1} \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} = U_2\Lambda_2^{-1}U_2^T.$$

Note that  $CU_1 = 0$ . We have

$$\begin{aligned} H(s) &= CU(sI - \Lambda)^{-1}U^TB \\ &= C [U_1 \ U_2] \begin{bmatrix} s^{-1}I & 0 \\ 0 & (sI - \Lambda_2)^{-1} \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} B \\ &= CU_2(sI - \Lambda_2)^{-1}U_2^TB. \end{aligned}$$

Hence,  $H(s)$  is defined at  $s = 0$  and  $H(0) = -CU_2\Lambda_2^{-1}U_2^TB = -CA^+B$ .  $\square$

### 4.2.3 Problem formulation

We consider leader-follower multi-agent systems of the form

$$\dot{x}(t) = (I_n \otimes A - \mathfrak{L} \otimes B)x(t) + (\mathfrak{B} \otimes F)u(t), \quad (4.5a)$$

$$y(t) = (\mathfrak{L} \otimes I_n)x(t), \quad (4.5b)$$

defined on an undirected, weighted, connected graph  $\mathfrak{G} = (\mathfrak{V}, \mathfrak{E}, \mathfrak{w})$  with vertex set  $\mathfrak{V} = \{1, 2, \dots, n\}$ , adjacency matrix  $\mathfrak{A}$ , Laplacian matrix  $\mathfrak{L}$ , leader set  $\mathfrak{V}_L = \{v_1, v_2, \dots, v_m\} \subseteq \mathfrak{V}$ , and follower set  $\mathfrak{V}_F = \mathfrak{V} \setminus \mathfrak{V}_L$ . Furthermore,  $x_i(t) \in \mathbb{R}^n$  is the state of agent  $i$ , and  $u_\ell(t) \in \mathbb{R}^m$  is the external input to the leader  $v_\ell$ . Finally,  $A \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times n}$ , and  $F \in \mathbb{R}^{n \times m}$  are real matrices. Denote  $x(t) = \text{col}(x_1(t), x_2(t), \dots, x_n(t))$  and  $u(t) = \text{col}(u_1(t), u_2(t), \dots, u_m(t))$ .

The goal of this section is to find a reduced order networked system, whose dynamics is a good approximation of the networked system (4.5a). Following [MTC14], the idea to obtain such an approximation is to *cluster* groups of agents in the network, and to treat each of the resulting clusters as a vertex in a new, reduced order, network. The reduced order network will again be a leader-follower network, and by the clustering procedure, essential interconnection features of the network will be preserved. We will also require that the *synchronization* properties of the network are preserved after reduction. We assume that the original network is synchronized, meaning that if the external inputs satisfy  $u_\ell = 0$  for  $\ell = 1, 2, \dots, m$ , then for all  $i, j \in \mathfrak{V}$ , we have

$$x_i(t) - x_j(t) \rightarrow 0$$

as  $t \rightarrow \infty$ . We impose that the reduction procedure preserves this property.

Being a measure for the disagreement between the states of the agents in (4.5a), we choose  $y(t) = (\mathfrak{L} \otimes I_n)x(t)$  as the output of the original network. Indeed, this output  $y$  can be considered a measure of the disagreement in the network, in the sense that  $y(t)$  is small if and only if the network is close to being synchronized.

Note that the state space dimension of (4.5) is equal to  $nm$ , its number of inputs equals  $mm$ , and the number of outputs is  $nm$ .

Here, we will use clustering to obtain a reduced order network, i.e., a network with a reduced number of agents, as an approximation of the original network (4.5).

### 4.2.4 Graph partitions and reduction by clustering

We consider networks whose interaction topologies are represented by weighted graphs  $\mathfrak{G}$  with vertex set  $\mathfrak{V}$ . The graph of the original network (4.5a) is undirected, however, our reduction procedure will lead to networks on directed graphs.

Next, we will construct a reduced order approximation of (4.5) by clustering the agents in the network using a partition of  $\mathfrak{G}$ . Extending the main idea in [MTC14], we take as reduced order system the Petrov-Galerkin projection of the original system (4.5), with the following choice for the matrices  $V$  and  $W$ :

$$W = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \otimes I_n \in \mathbb{R}^{nm \times vn}, \quad V = \mathfrak{P} \otimes I_n \in \mathbb{R}^{nm \times vn}.$$



The dynamics of the resulting reduced order model is then given by

$$\begin{aligned}\dot{\hat{x}}(t) &= (I_{\mathfrak{r}} \otimes A - \widehat{\mathfrak{L}} \otimes B)\hat{x}(t) + (\widehat{\mathfrak{B}} \otimes F)u(t), \\ \hat{y}(t) &= (\mathfrak{L}\mathfrak{P} \otimes I_n)\hat{x}(t),\end{aligned}\tag{4.6}$$

where

$$\begin{aligned}\widehat{\mathfrak{L}} &= (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{L} \mathfrak{P} \in \mathbb{R}^{\mathfrak{r} \times \mathfrak{r}}, \\ \widehat{\mathfrak{B}} &= (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \in \mathbb{R}^{\mathfrak{r} \times \mathfrak{m}}.\end{aligned}$$

It can be seen by inspection that the matrix  $\widehat{\mathfrak{L}}$  is the Laplacian of a weighted *directed* graph with vertex set  $\{1, 2, \dots, \mathfrak{r}\}$ , with  $\mathfrak{r}$  equal to the number of clusters in the partition  $\pi$ , and adjacency matrix  $\widehat{\mathfrak{A}} = [\widehat{\mathfrak{a}}_{pq}]$ , with

$$\widehat{\mathfrak{a}}_{pq} = \frac{1}{|\mathfrak{C}_p|} \sum_{i \in \mathfrak{C}_p, j \in \mathfrak{C}_q} \mathfrak{a}_{ij},$$

where  $|\mathfrak{C}_p|$  the cardinality of  $\mathfrak{C}_p$ . In other words: in the reduced graph, the edge from vertex  $q$  to vertex  $p$  is obtained by summing over all  $j \in \mathfrak{C}_q$  the weights of all edges to  $i \in \mathfrak{C}_p$  and dividing this sum by the cardinality of  $\mathfrak{C}_p$ . The row sums of  $\widehat{\mathfrak{L}}$  are indeed equal to zero since  $\widehat{\mathfrak{L}}\mathbf{1}_{\mathfrak{r}} = 0$ . The matrix  $\widehat{\mathfrak{B}} \in \mathbb{R}^{\mathfrak{r} \times \mathfrak{m}}$  satisfies

$$[\widehat{\mathfrak{B}}]_{pj} = \begin{cases} \frac{1}{|\mathfrak{C}_p|} & \text{if } \mathfrak{v}_j \in \mathfrak{C}_p, \\ 0 & \text{otherwise,} \end{cases}$$

where  $\mathfrak{v}_1, \mathfrak{v}_2, \dots, \mathfrak{v}_m$  are the leader vertices,  $p = 1, 2, \dots, \mathfrak{r}$ , and  $j = 1, 2, \dots, \mathfrak{m}$ .

Clearly, the state space dimension of the reduced order network (4.6) is equal to  $\mathfrak{r}n$ , whereas the dimensions  $\mathfrak{m}n$  and  $\mathfrak{n}n$  of the input and output have remained unchanged. Thus we can investigate the error between the original and reduced order network by looking at the difference of their transfer functions. In the sequel, we will investigate both the  $\mathcal{H}_2$ -norm as well as the  $\mathcal{H}_\infty$ -norm of this difference.

Before doing this, we will now first study the question whether our reduction procedure preserves synchronization. It is important to note that since, by assumption, the original undirected graph is connected, it has a directed spanning tree. It is easily verified that this property is preserved by our clustering procedure. Then, since the property of having a directed spanning tree is equivalent with 0 being a simple eigenvalue of the Laplacian (see [ME10, Proposition 3.8]), the reduced order Laplacian  $\widehat{\mathfrak{L}}$  has again 0 as a simple eigenvalue.

Now assume that the original network (4.5) is synchronized. It is well known, see e.g. [TTM13], that this is equivalent with the condition that for each nonzero eigenvalue  $\lambda$  of the Laplacian  $\mathfrak{L}$  the matrix  $A - \lambda B$  is Hurwitz. Thus, synchronization is preserved if and only if for each nonzero eigenvalue  $\widehat{\lambda}$  of the reduced order Laplacian  $\widehat{\mathfrak{L}}$  the matrix  $A - \widehat{\lambda}B$  is Hurwitz.

Unfortunately, in general  $A - \lambda B$  Hurwitz for all nonzero  $\lambda \in \sigma(\mathfrak{L})$  does *not* imply that  $A - \widehat{\lambda}B$  Hurwitz for all nonzero  $\lambda \in \sigma(\widehat{\mathfrak{L}})$ . An exception is the ‘single integrator’

case  $A = 0$  and  $B = 1$ , where this condition is trivially satisfied, so in this special case synchronization is preserved. Also if we restrict ourselves to a special type of graph partitions, namely *almost equitable partitions*, then synchronization turns out to be preserved.

As an immediate consequence, the reduced Laplacian  $\widehat{\mathcal{L}}$  resulting from an AEP satisfies  $\mathcal{L}\mathfrak{P} = \mathfrak{P}\widehat{\mathcal{L}}$ . Indeed, since  $\text{im}(\mathfrak{P})$  is  $\mathcal{L}$ -invariant we have  $\mathcal{L}\mathfrak{P} = \mathfrak{P}X$  for some matrix  $X$ . Obviously, we must then have  $X = (\mathfrak{P}^T\mathfrak{P})^{-1}\mathfrak{P}^T\mathcal{L}\mathfrak{P} = \widehat{\mathcal{L}}$ . From this, it follows that  $\sigma(\widehat{\mathcal{L}}) \subseteq \sigma(\mathcal{L})$ . It then readily follows that synchronization is preserved if we cluster according to an AEP:

**Theorem 4.4:**

Assume that the network (4.5) is synchronized. Let  $\pi$  be an AEP. Then the reduced order network (4.6) obtained by clustering according to  $\pi$  is synchronized.  $\diamond$

### 4.2.5 $\mathcal{H}_2$ -error bounds

In this section, we will formulate the first main theorem. The theorem gives an a priori upper bound for the  $\mathcal{H}_2$ -norm of the approximation error in the case that we cluster according to an AEP. After formulating the theorem, in the remainder of this section we will establish a proof.

Before stating the theorem, we will now first discuss some important ingredients. Let  $H$  and  $\widehat{H}$  denote the transfer functions of the original (4.5) and reduced order network (4.6), respectively. We will measure the approximation error by the  $\mathcal{H}_2$ -norm  $\|H - \widehat{H}\|_{\mathcal{H}_2}$  of these transfer functions. An important role will be played by the  $\mathbf{n} - 1$  auxiliary input-state-output systems

$$\begin{aligned} \dot{x}(t) &= (A - \lambda B)x(t) + Fd(t), \\ z(t) &= \lambda x(t), \end{aligned} \tag{4.7}$$

where  $\lambda$  ranges over the  $\mathbf{n} - 1$  nonzero eigenvalues of the Laplacian  $\mathcal{L}$ . Let  $H_\lambda(s) = \lambda(sI_n - A + \lambda B)^{-1}F$  be the transfer matrices of these systems. We assume that the original network (4.5) is synchronized, so that all of the  $A - \lambda B$  are Hurwitz. Let  $\|H_\lambda\|_{\mathcal{H}_2}$  denote the  $\mathcal{H}_2$ -norm of  $H_\lambda$ . Recall that the set of leader vertices is  $\mathfrak{V}_L = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ . Vertex  $\mathbf{v}_i$  will be called leader  $i$ . This leader is an element of cluster  $\mathfrak{C}_{k_i}$  for some  $k_i \in \{1, 2, \dots, \mathbf{r}\}$ . We now have the following theorem.

**Theorem 4.5:**

Assume that the network (4.5) is synchronized. Let  $\pi$  be an AEP of the graph  $\mathfrak{G}$ . The absolute approximation error when clustering  $\mathfrak{G}$  according to  $\pi$  then satisfies

$$\|H - \widehat{H}\|_{\mathcal{H}_2}^2 \leq H_{\max, \mathcal{H}_2}^2 \sum_{i=1}^m \left(1 - \frac{1}{|\mathfrak{C}_{k_i}|}\right),$$

where  $\mathfrak{C}_{k_i}$  is the cluster containing the leader  $i$ , and

$$H_{\max, \mathcal{H}_2} := \max_{\lambda \in \sigma(\mathcal{L}) \setminus \sigma(\widehat{\mathcal{L}})} \|H_\lambda\|_{\mathcal{H}_2}.$$

Furthermore, the relative approximation error satisfies

$$\frac{\|H - \widehat{H}\|_{\mathcal{H}_2}^2}{\|H\|_{\mathcal{H}_2}^2} \leq \left( \frac{H_{\max, \mathcal{H}_2}}{H_{\min, \mathcal{H}_2}} \right)^2 \frac{\sum_{i=1}^m \left( 1 - \frac{1}{|\mathbf{c}_{k_i}|} \right)}{\mathbf{m} \left( 1 - \frac{1}{n} \right)},$$

where

$$H_{\min, \mathcal{H}_2} := \min_{\lambda \in \sigma(\mathcal{L}) \setminus \{0\}} \|H_\lambda\|_{\mathcal{H}_2}. \quad \diamond$$

**Remark 4.6:**

We see that, with fixed number of agents and fixed number of leaders, the approximation error is equal to 0 if in each cluster that contains a leader, the leader is the only vertex in that cluster. In general, the upper bound increases if the numbers of cellmates of the leaders increase. The upper bound also depends multiplicatively on the maximal  $\mathcal{H}_2$ -norm of the auxiliary systems (4.7) over all Laplacian eigenvalues in the complement of the spectrum of the reduced Laplacian  $\widehat{\mathcal{L}}$ . The relative error in addition depends on the minimal  $\mathcal{H}_2$ -norm of the auxiliary systems (4.7) over all nonzero eigenvalues of the Laplacian  $\mathcal{L}$ .  $\diamond$

**Remark 4.7:**

For the special case that the agents are single integrators (so  $n = 1$ ,  $A = 0$ ,  $B = 1$ , and  $F = 1$ ) it is easily seen that  $H_{\max, \mathcal{H}_2} = \frac{1}{2} \max\{\lambda : \lambda \in \sigma(L) \setminus \sigma(\widehat{L})\}$  and  $H_{\min, \mathcal{H}_2} = \frac{1}{2} \min\{\lambda : \lambda \in \sigma(L), \lambda \neq 0\}$ . Thus, in the single integrator case the corresponding a priori upper bounds explicitly involve the Laplacian eigenvalues. As already noted in Section 4.2.1, the single integrator case was also studied in [MTC14] for the slightly different set up that the output equation in the original network (4.5) is taken as  $y(t) = (\mathfrak{W}^{\frac{1}{2}} \mathfrak{R}^T \otimes I_n)x(t)$  instead of  $y(t) = (\mathcal{L} \otimes I_n)x(t)$ . Here,  $\mathfrak{R}$  is the incidence matrix of the graph and  $\mathfrak{W}$  the diagonal matrix with the edge weights on the diagonal (in other words,  $\mathcal{L} = \mathfrak{R}\mathfrak{W}\mathfrak{R}^T$ ). It was shown in [MTC14] that in that case the absolute and relative approximation errors even admit the explicit formulas

$$\|H - \widehat{H}\|_{\mathcal{H}_2}^2 = \frac{1}{2} \sum_{i=1}^m \left( 1 - \frac{1}{|\mathbf{c}_{k_i}|} \right),$$

and

$$\frac{\|H - \widehat{H}\|_{\mathcal{H}_2}^2}{\|H\|_{\mathcal{H}_2}^2} = \frac{\sum_{i=1}^m \left( 1 - \frac{1}{|\mathbf{c}_{k_i}|} \right)}{\mathbf{m} \left( 1 - \frac{1}{n} \right)}. \quad \diamond$$

In the remainder of this section, we will establish a proof of Theorem 4.5. As a first step, we establish the following lemma (see also [MTC14], where only the single integrator case was treated).

**Lemma 4.8:**

Let  $\pi$  be an AEP of the graph  $\mathfrak{G}$ . The approximation error when clustering  $\mathfrak{G}$  according to  $\pi$  then satisfies

$$\|H - \widehat{H}\|_{\mathcal{H}_2}^2 = \|H\|_{\mathcal{H}_2}^2 - \|\widehat{H}\|_{\mathcal{H}_2}^2. \quad \diamond$$

*Proof.* First, note that the columns of  $\mathfrak{P}$  are orthogonal. We construct a matrix  $\mathfrak{T} = \begin{bmatrix} \mathfrak{P} & \mathfrak{P}_\perp \end{bmatrix}$ , where the  $\mathbf{n} \times (\mathbf{n} - \mathbf{r})$  matrix  $\mathfrak{P}_\perp$  is chosen such that the columns of  $\mathfrak{T}$  form an orthogonal basis for  $\mathbb{R}^{\mathbf{n}}$ . In this case, we have  $\mathfrak{P}^T \mathfrak{P}_\perp = 0$ . Next, we apply the state space transformation  $x(t) = (\mathfrak{T} \otimes I_n) \tilde{x}(t)$  to system (4.5). We obtain

$$\begin{aligned} \begin{bmatrix} \dot{\tilde{x}}_1(t) \\ \dot{\tilde{x}}_2(t) \end{bmatrix} &= \tilde{A} \begin{bmatrix} \tilde{x}_1(t) \\ \tilde{x}_2(t) \end{bmatrix} + \tilde{B}u(t), \\ y(t) &= \tilde{C} \begin{bmatrix} \tilde{x}_1(t) \\ \tilde{x}_2(t) \end{bmatrix}, \end{aligned} \quad (4.8)$$

where the matrices  $\tilde{A}$ ,  $\tilde{B}$ , and  $\tilde{C}$  are given by

$$\begin{aligned} \tilde{A} &= \begin{bmatrix} I_{\mathbf{r}} \otimes A - (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{L} \mathfrak{P} \otimes B & - (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{L} \mathfrak{P}_\perp \otimes B \\ - (\mathfrak{P}_\perp^T \mathfrak{P}_\perp)^{-1} \mathfrak{P}_\perp^T \mathfrak{L} \mathfrak{P} \otimes B & I_{\mathbf{n}-\mathbf{r}} \otimes A - (\mathfrak{P}_\perp^T \mathfrak{P}_\perp)^{-1} \mathfrak{P}_\perp^T \mathfrak{L} \mathfrak{P}_\perp \otimes B \end{bmatrix}, \\ \tilde{B} &= \begin{bmatrix} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \otimes F \\ (\mathfrak{P}_\perp^T \mathfrak{P}_\perp)^{-1} \mathfrak{P}_\perp^T \mathfrak{B} \otimes F \end{bmatrix}, \quad \tilde{C} = [\mathfrak{L} \mathfrak{P} \otimes I_n \quad \mathfrak{L} \mathfrak{P}_\perp \otimes I_n]. \end{aligned}$$

Obviously, in (4.8) the transfer function from  $u$  to  $y$  is equal to  $H$ . Furthermore, if the state component  $\tilde{x}_2$  is truncated from (4.8), what we are left with is the reduced order model (4.6). Since  $\pi$  is an AEP of  $\mathfrak{G}$ , by Lemma 2.48,  $\text{im}(\mathfrak{P})$  is invariant under  $\mathfrak{L}$ . From this, it follows that not only  $\mathfrak{P}_\perp^T \mathfrak{P} = 0$ , but also

$$\mathfrak{P}_\perp^T \mathfrak{L} \mathfrak{P} = 0 \text{ and } \mathfrak{P}_\perp^T \mathfrak{L}^2 \mathfrak{P} = 0. \quad (4.9)$$

It is easily checked that

$$H(s) = \widehat{H}(s) + H_{\text{err}}(s),$$

where  $H_{\text{err}}(s)$  is given by

$$\begin{aligned} H_{\text{err}}(s) &= (\mathfrak{L} \mathfrak{P}_\perp \otimes I_n) \left( sI - \left( I_{\mathbf{n}-\mathbf{r}} \otimes A - (\mathfrak{P}_\perp^T \mathfrak{P}_\perp)^{-1} \mathfrak{P}_\perp^T \mathfrak{L} \mathfrak{P}_\perp \otimes B \right) \right)^{-1} \\ &\quad \cdot \left( (\mathfrak{P}_\perp^T \mathfrak{P}_\perp)^{-1} \mathfrak{P}_\perp^T \mathfrak{B} \otimes F \right). \end{aligned} \quad (4.10)$$

From (4.9) and (4.10), we have  $\widehat{H}(-s)^T H_{\text{err}}(s) = 0$ . Thus, we find that

$$\|H\|_{\mathcal{H}_2}^2 = \|\widehat{H}\|_{\mathcal{H}_2}^2 + \|H_{\text{err}}\|_{\mathcal{H}_2}^2,$$

which concludes the proof.  $\square$

Recall that, since  $\pi$  is an AEP, we have  $\sigma(\widehat{\mathfrak{L}}) \subseteq \sigma(\mathfrak{L})$ . Label the eigenvalues of  $\mathfrak{L}$  as  $0, \lambda_2, \lambda_3, \dots, \lambda_n$  in such a way that  $0, \lambda_2, \lambda_3, \dots, \lambda_r$  are the eigenvalues of  $\widehat{\mathfrak{L}}$ . Also, without loss of generality, we assume that  $\pi$  is *regularly formed*, i.e., all ones in each of the columns of  $\mathfrak{P}$  are consecutive. One can always relabel the agents in the graph in such a way that this is achieved. Consider now the symmetric matrix

$$\overline{\mathfrak{L}} := (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{\mathfrak{L}} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} = (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \mathfrak{P}^T \mathfrak{L} \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}}. \quad (4.11)$$

Note that the eigenvalues of  $\overline{\mathfrak{L}}$  and  $\widehat{\mathfrak{L}}$  coincide. Let  $\widehat{U}$  be an orthogonal matrix that diagonalizes  $\overline{\mathfrak{L}}$ . We then have

$$\widehat{U}^T \overline{\mathfrak{L}} \widehat{U} = \text{diag}(0, \lambda_2, \dots, \lambda_r) =: \widehat{\Lambda}. \quad (4.12)$$

Next, take  $U_1 = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \widehat{U}$ . The columns of  $U_1$  are orthonormal:

$$U_1^T U_1 = \widehat{U}^T (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \mathfrak{P}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \widehat{U} = \widehat{U}^T \widehat{U} = I.$$

Furthermore, we have that

$$U_1^T \mathfrak{L} U_1 = \widehat{U}^T \overline{\mathfrak{L}} \widehat{U} = \widehat{\Lambda}.$$

Now choose  $U_2$  such that  $U = [U_1 \ U_2]$  is an orthogonal matrix and

$$\Lambda := U^T \mathfrak{L} U = \begin{bmatrix} \widehat{\Lambda} & 0 \\ 0 & \overline{\Lambda} \end{bmatrix}, \quad (4.13)$$

where  $\overline{\Lambda} = \text{diag}(\lambda_{r+1}, \dots, \lambda_n)$ . It is easily verified that the first column of  $U_1$ , and thus the first column of  $U$ , is given by  $\frac{1}{\sqrt{n}} \mathbf{1}_n$ , a fact that we will use in the remainder of this section.

Using the above, we will now first establish explicit formulas for the  $\mathcal{H}_2$ -norms of  $H$  and  $\widehat{H}$  separately. The following lemma gives a formula for the  $\mathcal{H}_2$ -norm of the original transfer function  $H$ .

**Lemma 4.9:**

Let  $U$  be as in (4.13). For  $i = 2, \dots, n$ , let  $X_i$  be the observability Gramian of the auxiliary system  $(A - \lambda_i B, F, \lambda_i I)$  in (4.7), i.e., the unique solution of the Lyapunov equation  $(A - \lambda_i B)^T X_i + X_i (A - \lambda_i B) + \lambda_i^2 I = 0$ . Then the  $\mathcal{H}_2$ -norm of  $H$  is given by:

$$\|H\|_{\mathcal{H}_2}^2 = \text{tr}((U^T \mathfrak{B} \mathfrak{B}^T U \otimes I) \text{diag}(0, F^T X_2 F, \dots, F^T X_n F)). \quad (4.14)$$

◇

*Proof.* It can be verified, using the fact that  $A - \lambda_i B$  is Hurwitz for  $i = 2, 3, \dots, n$ , that

$$\mathcal{X}_+(I \otimes A - \mathfrak{L} \otimes B) = \mathbf{1}_n \otimes \mathcal{X}_+(A).$$

This immediately implies that  $\mathcal{X}_+(I \otimes A - \mathfrak{L} \otimes B) \subseteq \ker(\mathfrak{L} \otimes I)$ . As a consequence of Proposition 4.1, we have

$$\|H\|_{\mathcal{H}_2}^2 = \text{tr}((\mathfrak{B}^T \otimes F^T) X (\mathfrak{B} \otimes F)),$$

where  $X$  is the unique positive semi-definite solution to the Lyapunov equation

$$(I \otimes A^T - \mathfrak{L} \otimes B^T) X + X(I \otimes A - \mathfrak{L} \otimes B) + \mathfrak{L}^2 \otimes I = 0 \quad (4.15)$$

with the property that  $\mathcal{X}_+(I \otimes A - \mathfrak{L} \otimes B) \subseteq \ker(X)$ . In order to compute this solution  $X$ , premultiply (4.15) by  $U^T \otimes I$  and postmultiply by  $U \otimes I$ , and substitute  $Z = (U^T \otimes I)X(U \otimes I)$  to obtain

$$(I \otimes A^T - \Lambda \otimes B^T) Z + Z(I \otimes A - \Lambda \otimes B) + \Lambda^2 \otimes I = 0. \quad (4.16)$$

Solving (4.16), we take  $Z$  as

$$Z = \text{diag}(0, X_2, \dots, X_n),$$

where  $X_i$ , for  $i = 2, \dots, n$ , is the observability Gramian of the auxiliary system  $(A - \lambda_i B, F, \lambda_i I)$  in (4.7). Next,  $X := (U \otimes I)Z(U^T \otimes I)$  is a solution of the original Lyapunov equation, and it is easily verified that indeed  $\mathcal{X}_+(I \otimes A - \mathfrak{L} \otimes B) \subseteq \ker(X)$ . Thus, we obtain the following expression for the  $\mathcal{H}_2$ -norm of  $H$ :

$$\begin{aligned} \|H\|_{\mathcal{H}_2}^2 &= \text{tr}((\mathfrak{B}^T U \otimes F^T) \text{diag}(0, X_2, \dots, X_n) (U^T \mathfrak{B} \otimes F)) \\ &= \text{tr}((U^T \mathfrak{B} \mathfrak{B}^T U \otimes I) \text{diag}(0, F^T X_2 F, \dots, F^T X_n F)). \quad \square \end{aligned}$$

We proceed with finding a formula for the  $\mathcal{H}_2$ -norm for the reduced system. This will be dealt with in the following lemma.

**Lemma 4.10:**

Let  $\widehat{U}$  be as in (4.12) above. For  $i = 2, \dots, \mathfrak{r}$ , let  $X_i$  be the observability Gramian of the auxiliary system  $(A - \lambda_i B, F, \lambda_i I)$  in (4.7), i.e., the unique solution of the Lyapunov equation  $(A - \lambda_i B)^T X_i + X_i(A - \lambda_i B) + \lambda_i^2 I = 0$ . Then the  $\mathcal{H}_2$ -norm of  $\widehat{H}$  is given by

$$\|\widehat{H}\|_{\mathcal{H}_2}^2 = \text{tr}\left(\left(\widehat{U}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{\mathfrak{B}} \widehat{\mathfrak{B}}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{U} \otimes I\right) \text{diag}(0, F^T X_2 F, \dots, F^T X_{\mathfrak{r}} F)\right). \quad (4.17)$$

◇

*Proof.* Firstly, it can be verified that

$$\mathcal{X}_+(I \otimes A - \widehat{\mathfrak{L}} \otimes B) = \mathbf{1}_{\mathfrak{r}} \otimes \mathcal{X}_+(A).$$

This implies that  $\mathcal{X}_+(I \otimes A - \widehat{\mathfrak{L}} \otimes B) \subseteq \ker(\mathfrak{L} \mathfrak{P} \otimes I)$ . By Proposition 4.1, we then have

$$\|\widehat{H}\|_{\mathcal{H}_2}^2 = \text{tr}\left(\left(\widehat{\mathfrak{B}}^T \otimes F^T\right) \widehat{X} \left(\widehat{\mathfrak{B}} \otimes F\right)\right),$$

where  $\widehat{X}$  is the unique positive semi-definite solution to the Lyapunov equation

$$\left(I \otimes A^T - \widehat{\mathcal{L}}^T \otimes B^T\right) \widehat{X} + \widehat{X} \left(I \otimes A - \widehat{\mathcal{L}} \otimes B\right) + \mathfrak{P}^T \mathcal{L}^2 \mathfrak{P} \otimes I = 0 \quad (4.18)$$

satisfying the property that  $\mathcal{X}_+(I \otimes A - \widehat{\mathcal{L}} \otimes B) \subseteq \ker(\widehat{X})$ . In order to compute this solution, pre- and postmultiply (4.18) by  $(\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes I$  and substitute

$$\widehat{Y} = \left((\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes I\right) \widehat{X} \left((\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes I\right)$$

to obtain

$$\left(I \otimes A^T - \overline{\mathcal{L}} \otimes B^T\right) \widehat{Y} + \widehat{Y} \left(I \otimes A - \overline{\mathcal{L}} \otimes B\right) + (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \mathfrak{P}^T \mathcal{L}^2 \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes I = 0. \quad (4.19)$$

Recall from Section 4.2.4 that  $\mathcal{L} \mathfrak{P} = \mathfrak{P} \widehat{\mathcal{L}}$ . From this it follows that

$$(\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \mathfrak{P}^T \mathcal{L}^2 \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} = \overline{\mathcal{L}}^2.$$

Consequently, we can diagonalize the corresponding term in (4.19) by premultiplying by  $\widehat{U}^T \otimes I$  and postmultiplying by  $\widehat{U} \otimes I$ , where  $\widehat{U}$  is as in (4.12). Next, we denote  $\widehat{Z} = (\widehat{U}^T \otimes I) \widehat{Y} (\widehat{U} \otimes I)$  so that (4.19) reduces to

$$\left(I \otimes A^T - \widehat{\Lambda} \otimes B^T\right) \widehat{Z} + \widehat{Z} \left(I \otimes A - \widehat{\Lambda} \otimes B\right) + \widehat{\Lambda}^2 \otimes I = 0,$$

which can be solved by taking

$$\widehat{Z} = \text{diag}(0, X_2, \dots, X_{\mathfrak{r}}),$$

where again  $X_i$ , for  $i = 2, \dots, \mathfrak{r}$ , is the observability Gramian of the auxiliary system  $(A - \lambda_i B, F, \lambda_i I)$  in (4.7). Next,

$$\widehat{X} = \left((\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{U} \otimes I\right) \widehat{Z} \left(\widehat{U}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \otimes I\right)$$

then satisfies (4.18), and it can be verified that  $\mathcal{X}_+(I \otimes A - \widehat{\mathcal{L}} \otimes B) \subseteq \ker(\widehat{X})$ . Thus, the  $\mathcal{H}_2$ -norm of  $\widehat{H}$  is given by

$$\begin{aligned} \|\widehat{H}\|_{\mathcal{H}_2}^2 &= \text{tr} \left( \left( \widehat{\mathfrak{B}}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{U} \otimes F^T \right) \text{diag}(0, X_2, \dots, X_{\mathfrak{r}}) \left( \widehat{U}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{\mathfrak{B}} \otimes F \right) \right) \\ &= \text{tr} \left( \left( \widehat{U}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{\mathfrak{B}} \widehat{\mathfrak{B}}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{U} \otimes I \right) \text{diag}(0, F^T X_2 F, \dots, F^T X_{\mathfrak{r}} F) \right). \quad \square \end{aligned}$$

We will now combine the previous lemmas, and give a proof of Theorem 4.5.

*Proof of Theorem 4.5.* Using Lemma 4.8, and formulas (4.14) and (4.17), we compute

$$\begin{aligned}
\|H - \widehat{H}\|_{\mathcal{H}_2}^2 &= \text{tr}((U^T \mathfrak{B} \mathfrak{B}^T U \otimes I) \text{diag}(0, F^T X_2 F, \dots, F^T X_n F)) \\
&\quad - \text{tr}\left(\left(\widehat{U}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{\mathfrak{B}} \widehat{\mathfrak{B}}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{U} \otimes I\right) \text{diag}(0, F^T X_2 F, \dots, F^T X_r F)\right) \\
&= \text{tr}\left(\left(\begin{bmatrix} U_1^T \mathfrak{B} \mathfrak{B}^T U_1 & U_1^T \mathfrak{B} \mathfrak{B}^T U_2 \\ U_2^T \mathfrak{B} \mathfrak{B}^T U_1 & U_2^T \mathfrak{B} \mathfrak{B}^T U_2 \end{bmatrix} \otimes I\right) \text{diag}(0, F^T X_2 F, \dots, F^T X_n F)\right) \\
&\quad - \text{tr}((U_1^T \mathfrak{B} \mathfrak{B}^T U_1 \otimes I) \text{diag}(0, F^T X_2 F, \dots, F^T X_r F)) \\
&= \text{tr}((U_2^T \mathfrak{B} \mathfrak{B}^T U_2 \otimes I) \text{diag}(F^T X_{r+1} F, \dots, F^T X_n F)), \tag{4.20}
\end{aligned}$$

where the second equality follows from the fact that

$$\begin{aligned}
\widehat{\mathfrak{B}}^T (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{U} &= \mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{U} \\
&= \mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \widehat{U} \\
&= \mathfrak{B}^T U_1.
\end{aligned}$$

Next, observe that (4.20) can be rewritten as

$$\begin{aligned}
\|H - \widehat{H}\|_{\mathcal{H}_2}^2 &= \text{tr}((U_2^T \mathfrak{B} \mathfrak{B}^T U_2 \otimes I) \text{diag}(F^T X_{r+1} F, \dots, F^T X_n F)) \\
&= \text{tr}((U_2^T \mathfrak{B} \mathfrak{B}^T U_2) \text{diag}(\text{tr}(F^T X_{r+1} F), \dots, \text{tr}(F^T X_n F))) \\
&= \text{tr}\left((U_2^T \mathfrak{B} \mathfrak{B}^T U_2) \text{diag}\left(\|H_{\lambda_{r+1}}\|_{\mathcal{H}_2}^2, \dots, \|H_{\lambda_n}\|_{\mathcal{H}_2}^2\right)\right),
\end{aligned}$$

where  $H_{\lambda_i}$  for  $i = r+1, \dots, n$  is the transfer function of the auxiliary system (4.7). An upper bound for this expression is given by

$$\text{tr}\left((U_2^T \mathfrak{B} \mathfrak{B}^T U_2) \text{diag}\left(\|H_{\lambda_{r+1}}\|_{\mathcal{H}_2}^2, \dots, \|H_{\lambda_n}\|_{\mathcal{H}_2}^2\right)\right) \leq H_{\max, \mathcal{H}_2}^2 \text{tr}(U_2^T \mathfrak{B} \mathfrak{B}^T U_2),$$

where  $H_{\max, \mathcal{H}_2} = \max_{r+1 \leq j \leq n} \|H_{\lambda_j}\|_{\mathcal{H}_2}$ . Furthermore, we have

$$\begin{aligned}
\text{tr}(U_2^T \mathfrak{B} \mathfrak{B}^T U_2) &= \text{tr}(U^T \mathfrak{B} \mathfrak{B}^T U) - \text{tr}(U_1^T \mathfrak{B} \mathfrak{B}^T U_1) \\
&= m - \text{tr}\left(\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \mathfrak{B}^T\right).
\end{aligned}$$

Since, by assumption, the partition  $\pi$  is regularly formed,  $\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T$  is a block diagonal matrix of the form

$$\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T = \text{diag}(\mathfrak{P}_1, \mathfrak{P}_2, \dots, \mathfrak{P}_k).$$

It is easily verified that each  $\mathfrak{P}_i$  is a  $|\mathcal{C}_i| \times |\mathcal{C}_i|$  matrix whose elements are all equal to  $\frac{1}{|\mathcal{C}_i|}$ . The matrix  $\mathfrak{B} \mathfrak{B}^T$  is a diagonal matrix whose diagonal entries are either 0 or 1. We then have that the  $i$ th column of  $\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \mathfrak{B}^T$  is either equal to the  $i$ th



column of  $\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T$  if agent  $i$  is a leader, or zero otherwise. It then follows that the diagonal elements of  $\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \mathfrak{B}^T$  are either zero or  $\frac{1}{|\mathfrak{C}_{k_i}|}$  if  $i$  is part of the leader set, where  $\mathfrak{C}_{k_i}$  is the cluster containing agent  $i$ . Hence, we have

$$\text{tr}(U_1^T \mathfrak{B} \mathfrak{B}^T U_1) = \sum_{i=1}^m \frac{1}{|\mathfrak{C}_{k_i}|},$$

and consequently

$$\text{tr}(U_2^T \mathfrak{B} \mathfrak{B}^T U_2) = m - \sum_{i=1}^m \frac{1}{|\mathfrak{C}_{k_i}|}.$$

In conclusion, we have

$$\|H - \widehat{H}\|_{\mathcal{H}_2}^2 \leq H_{\max, \mathcal{H}_2}^2 \sum_{i=1}^m \left(1 - \frac{1}{|\mathfrak{C}_{k_i}|}\right),$$

which completes the proof of the first part of the theorem.

We now prove the statement about the relative error. For this, we will establish a lower bound for  $\|H\|_{\mathcal{H}_2}^2$ . By (4.14), we have

$$\begin{aligned} \|H\|_{\mathcal{H}_2}^2 &= \text{tr}((U^T \mathfrak{B} \mathfrak{B}^T U \otimes I) \text{diag}(0, F^T X_2 F, \dots, F^T X_n F)) \\ &= \text{tr}((U^T \mathfrak{B} \mathfrak{B}^T U) \text{diag}(0, \text{tr}(F^T X_2 F), \dots, \text{tr}(F^T X_n F))). \end{aligned} \quad (4.21)$$

The first column of  $U$  spans the eigenspace corresponding to the eigenvalue 0 of  $\mathfrak{L}$  and hence must be equal to  $u_1 = \frac{1}{\sqrt{n}} \mathbf{1}_n$ . Let  $\bar{U}$  be such that  $U = [u_1 \ \bar{U}]$ . It is then easily verified using (4.21) that

$$\begin{aligned} \|H\|_{\mathcal{H}_2}^2 &= \text{tr}\left(\left(\bar{U}^T \mathfrak{B} \mathfrak{B}^T \bar{U}\right) \text{diag}(\text{tr}(F^T X_2 F), \dots, \text{tr}(F^T X_n F))\right) \\ &= \text{tr}\left(\left(\bar{U}^T \mathfrak{B} \mathfrak{B}^T \bar{U}\right) \text{diag}(\|H_{\lambda_2}\|_{\mathcal{H}_2}^2, \dots, \|H_{\lambda_n}\|_{\mathcal{H}_2}^2)\right). \end{aligned}$$

Finally, since

$$\text{tr}(\bar{U}^T \mathfrak{B} \mathfrak{B}^T \bar{U}) = \text{tr}(\mathfrak{B}^T \bar{U} \bar{U}^T \mathfrak{B}) = \text{tr}(\mathfrak{B}^T (U U^T - u_1 u_1^T) \mathfrak{B}) = m - \frac{m}{n},$$

we obtain that  $\|H\|_{\mathcal{H}_2}^2 \geq m \left(1 - \frac{1}{n}\right) (H_{\min, \mathcal{H}_2})^2$ . This then yields the upper bound for the relative error as claimed.  $\square$

**Remark 4.11:**

Note that by our labeling of the eigenvalues of  $\mathfrak{L}$ , in the formulation of Theorem 4.5, we have that  $\sigma(\mathfrak{L}) \setminus \sigma(\widehat{\mathfrak{L}})$  is equal to  $\{\lambda_{\tau+1}, \dots, \lambda_n\}$  used in the proof. We stress that this should not be confused with the notation often used in the literature, where the  $\lambda_i$ 's are labeled in increasing order.  $\diamond$

### 4.2.6 $\mathcal{H}_\infty$ -error bounds

Whereas in the previous section we studied a priori upper bounds for the approximation error in terms of the  $\mathcal{H}_2$ -norm, the present section aims at expressing the approximation error in terms of the  $\mathcal{H}_\infty$ -norm. This section consists of two subsections. In the first subsection, we consider the special case that the agent dynamics is a single integrator system. Here we obtain an explicit formula for the  $\mathcal{H}_\infty$ -norm of the error. In the second subsection, we find an upper bound for the  $\mathcal{H}_\infty$ -error for symmetric systems.

#### 4.2.6.1 The single integrator case

Here we consider the special case that the agent dynamics is a single integrator system. In this case, we have  $A = 0$ ,  $B = 1$ , and  $F = 1$  and the original system (4.5) reduces to

$$\begin{aligned} \dot{x}(t) &= -\mathfrak{L}x(t) + \mathfrak{B}u(t), \\ y(t) &= \mathfrak{L}x(t). \end{aligned} \tag{4.22}$$

The state space dimension of (4.22) is then simply  $\mathbf{n}$ , the number of agents. For a given partition  $\pi = \{\mathfrak{C}_1, \mathfrak{C}_2, \dots, \mathfrak{C}_\tau\}$ , the reduced system (4.6) is now given by

$$\begin{aligned} \dot{\hat{x}}(t) &= -\hat{\mathfrak{L}}\hat{x}(t) + \hat{\mathfrak{B}}u(t), \\ \hat{y}(t) &= \mathfrak{L}\mathfrak{P}\hat{x}(t), \end{aligned}$$

where  $\mathfrak{P}$  is again the characteristic matrix of  $\pi$  and  $\hat{x}(t) \in \mathbb{R}^\tau$ . The transfer functions  $H$  and  $\hat{H}$ , of the original and reduced system respectively, are given by

$$\begin{aligned} H(s) &= \mathfrak{L}(sI_n + \mathfrak{L})^{-1}\mathfrak{B}, \\ \hat{H}(s) &= \mathfrak{L}\mathfrak{P}(sI_\tau + \hat{\mathfrak{L}})^{-1}\hat{\mathfrak{B}}. \end{aligned}$$

The first main result of this section is the following explicit formula for the  $\mathcal{H}_\infty$  model reduction error. It complements the formula for the  $\mathcal{H}_2$ -error obtained in [MTC14] (see also Remark 4.7).

**Theorem 4.12:**

Let  $\pi$  be an AEP of the graph  $\mathfrak{G}$ . If the network with single integrator agent dynamics (4.22) is clustered according to  $\pi$ , then the  $\mathcal{H}_\infty$ -error is given by

$$\|H - \hat{H}\|_{\mathcal{H}_\infty}^2 = \begin{cases} \max_{1 \leq i \leq \mathbf{m}} \left(1 - \frac{1}{|\mathfrak{C}_{k_i}|}\right) & \text{if the leaders are in different clusters,} \\ 1 & \text{otherwise,} \end{cases}$$

where, for some  $k_i \in \{1, 2, \dots, \mathbf{m}\}$ ,  $\mathfrak{C}_{k_i}$  is the cluster containing the leader  $i$ . Furthermore,  $\|H\|_{\mathcal{H}_\infty} = 1$ , hence the relative and absolute  $\mathcal{H}_\infty$ -errors coincide.  $\diamond$

**Remark 4.13:**

We see that the  $\mathcal{H}_\infty$ -error lies in the interval  $[0, 1]$ . The error is maximal ( $= 1$ ) if and only if two or more leader vertices occupy one and the same cluster. The error is minimal ( $= 0$ ) if and only if each leader vertex occupies a different cluster, and is the only vertex in this cluster. In general, the error increases if the number of cellmates of the leaders increases.  $\diamond$

*Proof of Theorem 4.12.* To simplify notation, denote  $H_{\text{err}}(s) = H(s) - \widehat{H}(s)$ . Note that both  $H$  and  $\widehat{H}$  have all poles in the open left half plane. We now first show that, since  $\pi$  is an AEP, we have

$$\|H_{\text{err}}\|_{\mathcal{H}_\infty} = \|H_{\text{err}}(0)\|_2. \quad (4.23)$$

First note that  $\widehat{H}(s) = \mathfrak{L}\mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}}(sI_\tau + \overline{\mathfrak{L}})^{-1}(\mathfrak{P}^T\mathfrak{P})^{\frac{1}{2}}\widehat{\mathfrak{B}}$ , where the symmetric matrix  $\overline{\mathfrak{L}}$  is given by (4.11). Thus, a state space representation for the error system is given by

$$\begin{aligned} \dot{x}_e(t) &= \begin{bmatrix} -\mathfrak{L} & 0 \\ 0 & -\overline{\mathfrak{L}} \end{bmatrix} x_e(t) + \begin{bmatrix} \mathfrak{B} \\ (\mathfrak{P}^T\mathfrak{P})^{\frac{1}{2}}\widehat{\mathfrak{B}} \end{bmatrix} u(t), \\ e(t) &= \begin{bmatrix} \mathfrak{L} & -\mathfrak{L}\mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}} \end{bmatrix} x_e(t). \end{aligned} \quad (4.24)$$

Next, we show that (4.23) holds by applying Lemma 4.2 to system (4.24). Indeed, with  $X = -\mathfrak{L}$ , we have

$$\begin{aligned} \begin{bmatrix} \mathfrak{L} & -\mathfrak{L}\mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}} \end{bmatrix} \begin{bmatrix} -\mathfrak{L} & 0 \\ 0 & -\overline{\mathfrak{L}} \end{bmatrix} &= \begin{bmatrix} -\mathfrak{L}^2 & \mathfrak{L}\mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}}\overline{\mathfrak{L}} \end{bmatrix} \\ &= \begin{bmatrix} -\mathfrak{L}^2 & \mathfrak{L}\widehat{\mathfrak{L}}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}} \end{bmatrix} \\ &= \begin{bmatrix} -\mathfrak{L}^2 & \mathfrak{L}^2\mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}} \end{bmatrix} \\ &= X \begin{bmatrix} \mathfrak{L} & -\mathfrak{L}\mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}} \end{bmatrix}, \end{aligned}$$

and from Lemma 4.2 it then immediately follows that  $\|H_{\text{err}}\|_{\mathcal{H}_\infty} = \|H_{\text{err}}(0)\|_2$ . To compute  $\|H_{\text{err}}(0)\|_2$ , we apply Lemma 4.3 to system (4.24). First, it is easily verified that

$$\ker \begin{bmatrix} -\mathfrak{L} & 0 \\ 0 & -\overline{\mathfrak{L}} \end{bmatrix} \subseteq \ker \begin{bmatrix} \mathfrak{L} & -\mathfrak{L}\mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}} \end{bmatrix}.$$

By applying Lemma 4.3 we then obtain

$$\begin{aligned} H_{\text{err}}(0) &= \begin{bmatrix} \mathfrak{L} & -\mathfrak{L}\mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \mathfrak{L} & 0 \\ 0 & \overline{\mathfrak{L}} \end{bmatrix}^+ \begin{bmatrix} \mathfrak{B} \\ (\mathfrak{P}^T\mathfrak{P})^{\frac{1}{2}}\widehat{\mathfrak{B}} \end{bmatrix} \\ &= \mathfrak{L} \left( \mathfrak{L}^+ - \mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}}\overline{\mathfrak{L}}^+(\mathfrak{P}^T\mathfrak{P})^{-\frac{1}{2}}\mathfrak{P}^T \right) \mathfrak{B}. \end{aligned} \quad (4.25)$$

Recall that  $\widehat{U}$  in (4.12) is an orthogonal matrix that diagonalizes  $\overline{\mathfrak{L}}$  and that  $U_1 = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \widehat{U}$ . Then  $\overline{\mathfrak{L}}^+ = \widehat{U} \widehat{\Lambda}^+ \widehat{U}^T$ . Thus, we have

$$\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \overline{\mathfrak{L}}^+ (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \mathfrak{P}^T = U_1 \widehat{\Lambda}^+ U_1^T.$$

Next, we compute

$$\begin{aligned} \mathfrak{L} \mathfrak{L}^+ &= U \Lambda U^T U \Lambda^+ U^T \\ &= U \Lambda \Lambda^+ U^T \\ &= I_n - \frac{1}{\mathbf{n}} \mathbf{1}_n \mathbf{1}_n^T, \end{aligned} \tag{4.26}$$

where the last equality follows from the fact that the first column of  $U$  is  $\frac{1}{\sqrt{\mathbf{n}}} \mathbf{1}_n$ . Now observe that

$$\begin{aligned} \mathfrak{L} U_1 \widehat{\Lambda}^+ U_1^T &= U \Lambda U^T U_1 \widehat{\Lambda}^+ U_1^T \\ &= U_1 \widehat{\Lambda} \widehat{\Lambda}^+ U_1^T \\ &= U_1 U_1^T - \frac{1}{\mathbf{n}} \mathbf{1}_n \mathbf{1}_n^T \\ &= \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T - \frac{1}{\mathbf{n}} \mathbf{1}_n \mathbf{1}_n^T. \end{aligned} \tag{4.27}$$

Combining (4.26) and (4.27) with (4.25), we obtain

$$H_{\text{err}}(0) = \left( I_n - \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \right) \mathfrak{B}.$$

From (4.23) then, we have that the  $\mathcal{H}_\infty$ -error is given by

$$\begin{aligned} \|H - \widehat{H}\|_{\mathcal{H}_\infty}^2 &= \lambda_{\max} \left( H_{\text{err}}(0)^T H_{\text{err}}(0) \right) \\ &= \lambda_{\max} \left( \mathfrak{B}^T \left( I_n - \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \right)^2 \mathfrak{B} \right) \\ &= \lambda_{\max} \left( I_m - \mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \right) \\ &= 1 - \lambda_{\min} \left( \mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \right). \end{aligned} \tag{4.28}$$

All that is left is to compute the minimal eigenvalue of  $\mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B}$ . Again, let  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$  be the set of leaders and note that  $\mathfrak{B}$  satisfies

$$\mathfrak{B} = [e_{v_1} \quad e_{v_2} \quad \cdots \quad e_{v_m}].$$

As before, without loss of generality, assume that  $\pi$  is regularly formed. Then the matrix  $\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T$  is block diagonal where each diagonal block  $\mathfrak{P}_i$  is a  $|\mathfrak{C}_i| \times |\mathfrak{C}_i|$  matrix whose entries are all  $\frac{1}{|\mathfrak{C}_i|}$ . Let  $k_i \in \{1, 2, \dots, \mathfrak{r}\}$  be such that  $\mathbf{v}_i \in \mathfrak{C}_{k_i}$ . If all the leaders are in different clusters, then

$$\mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} = \text{diag} \left( \frac{1}{|\mathfrak{C}_{k_1}|}, \frac{1}{|\mathfrak{C}_{k_2}|}, \dots, \frac{1}{|\mathfrak{C}_{k_m}|} \right),$$

and so

$$\lambda_{\min} \left( \mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \right) = \min_{1 \leq i \leq m} \frac{1}{|\mathfrak{C}_{k_i}|}. \quad (4.29)$$

Now suppose that two leaders  $\mathbf{v}_i$  and  $\mathbf{v}_j$  are cellmates. Then we have

$$\mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} (e_i - e_j) = \mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T (e_{\mathbf{v}_i} - e_{\mathbf{v}_j}) = 0.$$

which together with  $\mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \succcurlyeq 0$  implies

$$\lambda_{\min} \left( \mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathfrak{B} \right) = 0. \quad (4.30)$$

From (4.28), (4.29), and (4.30), we find the absolute  $\mathcal{H}_\infty$ -error. To find the relative  $\mathcal{H}_\infty$ -error, we compute  $\|H\|_{\mathcal{H}_\infty}$  by applying Lemma 4.2 and Lemma 4.3 to the original system (4.22). Combined with (4.26), this results in the  $\mathcal{H}_\infty$ -norm of the original system:

$$\|H\|_{\mathcal{H}_\infty}^2 = \lambda_{\max} \left( H(0)^T H(0) \right) = \lambda_{\max} \left( \mathfrak{B}^T \left( I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) \mathfrak{B} \right) = 1.$$

This completes the proof.  $\square$

#### 4.2.6.2 The general case with symmetric agent dynamics

In this subsection, we return to the general case that the agent dynamics is given by an arbitrary multivariable input-state-output system. Thus, the original and reduced networks are again given by (4.5) and (4.6), respectively. As in the proof of Theorem 4.12, we will rely heavily on Lemma 4.3 to compute the  $\mathcal{H}_\infty$ -error. Since Lemma 4.3 relies on a symmetry argument, we will need to assume that the matrices  $A$  and  $B$  are both symmetric, which will be a standing assumption in the remainder of this section.

We will now establish an a priori upper bound for the  $\mathcal{H}_\infty$ -norm of the approximation error in the case that we cluster according to an AEP. Again, an important role is played by the  $n - 1$  auxiliary systems (4.7) with  $\lambda$  ranging over the nonzero eigenvalues of the Laplacian  $\mathfrak{L}$ . Again, let  $H_\lambda(s) = \lambda(sI - A + \lambda B)^{-1} F$  be their transfer functions. We assume that the original network (4.5) is synchronized, so that all of the  $A - \lambda B$  are Hurwitz. We again use  $H$ ,  $\hat{H}$ , and  $H_{\text{err}}$  to denote the relevant transfer functions.

The following is the second main theorem.

**Theorem 4.14:**

Assume the network (4.5) is synchronized and that  $A$  and  $B$  are symmetric matrices. Let  $\pi$  be an AEP of the graph  $\mathfrak{G}$ . The  $\mathcal{H}_\infty$ -error when clustering  $\mathfrak{G}$  according to  $\pi$  then satisfies

$$\|H - \hat{H}\|_{\mathcal{H}_\infty}^2 \leq \begin{cases} H_{\max, \mathcal{H}_\infty}^2 \max_{1 \leq i \leq m} \left( 1 - \frac{1}{|\mathfrak{C}_{k_i}|} \right) & \text{if the leaders are in different clusters,} \\ H_{\max, \mathcal{H}_\infty}^2 & \text{otherwise,} \end{cases}$$

and

$$\frac{\|H - \widehat{H}\|_{\mathcal{H}_\infty}^2}{\|H\|_{\mathcal{H}_\infty}^2} \leq \begin{cases} \left(\frac{H_{\max, \mathcal{H}_\infty}}{H_{\min, \mathcal{H}_\infty}}\right)^2 \max_{1 \leq i \leq m} \left(1 - \frac{1}{|\mathfrak{e}_{k_i}|}\right) & \text{if the leaders are in different clusters,} \\ \left(\frac{H_{\max, \mathcal{H}_\infty}}{H_{\min, \mathcal{H}_\infty}}\right)^2 & \text{otherwise,} \end{cases}$$

where

$$H_{\max, \mathcal{H}_\infty} := \max_{\lambda \in \sigma(\mathfrak{L}) \setminus \sigma(\widehat{\mathfrak{L}})} \|H_\lambda\|_{\mathcal{H}_\infty}, \quad (4.31)$$

and

$$H_{\min, \mathcal{H}_\infty} := \min_{\lambda \in \sigma(\mathfrak{L}) \setminus \{0\}} \sigma_{\min}(H_\lambda(0)), \quad (4.32)$$

with  $H_\lambda$  the transfer functions of the auxiliary systems (4.7).  $\diamond$

**Remark 4.15:**

The absolute  $\mathcal{H}_\infty$ -error thus lies in the interval  $[0, H_{\max, \mathcal{H}_\infty}]$  with  $H_{\max, \mathcal{H}_\infty}$  the maximum over the  $\mathcal{H}_\infty$ -norms of the transfer functions  $H_\lambda$  with  $\lambda \in \sigma(\mathfrak{L}) \setminus \sigma(\widehat{\mathfrak{L}})$ . The error is minimal (= 0) if each leader vertex occupies a different cluster, and is the only vertex in this cluster. In general, the upper bound increases if the number of cellmates of the leaders increases.  $\diamond$

*Proof of Theorem 4.14.* First note that the transfer function  $\widehat{H}$  of the reduced network (4.6) is equal to

$$\widehat{H}(s) = \left(\mathfrak{L}\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes I_n\right) (sI - I_r \otimes A + \overline{\mathfrak{L}} \otimes B)^{-1} \left((\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{\mathfrak{B}} \otimes F\right), \quad (4.33)$$

with the symmetric matrix  $\overline{\mathfrak{L}}$  given by (4.11). Analogous to the proof of Theorem 4.12, we first apply Lemma 4.2 to the error system

$$\begin{aligned} \dot{x}_e(t) &= \begin{bmatrix} I_n \otimes A - \mathfrak{L} \otimes B & 0 \\ 0 & I_r \otimes A - \overline{\mathfrak{L}} \otimes B \end{bmatrix} x_e(t) + \begin{bmatrix} \mathfrak{B} \otimes F \\ (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \widehat{\mathfrak{B}} \otimes F \end{bmatrix} u(t), \\ e(t) &= \begin{bmatrix} \mathfrak{L} \otimes I_n & -\mathfrak{L}\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes I_n \end{bmatrix} x_e(t), \end{aligned}$$

with transfer function  $H_{\text{err}}$ . Take  $X = I_n \otimes A - \mathfrak{L} \otimes B$ . We then have

$$\begin{aligned} &\begin{bmatrix} \mathfrak{L} \otimes I_n & -\mathfrak{L}\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes I_n \end{bmatrix} \begin{bmatrix} I_n \otimes A - \mathfrak{L} \otimes B & 0 \\ 0 & I_r \otimes A - \overline{\mathfrak{L}} \otimes B \end{bmatrix} \\ &= X \begin{bmatrix} \mathfrak{L} \otimes I_n & -\mathfrak{L}\mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes I_n \end{bmatrix}. \end{aligned}$$

From Lemma 4.2, we thus obtain that

$$\|H_{\text{err}}\|_{\mathcal{H}_\infty} = \|H_{\text{err}}(0)\|_2 = \lambda_{\max} \left( H_{\text{err}}(0)^T H_{\text{err}}(0) \right)^{\frac{1}{2}}.$$

In the proof of Lemma 4.8, it was shown that

$$\widehat{H}(-s)^T H_{\text{err}}(s) = \widehat{H}(-s)^T (H(s) - \widehat{H}(s)) = 0.$$

Since all transfer functions involved are asymptotically stable, in particular this holds for  $s = 0$ . We then have that  $\widehat{H}(0)^T (H(0) - \widehat{H}(0)) = 0$ , i.e.,  $\widehat{H}(0)^T H(0) = \widehat{H}(0)^T \widehat{H}(0)$ . By transposing, we also have  $H(0)^T \widehat{H}(0) = \widehat{H}(0)^T \widehat{H}(0)$ . Therefore,

$$\begin{aligned} H_{\text{err}}(0)^T H_{\text{err}}(0) &= (H(0) - \widehat{H}(0))^T (H(0) - \widehat{H}(0)) \\ &= H(0)^T H(0) - H(0)^T \widehat{H}(0) - \widehat{H}(0)^T H(0) + \widehat{H}(0)^T \widehat{H}(0) \\ &= H(0)^T H(0) - \widehat{H}(0)^T \widehat{H}(0). \end{aligned}$$

By applying Lemma 4.3 to system (4.5), we obtain

$$\begin{aligned} H(0)^T H(0) &= (\mathfrak{B}^T \otimes F^T) (I_n \otimes A - \mathfrak{L} \otimes B)^+ (\mathfrak{L}^2 \otimes I_n) (I_n \otimes A - \mathfrak{L} \otimes B)^+ (\mathfrak{B} \otimes F) \\ &= (\mathfrak{B}^T \otimes F^T) (U \otimes I_n) (I_n \otimes A - \Lambda \otimes B)^+ (\Lambda^2 \otimes I_n) \\ &\quad \cdot (I_n \otimes A - \Lambda \otimes B)^+ (U^T \otimes I_n) (\mathfrak{B} \otimes F) \\ &= (\mathfrak{B}^T U \otimes F^T) \text{diag}(0, \lambda_2^2 (A - \lambda_2 B)^{-2}, \dots, \lambda_n^2 (A - \lambda_n B)^{-2}) (U^T \mathfrak{B} \otimes F) \\ &= (\mathfrak{B}^T U \otimes I_m) \text{diag}\left(0, H_{\lambda_2}(0)^T H_{\lambda_2}(0), \dots, H_{\lambda_n}(0)^T H_{\lambda_n}(0)\right) (U^T \mathfrak{B} \otimes I_m), \end{aligned} \tag{4.34}$$

where  $H_\lambda$  is again given by (4.7). Recall that  $\widehat{\mathfrak{B}} = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T M$  and  $U_1 = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \widehat{U}$ . Now apply Lemma 4.3 to the transfer function (4.33) of the system (4.6):

$$\begin{aligned} \widehat{H}(0)^T \widehat{H}(0) &= \left( \mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes F^T \right) (I_n \otimes A - \overline{\mathfrak{L}} \otimes B)^+ \\ &\quad \cdot \left( (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \mathfrak{P}^T \mathfrak{L}^2 \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes I_n \right) \\ &\quad \cdot (I_n \otimes A - \overline{\mathfrak{L}} \otimes B)^+ \left( (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \mathfrak{P}^T \mathfrak{B} \otimes F \right) \\ &= \left( \mathfrak{B}^T \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \otimes F^T \right) \left( \widehat{U} \otimes I_n \right) \left( I_n \otimes A - \widehat{\Lambda} \otimes B \right)^+ \\ &\quad \cdot \left( \widehat{\Lambda}^2 \otimes I_n \right) \left( I_n \otimes A - \widehat{\Lambda} \otimes B \right)^+ \\ &\quad \cdot \left( \widehat{U}^T \otimes I_n \right) \left( (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}} \mathfrak{P}^T \mathfrak{B} \otimes F \right) \\ &= (\mathfrak{B}^T U_1 \otimes F^T) \text{diag}(0, \lambda_2^2 (A - \lambda_2 B)^{-2}, \dots, \lambda_r^2 (A - \lambda_r B)^{-2}) (U_1^T \mathfrak{B} \otimes F) \\ &= (\mathfrak{B}^T U_1 \otimes I_m) \text{diag}\left(0, H_{\lambda_2}(0)^T H_{\lambda_2}(0), \dots, H_{\lambda_r}(0)^T H_{\lambda_r}(0)\right) (U_1^T \mathfrak{B} \otimes I_m). \end{aligned}$$

Combining the two expressions above, it immediately follows that

$$\begin{aligned} H_{\text{err}}(0)^T H_{\text{err}}(0) &= H(0)^T H(0) - \widehat{H}(0)^T \widehat{H}(0) \\ &= (\mathfrak{B}^T U_2 \otimes I_m) \text{diag}\left(H_{\lambda_{r+1}}(0)^T H_{\lambda_{r+1}}(0), \dots, H_{\lambda_N}(0)^T H_{\lambda_N}(0)\right) \\ &\quad \cdot (U_2^T \mathfrak{B} \otimes I_m). \end{aligned}$$

By taking  $H_{\max, \mathcal{H}_\infty}$  as defined by (4.31) it then follows that

$$\begin{aligned} H_{\text{err}}(0)^\top H_{\text{err}}(0) &\preceq (\mathfrak{B}^\top U_2 \otimes I_m) \text{diag}(H_{\max, \mathcal{H}_\infty}^2 I_m, \dots, H_{\max, \mathcal{H}_\infty}^2 I_m) (U_2^\top \mathfrak{B} \otimes I_m) \\ &= H_{\max, \mathcal{H}_\infty}^2 (\mathfrak{B}^\top U_2 U_2^\top \mathfrak{B} \otimes I_m) \\ &= H_{\max, \mathcal{H}_\infty}^2 (\mathfrak{B}^\top (I_n - U_1 U_1^\top) \mathfrak{B} \otimes I_m) \\ &= H_{\max, \mathcal{H}_\infty}^2 \left( \left( I_m - \mathfrak{B}^\top \mathfrak{P} (\mathfrak{P}^\top \mathfrak{P})^{-1} \mathfrak{P}^\top \mathfrak{B} \right) \otimes I_m \right). \end{aligned}$$

Continuing as in the proof of Theorem 4.12, we find an upper bound for the  $\mathcal{H}_\infty$ -error:

$$\|H_{\text{err}}\|_{\mathcal{H}_\infty}^2 \leq H_{\max, \mathcal{H}_\infty}^2 \lambda_{\max} \left( I_m - \mathfrak{B}^\top \mathfrak{P} (\mathfrak{P}^\top \mathfrak{P})^{-1} \mathfrak{P}^\top \mathfrak{B} \right).$$

To compute an upper bound for the relative  $\mathcal{H}_\infty$ -error, we bound the  $\mathcal{H}_\infty$ -norm of system (4.5) from below. Again, let  $\bar{U}$  be such that  $U = [u_1 \ \bar{U}]$  and let  $H_{\min, \mathcal{H}_\infty}$  be as defined by (4.32). From (4.34) it now follows that

$$\begin{aligned} H(0)^\top H(0) &= (\mathfrak{B}^\top \bar{U} \otimes I_m) \text{diag} \left( H_{\lambda_2}(0)^\top H_{\lambda_2}(0), \dots, H_{\lambda_n}(0)^\top H_{\lambda_n}(0) \right) \left( \bar{U}^\top \mathfrak{B} \otimes I_m \right) \\ &\succeq (\mathfrak{B}^\top \bar{U} \otimes I_m) \text{diag} (H_{\min, \mathcal{H}_\infty}^2 I_m, \dots, H_{\min, \mathcal{H}_\infty}^2 I_m) \left( \bar{U}^\top \mathfrak{B} \otimes I_m \right) \\ &= (H_{\min, \mathcal{H}_\infty})^2 (\mathfrak{B}^\top \bar{U} \bar{U}^\top \mathfrak{B} \otimes I_m) \\ &= (H_{\min, \mathcal{H}_\infty})^2 \left( \mathfrak{B}^\top \left( I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \mathfrak{B} \otimes I_m \right). \end{aligned}$$

Again using Lemma 4.3, we find a lower bound to the  $\mathcal{H}_\infty$ -norm of  $H$ :

$$\|H\|_{\mathcal{H}_\infty}^2 = \lambda_{\max} \left( H(0)^\top H(0) \right) \geq H_{\min, \mathcal{H}_\infty}^2,$$

which concludes the proof of the theorem.  $\square$

## 4.2.7 Towards a priori error bounds for general graph partitions

Up to now, we have only dealt with establishing error bounds for network reduction by clustering using almost equitable partitions of the network graph. Of course, we would also like to obtain error bounds for *arbitrary*, possibly non almost equitable, partitions. In this section, we present some ideas to address this more general problem. We will first study the single integrator case. Subsequently, we will look at the general case.

### 4.2.7.1 The single integrator case

Consider the multi-agent network

$$\begin{aligned} \dot{x}(t) &= -\mathfrak{L}x(t) + \mathfrak{B}u(t), \\ y(t) &= \mathfrak{L}x(t). \end{aligned} \tag{4.35}$$



As before, assume that the underlying graph  $\mathfrak{G}$  is connected. The network is then synchronized. Let  $\pi = \{\mathfrak{C}_1, \mathfrak{C}_2, \dots, \mathfrak{C}_r\}$  be a graph partition, not necessarily an AEP, and let  $\mathfrak{P} \in \mathbb{R}^{n \times r}$  be its characteristic matrix. As before, the reduced order network is taken to be the Petrov-Galerkin projection of (4.35), and is represented by

$$\begin{aligned}\dot{\hat{x}}(t) &= -\hat{\mathfrak{L}}\hat{x}(t) + \hat{\mathfrak{B}}u(t), \\ \hat{y}(t) &= \mathfrak{L}\mathfrak{P}\hat{x}(t),\end{aligned}\tag{4.36}$$

Again, let  $H$  and  $\hat{H}$  be the transfer functions of (4.35) and (4.36), respectively. We will address the problem of obtaining a priori upper bounds for  $\|H - \hat{H}\|_{\mathcal{H}_2}$  and  $\|H - \hat{H}\|_{\mathcal{H}_\infty}$ . We will pursue the following idea: as a first step we will approximate the original Laplacian matrix  $\mathfrak{L}$  (of the original network graph  $\mathfrak{G}$ ) by a new Laplacian matrix, denoted by  $\mathfrak{L}_{\text{AEP}}$  (corresponding to a ‘nearby’ graph  $\mathfrak{G}_{\text{AEP}}$ ) such that the given partition  $\pi$  is an AEP for this new graph  $\mathfrak{G}_{\text{AEP}}$ . This new graph  $\mathfrak{G}_{\text{AEP}}$  defines a new multi-agent system with transfer function  $H_{\text{AEP}}(s) = \mathfrak{L}_{\text{AEP}}(sI_n + \mathfrak{L}_{\text{AEP}})^{-1}\mathfrak{B}$ . The reduced order network of  $H_{\text{AEP}}$  (using the AEP  $\pi$ ) has transfer function  $\hat{H}_{\text{AEP}}(s) = \mathfrak{L}_{\text{AEP}}\mathfrak{P}(sI_r + \hat{\mathfrak{L}}_{\text{AEP}})^{-1}\hat{\mathfrak{B}}$ . Then using the triangle inequality, both for  $p = 2$  and  $p = \infty$ , we have

$$\begin{aligned}\|H - \hat{H}\|_{\mathcal{H}_p} &= \|H - H_{\text{AEP}} + H_{\text{AEP}} - \hat{H}_{\text{AEP}} + \hat{H}_{\text{AEP}} - \hat{H}\|_{\mathcal{H}_p} \\ &\leq \|H - H_{\text{AEP}}\|_{\mathcal{H}_p} + \|H_{\text{AEP}} - \hat{H}_{\text{AEP}}\|_{\mathcal{H}_p} + \|\hat{H}_{\text{AEP}} - \hat{H}\|_{\mathcal{H}_p}.\end{aligned}\tag{4.37}$$

The idea is to obtain a priori upper bounds for all three terms in (4.37). We first propose an approximating Laplacian matrix  $\mathfrak{L}_{\text{AEP}}$ , and subsequently study the problems of establishing upper bounds for the three terms in (4.37) separately.

In the following, denote  $\mathcal{P} := \mathfrak{P}(\mathfrak{P}^T\mathfrak{P})^{-1}\mathfrak{P}^T$ . Note that  $\mathcal{P}$  is the orthogonal projector onto  $\text{im}(\mathfrak{P})$ . As approximation for  $\mathfrak{L}$ , we compute the unique solution to the convex optimization problem

$$\begin{aligned}\underset{\mathfrak{L}_{\text{AEP}}}{\text{minimize}} \quad & \|\mathfrak{L} - \mathfrak{L}_{\text{AEP}}\|_{\text{F}}^2, \\ \text{subject to} \quad & (I_n - \mathcal{P})\mathfrak{L}_{\text{AEP}}\mathfrak{P} = 0, \\ & \mathfrak{L}_{\text{AEP}} = \mathfrak{L}_{\text{AEP}}^T, \\ & \mathfrak{L}_{\text{AEP}} \succcurlyeq 0, \\ & \mathfrak{L}_{\text{AEP}}\mathbf{1}_n = 0.\end{aligned}\tag{4.38}$$

In other words, we want to compute a positive semi-definite matrix  $\mathfrak{L}_{\text{AEP}}$  with row sums equal to zero, and with the property that  $\text{im}(\mathfrak{P})$  is invariant under  $\mathfrak{L}_{\text{AEP}}$  (equivalently, the given partition  $\pi$  is an AEP for the new graph). We will show that such an  $\mathfrak{L}_{\text{AEP}}$  may correspond to an undirected graph *with negative weights*. However, it is constrained to be positive semi-definite, so the results of Sections 4.2.4 to 4.2.6 in this section will remain valid.

**Theorem 4.16:**

The matrix  $\mathfrak{L}_{\text{AEP}} := \mathcal{P}\mathfrak{L}\mathcal{P} + (I_n - \mathcal{P})\mathfrak{L}(I_n - \mathcal{P})$  is the unique solution to the convex optimization problem (4.38). If  $\mathfrak{L}$  corresponds to a connected graph, then, in fact,  $\ker(\mathfrak{L}_{\text{AEP}}) = \text{im}(\mathbf{1}_n)$ .  $\diamond$

*Proof.* Clearly,  $\mathfrak{L}_{\text{AEP}}$  is symmetric and positive semi-definite since  $\mathfrak{L}$  is. Also,  $(I_n - \mathcal{P})\mathfrak{L}_{\text{AEP}}\mathfrak{P} = 0$  since  $(I_n - \mathcal{P})\mathfrak{P} = 0$ . It is also obvious that  $\mathfrak{L}_{\text{AEP}}\mathbf{1}_n = 0$  since  $\mathcal{P}\mathbf{1}_n = \mathbf{1}_n$ . We now show that  $\mathfrak{L}_{\text{AEP}}$  uniquely minimizes the distance to  $\mathfrak{L}$ . Let  $X$  satisfy the constraints and define  $\Delta = \mathfrak{L}_{\text{AEP}} - X$ . Then we have

$$\|\mathfrak{L} - X\|_{\text{F}}^2 = \|\mathfrak{L} - \mathfrak{L}_{\text{AEP}}\|_{\text{F}}^2 + \|\Delta\|_{\text{F}}^2 + 2\text{tr}((\mathfrak{L} - \mathfrak{L}_{\text{AEP}})\Delta).$$

It can be verified that  $\mathfrak{L} - \mathfrak{L}_{\text{AEP}} = (I_n - \mathcal{P})\mathfrak{L}\mathcal{P} + \mathcal{P}\mathfrak{L}(I_n - \mathcal{P})$ . Thus,

$$\text{tr}((\mathfrak{L} - \mathfrak{L}_{\text{AEP}})\Delta) = \text{tr}((I_n - \mathcal{P})\mathfrak{L}\mathcal{P}\Delta) + \text{tr}(\mathcal{P}\mathfrak{L}(I_n - \mathcal{P})\Delta).$$

Now, since both  $X$  and  $\mathfrak{L}_{\text{AEP}}$  satisfy the first constraint, we have  $(I_n - \mathcal{P})\Delta\mathcal{P} = 0$ . Using this we have

$$\text{tr}((I_n - \mathcal{P})\mathfrak{L}\mathcal{P}\Delta) = \text{tr}(\mathcal{P}\Delta(I_n - \mathcal{P})\mathfrak{L}) = \text{tr}(\mathfrak{L}(I_n - \mathcal{P})\Delta\mathcal{P}) = 0.$$

Also,

$$\text{tr}(\mathcal{P}\mathfrak{L}(I_n - \mathcal{P})\Delta) = \text{tr}(\mathfrak{L}(I_n - \mathcal{P})\Delta\mathcal{P}) = 0.$$

Thus, we obtain

$$\|\mathfrak{L} - X\|_{\text{F}}^2 = \|\mathfrak{L} - \mathfrak{L}_{\text{AEP}}\|_{\text{F}}^2 + \|\Delta\|_{\text{F}}^2,$$

from which it follows that  $\|\mathfrak{L} - X\|_{\text{F}}$  is minimal if and only if  $\Delta = 0$ , equivalently,  $X = \mathfrak{L}_{\text{AEP}}$ .

To prove the second statement, let  $x \in \ker(\mathfrak{L}_{\text{AEP}})$ , so  $x^T\mathfrak{L}_{\text{AEP}}x = 0$ . Then both  $x^T\mathcal{P}\mathfrak{L}\mathcal{P}x = 0$  and  $x^T(I_n - \mathcal{P})\mathfrak{L}(I_n - \mathcal{P})x = 0$ . This clearly implies  $\mathfrak{L}\mathcal{P}x = 0$  and  $\mathfrak{L}(I_n - \mathcal{P})x = 0$ . Since  $\mathfrak{L}$  corresponds to a connected graph, we must have  $\mathcal{P}x \in \text{im}(\mathbf{1}_n)$  and  $(I_n - \mathcal{P})x \in \text{im}(\mathbf{1}_n)$ . We conclude that  $x \in \text{im}(\mathbf{1}_n)$ , as desired.  $\square$

As announced above,  $\mathfrak{L}_{\text{AEP}}$  may have positive off-diagonal elements, corresponding to a graph with some of its edge weights being negative. For example, for

$$\mathfrak{L} = \begin{bmatrix} 1 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix}, \quad \mathfrak{P} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix},$$

we have

$$\mathfrak{L}_{\text{AEP}} = \begin{bmatrix} \frac{11}{9} & -\frac{7}{9} & -\frac{1}{9} & 0 & -\frac{1}{3} \\ -\frac{7}{9} & \frac{20}{9} & -\frac{10}{9} & 0 & -\frac{1}{3} \\ -\frac{1}{9} & -\frac{10}{9} & \frac{14}{9} & -\frac{1}{2} & \frac{1}{6} \\ 0 & 0 & -\frac{1}{2} & \frac{3}{2} & -1 \\ -\frac{1}{3} & -\frac{1}{3} & \frac{1}{6} & -1 & \frac{3}{2} \end{bmatrix},$$

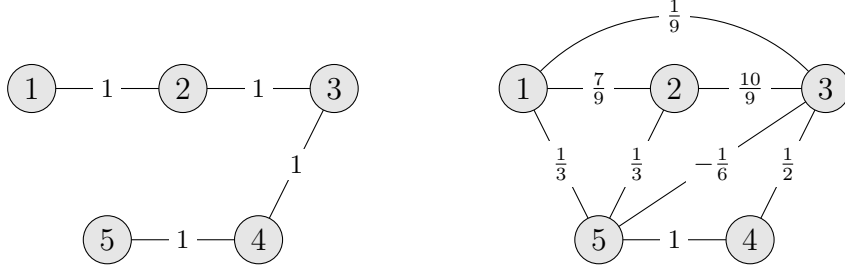


Figure 4.1: A path graph on 5 vertices and its closest graph such that the partition  $\{\{1, 2, 3\}, \{4, 5\}\}$  is almost equitable.

so the edge between vertices 3 and 5 has a negative weight. Figure 4.1 shows the graphs corresponding to  $\mathfrak{L}$  and  $\mathfrak{L}_{\text{AEP}}$ . Although  $\mathfrak{L}_{\text{AEP}}$  is not necessarily a Laplacian matrix with only nonpositive off-diagonal elements, it has all the properties we associate with a Laplacian matrix. Specifically, it can be checked that all results in this section remain valid, since they only depend on the symmetric positive semi-definiteness of the Laplacian matrix.

Using the approximating Laplacian  $\mathfrak{L}_{\text{AEP}} = \mathcal{P}\mathfrak{L}\mathcal{P} + (I_n - \mathcal{P})\mathfrak{L}(I_n - \mathcal{P})$  as above, we will now deal with establishing upper bounds for the three terms in (4.37). We start off with the middle term  $\|H_{\text{AEP}} - \hat{H}_{\text{AEP}}\|_{\mathcal{H}_p}$  in (4.37).

According to Remark 4.7, for  $p = 2$  this term has an upper bound depending on the maximal  $\lambda \in \sigma(\mathfrak{L}_{\text{AEP}}) \setminus \sigma(\hat{\mathfrak{L}}_{\text{AEP}})$ , and on the number of cellmates of the leaders with respect to the partitioning  $\pi$ . For  $p = \infty$ , in Theorem 4.12 this term was expressed in terms of the maximal number of cellmates with respect to the partitioning  $\pi$  (noting that it is equal to 1 in case two or more leaders share the same cluster).

Next, we will take a look at the first and third term in (4.37), i.e.,  $\|H - H_{\text{AEP}}\|_{\mathcal{H}_p}$  and  $\|\hat{H} - \hat{H}_{\text{AEP}}\|_{\mathcal{H}_p}$ . Let us denote  $\Delta\mathfrak{L} = \mathfrak{L} - \mathfrak{L}_{\text{AEP}}$ . We find

$$\begin{aligned}
 H(s) - H_{\text{AEP}}(s) &= \mathfrak{L}(sI_n + \mathfrak{L})^{-1}\mathfrak{B} - \mathfrak{L}_{\text{AEP}}(sI_n + \mathfrak{L}_{\text{AEP}})^{-1}\mathfrak{B} \\
 &= \mathfrak{L}(sI_n + \mathfrak{L})^{-1}\mathfrak{B} \\
 &\quad - \mathfrak{L}_{\text{AEP}} \left[ (sI_n + \mathfrak{L})^{-1} + (sI_n + \mathfrak{L}_{\text{AEP}})^{-1}\Delta\mathfrak{L}(sI_n + \mathfrak{L})^{-1} \right] \mathfrak{B} \\
 &= \mathfrak{L}(sI_n + \mathfrak{L})^{-1}\mathfrak{B} - \mathfrak{L}_{\text{AEP}}(sI_n + \mathfrak{L})^{-1}\mathfrak{B} \\
 &\quad - \mathfrak{L}_{\text{AEP}}(sI_n + \mathfrak{L}_{\text{AEP}})^{-1}\Delta\mathfrak{L}(sI_n + \mathfrak{L})^{-1}\mathfrak{B} \\
 &= \Delta\mathfrak{L}(sI_n + \mathfrak{L})^{-1}\mathfrak{B} - \mathfrak{L}_{\text{AEP}}(sI_n + \mathfrak{L}_{\text{AEP}})^{-1}\Delta\mathfrak{L}(sI_n + \mathfrak{L})^{-1}\mathfrak{B} \\
 &= \left[ I_n - \mathfrak{L}_{\text{AEP}}(sI_n + \mathfrak{L}_{\text{AEP}})^{-1} \right] \Delta\mathfrak{L}(sI_n + \mathfrak{L})^{-1}\mathfrak{B}.
 \end{aligned}$$

Thus, both for  $p = 2$  and  $p = \infty$ , we have

$$\|H - H_{\text{AEP}}\|_{\mathcal{H}_p} \leq \|I_n - \mathfrak{L}_{\text{AEP}}(sI_n + \mathfrak{L}_{\text{AEP}})^{-1}\|_{\mathcal{H}_\infty} \|\Delta\mathfrak{L}(sI_n + \mathfrak{L})^{-1}\mathfrak{B}\|_{\mathcal{H}_p}$$

Using

$$\|I_n - \mathfrak{L}_{\text{AEP}}(sI_n + \mathfrak{L}_{\text{AEP}})^{-1}\|_{\mathcal{H}_\infty} = \|s(sI_n + \mathfrak{L}_{\text{AEP}})^{-1}\|_{\mathcal{H}_\infty} = \sup_{\omega \in \mathbb{R}} \|\mathfrak{L}_{\text{AEP}}(\mathfrak{L}_{\text{AEP}} + \omega^2 I_n)^{-1}\|_2$$

$$= \sup_{\omega \in \mathbb{R}} \frac{|\omega|}{\sqrt{\lambda_{\min}(\omega^2 I_n + \mathcal{L}_{\text{AEP}}^2)}} = 1,$$

we get

$$\|H - H_{\text{AEP}}\|_{\mathcal{H}_p} \leq \|\Delta \mathcal{L}(sI_n + \mathcal{L})^{-1} \mathfrak{B}\|_{\mathcal{H}_p}. \quad (4.39)$$

It is also easily seen that  $\widehat{\mathcal{L}}_{\text{AEP}} = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathcal{L}_{\text{AEP}} \mathfrak{P} = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathcal{L} \mathfrak{P} = \widehat{\mathcal{L}}$  and  $\mathcal{L}_{\text{AEP}} \mathfrak{P} = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \mathcal{L} \mathfrak{P} = \mathfrak{P} \widehat{\mathcal{L}}$ . Therefore,

$$\begin{aligned} \widehat{H}(s) - \widehat{H}_{\text{AEP}}(s) &= \mathcal{L} \mathfrak{P} (sI_n + \widehat{\mathcal{L}})^{-1} \widehat{\mathfrak{B}} - \mathcal{L}_{\text{AEP}} \mathfrak{P} (sI_n + \widehat{\mathcal{L}}_{\text{AEP}})^{-1} \widehat{\mathfrak{B}} \\ &= \mathcal{L} \mathfrak{P} (sI_n + \widehat{\mathcal{L}})^{-1} \widehat{\mathfrak{B}} - \mathfrak{P} \widehat{\mathcal{L}} (sI_n + \widehat{\mathcal{L}})^{-1} \widehat{\mathfrak{B}} \\ &= (\mathcal{L} \mathfrak{P} - \mathfrak{P} \widehat{\mathcal{L}}) (sI_n + \widehat{\mathcal{L}})^{-1} \widehat{\mathfrak{B}}. \end{aligned}$$

Since  $(\mathcal{L} \mathfrak{P} - \mathfrak{P} \widehat{\mathcal{L}})^T (\mathcal{L} \mathfrak{P} - \mathfrak{P} \widehat{\mathcal{L}}) = \mathfrak{P}^T (\Delta \mathcal{L})^2 \mathfrak{P}$ , for  $p = 2$  and  $p = \infty$ , we obtain

$$\|\widehat{H} - \widehat{H}_{\text{AEP}}\|_{\mathcal{H}_p} = \left\| \Delta \mathcal{L} \mathfrak{P} (sI_n + \widehat{\mathcal{L}})^{-1} \widehat{\mathfrak{B}} \right\|_{\mathcal{H}_p}. \quad (4.40)$$

Finally, combining (4.37), (4.39) and (4.40), we get

$$\|H - \widehat{H}\|_{\mathcal{H}_p} \leq \|\Delta \mathcal{L}(sI_n + \mathcal{L})^{-1} \mathfrak{B}\|_{\mathcal{H}_p} + \|H_{\text{AEP}} - \widehat{H}_{\text{AEP}}\|_{\mathcal{H}_p} + \left\| \Delta \mathcal{L} \mathfrak{P} (sI_n + \widehat{\mathcal{L}})^{-1} \widehat{\mathfrak{B}} \right\|_{\mathcal{H}_p}.$$

Thus, both in (4.39) and (4.40) the upper bound involves the difference  $\Delta \mathcal{L} = \mathcal{L} - \mathcal{L}_{\text{AEP}}$  between the original Laplacian and its optimal approximation in the set of Laplacian matrices for which the given partition  $\pi$  is an AEP. In a sense, the difference  $\Delta \mathcal{L}$  measures how far  $\pi$  is away from being an AEP for the original graph  $\mathfrak{G}$ . Obviously,  $\Delta \mathcal{L} = 0$  if and only if  $\pi$  is an AEP for  $\mathfrak{G}$ . In that case only the middle term in (4.37) is present.

#### 4.2.7.2 The general case

In this final subsection, we will put forward some ideas to deal with the case that the agent dynamics is a general linear input-state-output system and the given graph partition  $\pi$ , with characteristic matrix  $\mathfrak{P}$ , is not almost equitable. In this case, the original network is given by (4.5) and the reduced network by (4.6). Their transfer functions are  $H$  and  $\widehat{H}$ , respectively. Let  $\mathcal{L}_{\text{AEP}}$  and  $\widehat{\mathcal{L}}_{\text{AEP}}$  as in the previous subsection and let

$$H_{\text{AEP}}(s) = (\mathcal{L}_{\text{AEP}} \otimes I_n)(sI - I_n \otimes A + \mathcal{L}_{\text{AEP}} \otimes B)^{-1} (\mathfrak{B} \otimes F)$$

and

$$\widehat{H}_{\text{AEP}}(s) = (\mathcal{L}_{\text{AEP}} \mathfrak{P} \otimes I_n) (sI - I_n \otimes A + \widehat{\mathcal{L}}_{\text{AEP}} \otimes B)^{-1} (\widehat{\mathfrak{B}} \otimes F).$$

As before, we assume that (4.5) is synchronized, so  $H$  is asymptotically stable. However, since the partition  $\pi$  is no longer assumed to be an AEP, the reduced transfer function  $\widehat{H}$  need not be asymptotically stable anymore. Also,  $H_{\text{AEP}}$  and  $\widehat{H}_{\text{AEP}}$  need not be asymptotically stable. We will now first study under what conditions these are asymptotically stable. First note that, using Proposition 2.51,  $\widehat{H}$  is asymptotically stable if and only if  $A - \widehat{\lambda}B$  is Hurwitz for all nonzero eigenvalues  $\widehat{\lambda}$  of  $\widehat{\mathfrak{L}}$ . Moreover,  $H_{\text{AEP}}$  and  $\widehat{H}_{\text{AEP}}$  are asymptotically stable if and only if  $A - \lambda B$  is Hurwitz for all nonzero eigenvalues  $\lambda$  of  $\mathfrak{L}_{\text{AEP}}$ . In the following, let  $\lambda_{\min}(\mathfrak{L})$  and  $\lambda_{\max}(\mathfrak{L})$  denote the smallest nonzero and largest eigenvalue of  $\mathfrak{L}$ , respectively. We have the following lemma about the location of the nonzero eigenvalues of  $\widehat{\mathfrak{L}}$  and  $\mathfrak{L}_{\text{AEP}}$ .

**Lemma 4.17:**

All nonzero eigenvalues of  $\widehat{\mathfrak{L}}$  and of  $\mathfrak{L}_{\text{AEP}}$  lie in the closed interval  $[\lambda_{\min}(\mathfrak{L}), \lambda_{\max}(\mathfrak{L})]$ .  $\diamond$

*Proof.* The claim about the eigenvalues of  $\widehat{\mathfrak{L}}$  follows from Theorem 2.6. Next, note that  $\mathcal{P} = Q_1 Q_1^T$ , with  $Q_1 = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-\frac{1}{2}}$ . Since the columns of  $Q_1$  are orthonormal, there exists a matrix  $Q_2 \in \mathbb{R}^{n \times (n-r)}$  such that  $[Q_1 \ Q_2]$  is an orthogonal matrix. Then, we have  $I_n - \mathcal{P} = Q_2 Q_2^T$  and we find

$$\begin{aligned} \mathfrak{L}_{\text{AEP}} &= \mathcal{P} \mathfrak{L} \mathcal{P} + (I_n - \mathcal{P}) \mathfrak{L} (I_n - \mathcal{P}) \\ &= Q_1 Q_1^T \mathfrak{L} Q_1 Q_1^T + Q_2 Q_2^T \mathfrak{L} Q_2 Q_2^T \\ &= [Q_1 \ Q_2] \begin{bmatrix} Q_1^T \mathfrak{L} Q_1 & 0 \\ 0 & Q_2^T \mathfrak{L} Q_2 \end{bmatrix} \begin{bmatrix} Q_1^T \\ Q_2^T \end{bmatrix}. \end{aligned}$$

It follows that  $\sigma(\mathfrak{L}_{\text{AEP}}) = \sigma(Q_1^T \mathfrak{L} Q_1) \cup \sigma(Q_2^T \mathfrak{L} Q_2)$ . By Theorem 2.6, both the eigenvalues of  $Q_1^T \mathfrak{L} Q_1$  and  $Q_2^T \mathfrak{L} Q_2$  are interlaced with the eigenvalues of  $\mathfrak{L}$ , so in particular we have that all eigenvalues  $\lambda$  of  $\mathfrak{L}_{\text{AEP}}$  satisfy  $\lambda \leq \lambda_{\max}(\mathfrak{L})$ . In order to prove the lower bound, note that  $Q_1^T \mathfrak{L} Q_1$  is similar to  $\widehat{\mathfrak{L}}$ , for which we know that its nonzero eigenvalues are between the nonzero eigenvalues of  $\mathfrak{L}$ . As for the eigenvalues of  $Q_2^T \mathfrak{L} Q_2$ , note that  $\|Q_2 x\|_2 = \|x\|_2$  for all  $x \in \mathbb{R}^{n-r}$  and  $\mathbf{1}^T Q_2 = 0$  (since  $Q_1 (\mathfrak{P}^T \mathfrak{P})^{\frac{1}{2}} \mathbf{1} = \mathbf{1}$ ). Thus, we find

$$\min_{\|x\|_2=1} x^T Q_2^T \mathfrak{L} Q_2 x \geq \min_{\substack{\mathbf{1}^T y=0 \\ \|y\|_2=1}} y^T \mathfrak{L} y = \lambda_{\min}(\mathfrak{L}).$$

Therefore, the smallest eigenvalue of  $Q_2^T \mathfrak{L} Q_2$  is larger than the smallest positive eigenvalue of  $\mathfrak{L}$ . We conclude that indeed  $\lambda \geq \lambda_{\min}(\mathfrak{L})$  for all nonzero eigenvalues  $\lambda$  of  $\mathfrak{L}_{\text{AEP}}$ .  $\square$

Using this lemma, we see that a sufficient condition for  $\widehat{H}$ ,  $H_{\text{AEP}}$ , and  $\widehat{H}_{\text{AEP}}$  to be asymptotically stable is that for each  $\lambda \in [\lambda_{\min}(\mathfrak{L}), \lambda_{\max}(\mathfrak{L})]$ , the strict Lyapunov inequality

$$(A - \lambda B)X + X(A - \lambda B)^T \prec 0$$

has a positive definite solution  $X$ . This sufficient condition can be checked by verifying solvability of a *single* linear matrix inequality, whose size does not depend on the number of agents, see [OP05]. After having checked this, it would then remain to establish upper bounds for the first and third term in (4.37). This can be done in an analogous way as in the previous subsection. Specifically, it can be shown that for  $p = 2$  and  $p = \infty$  we have

$$\|H - H_{\text{AEP}}\|_{\mathcal{H}_p} \leq \left(1 + \|(\mathfrak{L}_{\text{AEP}} \otimes I_n)(sI - I_n \otimes A + \mathfrak{L}_{\text{AEP}} \otimes B)^{-1}(I_n \otimes B)\|_{\mathcal{H}_\infty}\right) \cdot \|(\Delta \mathfrak{L} \otimes I_n)(sI - I_n \otimes A + \mathfrak{L} \otimes B)^{-1}(\mathfrak{B} \otimes F)\|_{\mathcal{H}_p}$$

and

$$\|\widehat{H} - \widehat{H}_{\text{AEP}}\|_{\mathcal{H}_p} = \left\| (\Delta \mathfrak{L} \mathfrak{P} \otimes I_n) \left( sI - I_n \otimes A + \widehat{\mathfrak{L}} \otimes B \right)^{-1} (\widehat{\mathfrak{B}} \otimes F) \right\|_{\mathcal{H}_p}.$$

#### 4.2.8 Numerical examples

To illustrate the error bounds we have established in this section, consider the graph with 10 vertices taken from [MTC14], as shown in Figure 3.1. Its Laplacian matrix is

$$\mathfrak{L} = \begin{bmatrix} 5 & 0 & 0 & 0 & 0 & -5 & 0 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 & -3 & -2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 6 & -1 & -2 & -3 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 6 & -5 & 0 & 0 & 0 & 0 & 0 \\ 0 & -3 & -2 & -5 & 25 & -2 & -6 & -7 & 0 & 0 \\ -5 & -2 & -3 & 0 & -2 & 25 & -6 & -7 & 0 & 0 \\ 0 & 0 & 0 & 0 & -6 & -6 & 15 & -1 & -1 & -1 \\ 0 & 0 & 0 & 0 & -7 & -7 & -1 & 15 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -1 & 0 & 0 & 1 \end{bmatrix},$$

with spectrum (rounded to three significant digits)

$$\sigma(\mathfrak{L}) \approx \{0, 1, 1.08, 4.14, 5, 6.7, 8.36, 16.1, 28.2, 33.5\}.$$

First, we illustrate the  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  error bounds from Theorems 4.5 and 4.14. We take  $\pi = \{\{1, 2, 3, 4\}, \{5, 6\}, \{7\}, \{8\}, \{9, 10\}\}$  and

$$A = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}, \quad B = F = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

Note that, indeed,  $\pi$  is an AEP. Also, in order to satisfy the assumptions of Theorem 4.14, we have taken  $A$  and  $B$  symmetric. Note that  $A - \lambda B$  is Hurwitz for all nonzero eigenvalues  $\lambda$  of the Laplacian matrix  $\mathfrak{L}$ . Therefore, the multi-agent system is synchronized. It remains to choose the set of leaders  $\mathfrak{V}_L$ . For demonstration, we compute the  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  upper bounds and the true errors for all possible choices of  $\mathfrak{V}_L$ . Since the sets of leaders are nonempty subsets of  $\mathfrak{V}$ , it follows that there are  $2^{10} - 1 = 1023$  possible sets of leaders. Figure 4.2 shows all the ratios of upper bounds and corresponding true errors, where we define  $\frac{0}{0} := 1$ . We see that in this example,

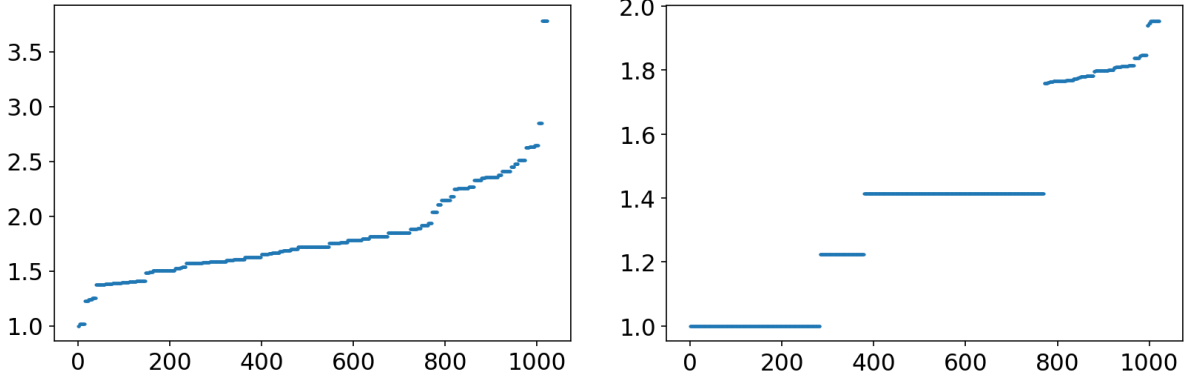


Figure 4.2: Ratios of  $\mathcal{H}_2$  (left) and  $\mathcal{H}_\infty$  (right) upper bounds and corresponding true errors, for a fixed almost equitable partition and all possible sets of leaders. In both figures, the sets of leaders are sorted such that the ratio is increasing (in particular, the ordering of the sets of leaders is not the same).

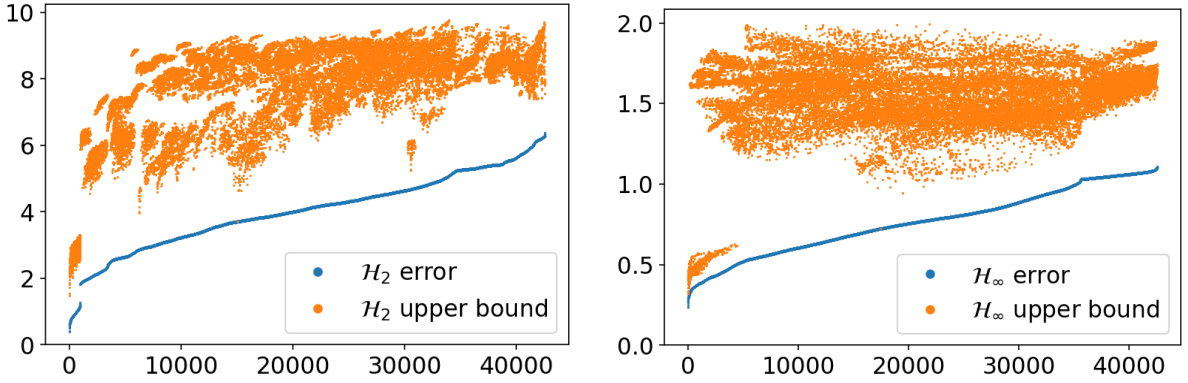


Figure 4.3: True  $\mathcal{H}_2$  (left) and  $\mathcal{H}_\infty$  (right) errors and upper bounds, for a fixed set of leaders and all partitions with five clusters. In each figure, partitions were sorted such that the true errors are increasing.

all true errors and upper bounds are within one order of magnitude, and that in most cases the ratio is below 2.

Next, we compare the true errors with the triangle inequality-based error bounds from (4.37) for a fixed set of leaders and all possible partitions consisting of five clusters. For the set of leaders, we take  $\mathfrak{V}_L = \{6, 7\}$ , as was also used in [MTC14]. With this choice of leaders, the systems norms are  $\|S\|_{\mathcal{H}_2} \approx 6.4$  and  $\|S\|_{\mathcal{H}_\infty} \approx 1.03$  (rounded to three significant digits). Figure 4.3 shows true errors and upper bounds for all partitions of  $\mathfrak{V}$  with five clusters (there are 42 525 such partitions). We observe that the upper bounds vary significantly as the true error increases, but the ratio is still less than one order of magnitude. Additionally, we notice that partitions giving small  $\mathcal{H}_2$  errors give smaller upper bounds, as seen more clearly in the left subfigure of Figure 4.4. Furthermore, we observe a jump after the 966th partition. In fact, the 966

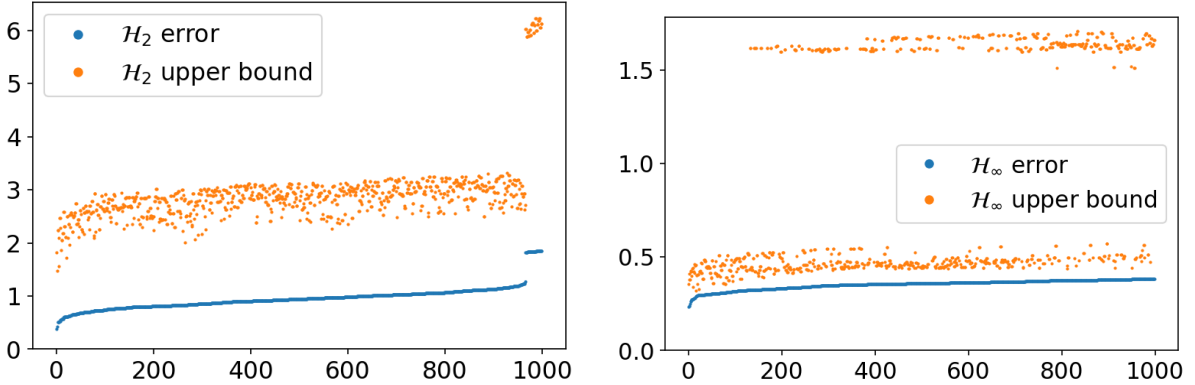


Figure 4.4: First 1000 true errors and upper bounds from Figure 4.3.

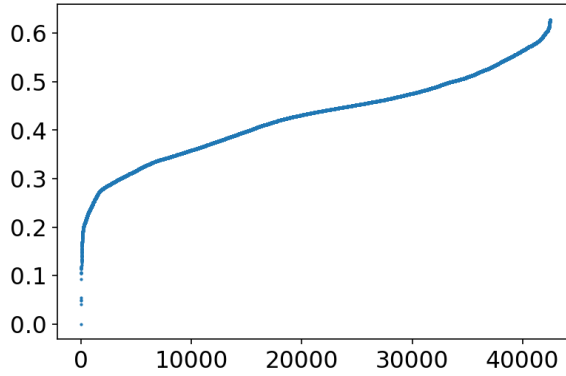


Figure 4.5: Relative error of  $\mathcal{L}$  by  $\mathcal{L}_{\text{AEP}}$  in Frobenius norm for all partitions with five clusters. The partitions are ordered such that the errors are increasing.

partitions giving the smallest  $\mathcal{H}_2$  error are all those partitions where the leaders are the only members in their cluster. For the  $\mathcal{H}_\infty$  error this is not the case, i.e., there are partitions with leaders sharing a cluster with more agents that give a smaller  $\mathcal{H}_\infty$  error than a partition with leaders not sharing a cluster. On the other hand, partitions with the smallest  $\mathcal{H}_2$  or  $\mathcal{H}_\infty$  upper bound are close to the optimal true error.

In the following, we also compute the errors  $\|\mathcal{L} - \mathcal{L}_{\text{AEP}}\|_{\text{F}}$  for all partitions with five clusters. Figure 4.5 shows the relative approximation errors  $\frac{\|\mathcal{L} - \mathcal{L}_{\text{AEP}}\|_{\text{F}}}{\|\mathcal{L}\|_{\text{F}}}$ . We see that only a few (six, to be precise) partitions give a relative error less than 0.1. Irrespective of this, a small triangle inequality-based error bound (4.37) seems to indicate good partitions.

Finally, we compare the bound (4.37) with those from Ishizaki et al. [IKIA14, IKG<sup>+</sup>15, IKI16a]. There are also error bounds developed in [CKS16] and [BSJ16], but they depend on the proposed MOR methods and cannot be evaluated for an arbitrary partition. The  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  error bounds from Ishizaki et al. are based on the decomposition (see



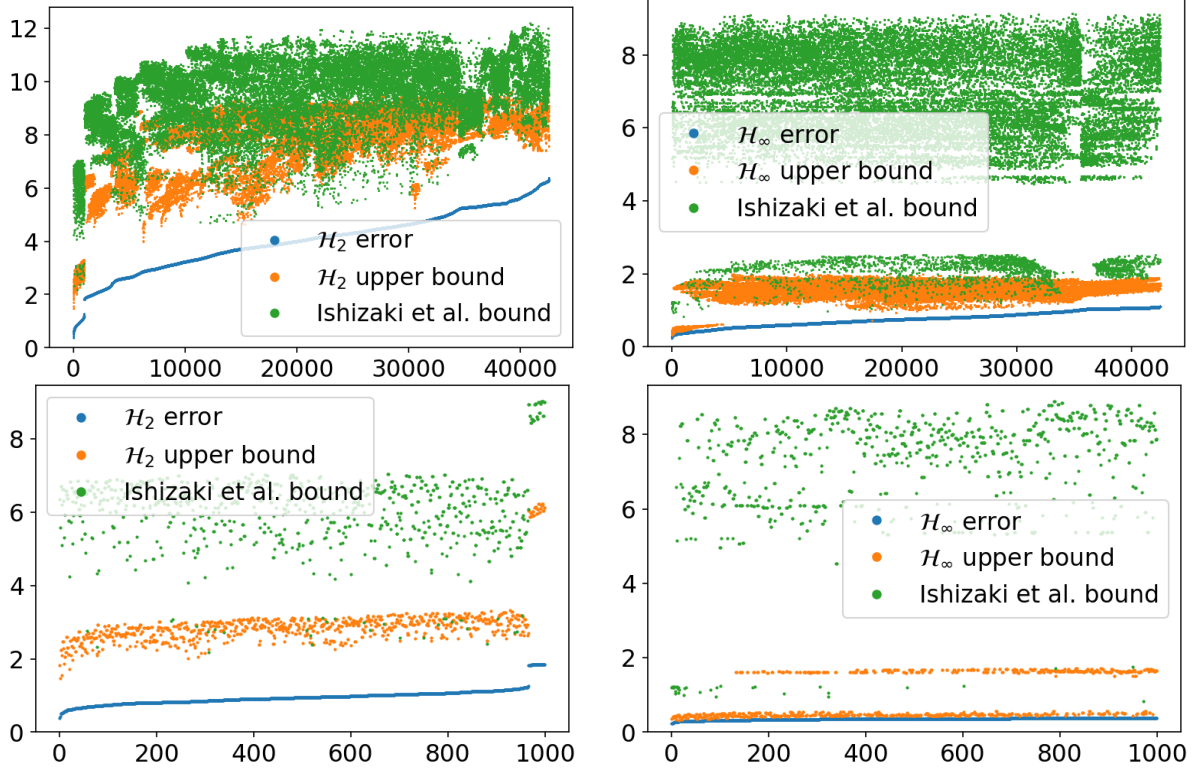


Figure 4.6: Comparison with error bounds from Ishizaki et al. [IKIA14, IKG<sup>+</sup>15, IKI16a]. The first column shows the  $\mathcal{H}_2$  errors and bounds, the second column the  $\mathcal{H}_\infty$  errors and bounds. The first row contains values for all partitions with five clusters, the second row only the first 1000 best ones.

equation (31) in [IKIA14], (20) in [II15], or (17) in [IKI16a])

$$H(s) - \hat{H}(s) = \Xi(s)QQ^T X(s),$$

where

$$X(s) = (sI - I_n \otimes A + P^T \mathcal{L}P \otimes B)^{-1} (\mathfrak{B} \otimes F),$$

$$\Xi(s) = (\mathcal{L}P \otimes I_n) (sI - I_n \otimes A + P^T \mathcal{L}P \otimes B)^{-1} (P^T \otimes A - P^T \mathcal{L} \otimes B) + \mathcal{L} \otimes I_n,$$

$P = \mathfrak{P} (\mathfrak{P}^T \mathfrak{P})^{-1}$ , and  $Q$  is such that  $[P \ Q]$  is orthogonal. The error bounds are then

$$\|H - \hat{H}\|_{\mathcal{H}_p} \leq \|\Xi\|_{\mathcal{H}_\infty} \|QQ^T X\|_{\mathcal{H}_p} = \|\Xi\|_{\mathcal{H}_\infty} \|Q^T X\|_{\mathcal{H}_p},$$

for  $p = 2$  and  $p = \infty$ . Figure 4.6 shows the comparison between these bounds, the triangle inequality-based bound (4.37), and the true errors. In this example, our bounds are, for most partitions, lower than those from Ishizaki et al. Yet, they do share some qualitative properties: both vary significantly as the true error increases and those partitions with the small bounds are close to the optimal.

## 4.3 Exact clustering-based model order reduction for nonlinear power systems

### 4.3.1 Introduction

A power system is a network of electrical generators, loads, and their associated control elements. Each of these components may be thought of as vertex of a graph, while the transmission lines connecting them can be regarded as the edges of the graph. The vertices are modeled by physical laws that typically lead to a set of differential equations. These differential equations are coupled to each other across the edges. One question that has been of interest to power engineers over many years is how do the graph-theoretic properties of these types of electrical networks impact system-theoretic properties of the grid model [AA13].

Here, we study synchronization properties of power systems (see [DB14] for an overview) using graph-theoretic tools. Specifically, we show relations to graph symmetry and equitable partitions [RJME09], extending the work in [IKI16b] for linear systems to nonlinear power systems. Additionally, based on our results about synchronization, we propose a structure-preserving, clustering-based MOR framework for nonlinear power systems. Further, we show that for certain partitions this reduction is exact. In general, the dynamics of the reduced system can be used to approximate the dynamics of the original power system.

The motivation for clustering, in addition to reducing simulation time, is the possibility to simulate or control only a certain part of the grid, or a certain phenomenon that happens only over a certain time-scale. Some recent work on clustering of linear network systems can be found in [IKIA14, IKG<sup>+</sup>15, MGB15, CKS16, XC16, CKS17].

In Section 4.3.2, we describe the system we analyze. Next, we introduce synchronization for a pair of generators and prove necessary and sufficient conditions in Section 4.3.3. In Section 4.3.4, we continue in a similar way with two notions of synchronization with respect to a partition. We discuss clustering-based MOR in Section 4.3.5. Finally, we demonstrate our results in Section 4.3.6.

### 4.3.2 System description

We use the power system example in Figure 4.7 to introduce the type of system we analyze and to illustrate our results. As in the example in Figure 4.7, we consider power systems consisting of generators and buses, where each generator is connected to exactly one bus and buses can be classified into *generator buses* (those connected to one generator and some buses) and *non-generator buses* (those connected only to other buses). We follow the classical model of a synchronous generator [Kun94], which means that the generators' voltage amplitude is constant over time  $t$ .

Let  $\mathcal{G} := \{1, 2, \dots, n\}$  and  $\bar{\mathcal{G}} := \{n + 1, n + 2, \dots, n + \bar{n}\}$  denote the label sets of generator and non-generator buses. In the example in Figure 4.7, we have  $n = 5$  and

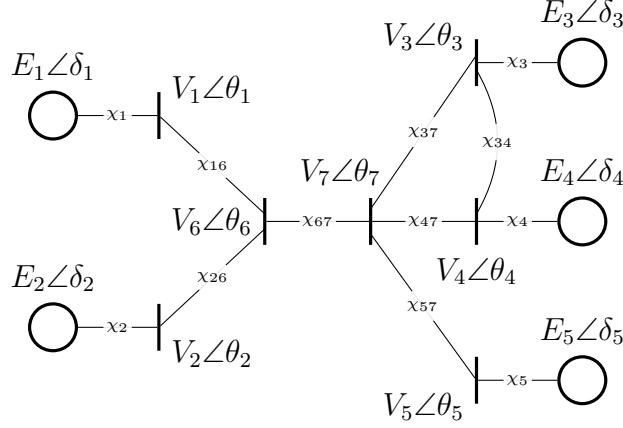


Figure 4.7: Power system consisting of generators (circles) and buses (vertical bars), where the  $i$ th generator is only connected to the  $i$ th bus. See Table 4.1 for the notation.

$\bar{n} = 2$ . The vector of currents from generators to generator buses is given as

$$\mathbf{I}_g(t) = \frac{1}{\mathbf{l}} L_D (\mathbf{E}_g(t) - \mathbf{V}_g(t)), \quad (4.41)$$

where the vectors of voltages of generators and generator buses are denoted as

$$\begin{aligned} \mathbf{E}_g(t) &:= [E_i(\cos \delta_i(t) + \mathbf{l} \sin \delta_i(t))]_{i \in \mathcal{G}} \in \mathbb{C}^n, \\ \mathbf{V}_g(t) &:= [V_i(t)(\cos \theta_i(t) + \mathbf{l} \sin \theta_i(t))]_{i \in \mathcal{G}} \in \mathbb{C}^n, \end{aligned}$$

and  $L_D$  is a positive diagonal reactance matrix given as

$$L_D := \text{diag}([\chi_i^{-1}]_{i \in \mathcal{G}}),$$

where  $\chi_i$  is the reactance between the  $i$ th generator and its bus (see Figure 4.7). We assume the generator voltage amplitudes  $E_i$  and reactances  $\chi_i$  are given constants. Additionally, we assume the line resistances to be negligible.

The relation between the currents and voltages is given as

$$\begin{bmatrix} \mathbf{I}_g(t) \\ 0 \end{bmatrix} = \frac{1}{\mathbf{l}} \begin{bmatrix} \mathfrak{L}_{11} & \mathfrak{L}_{12} \\ \mathfrak{L}_{12}^T & \mathfrak{L}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{V}_g(t) \\ \mathbf{V}_{\bar{\mathcal{G}}}(t) \end{bmatrix}, \quad (4.42)$$

where the voltage vector of non-generator buses is denoted as

$$\mathbf{V}_{\bar{\mathcal{G}}}(t) := [V_i(t)(\cos \theta_i(t) + \mathbf{l} \sin \theta_i(t))]_{i \in \bar{\mathcal{G}}} \in \mathbb{C}^{\bar{n}}$$

and  $\mathfrak{L} = [\mathfrak{L}_{ij}] \in \mathbb{R}^{(n+\bar{n}) \times (n+\bar{n})}$  denotes the weighted graph Laplacian of the reactance network. In particular, the  $(i, j)$ -th element of  $\mathfrak{L}$  is  $-\chi_{ij}^{-1}$  if the  $i$ th and  $j$ th buses are

Table 4.1: Notation

Symbol	Description
$[a_i]_{i \in S}$	vector $(a_{i_1}, a_{i_2}, \dots, a_{i_n})$ , if $S = \{i_1, i_2, \dots, i_n\}$
sin, cos	functions applied element-wise to a vector or a matrix
$\mathcal{G}$	label set of generator buses
$\bar{\mathcal{G}}$	label set of non-generator buses
$\mathbf{E}_{\mathcal{G}}(t)$	voltages of the generators at time $t$
$E_i$	voltage amplitude of the $i$ th generator
$\delta_i(t)$	voltage phase of the $i$ th generator at time $t$
$\mathbf{V}_{\mathcal{G}}(t)$	voltages of the generator buses at time $t$
$\mathbf{V}_{\bar{\mathcal{G}}}(t)$	voltages of the non-generator buses at time $t$
$V_i(t)$	voltage amplitude of the $i$ th bus at time $t$
$\theta_i(t)$	voltage phase of the $i$ th bus at time $t$
$\mathbf{I}_{\mathcal{G}}(t)$	currents from generators to generator buses at time $t$
$\chi_i$	reactance between the $i$ th generator and its bus
$\chi_{ij}$	reactance between the $i$ th and $j$ th bus
$L_D$	reactance matrix, $\text{diag}([\chi_i^{-1}]_{i \in \mathcal{G}})$
$\mathfrak{L}$	$[\mathfrak{L}_{ij}]_{i,j \in \{1,2\}}$ , weighted graph Laplacian of the reactance network
$\delta(t)$	$[\delta_i(t)]_{i \in \mathcal{G}}$
$M$	diagonal matrix of inertias $M_i$ of the generators
$D$	diagonal matrix of dissipatives $D_i$ of the generators
$f$	vector of powers $f_i$ to the generators
$X$	$(L_D + \mathfrak{L}_{11} - \mathfrak{L}_{12} \mathfrak{L}_{22}^{-1} \mathfrak{L}_{12}^T)^{-1} L_D$
$\Gamma$	$L_D (L_D + \mathfrak{L}_{11} - \mathfrak{L}_{12} \mathfrak{L}_{22}^{-1} \mathfrak{L}_{12}^T)^{-1} L_D$
$\gamma_{ij}$	$[\Gamma]_{ij}^{-1}$
$E$	$[E_i]_{i \in \mathcal{G}}$
$V_{\mathcal{G}}(t)$	$[V_i(t)]_{i \in \mathcal{G}}$
$\theta_{\mathcal{G}}(t)$	$[\theta_i(t)]_{i \in \mathcal{G}}$
$\mathcal{X}_{ij}$	subspace of synchronism $\{x \in \mathbb{R}^n : x_i = x_j\}$
$\Pi_{ij}$	permutation matrix that swaps $i$ th and $j$ th components
$\mathcal{S}_{ij}$	set of symmetrical matrices $\{A \in \mathbb{R}^{n \times n} : A \Pi_{ij} = \Pi_{ij} A\}$
$\mathcal{X}_{\text{cl}}$	$\bigcap_{\ell \in \bar{\mathcal{G}}} \bigcap_{i,j \in \mathcal{C}_{\ell}} \mathcal{X}_{ij}$
$\mathcal{S}_{\text{cl}}$	$\bigcap_{\ell \in \bar{\mathcal{G}}} \bigcap_{i,j \in \mathcal{C}_{\ell}} \mathcal{S}_{ij}$

connected (see Figure 4.7) and the  $i$ th diagonal element is  $\sum_{j \neq i} \chi_{ij}^{-1}$ . In the following, we assume that the reactance network is connected, i.e.,  $\mathfrak{L}$  is irreducible. This assumption can be made without loss of generality because the same arguments can be applied to each connected component. For the example in Figure 4.7 with  $\chi_{ij} = 1$  for all  $i, j$ , we have

$$\mathfrak{L} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 0 & 2 & -1 & 0 & 0 & -1 \\ 0 & 0 & -1 & 2 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 1 & 0 & -1 \\ -1 & -1 & 0 & 0 & 0 & 3 & -1 \\ 0 & 0 & -1 & -1 & -1 & -1 & 4 \end{bmatrix}.$$

The dynamics of generators is given by

$$M\ddot{\delta}(t) + D\dot{\delta}(t) = f - \left[ \frac{E_i V_i(t)}{\chi_i} \sin(\delta_i(t) - \theta_i(t)) \right]_{i \in \mathcal{G}}, \quad (4.43a)$$

with voltage phases  $\delta(t) := [\delta_i(t)]_{i \in \mathcal{G}}$ , inertia constants  $M := \text{diag}([M_i]_{i \in \mathcal{G}})$ ,  $M_i > 0$ , damping constants  $D := \text{diag}([D_i]_{i \in \mathcal{G}})$ ,  $D_i \geq 0$ , and input powers  $f \in \mathbb{R}^n$  [Kun94]. Eliminating  $\mathbf{I}_{\mathcal{G}}(t)$  from (4.41) and (4.42), we obtain

$$\begin{bmatrix} L_D (\mathbf{E}_{\mathcal{G}}(t) - \mathbf{V}_{\mathcal{G}}(t)) \\ 0 \end{bmatrix} = \begin{bmatrix} \mathfrak{L}_{11} & \mathfrak{L}_{12} \\ \mathfrak{L}_{12}^T & \mathfrak{L}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{V}_{\mathcal{G}}(t) \\ \mathbf{V}_{\overline{\mathcal{G}}}(t) \end{bmatrix}, \quad (4.43b)$$

The set of equations (4.43) forms a differential-algebraic system. We can remove the algebraic constraints to find an equivalent set of differential equations using Kron reduction [Kro39]. First, from (4.43b), we find

$$\begin{aligned} \mathbf{V}_{\overline{\mathcal{G}}}(t) &= -\mathfrak{L}_{22}^{-1} \mathfrak{L}_{12}^T \mathbf{V}_{\mathcal{G}}(t), \\ \mathbf{V}_{\mathcal{G}}(t) &= X \mathbf{E}_{\mathcal{G}}(t), \end{aligned} \quad (4.44)$$

where

$$X := (L_D + \mathfrak{L}_{11} - \mathfrak{L}_{12} \mathfrak{L}_{22}^{-1} \mathfrak{L}_{12}^T)^{-1} L_D. \quad (4.45)$$

It follows that

$$\Gamma := L_D (L_D + \mathfrak{L}_{11} - \mathfrak{L}_{12} \mathfrak{L}_{22}^{-1} \mathfrak{L}_{12}^T)^{-1} L_D = L_D X$$

is a positive definite matrix with positive elements, since  $L_D + \mathfrak{L}_{11} - \mathfrak{L}_{12} \mathfrak{L}_{22}^{-1} \mathfrak{L}_{12}^T$  is positive definite and an irreducible  $M$ -matrix [BP94, pp. 141]. We denote its elements by  $\gamma_{ij}^{-1} := [\Gamma]_{ij}$ . Then, multiplying (4.44) from the left by  $L_D$ , we find

$$\begin{aligned} \left[ \frac{V_i(t)}{\chi_i} \cos \theta_i(t) \right]_{i \in \mathcal{G}} &= \Gamma [E_i \cos \delta_i(t)]_{i \in \mathcal{G}}, \\ \left[ \frac{V_i(t)}{\chi_i} \sin \theta_i(t) \right]_{i \in \mathcal{G}} &= \Gamma [E_i \sin \delta_i(t)]_{i \in \mathcal{G}}, \end{aligned}$$

which, when inserted in (4.43a), gives us

$$\begin{aligned} M\ddot{\delta}(t) + D\dot{\delta}(t) &= f - \left[ \frac{E_i V_i(t)}{\chi_i} (\sin \delta_i(t) \cos \theta_i(t) - \cos \delta_i(t) \sin \theta_i(t)) \right]_{i \in \mathcal{G}} \\ &= f - \left( \text{diag}([E_i \sin \delta_i(t)]_{i \in \mathcal{G}}) \left[ \frac{V_i(t)}{\chi_i} \cos \theta_i(t) \right]_{i \in \mathcal{G}} \right. \\ &\quad \left. - \text{diag}([E_i \cos \delta_i(t)]_{i \in \mathcal{G}}) \left[ \frac{V_i(t)}{\chi_i} \sin \theta_i(t) \right]_{i \in \mathcal{G}} \right) \\ &= f - \left( \text{diag}([E_i \sin \delta_i(t)]_{i \in \mathcal{G}}) \Gamma [E_i \cos \delta_i(t)]_{i \in \mathcal{G}} \right. \\ &\quad \left. - \text{diag}([E_i \cos \delta_i(t)]_{i \in \mathcal{G}}) \Gamma [E_i \sin \delta_i(t)]_{i \in \mathcal{G}} \right). \end{aligned}$$

Thus, now by using  $\sin \delta_i(t) \cos \delta_j(t) - \cos \delta_i(t) \sin \delta_j(t) = \sin(\delta_i(t) - \delta_j(t))$ , the Kron-reduced system of (4.43) is given as

$$M_i \ddot{\delta}_i(t) + D_i \dot{\delta}_i(t) = f_i - \sum_{k=1}^n \frac{E_i E_k}{\gamma_{ik}} \sin(\delta_i(t) - \delta_k(t)), \quad (4.46a)$$

with generator buses' voltages and phases satisfying

$$L_D \mathbf{V}_g(t) = \Gamma \mathbf{E}_g(t). \quad (4.46b)$$

Denoting  $E := [E_i]_{i \in \mathcal{G}}$ ,  $V_g(t) := [V_i(t)]_{i \in \mathcal{G}}$ , and  $\theta_g(t) := [\theta_i(t)]_{i \in \mathcal{G}}$ , we can write (4.43a) and (4.46a) more compactly as

$$M \ddot{\delta}(t) + D \dot{\delta}(t) = f - L_D (E \circ V_g(t) \circ \sin(\delta(t) - \theta_g(t))), \quad (4.47)$$

and

$$M \ddot{\delta}(t) + D \dot{\delta}(t) = f - \left( \Gamma \circ E E^T \circ \sin\left(\delta(t) \mathbf{1}_n^T - \mathbf{1}_n \delta(t)^T\right) \right) \mathbf{1}_n.$$

### 4.3.3 Synchronization of generator pair

Let us denote the subspace of the synchronism between the  $i$ th and  $j$ th elements by

$$\mathcal{X}_{ij} := \{x \in \mathbb{R}^n : x_i = x_j\}.$$

In this notation, we introduce the following notion of synchronism for the power system (4.43).

**Definition 4.18:**

Consider the power system (4.43). The  $i$ th and  $j$ th generators are said to be *synchronized* if

$$\delta(t) \in \mathcal{X}_{ij} \text{ and } \mathbf{V}_g(t) \in \mathcal{X}_{ij}, \text{ for all } t \geq 0$$

and for any initial condition  $\delta(0), \dot{\delta}(0) \in \mathcal{X}_{ij}$ .  $\diamond$

To characterize this generator synchronism in an algebraic manner, let us define a set of symmetrical matrices with respect to the permutation of the  $i$ th and  $j$ th columns and rows by

$$\mathcal{S}_{ij} := \{A \in \mathbb{R}^{n \times n} : A \Pi_{ij} = \Pi_{ij} A\}, \quad (4.48)$$

where  $\Pi_{ij}$  denotes the permutation matrix associated with the  $i$ th and  $j$ th elements, i.e., all diagonal elements of  $\Pi_{ij}$  other than the  $i$ th and  $j$ th elements are 1, the  $(i, j)$ -th and  $(j, i)$ -th elements are 1, and the others are zero. Note that  $\mathcal{S}_{ij}$  is not the set of usual symmetric (Hermitian) matrices; the condition in (4.48) represents the invariance with respect to the permutation of the  $i$ th and  $j$ th columns and rows, i.e.,  $\Pi_{ij}^T A \Pi_{ij} = A$ . The following lemma gives necessary and sufficient conditions for a symmetric matrix to be symmetrical.

**Lemma 4.19:**

Let  $A \in \mathbb{R}^{n \times n}$  be a symmetric matrix and  $i, j \in \{1, 2, \dots, n\}$  such that  $i \neq j$ . Then  $A \in \mathcal{S}_{ij}$  if and only if  $a_{ii} = a_{jj}$  and  $a_{ik} = a_{jk}$  for all  $k \neq i, j$ .  $\diamond$

*Proof.* From the definition, it can be seen that  $A \in \mathcal{S}_{ij}$  is equivalent to  $a_{ii} = a_{jj}$ ,  $a_{ij} = a_{ji}$ ,  $a_{ik} = a_{jk}$ , and  $a_{ki} = a_{kj}$  for all  $k \neq i, j$ . Using that  $A$  is symmetric, the conditions of the lemma follow.  $\square$

Some basic properties of symmetrical matrices are given in the following lemma.

**Lemma 4.20:**

Let  $A, B \in \mathcal{S}_{ij}$  for some  $i, j \in \{1, 2, \dots, n\}$  such that  $i \neq j$  and  $\alpha, \beta \in \mathbb{R}$ . Then,

1.  $\alpha A + \beta B \in \mathcal{S}_{ij}$ ,
2.  $AB \in \mathcal{S}_{ij}$ , and
3. if  $A$  is invertible, then  $A^{-1} \in \mathcal{S}_{ij}$ .  $\diamond$

*Proof.* Follows directly from the definition of  $\mathcal{S}_{ij}$  in (4.48).  $\square$

We state the main result about synchronization of a pair of generators and prove it in the remainder of this section.

**Theorem 4.21:**

Consider the power system (4.43). The following two statements hold.

1. Let  $n = 2$  and  $M_1 = M_2$ . Then the two generators are synchronized if and only if  $D_1 = D_2$ ,  $f_1 = f_2$ , and  $E_1 = E_2$ .
2. Let  $n \geq 3$  and  $M \in \mathcal{S}_{ij}$ . Then the  $i$ th and  $j$ th generators are synchronized if and only if  $D \in \mathcal{S}_{ij}$ ,  $f \in \mathcal{X}_{ij}$ ,  $E \in \mathcal{X}_{ij}$ , and  $\Gamma \in \mathcal{S}_{ij}$ .  $\diamond$

**Remark 4.22:**

Essentially, this result shows that the  $i$ th and  $j$ th generators are synchronized when the system equation are invariant under swapping the  $i$ th and  $j$ th label.  $\diamond$

We arrange the proof of Theorem 4.21 into a sequence of propositions in this section. We begin by analyzing the equations of the system (4.43) without assumptions on  $n$  and  $M$ .

**Proposition 4.23:**

The  $i$ th and  $j$ th generators are synchronized if and only if

$$\frac{D_i}{M_i} = \frac{D_j}{M_j}, \quad (4.49a)$$

$$\frac{f_i}{M_i} = \frac{f_j}{M_j}, \quad (4.49b)$$

$$\frac{E_i}{M_i \gamma_{ik}} = \frac{E_j}{M_j \gamma_{jk}}, \text{ for } k \neq i, j, \quad (4.49c)$$

$$\frac{\chi_i}{\gamma_{ik}} = \frac{\chi_j}{\gamma_{jk}}, \text{ for } k \neq i, j, \text{ and} \quad (4.49d)$$

$$\frac{\chi_i E_i}{\gamma_{ii}} + \frac{\chi_i E_j}{\gamma_{ij}} = \frac{\chi_j E_i}{\gamma_{ji}} + \frac{\chi_j E_j}{\gamma_{jj}}. \quad (4.49e)$$

◇

*Proof.* From (4.46a), we get

$$\ddot{\delta}_i - \ddot{\delta}_j = -\frac{D_i}{M_i} \dot{\delta}_i + \frac{D_j}{M_j} \dot{\delta}_j + \frac{f_i}{M_i} - \frac{f_j}{M_j} - \sum_{k=1}^n \left( \frac{E_i E_k}{M_i \gamma_{ik}} \sin(\delta_i - \delta_k) - \frac{E_j E_k}{M_j \gamma_{jk}} \sin(\delta_j - \delta_k) \right).$$

It is clear that, if (4.49a), (4.49b), and (4.49c) are true, then  $\delta, \dot{\delta} \in \mathcal{X}_{ij}$  implies  $\ddot{\delta} \in \mathcal{X}_{ij}$ . For the other direction, let us assume that the  $i$ th and  $j$ th generators are synchronized. Then we necessarily have

$$-\left( \frac{D_i}{M_i} - \frac{D_j}{M_j} \right) \dot{\delta}_i + \left( \frac{f_i}{M_i} - \frac{f_j}{M_j} \right) - \sum_{k=1}^n \left( \left( \frac{E_i E_k}{M_i \gamma_{ik}} - \frac{E_j E_k}{M_j \gamma_{jk}} \right) \sin(\delta_i - \delta_k) \right) = 0,$$

for any  $\delta_i, \dot{\delta}_i$ , and  $\delta_k, k \neq i, j$ . Choosing  $\dot{\delta}_i = 0$  and  $\delta_k = \delta_i$ , condition (4.49b) follows. Taking  $\dot{\delta}_i = 1$  and  $\delta_k = \delta_i$ , we find condition (4.49a). Lastly, with  $\delta_i - \delta_k = \frac{\pi}{2}$  for some  $k \neq i, j$  and  $\delta_i - \delta_\ell = 0$  for  $\ell \neq i, j, k$ , condition (4.49c) follows for the chosen  $k$ .

From (4.46b), we have

$$\begin{aligned} V_i \cos \theta_i - V_j \cos \theta_j &= \left( \frac{\chi_i E_i}{\gamma_{ii}} - \frac{\chi_j E_i}{\gamma_{ji}} \right) \cos \delta_i + \left( \frac{\chi_i E_j}{\gamma_{ij}} - \frac{\chi_j E_j}{\gamma_{jj}} \right) \cos \delta_j \\ &\quad + \sum_{\substack{k=1 \\ k \neq i, j}}^n \left( \frac{\chi_i}{\gamma_{ik}} - \frac{\chi_j}{\gamma_{jk}} \right) E_k \cos \delta_k, \\ V_i \sin \theta_i - V_j \sin \theta_j &= \left( \frac{\chi_i E_i}{\gamma_{ii}} - \frac{\chi_j E_i}{\gamma_{ji}} \right) \sin \delta_i + \left( \frac{\chi_i E_j}{\gamma_{ij}} - \frac{\chi_j E_j}{\gamma_{jj}} \right) \sin \delta_j \\ &\quad + \sum_{\substack{k=1 \\ k \neq i, j}}^n \left( \frac{\chi_i}{\gamma_{ik}} - \frac{\chi_j}{\gamma_{jk}} \right) E_k \sin \delta_k. \end{aligned}$$

Similarly, if we assume conditions (4.49d) and (4.49e) to be true, then  $\delta_i = \delta_j$  implies  $V_i \cos \theta_i = V_j \cos \theta_j$  and  $V_i \sin \theta_i = V_j \sin \theta_j$ , which in turn implies that  $\mathbf{V}_g \in \mathcal{X}_{ij}$ .



Conversely, we have

$$0 = \left( \frac{\chi_i E_i}{\gamma_{ii}} + \frac{\chi_i E_j}{\gamma_{ij}} - \frac{\chi_j E_i}{\gamma_{ji}} - \frac{\chi_j E_j}{\gamma_{jj}} \right) \cos \delta_i + \sum_{\substack{k=1 \\ k \neq i, j}}^n \left( \frac{\chi_i}{\gamma_{ik}} - \frac{\chi_j}{\gamma_{jk}} \right) E_k \cos \delta_k,$$

$$0 = \left( \frac{\chi_i E_i}{\gamma_{ii}} + \frac{\chi_i E_j}{\gamma_{ij}} - \frac{\chi_j E_i}{\gamma_{ji}} - \frac{\chi_j E_j}{\gamma_{jj}} \right) \sin \delta_i + \sum_{\substack{k=1 \\ k \neq i, j}}^n \left( \frac{\chi_i}{\gamma_{ik}} - \frac{\chi_j}{\gamma_{jk}} \right) E_k \sin \delta_k,$$

for arbitrary  $\delta_i$  and  $\delta_k$  for  $k \neq i, j$ . By appropriate choices of  $\delta_i$  and  $\delta_k$ , conditions (4.49d) and (4.49e) follow.  $\square$

The following lemma gives an important property of the matrix  $X$ .

**Lemma 4.24:**

For  $X$  as in (4.45), we have  $X\mathbf{1} = \mathbf{1}$ .  $\diamond$

*Proof.* Recalling the definition of  $X$  from (4.45), after some algebraic manipulation, it is clear that  $X\mathbf{1} = \mathbf{1}$  is equivalent to

$$(\mathfrak{L}_{11} - \mathfrak{L}_{12}\mathfrak{L}_{22}^{-1}\mathfrak{L}_{12}^T)\mathbf{1} = 0,$$

which follows from  $\mathfrak{L}\mathbf{1} = 0$ .  $\square$

Let us now assume that  $E_i \neq E_j$  and find what follows from conditions of Proposition 4.23. From (4.49d) and Lemma 4.24, it follows that  $\frac{\chi_i}{\gamma_{ii}} + \frac{\chi_i}{\gamma_{ij}} = \frac{\chi_j}{\gamma_{ji}} + \frac{\chi_j}{\gamma_{jj}}$ . Then, by (4.49e) and  $E_i \neq E_j$ , it is necessary that  $\frac{\chi_i}{\gamma_{ii}} = \frac{\chi_j}{\gamma_{ji}}$  and  $\frac{\chi_i}{\gamma_{ij}} = \frac{\chi_j}{\gamma_{jj}}$ . This, together with (4.49d), means that the  $i$ th and  $j$ th rows in  $X$  are equal, which is a contradiction with  $X$  being invertible. Therefore, for  $i$ th and  $j$ th generators to be synchronized, it is necessary that  $E_i = E_j$ . This allows us to simplify the statement of Proposition 4.23. We can simplify it further by assuming  $M_i = M_j$ , which gives us the following corollary.

**Corollary 4.25:**

Let  $M_i = M_j$ . Then the  $i$ th and  $j$ th generators are synchronized if and only if

$$\begin{aligned} D_i &= D_j, \\ f_i &= f_j, \\ E_i &= E_j, \\ \gamma_{ik} &= \gamma_{jk}, \text{ for } k \neq i, j, \end{aligned} \tag{4.50a}$$

$$\frac{\chi_i}{\gamma_{ik}} = \frac{\chi_j}{\gamma_{jk}}, \text{ for } k \neq i, j, \text{ and} \tag{4.50b}$$

$$\frac{\chi_i}{\gamma_{ii}} + \frac{\chi_i}{\gamma_{ij}} = \frac{\chi_j}{\gamma_{ji}} + \frac{\chi_j}{\gamma_{jj}}. \tag{4.50c}$$

$\diamond$

In the following, we separate the  $n = 2$  and  $n \geq 3$  cases. First, we use Corollary 4.25 to prove part 1 of Theorem 4.21.

*Proof of Theorem 4.21, part 1.* This is true since (4.50a) and (4.50b) are empty statements, while (4.50c) follows immediately from Lemma 4.24.  $\square$

Finally, to prove part 2 of Theorem 4.21, we simplify the statement of Corollary 4.25 for the case of  $n \geq 3$ . This gives us the following corollary.

**Corollary 4.26:**

Let  $n \geq 3$  and  $M_i = M_j$ . Then the  $i$ th and  $j$ th generators are synchronized if and only if

$$\begin{aligned} D_i &= D_j, \\ f_i &= f_j, \\ E_i &= E_j, \\ \gamma_{ik} &= \gamma_{jk}, \text{ for } k \neq i, j, \end{aligned} \tag{4.51a}$$

$$\chi_i = \chi_j, \text{ and} \tag{4.51b}$$

$$\gamma_{ii} = \gamma_{jj}. \tag{4.51c}$$

$\diamond$

*Proof.* Condition (4.51b) follows from (4.50a) and (4.50b), using that there are at least three generators. Then (4.51c) follows from (4.50c), (4.51b), and symmetry  $\gamma_{ij} = \gamma_{ji}$ .  $\square$

Corollary 4.26, together with the following lemma allows us to complete the proof of Theorem 4.21.

**Lemma 4.27:**

Let  $i, j \in \{1, 2, \dots, n\}$  be such that  $i \neq j$ . We have  $\Gamma \in \mathcal{S}_{ij}$  if and only if  $L_D \in \mathcal{S}_{ij}$  and  $\mathfrak{L}_{11} - \mathfrak{L}_{12}\mathfrak{L}_{22}^{-1}\mathfrak{L}_{12}^T \in \mathcal{S}_{ij}$ .  $\diamond$

*Proof.*  $\Leftarrow$  Follows from Lemma 4.20.

$\Rightarrow$  First we show that  $L_D \in \mathcal{S}_{ij}$ . Using  $\Gamma = L_D X$ ,  $\Pi_{ij}\mathbf{1} = \mathbf{1}$ , and  $X\mathbf{1} = \mathbf{1}$ , from  $\Gamma\Pi_{ij}\mathbf{1} = \Pi_{ij}\Gamma\mathbf{1}$  it follows that  $L_D\mathbf{1} = \Pi_{ij}L_D\mathbf{1}$ . Since  $L_D$  is a diagonal matrix, from this we see that  $L_D \in \mathcal{S}_{ij}$ . Now  $\mathfrak{L}_{11} - \mathfrak{L}_{12}\mathfrak{L}_{22}^{-1}\mathfrak{L}_{12}^T \in \mathcal{S}_{ij}$  follows from Lemma 4.20.  $\square$

Now we can complete the proof of Theorem 4.21.

*Proof of Theorem 4.21, part 2.* Conditions (4.51a) and (4.51c), by Lemma 4.19, are equivalent to  $\Gamma \in \mathcal{S}_{ij}$ , which, by Lemma 4.27, is in turn equivalent to  $L_D \in \mathcal{S}_{ij}$  and  $\mathfrak{L}_{11} - \mathfrak{L}_{12}\mathfrak{L}_{22}^{-1}\mathfrak{L}_{12}^T \in \mathcal{S}_{ij}$ . Therefore, (4.51a) and (4.51c) imply (4.51b).  $\square$

### 4.3.4 Synchronization of generator partition

Let  $\pi = \{\mathcal{C}_\ell\}_{\ell \in \widehat{\mathcal{G}}}$  be a partition of the set  $\mathcal{G}$ , where  $\widehat{\mathcal{G}} = \{1, 2, \dots, \widehat{n}\}$  and  $\widehat{n} \leq n$ . Let us denote

$$\mathcal{X}_{\text{cl}} := \bigcap_{\ell \in \widehat{\mathcal{G}}} \bigcap_{i, j \in \mathcal{C}_\ell} \mathcal{X}_{ij}, \quad \mathcal{S}_{\text{cl}} := \bigcap_{\ell \in \widehat{\mathcal{G}}} \bigcap_{i, j \in \mathcal{C}_\ell} \mathcal{S}_{ij}.$$

and  $\mathfrak{P}$  the characteristic matrix of the partition  $\pi$ . Notice that  $\mathcal{X}_{\text{cl}} = \text{im}(\mathfrak{P})$ .

We define two notions generalizing the synchronization of two generators to a partition of generators.

**Definition 4.28:**

The system (4.43) is said to be *strongly synchronized with respect to partition  $\pi$*  if the  $i$ th and  $j$ th generators are synchronized for all  $i, j \in \mathfrak{C}_\ell$  and all  $\ell \in \widehat{\mathfrak{G}}$ , i.e.,  $\delta(t) \in \mathcal{X}_{ij}$  and  $\mathbf{V}_{\mathfrak{G}}(t) \in \mathcal{X}_{ij}$  for all  $t \geq 0$  and for any  $\delta(0), \dot{\delta}(0) \in \mathcal{X}_{ij}$ ,  $i, j \in \mathfrak{C}_\ell$ , and  $\ell \in \widehat{\mathfrak{G}}$ .

The system (4.43) is said to be *weakly synchronized with respect to partition  $\pi$*  if, for arbitrary  $\delta(0), \dot{\delta}(0) \in \mathcal{X}_{\text{cl}}$ , there exist functions  $\widehat{\delta}: [0, \infty) \rightarrow \mathbb{R}^{\widehat{n}}$  and  $\widehat{\mathbf{V}}_{\widehat{\mathfrak{G}}}: [0, \infty) \rightarrow \mathbb{C}^{\widehat{n}}$  such that  $\delta(t) = \mathfrak{P}\widehat{\delta}(t)$  and  $\mathbf{V}_{\mathfrak{G}}(t) = \mathfrak{P}\widehat{\mathbf{V}}_{\widehat{\mathfrak{G}}}(t)$ , i.e.,  $\delta(t) \in \mathcal{X}_{\text{cl}}$  and  $\mathbf{V}_{\mathfrak{G}}(t) \in \mathcal{X}_{\text{cl}}$  for all  $t \geq 0$  and for any  $\delta(0), \dot{\delta}(0) \in \mathcal{X}_{\text{cl}}$ .  $\diamond$

**Remark 4.29:**

Notice that strong synchronization is equivalent to  $\mathcal{X}_{ij} \times \mathcal{X}_{ij} \times \mathcal{X}_{ij}$  being an invariant set for  $(\delta, \dot{\delta}, \widehat{\mathbf{V}}_{\widehat{\mathfrak{G}}})$  for any  $i, j \in \mathfrak{C}_\ell$  and  $\ell \in \widehat{\mathfrak{G}}$ , while weak synchronization is equivalent to an invariant set being  $\mathcal{X}_{\text{cl}} \times \mathcal{X}_{\text{cl}} \times \mathcal{X}_{\text{cl}}$ . This means that, if the power system is strongly synchronized, when two generators and their buses in the same cluster have equal state, they will remain equal. If the power system is weakly synchronized, then when the states of every generator and its bus are equal to all others in the same cluster, they will stay equal. From this, we see that that if the system (4.43) is strongly synchronized with respect to  $\pi$ , then it is also weakly synchronized with respect to  $\pi$ , since  $\mathcal{X}_{\text{cl}} \times \mathcal{X}_{\text{cl}} \times \mathcal{X}_{\text{cl}} \subseteq \mathcal{X}_{ij} \times \mathcal{X}_{ij} \times \mathcal{X}_{ij}$ , for all  $i, j \in \mathfrak{C}_\ell$  and all  $\ell \in \widehat{\mathfrak{G}}$ .

Further, the  $i$ th and  $j$ th generators are synchronized if and only if (4.43) is either strongly or weakly synchronized with respect to  $\{\{i, j\}\} \cup \{\{k\} : k \neq i, j\}$ .

Finally, notice that (4.43) is always both strongly and weakly synchronized with respect to  $\{\{i\} : i \in \mathfrak{G}\}$ .  $\diamond$

In the following, we show necessary and sufficient conditions for the two synchronization notions. To start, in the next proposition, we present cases when the structure of  $\Gamma$  has no influence. It also illustrates the relation between strong and weak synchronization.

**Proposition 4.30:**

Let  $\pi = \{\mathfrak{G}\}$ ,  $M, D \in \mathcal{S}_{\text{cl}}$ , and  $f, E \in \mathcal{X}_{\text{cl}}$ . Then the system (4.43) is weakly synchronized with respect to  $\{\mathfrak{G}\}$ . If additionally  $n = 2$ , then (4.43) is also strongly synchronized with respect to  $\{\mathfrak{G}\}$ .  $\diamond$

*Proof.* From the assumptions, it follows that  $M = \widehat{m}I$ ,  $D = \widehat{d}I$ ,  $f = \widehat{f}\mathbf{1}$ , and  $E = \widehat{E}\mathbf{1}$ , for some  $\widehat{m} > 0$ ,  $\widehat{d} \geq 0$ , and  $\widehat{f}, \widehat{E} \in \mathbb{R}$ . Notice that for  $\pi = \{\mathfrak{G}\}$ , we have  $\mathfrak{P} = \mathbf{1}$ .

Let us assume that  $\delta(0), \dot{\delta}(0) \in \text{im}(\mathbf{1})$ . To prove weak synchronization, we need to show that  $\delta(t) \in \text{im}(\mathbf{1})$  and  $\mathbf{V}_{\mathfrak{G}}(t) \in \text{im}(\mathbf{1})$ . For the former, it is enough to show that  $\ddot{\delta}(t) \in \text{im}(\mathbf{1})$  if  $\delta(t), \dot{\delta}(t) \in \text{im}(\mathbf{1})$ , which is clear, since then  $\ddot{\delta}(t) = -M^{-1}D\dot{\delta}(t) + M^{-1}f = -\frac{\widehat{d}}{\widehat{m}}\dot{\delta}(t) + \frac{\widehat{f}}{\widehat{m}}\mathbf{1}$ . For the latter, we see that  $\mathbf{V}_{\mathfrak{G}} = L_{\text{D}}^{-1}\Gamma\mathbf{E}_{\mathfrak{G}} \in \text{im}(\mathbf{1})$  whenever  $\mathbf{E}_{\mathfrak{G}} \in \text{im}(\mathbf{1})$ , which is equivalent to  $\delta \in \text{im}(\mathbf{1})$ .

The second part follows from part 1 of Theorem 4.21.  $\square$

We continue with the first main result of this section—the necessary and sufficient conditions for strong synchronization. Here, symmetrical conditions for  $\Gamma$  are relevant.

**Theorem 4.31:**

Let  $n \geq 3$ ,  $\pi$  arbitrary, and  $M \in \mathcal{S}_{\text{cl}}$ . Then the system (4.43) is strongly synchronized with respect to  $\pi$  if and only if  $D \in \mathcal{S}_{\text{cl}}$ ,  $f \in \mathcal{X}_{\text{cl}}$ ,  $E \in \mathcal{X}_{\text{cl}}$ , and  $\Gamma \in \mathcal{S}_{\text{cl}}$ .  $\diamond$

*Proof.* Follows from applying part 2 of Theorem 4.21 for every  $i$ th and  $j$ th generator where  $i, j \in \mathfrak{C}_\ell$  and  $\ell \in \widehat{\mathfrak{G}}$ .  $\square$

We conclude this section with the second main result—the necessary and sufficient conditions for weak synchronization. Instead of symmetrical conditions,  $\mathcal{X}_{\text{cl}}$  being  $\Gamma$ -invariant is one of the conditions. Since  $\mathcal{X}_{\text{cl}} = \text{im}(\mathfrak{P})$ , this actually means that  $\pi$  is an equitable partition for a graph whose adjacency matrix is  $\Gamma$  (see Lemma 2.48).

**Theorem 4.32:**

Let  $|\pi| \geq 2$ ,  $M, D \in \mathcal{S}_{\text{cl}}$ , and  $f, E \in \mathcal{X}_{\text{cl}}$ . Then the system (4.43) is weakly synchronized with respect to  $\pi$  if and only if

$$L_D \in \mathcal{S}_{\text{cl}} \quad \text{and} \quad \mathcal{X}_{\text{cl}} \text{ is } \Gamma\text{-invariant.} \quad (4.52)$$

$\diamond$

*Proof.* From the definition, we see that (4.43) is weakly synchronized with respect to  $\pi$  if and only if

$$\left( \forall \delta, \dot{\delta} \in \mathcal{X}_{\text{cl}} \right) M^{-1} \left( -D\dot{\delta} + f - (\Gamma \circ EE^T \circ \sin(\delta \mathbf{1}_n^T - \mathbf{1}_n \delta^T)) \mathbf{1}_n \right) \in \mathcal{X}_{\text{cl}} \quad (4.53)$$

and

$$(\forall \delta \in \mathcal{X}_{\text{cl}}) L_D^{-1} \Gamma \mathbf{E}_{\mathfrak{G}} \in \mathcal{X}_{\text{cl}}. \quad (4.54)$$

Since  $M, D \in \mathcal{S}_{\text{cl}}$  and  $f \in \mathcal{X}_{\text{cl}}$ , condition (4.53) is equivalent to

$$(\forall \delta \in \mathcal{X}_{\text{cl}}) (\Gamma \circ EE^T \circ \sin(\delta \mathbf{1}_n^T - \mathbf{1}_n \delta^T)) \mathbf{1}_n \in \mathcal{X}_{\text{cl}}.$$

Using  $\delta = \mathfrak{P}\widehat{\delta}$ ,  $E = \mathfrak{P}\widehat{E}$ ,  $\mathbf{1}_n = \mathfrak{P}\mathbf{1}_{\widehat{n}}$ , and that  $v \in \mathcal{X}_{\text{cl}}$  is equivalent to  $\Pi_{ij}v = v$  for all  $i, j \in \mathfrak{C}_\ell$  and  $\ell \in \widehat{\mathfrak{G}}$ , we find that the above condition is equivalent to

$$\begin{aligned} & \left( \forall \widehat{\delta} \in \mathbb{R}^{\widehat{n}} \right) \left( \forall \ell \in \widehat{\mathfrak{G}} \right) \left( \forall i, j \in \mathfrak{C}_\ell \right) \\ & \left( (\Gamma \mathfrak{P} - \Pi_{ij} \Gamma \mathfrak{P}) \circ \mathfrak{P} \widehat{E} \widehat{E}^T \circ \mathfrak{P} \sin \left( \widehat{\delta} \mathbf{1}_{\widehat{n}}^T - \mathbf{1}_{\widehat{n}} \widehat{\delta}^T \right) \right) \mathbf{1}_{\widehat{n}} = 0. \end{aligned} \quad (4.55)$$

In a similar way, we find that the condition (4.54) is equivalent to

$$\left( \forall \widehat{\mathbf{E}}_{\widehat{\mathfrak{G}}} \in \mathbb{C}^{\widehat{n}} \right) \left( \forall \ell \in \widehat{\mathfrak{G}} \right) \left( \forall i, j \in \mathfrak{C}_\ell \right) L_D^{-1} \Gamma \mathfrak{P} \widehat{\mathbf{E}}_{\widehat{\mathfrak{G}}} = \Pi_{ij} L_D^{-1} \Gamma \mathfrak{P} \widehat{\mathbf{E}}_{\widehat{\mathfrak{G}}},$$

or, more simply,

$$\left(\forall \ell \in \widehat{\mathcal{G}}\right)(\forall i, j \in \mathfrak{C}_\ell) L_D^{-1} \Gamma \mathfrak{P} = \Pi_{ij} L_D^{-1} \Gamma \mathfrak{P}. \quad (4.56)$$

It is straightforward to check that (4.52) implies (4.55) and (4.56). For the other direction, choosing  $\widehat{\delta} = e_{\ell_2}$  for  $\ell_2 \neq \ell$  in (4.55), we find from the  $i$ th row that

$$\left(\forall \ell, \ell_2 \in \widehat{\mathcal{G}}, \ell_2 \neq \ell\right)(\forall i, j \in \mathfrak{C}_\ell) \sum_{k \in \mathfrak{C}_{\ell_2}} \frac{1}{\gamma_{ik}} = \sum_{k \in \mathfrak{C}_{\ell_2}} \frac{1}{\gamma_{jk}}. \quad (4.57)$$

The  $i$ th row and  $\ell_2$ th column in condition (4.56) gives

$$\left(\forall \ell, \ell_2 \in \widehat{\mathcal{G}}\right)(\forall i, j \in \mathfrak{C}_\ell) \chi_i \sum_{k \in \mathfrak{C}_{\ell_2}} \frac{1}{\gamma_{ik}} = \chi_j \sum_{k \in \mathfrak{C}_{\ell_2}} \frac{1}{\gamma_{jk}}. \quad (4.58)$$

Since the assumption is that there are at least two clusters in  $\pi$ , from (4.57) and (4.58) we find that  $\chi_i = \chi_j$ , for all  $i, j \in \mathfrak{C}_\ell$  and all  $\ell \in \widehat{\mathcal{G}}$ , i.e.,  $L_D \in \mathcal{S}_{\text{cl}}$ . This, together with (4.58), gives

$$\left(\forall \ell, \ell_2 \in \widehat{\mathcal{G}}\right)(\forall i, j \in \mathfrak{C}_\ell) \sum_{k \in \mathfrak{C}_{\ell_2}} \frac{1}{\gamma_{ik}} = \sum_{k \in \mathfrak{C}_{\ell_2}} \frac{1}{\gamma_{jk}},$$

which is equivalent to  $\text{im}(\Gamma \mathfrak{P}) \subseteq \mathcal{X}_{\text{cl}}$ , i.e.,  $\Gamma \mathcal{X}_{\text{cl}} \subseteq \mathcal{X}_{\text{cl}}$ .  $\square$

### 4.3.5 Clustering of power systems

Let us assume that the system (4.43) is weakly synchronized with respect to a partition  $\pi$ . Let also the initial condition satisfy  $\delta(0), \dot{\delta}(0) \in \mathcal{X}_{\text{cl}}$ . Then there exist  $\widehat{\delta}$  and  $\widehat{\mathbf{V}}_{\widehat{\mathcal{G}}}$  such that  $\delta(t) = \mathfrak{P}\widehat{\delta}(t)$  and  $\mathbf{V}_{\widehat{\mathcal{G}}}(t) = \mathfrak{P}\widehat{\mathbf{V}}_{\widehat{\mathcal{G}}}(t)$ , which also gives us  $V_{\widehat{\mathcal{G}}}(t) = \mathfrak{P}\widehat{V}_{\widehat{\mathcal{G}}}(t)$  and  $\theta_{\widehat{\mathcal{G}}}(t) = \mathfrak{P}\widehat{\theta}_{\widehat{\mathcal{G}}}(t)$ . Inserting this into (4.43) with dynamics rewritten as in (4.47), we find

$$\begin{aligned} M\ddot{\mathfrak{P}}\widehat{\delta}(t) + D\dot{\mathfrak{P}}\widehat{\delta}(t) &= f - L_D \left( E \circ \mathfrak{P}\widehat{\mathbf{V}}_{\widehat{\mathcal{G}}}(t) \circ \sin\left(\mathfrak{P}\widehat{\delta}(t) - \mathfrak{P}\widehat{\theta}_{\widehat{\mathcal{G}}}(t)\right) \right), \\ \begin{bmatrix} L_D \left( \mathbf{E}_{\widehat{\mathcal{G}}}(t) - \mathfrak{P}\widehat{\mathbf{V}}_{\widehat{\mathcal{G}}}(t) \right) \\ 0 \end{bmatrix} &= \begin{bmatrix} \mathfrak{L}_{11} & \mathfrak{L}_{12} \\ \mathfrak{L}_{12}^T & \mathfrak{L}_{22} \end{bmatrix} \begin{bmatrix} \mathfrak{P}\widehat{\mathbf{V}}_{\widehat{\mathcal{G}}}(t) \\ \mathbf{V}_{\widehat{\mathcal{G}}}(t) \end{bmatrix}. \end{aligned}$$

Assuming additionally that  $E \in \mathcal{X}_{\text{cl}}$ , i.e.,  $E = \mathfrak{P}\widehat{E}$  for some  $\widehat{E} \in \mathbb{R}^{\widehat{n}}$ , and premultiplying the above dynamics and first block-row of the constraint by  $\mathfrak{P}^T$ , we obtain

$$\widehat{M}\ddot{\widehat{\delta}}(t) + \widehat{D}\dot{\widehat{\delta}}(t) = \widehat{f} - \widehat{L}_D \left( \widehat{E} \circ \widehat{\mathbf{V}}_{\widehat{\mathcal{G}}}(t) \circ \sin\left(\widehat{\delta}(t) - \widehat{\theta}_{\widehat{\mathcal{G}}}(t)\right) \right), \quad (4.59a)$$

$$\begin{bmatrix} \widehat{L}_D \left( \widehat{\mathbf{E}}_{\widehat{\mathcal{G}}}(t) - \widehat{\mathbf{V}}_{\widehat{\mathcal{G}}}(t) \right) \\ 0 \end{bmatrix} = \begin{bmatrix} \widehat{\mathfrak{L}}_{11} & \widehat{\mathfrak{L}}_{12} \\ \widehat{\mathfrak{L}}_{12}^T & \widehat{\mathfrak{L}}_{22} \end{bmatrix} \begin{bmatrix} \widehat{\mathbf{V}}_{\widehat{\mathcal{G}}}(t) \\ \mathbf{V}_{\widehat{\mathcal{G}}}(t) \end{bmatrix}, \quad (4.59b)$$

where  $\widehat{M} = \mathfrak{P}^T M \mathfrak{P}$ ,  $\widehat{D} = \mathfrak{P}^T D \mathfrak{P}$ ,  $\widehat{f} = \mathfrak{P}^T f$ ,  $\widehat{L}_D = \mathfrak{P}^T L_D \mathfrak{P}$ ,  $\widehat{\mathcal{L}}_{11} = \mathfrak{P}^T \mathcal{L}_{11} \mathfrak{P}$ ,  $\widehat{\mathcal{L}}_{12} = \mathfrak{P}^T \mathcal{L}_{12}$ . Moreover, from  $\delta(t) = \mathfrak{P} \widehat{\delta}(t)$  follows that  $\widehat{\delta}(0) = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \delta(0)$  and  $\dot{\widehat{\delta}}(0) = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \dot{\delta}(0)$ .

Notice that the reduced model (4.59) is again a power system of the same form as (4.43). In particular, we have that  $\widehat{M}$ ,  $\widehat{D}$ , and  $\widehat{L}_D$  are positive definite diagonal matrices and that  $\widehat{\mathcal{L}}$  is a Laplacian matrix. Additionally, note that this projection-based MOR can be done for arbitrary power system and arbitrary partition. In general, we can take (4.59) with  $\widehat{\delta}(0) = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \delta(0)$ ,  $\dot{\widehat{\delta}}(0) = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \dot{\delta}(0)$ , and  $\widehat{E} = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T E$ . We can also apply Kron reduction to this reduced model.

### 4.3.6 Illustrative example

For the example in Figure 4.7, let  $\chi_i = 1$  and  $\chi_{ij} = 1$  for all  $i, j$ . Then we have

$$\Gamma = \frac{1}{32} \begin{bmatrix} 21 & 5 & 2 & 2 & 2 \\ 5 & 21 & 2 & 2 & 2 \\ 2 & 2 & 16 & 8 & 4 \\ 2 & 2 & 8 & 16 & 4 \\ 2 & 2 & 4 & 4 & 20 \end{bmatrix}.$$

Additionally, let  $M = D = I_5$ ,  $f = 0$ , and  $E = \mathbb{1}_5$ . Then, using Theorem 4.21, we see that the first and second generators are synchronized, and that the same is true for the third and fourth. By definition, this implies that the system is strongly synchronized with respect to  $\{\{1, 2\}, \{3, 4\}, \{5\}\}$ . On the other hand, from Theorem 4.32 and

$$\Gamma \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} = \frac{1}{16} \begin{bmatrix} 13 & 3 \\ 13 & 3 \\ 2 & 14 \\ 2 & 14 \\ 2 & 14 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \left( \frac{1}{16} \begin{bmatrix} 13 & 3 \\ 2 & 14 \end{bmatrix} \right),$$

we see that the system is weakly synchronized with respect to  $\{\{1, 2\}, \{3, 4, 5\}\}$ , but not strongly. Using the partition  $\pi = \{\{1, 2\}, \{3, 4, 5\}\}$  for clustering, we find the following reduced quantities:  $\widehat{M} = \widehat{D} = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}$ ,  $\widehat{f} = 0$ ,  $\widehat{E} = \mathbb{1}_2$ ,  $\widehat{L}_D = \widehat{\mathcal{L}}_{11} = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}$ ,  $\widehat{\mathcal{L}}_{12} = \begin{bmatrix} -2 & 0 \\ 0 & -3 \end{bmatrix}$ ,  $\widehat{\Gamma} = \frac{1}{8} \begin{bmatrix} 13 & 3 \\ 3 & 21 \end{bmatrix}$ . The Figure 4.8 shows the partition and Figure 4.9 the associated reduced power system. From the definition of weak synchronization, we know that this reduced power system exactly reproduces the initial value response of the original system for any initial condition  $\delta(0), \dot{\delta}(0) \in \mathcal{X}_{\text{cl}}$ , taking the initial condition of the reduced model to be  $\widehat{\delta}(0) = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \delta(0)$  and  $\dot{\widehat{\delta}}(0) = (\mathfrak{P}^T \mathfrak{P})^{-1} \mathfrak{P}^T \dot{\delta}(0)$ .

To demonstrate the possibility to cluster using any partition, including those with respect to which the power system is not weakly synchronized, and any initial condition, we show simulation result for partition  $\{\{1, 2, 3\}, \{4, 5\}\}$  in Figure 4.10. We see that, in this case, the reduced model matches the steady state and approximates the transient behavior. Finding sufficient conditions for matching the steady state and deriving error bounds is a possible topic of future research.

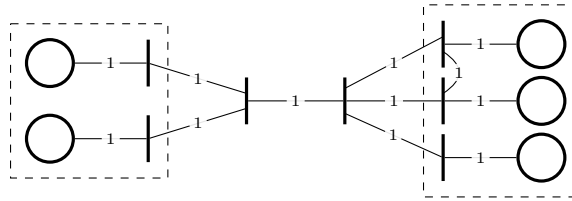


Figure 4.8: Partition  $\{\{1, 2\}, \{3, 4, 5\}\}$  applied to the original power system in Figure 4.7 with  $\chi_i = \chi_{ij} = 1$  for all  $i, j$ .

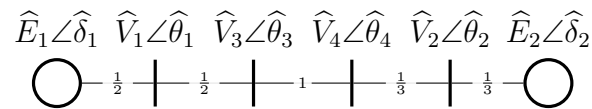


Figure 4.9: Reduced power system obtained by clustering the system in Figure 4.8 with  $M = D = I_5$ ,  $f = 0$ , and  $E = \mathbb{1}_5$ .

## 4.4 Conclusion

In Section 4.2, we have extended results on MOR of leader-follower networks with single integrator agent dynamics from [MTC14] to leader-follower networks with arbitrary linear multivariable agent dynamics. We have also extended these results to the case that the approximation error is measured in the  $\mathcal{H}_\infty$ -norm. The proposed MOR technique reduces the complexity of the network topology by clustering the agents. We have shown that clustering amounts to applying a specific Petrov-Galerkin projection associated with the graph partition. The resulting reduced order model can be interpreted as a networked multi-agent system with a weighted, directed network graph. If the original network is clustered using an almost equitable graph partition, then its consensus properties are preserved. We have provided a priori upper bounds on the  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  model reduction errors in this case. These error bounds depend on an auxiliary system related to the agent dynamics, the eigenvalues of the Laplacian matrices of the original and the reduced network, and on the number of cellmates of the leaders in the network. Finally, we have provided some insight into the general case of clustering according to arbitrary, not necessarily almost equitable, partitions. Here, direct computation of a priori upper bounds on the error is not as straightforward as in the case of almost equitable partitions. We have shown that in this more general case, one can bound the model reduction errors by first optimally approximating the original network by a new network for which the chosen partition is almost equitable, and then bounding the  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  errors using the triangle inequality.

In Section 4.3, we analyzed power systems consisting of generators and buses. We introduced a notion of synchronization for a pair of generators and two for a partition of the set of generators. We proved equivalent conditions depending on the Kron-reduced system being symmetrical or equitable. This additionally gives a relation between symmetrical matrices and equitable partitions. We showed how a synchronized power

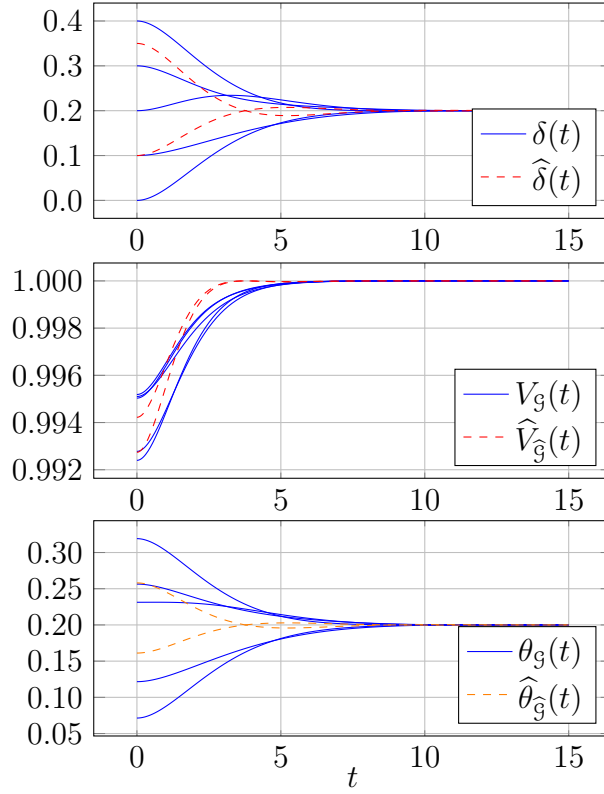


Figure 4.10: Initial value response of the original power system from Figure 4.7 and a reduced system obtained by clustering with partition  $\{\{1, 2, 3\}, \{4, 5\}\}$ . Original system's parameters are  $\chi_i = \chi_{ij} = 1$  for all  $i, j$ ,  $M = D = I_5$ ,  $f = 0$ , and  $E = \mathbf{1}_5$ . The initial value is  $\delta(0) = (0, 0.1, 0.2, 0.3, 0.4)$  and  $\hat{\delta}(0) = 0$ .

system can be exactly approximated with a reduced system by clustering generators and their buses. Furthermore, this provides a clustering-based MOR method for arbitrary power systems, although finding bounds for the approximation error remains an open problem.

These results give motivation for further research into efficient clustering methods for linear and nonlinear multi-agent systems.



# CHAPTER 5

## SUBSYSTEM REDUCTION FOR INTERCONNECTED SYSTEMS

### Contents

---

5.1	Introduction . . . . .	99
5.2	Stability-preserving balancing-based model order reduction . . . . .	100
5.2.1	Preliminaries . . . . .	100
5.2.2	Bounded real balanced truncation . . . . .	101
5.2.3	Stability-preserving model order reduction . . . . .	102
5.2.4	Numerical example . . . . .	104
5.3	$\mathcal{H}_2$ -optimal subsystem reduction . . . . .	105
5.3.1	Interconnected systems . . . . .	106
5.3.2	Multi-agent systems . . . . .	109
5.4	Conclusion . . . . .	112

---

## 5.1 Introduction

Here, we consider MOR of interconnected systems, in particular of LTI systems linearly interconnected through their inputs and outputs. Moreover, we want to preserve the interconnection structure by only reducing the subsystems, while also preserving the stability of the coupled system. This is a complementary approach to clustering-based methods discussed in previous chapters.

Several approaches were investigated in the literature. Reis and Stykel [RS07, RS08b] proposed using BT for each subsystem and proved a sufficient condition for stability of the interconnected system. Furthermore, they developed an a priori error bound of the network system based on the error bounds for the individual subsystems. Vandendorpe and Van Dooren [VVD08] showed sufficient conditions for transfer function interpolation of the interconnected system while preserving the interconnection structure. Sandberg and Murray [SM09] used block-diagonal generalized Gramians to extend BT

for unstructured LTI systems to interconnected systems. Monshizadeh et al. [MTC13] proposed methods preserving stability or synchronization for multi-agent systems.

We consider balancing-based and  $\mathcal{H}_2$ -optimal MOR of interconnected systems in Section 5.2 and Section 5.3, respectively.

## 5.2 Stability-preserving balancing-based model order reduction

### 5.2.1 Preliminaries

We extend the stability-preserving MOR method, based on bounded real balanced truncation (BRBT), for multi-agent systems from [MTC13] to coupled systems considered in [RS07]. More information about BRBT is in the following section.

Reis and Stykel [RS07] study systems of  $\mathbf{n}$  coupled LTI subsystems

$$\begin{aligned} E_i \dot{x}_i(t) &= A_i x_i(t) + B_i u_i(t), \\ y_i(t) &= C_i x_i(t), \end{aligned} \quad (5.1a)$$

with interconnections and external input

$$u_i(t) = K_{i1} y_1(t) + \cdots + K_{in} y_n(t) + F_i u(t), \quad (5.1b)$$

and external output

$$y(t) = G_1 y_1(t) + \cdots + G_n y_n(t), \quad (5.1c)$$

where  $i \in \{1, 2, \dots, \mathbf{n}\}$ ,  $E_i, A_i \in \mathbb{R}^{n_i \times n_i}$ ,  $B_i \in \mathbb{R}^{n_i \times m_i}$ ,  $C_i \in \mathbb{R}^{p_i \times n_i}$ ,  $x_i(t) \in \mathbb{R}^{n_i}$  is the state,  $u_i(t) \in \mathbb{R}^{m_i}$  is the internal input,  $y_i(t) \in \mathbb{R}^{p_i}$  is the internal output,  $K_{ij} \in \mathbb{R}^{m_i \times p_j}$ ,  $F_i \in \mathbb{R}^{m_i \times m}$ ,  $u(t) \in \mathbb{R}^m$  is the external input,  $G_i \in \mathbb{R}^{p \times p_i}$ , and  $y(t) \in \mathbb{R}^p$  is the external output. Here, we assume the matrices  $E_i$  are invertible and that the matrix pencils  $A_i - \lambda E_i$  are asymptotically stable. The LTI system (5.1) can be written (see [RS07]) as

$$\begin{aligned} \mathcal{E} \dot{x}(t) &= \mathcal{A} x(t) + \mathcal{B} u(t), \\ y(t) &= \mathcal{C} x(t), \end{aligned} \quad (5.2)$$

with  $\mathcal{E}, \mathcal{A} \in \mathbb{R}^{n \times n}$ ,  $\mathcal{B} \in \mathbb{R}^{n \times m}$ ,  $\mathcal{C} \in \mathbb{R}^{p \times n}$ , and  $x(t) = \text{col}(x_1(t), x_2(t), \dots, x_n(t)) \in \mathbb{R}^n$ , where

$$\begin{aligned} \mathcal{E} &= E_D, \quad \mathcal{A} = A_D + B_D K C_D, \quad \mathcal{B} = B_D F, \quad \mathcal{C} = G C_D, \\ E_D &= \text{diag}(E_1, E_2, \dots, E_n), \quad A_D = \text{diag}(A_1, A_2, \dots, A_n), \\ B_D &= \text{diag}(B_1, B_2, \dots, B_n), \quad C_D = \text{diag}(C_1, C_2, \dots, C_n), \\ K &= \begin{bmatrix} K_{11} & K_{12} & \cdots & K_{1n} \\ K_{21} & K_{22} & \cdots & K_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ K_{n1} & K_{n2} & \cdots & K_{nn} \end{bmatrix}, \quad F = \begin{bmatrix} F_1 \\ F_2 \\ \vdots \\ F_n \end{bmatrix}, \quad G = [G_1 \quad G_2 \quad \cdots \quad G_n]. \end{aligned} \quad (5.3)$$

---

**Algorithm 5.1:** Bounded real balanced truncation [OJ88]

---

**Input:** Asymptotically stable system  $(E; A, B, C)$  and  $\gamma > 0$  such that

$$\|H\|_{\mathcal{H}_\infty} < \gamma.$$

**Output:** ROM  $(\widehat{E}; \widehat{A}, \widehat{B}, \widehat{C})$ .

1 Compute the maximal solutions  $P$  and  $Q$  of

$$APE^T + EPA^T + BB^T + \frac{1}{\gamma^2}EPC^T CPE^T = 0,$$

$$A^TQE + E^TQA + C^TC + \frac{1}{\gamma^2}E^TQBB^TQE = 0.$$

2 Proceed as in Algorithm 2.1 using the obtained  $P$  and  $Q$ .

---

Denote  $H_i(s) = C_i(sE_i - A_i)^{-1}B_i$  and  $\mathcal{H}(s) = \mathcal{C}(s\mathcal{E} - \mathcal{A})^{-1}\mathcal{B}$  the transfer functions of the subsystems and the interconnected system respectively. Furthermore, let  $H_D(s) = C_D(sE_D - A_D)^{-1}B_D = \text{diag}(H_1(s), H_2(s), \dots, H_n(s))$ . Then, we have

$$\mathcal{H}(s) = G(I - H_D(s)K)^{-1}H_D(s)F = GH_D(s)(I - KH_D(s))^{-1}F.$$

## 5.2.2 Bounded real balanced truncation

BRBT is a modification of the standard BT method, originally developed to preserve strict bounded realness, or equivalently the bound  $\|H\|_{\mathcal{H}_\infty} < 1$ , of the original system [OJ88]. Algorithm 5.1 describes a straightforward extension of BRBT which preserves the bound  $\|H\|_{\mathcal{H}_\infty} < \gamma$ , for any given  $\gamma > 0$ . The following theorem presents some properties of BRBT and an  $\mathcal{H}_\infty$ -error bound, and is also a straightforward extension of results in [OJ88].

**Theorem 5.1:**

Let  $H(s) = C(sE - A)^{-1}B$  be an LTI system of order  $n$  such that  $\|H\|_{\mathcal{H}_\infty} < \gamma$  and  $\widehat{H}(s) = \widehat{C}(s\widehat{E} - \widehat{A})^{-1}\widehat{B}$  its ROM of order  $r < n$  obtained by BRBT. Then  $\widehat{H}$  is asymptotically stable,  $\|\widehat{H}\|_{\mathcal{H}_\infty} < \gamma$ , and

$$\|H - \widehat{H}\|_{\mathcal{H}_\infty} \leq 2 \sum_{i=r+1}^n \xi_i,$$

where  $\xi_i$  are the *bounded real characteristic values* (square roots of eigenvalues of  $PE^TQE$  in Algorithm 5.1).  $\diamond$

### 5.2.3 Stability-preserving model order reduction

In the above setting, the multi-agent systems studied by Monshizadeh et al. [MTC13] can be represented with

$$\begin{aligned} E_i &= I_n, & A_i &= A_1, & B_i &= B_1, & C_i &= C_1, & \text{for } i &= 1, 2, \dots, \mathbf{n}, \\ K &= -\mathcal{L} \otimes I_{m_1}, & m_1 &= p_1, \end{aligned} \quad (5.4)$$

where  $\mathcal{L}$  is the Laplacian matrix of an undirected interconnection graph. For these systems, the following sufficient condition for asymptotic stability was proven.

**Lemma 5.2** ([MTC13, Lemma 3.1]):

If  $\|\mathcal{L}\|_2 \|H_1\|_{\mathcal{H}_\infty} < 1$ , then the multi-agent system (5.2) with (5.4) is asymptotically stable.  $\diamond$

Based on this, they propose using BRBT to reduce the subsystem  $(A_1, B_1, C_1)$  to  $(\hat{A}_1, \hat{B}_1, \hat{C}_1)$  with transfer function  $\hat{H}_1$  such that  $\|\hat{H}_1\|_{\mathcal{H}_\infty} < \frac{1}{\|\mathcal{L}\|_2}$ . Since the interconnection structure is preserved in the reduced network system, i.e., the Laplacian matrix  $\mathcal{L}$  remains the same for the reduced multi-agent system, asymptotic stability is also preserved.

On the other hand, Reis and Stykel [RS07] prove the following sufficient condition for asymptotic stability of the interconnected system (5.2).

**Theorem 5.3** ([RS07, Corollary 2.3]):

Let  $\Phi_2 \in \mathbb{R}^{\mathbf{n} \times \mathbf{n}}$  be given by

$$\Phi_2 = \begin{bmatrix} \|K_{11}\|_2 & \|K_{12}\|_2 & \cdots & \|K_{1\mathbf{n}}\|_2 \\ \|K_{21}\|_2 & \|K_{22}\|_2 & \cdots & \|K_{2\mathbf{n}}\|_2 \\ \vdots & \vdots & \ddots & \vdots \\ \|K_{\mathbf{n}1}\|_2 & \|K_{\mathbf{n}2}\|_2 & \cdots & \|K_{\mathbf{n}\mathbf{n}}\|_2 \end{bmatrix}$$

and  $\Psi \in \mathbb{R}^{\mathbf{n} \times \mathbf{n}}$  by

$$\Psi = \Phi_2 \text{diag}(\|H_1\|_{\mathcal{H}_\infty}, \|H_2\|_{\mathcal{H}_\infty}, \dots, \|H_{\mathbf{n}}\|_{\mathcal{H}_\infty}).$$

If  $\rho(\Psi) < 1$ , then the system (5.2) is asymptotically stable.  $\diamond$

The similarity between the conditions in Lemma 5.2 and Theorem 5.3 motivates us to extend the use of BRBT from multi-agent systems to general interconnected systems. Notice that in the case of multi-agent systems, we have  $\Phi_2 = |\mathcal{L}|$ ,  $\Psi = |\mathcal{L}| \|H_1\|_{\mathcal{H}_\infty}$ , and  $\rho(\Psi) = \rho(|\mathcal{L}|) \|H_1\|_{\mathcal{H}_\infty}$ . Next, from [HJ85, Theorem 8.1.18], we have  $\rho(|\mathcal{L}|) \geq \rho(\mathcal{L})$ . Since  $\mathcal{L}$  is positive semi-definite, we have  $\rho(\mathcal{L}) = \|\mathcal{L}\|_2$  and thus  $\rho(|\mathcal{L}|) \geq \|\mathcal{L}\|_2$ . Therefore,  $\rho(\Psi) \geq \|\mathcal{L}\|_2 \|H_1\|_{\mathcal{H}_\infty}$ .

The following result allows us to extend the idea from [MTC13] for preserving the sufficient condition for asymptotic stability by applying BRBT to subsystems of the interconnected system (5.1).

**Proposition 5.4:**

Let  $\rho(\Psi) < 1$  and

$$\|\widehat{H}_i\|_{\mathcal{H}_\infty} < \frac{1}{\rho(\Psi)} \|H_i\|_{\mathcal{H}_\infty}, \text{ for } i = 1, 2, \dots, \mathbf{n}.$$

Then  $\rho(\widehat{\Psi}) < 1$ . ◇

*Proof.* We see that there must exist  $\alpha \in (0, 1)$  such that

$$\|\widehat{H}_i\|_{\mathcal{H}_\infty} \leq \alpha \frac{1}{\rho(\Psi)} \|H_i\|_{\mathcal{H}_\infty}, \text{ for } i = 1, 2, \dots, \mathbf{n}.$$

It follows that

$$\begin{aligned} \rho(\widehat{\Psi}) &= \rho(\Phi_2 \text{diag}(\|\widehat{H}_1\|_{\mathcal{H}_\infty}, \|\widehat{H}_2\|_{\mathcal{H}_\infty}, \dots, \|\widehat{H}_\mathbf{n}\|_{\mathcal{H}_\infty})) \\ &\leq \rho\left(\Phi_2 \text{diag}\left(\frac{\alpha}{\rho(\Psi)} \|H_1\|_{\mathcal{H}_\infty}, \frac{\alpha}{\rho(\Psi)} \|H_2\|_{\mathcal{H}_\infty}, \dots, \frac{\alpha}{\rho(\Psi)} \|H_\mathbf{n}\|_{\mathcal{H}_\infty}\right)\right) \\ &= \frac{\alpha}{\rho(\Psi)} \rho(\Phi_2 \text{diag}(\|H_1\|_{\mathcal{H}_\infty}, \|H_2\|_{\mathcal{H}_\infty}, \dots, \|H_\mathbf{n}\|_{\mathcal{H}_\infty})) \\ &< 1, \end{aligned}$$

where we used that  $\rho$  is increasing on nonnegative matrices ([HJ85, Corollary 8.1.19]) and homogeneous. □

Clearly,  $\|H_i\|_{\mathcal{H}_\infty} < \frac{1}{\rho(\Psi)} \|H_i\|_{\mathcal{H}_\infty}$  when  $\rho(\Psi) < 1$ . Thus, the idea is to apply BRBT to  $H_i$  with  $\gamma_i = \frac{1}{\rho(\Psi)} \|H_i\|_{\mathcal{H}_\infty}$ . Using Theorem 5.1, we find that

$$\|H_i - \widehat{H}_i\|_{\mathcal{H}_\infty} \leq 2 \sum_{j=r_i+1}^{n_i} \xi_j^{(i)}(\gamma_i), \quad (5.5)$$

where  $r_i$  is the order of  $\widehat{H}_i$  and  $\xi_1^{(i)}(\gamma_i), \xi_2^{(i)}(\gamma_i), \dots, \xi_{n_i}^{(i)}(\gamma_i)$  are the bounded real characteristic values of  $H_i$ . From Corollary 4.2 in [RS07], we find the upper bound for the  $\mathcal{H}_\infty$ -error from reducing subsystems by BRBT in the interconnected system (5.1).

**Theorem 5.5:**

Denote

$$\delta := \max_{i=1,2,\dots,\mathbf{n}} 2 \sum_{j=r_i+1}^{n_i} \xi_j^{(i)}(\gamma_i),$$

the maximum of the bounds from (5.5). Furthermore, let

$$\begin{aligned} g &= \|K(I - HK)^{-1}\|_{\mathcal{H}_\infty}, \quad g_1 = \|G(I - HK)^{-1}\|_{\mathcal{H}_\infty}, \quad g_2 = \|(I - KH)^{-1}F\|_{\mathcal{H}_\infty}, \\ c_1 &= g_1(\|F\|_2 + g \|HF\|_{\mathcal{H}_\infty}), \quad \text{and } c_2 = g_2(\|R\|_2 + g \|GH\|_{\mathcal{H}_\infty}). \end{aligned} \quad (5.6)$$

If  $g\delta < 1$ , then

$$\|\mathcal{H} - \widehat{\mathcal{H}}\|_{\mathcal{H}_\infty} \leq \min\{c_1, c_2\} \frac{\delta}{1 - g\delta}. \quad (5.7)$$

◇

**Algorithm 5.2:** Balancing method for network systems preserving stability and structure

**Input:** Network system  $(\mathcal{E}; \mathcal{A}, \mathcal{B}, \mathcal{C})$  from (5.2) with  $\rho(\Psi) < 1$ , tolerance  $\varepsilon > 0$ .

**Output:** Asymptotically stable reduced network system  $(\widehat{\mathcal{E}}; \widehat{\mathcal{A}}, \widehat{\mathcal{B}}, \widehat{\mathcal{C}})$  with

$$\|\mathcal{H} - \widehat{\mathcal{H}}\|_{\mathcal{H}_\infty} \leq \varepsilon.$$

- 1 Compute  $\gamma_i = \frac{1}{\rho(\Psi)} \|H_i\|_{\mathcal{H}_\infty}$  for  $i = 1, 2, \dots, \mathbf{n}$ .
- 2 Compute bounded real characteristic values  $\xi_j^{(i)}(\gamma_i)$  for  $i = 1, 2, \dots, \mathbf{n}$ ,  
 $j = 1, 2, \dots, n_i$ .
- 3 Compute  $g$ ,  $c_1$ , and  $c_2$  from (5.6).
- 4 Find minimal  $r_i$  such that

$$2 \sum_{j=r_i+1}^{n_i} \xi_j^{(i)}(\gamma_i) \leq \frac{\varepsilon}{\min\{c_1, c_2\} + g\varepsilon}$$

for  $i = 1, 2, \dots, \mathbf{n}$ .

- 5 Reduce the  $i$ th subsystem to order  $r_i$  using BRBT with  $\gamma = \gamma_i$  (see Algorithm 5.1), for  $i = 1, 2, \dots, \mathbf{n}$ .
- 

Theorem 5.5 enables us to adaptively choose reduced orders  $r_1, r_2, \dots, r_n$  of the subsystems. To see this, let  $\varepsilon > 0$  be a tolerance for the error  $\|\mathcal{H} - \widehat{\mathcal{H}}\|_{\mathcal{H}_\infty}$ . From (5.7), we see that the bound  $\|\mathcal{H} - \widehat{\mathcal{H}}\|_{\mathcal{H}_\infty} \leq \varepsilon$  can be achieved by enforcing

$$\min\{c_1, c_2\} \frac{\delta}{1 - g\delta} \leq \varepsilon,$$

which is equivalent to

$$\delta \leq \frac{\varepsilon}{\min\{c_1, c_2\} + g\varepsilon}.$$

Finding minimal  $r_i$  such that the right hand side in (5.5) is less than or equal to  $\frac{\varepsilon}{\min\{c_1, c_2\} + g\varepsilon}$  guarantees that  $\|\mathcal{H} - \widehat{\mathcal{H}}\|_{\mathcal{H}_\infty} \leq \varepsilon$ . Note that  $g \cdot \frac{\varepsilon}{\min\{c_1, c_2\} + g\varepsilon} < 1$ , which is the assumption in Theorem 5.5.

## 5.2.4 Numerical example

We use the interconnected string-beam example from [RS07, Section 6], illustrated in Figure 5.1. After discretizing the associated partial differential equation by finite differences, the string subsystem has  $n_1 = 1006$  states,  $m_1 = 3$  inputs, and  $p_1 = 2$  outputs, while the beam subsystem has  $n_2 = 1006$  states,  $m_2 = 2$  inputs, and  $p_2 = 2$  outputs. The interconnected system has  $n = n_1 + n_2 = 2012$  states,  $m = 1$  input, and  $p = 2$  outputs. The left figure in Figure 5.2 shows magnitude plots of the subsystems and the interconnected system.

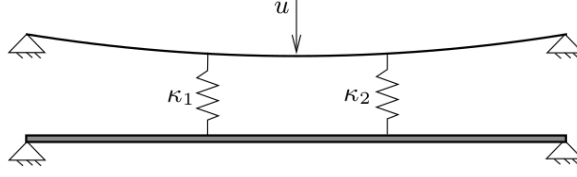


Figure 5.1: Interconnected string-beam example from [RS07]

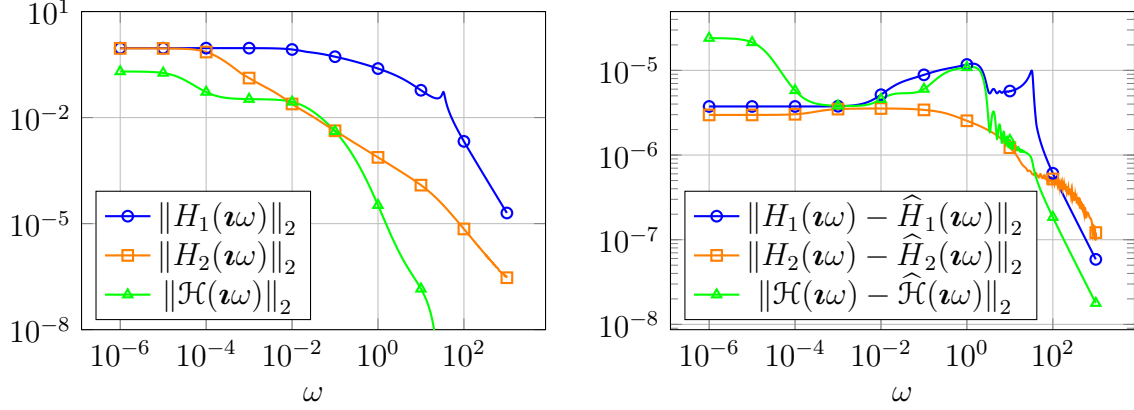


Figure 5.2: Magnitude plot of the full-order and error systems

For this system, we have  $\rho(\Psi) = 0.9157 < 1$ , which guarantees the interconnected system is asymptotically stable by Theorem 5.3. Furthermore, this allows us to use the method described in the previous section.

Let us set the reduction error tolerance for the coupled system to  $\varepsilon = 10^{-3}$ . From this, we find that the error tolerance for each of the subsystems is  $\frac{\varepsilon}{\min\{c_1, c_2\} + g\varepsilon} = 8.5821 \cdot 10^{-5}$ .

Applying BRBT from Algorithm 5.1, with  $\gamma_i = \frac{\|H_i\|_{\mathcal{H}_\infty}}{\rho(\Psi)}$  and adaptive reduced order decision, we determine reduced orders to be  $r_1 = 30$  and  $r_2 = 14$  for the string and beam subsystems respectively. We find that  $\|H_1 - \hat{H}_1\|_{\mathcal{H}_\infty} = 1.2001 \cdot 10^{-5}$ ,  $\|H_2 - \hat{H}_2\|_{\mathcal{H}_\infty} = 3.544 \cdot 10^{-6}$ , and  $\|H - \hat{H}\|_{\mathcal{H}_\infty} = 2.3893 \cdot 10^{-5} \leq 10^{-3}$ . The right figure in Figure 5.2 shows magnitude plots of error systems.

### 5.3 $\mathcal{H}_2$ -optimal subsystem reduction

Here, we are interested in subsystem reduction for interconnected systems (5.1) which minimizes the  $\mathcal{H}_2$ -error of the overall network system. We also consider multi-agent systems, where the additional constraint is that all subsystems are equal.

Notice that one special case of subsystem reduction for interconnected systems is weighted MOR [VVD08, Section 4], for which there are interpolatory  $\mathcal{H}_2$ -optimality conditions [ABGA13, BBG15]. We will focus on the general case.

### 5.3.1 Interconnected systems

We use the same notation as in the previous section. We want to find a reduced model with reduced subsystems

$$\begin{aligned}\widehat{E}_i \dot{\widehat{x}}_i(t) &= \widehat{A}_i \widehat{x}_i(t) + \widehat{B}_i \widehat{u}_i(t), \\ \widehat{y}_i(t) &= \widehat{C}_i \widehat{x}_i(t),\end{aligned}\tag{5.8a}$$

but preserved interconnections and external input

$$\widehat{u}_i(t) = K_{i1} \widehat{y}_1(t) + \cdots + K_{in} \widehat{y}_n(t) + F_i u(t),\tag{5.8b}$$

and external output

$$\widehat{y}(t) = G_1 \widehat{y}_1(t) + \cdots + G_n \widehat{y}_n(t),\tag{5.8c}$$

where  $\widehat{E}_i, \widehat{A}_i \in \mathbb{R}^{r_i \times r_i}$ ,  $\widehat{B} \in \mathbb{R}^{r_i \times m_i}$ ,  $\widehat{C} \in \mathbb{R}^{p_i \times r_i}$ , for some  $r_i < n_i$ ,  $i = 1, 2, \dots, n$ . Similarly as for the original system, the reduced system (5.8) can be represented by

$$\begin{aligned}\widehat{\mathcal{E}} \dot{\widehat{x}}(t) &= \widehat{\mathcal{A}} \widehat{x}(t) + \widehat{\mathcal{B}} u(t), \\ \widehat{y}(t) &= \widehat{\mathcal{C}} \widehat{x}(t),\end{aligned}\tag{5.9}$$

with  $\widehat{x}(t) = \text{col}(\widehat{x}_1(t), \widehat{x}_2(t), \dots, \widehat{x}_n(t))$  and

$$\begin{aligned}\widehat{\mathcal{E}} &= \widehat{E}, \quad \widehat{\mathcal{A}} = \widehat{A} + \widehat{B} K \widehat{C}, \quad \widehat{\mathcal{B}} = \widehat{B} F, \quad \widehat{\mathcal{C}} = G \widehat{C}, \\ \widehat{E} &= \text{diag}(\widehat{E}_1, \widehat{E}_2, \dots, \widehat{E}_n), \quad \widehat{A} = \text{diag}(\widehat{A}_1, \widehat{A}_2, \dots, \widehat{A}_n), \\ \widehat{B} &= \text{diag}(\widehat{B}_1, \widehat{B}_2, \dots, \widehat{B}_n), \quad \widehat{C} = \text{diag}(\widehat{C}_1, \widehat{C}_2, \dots, \widehat{C}_n).\end{aligned}\tag{5.10}$$

To find Gramian-based optimality conditions, we need to find gradients of the squared  $\mathcal{H}_2$ -error

$$\mathcal{J}(\widehat{E}_i, \widehat{A}_i, \widehat{B}_i, \widehat{C}_i) = \|H - \widehat{H}\|_{\mathcal{H}_2}^2.\tag{5.11}$$

With the Gramian-based formulation of the  $\mathcal{H}_2$ -norm, we can think of  $\mathcal{J}$  as a function  $f: X \rightarrow \mathbb{R}$  given implicitly by  $f(x) = g(x, y)$ , where  $h(x, y) = 0$  and  $D_y h(x, y)$  is bijective for all  $x$  and  $y$ . Here,  $x$  would represent the matrices of the ROM and  $y$  the Gramians. The following lemma states that we can use Lagrange multiplier method to compute  $Df(x)$ .

**Lemma 5.6:**

Let  $X, Y, Z$  be Banach spaces,  $U$  an open subset of  $X$ ,  $x_0$  an element of  $X$ , and  $f: U \rightarrow Z$ ,  $g: U \times Y \rightarrow Z$ ,  $h: U \times Y \rightarrow Y$  functions. Furthermore, suppose there exists  $y_0 \in Y$  such that  $h(x_0, y_0) = 0$ ,  $g$  is Fréchet differentiable at  $(x_0, y_0)$ ,  $D_y h$  exists on  $U \times Y$ ,  $D_y h(x_0, y_0)$  is bijective, and  $h$  is of class  $C^1$  in a neighborhood of  $(x_0, y_0)$ . Additionally, let  $\mathcal{L}: U \times Y \times B(Y, Z) \rightarrow Z$  be such that  $\mathcal{L}(x, y, \lambda) = g(x, y) - \lambda(h(x, y))$ .

Then  $f$  is Fréchet differentiable at  $x_0$  and

$$Df(x_0) = D_x \mathcal{L}(x_0, y_0, \lambda_0)$$

where  $\lambda_0$  is such that  $D_y \mathcal{L}(x_0, y_0, \lambda_0) = 0$ . ◇



*Proof.* Using Theorem 2.19 on  $h$ , we get that there exists an open subset  $U_{x_0}$  of  $U$  and a function  $k: U_{x_0} \rightarrow Y$  such that  $k(x_0) = y_0$ ,  $h(x, k(x)) = 0$  for  $x \in U_{x_0}$ , and  $Dk(x_0) = -D_y h(x_0, y_0)^{-1} D_x h(x_0, y_0)$ . Then

$$\begin{aligned} Df(x_0) &= D_x g(x_0, y_0) + D_y g(x_0, y_0) Dk(x_0) \\ &= D_x g(x_0, y_0) - D_y g(x_0, y_0) D_y h(x_0, y_0)^{-1} D_x h(x_0, y_0). \end{aligned}$$

On the other hand, we have

$$D_x \mathcal{L}(x_0, y_0, \lambda_0) = D_x g(x_0, y_0) - \lambda_0 D_x h(x_0, y_0), \quad (5.12a)$$

$$D_y \mathcal{L}(x_0, y_0, \lambda_0) = D_y g(x_0, y_0) - \lambda_0 D_y h(x_0, y_0). \quad (5.12b)$$

Therefore, from (5.12b) we get  $\lambda_0 = D_y g(x_0, y_0) D_y h(x_0, y_0)^{-1}$  and inserting it into (5.12a), we find

$$D_x \mathcal{L}(x_0, y_0, \lambda_0) = D_x g(x_0, y_0) - D_y g(x_0, y_0) D_y h(x_0, y_0)^{-1} D_x h(x_0, y_0),$$

which completes the proof.  $\square$

The following theorem gives the gradients of the squared  $\mathcal{H}_2$ -error (5.11).

**Theorem 5.7:**

Let (5.2) and (5.9) be an asymptotically stable systems. Then

$$\begin{aligned} \nabla_{\tilde{E}_i} \mathcal{J} &= 2 \sum_{j=1}^n \tilde{Q}_{ji}^T A_j \tilde{P}_{ji} + 2 \sum_{j,k=1}^n \tilde{Q}_{ji}^T B_j K_{jk} C_k \tilde{P}_{ki} \\ &\quad + 2 \sum_{j=1}^n \hat{Q}_{ij} \hat{A}_j \hat{P}_{ji} + 2 \sum_{j,k=1}^n \hat{Q}_{ij} \hat{B}_j K_{jk} \hat{C}_k \hat{P}_{ki}, \\ \nabla_{\hat{A}_i} \mathcal{J} &= 2 \sum_{j=1}^n \tilde{Q}_{ji}^T E_j \tilde{P}_{ji} + 2 \sum_{j=1}^n \hat{Q}_{ij} \hat{E}_j \hat{P}_{ji}, \\ \nabla_{\hat{B}_i} \mathcal{J} &= 2 \sum_{j,k=1}^n \tilde{Q}_{ji}^T E_j \tilde{P}_{jk} \hat{C}_k^T K_{ik}^T + 2 \sum_{j=1}^n \tilde{Q}_{ji}^T B_j F_j F_i^T \\ &\quad + 2 \sum_{j,k=1}^n \hat{Q}_{ij} \hat{E}_j \hat{P}_{jk} \hat{C}_k^T K_{ik}^T + 2 \sum_{j=1}^n \hat{Q}_{ji}^T \hat{B}_j F_j F_i^T, \\ \nabla_{\hat{C}_i} \mathcal{J} &= -2 \sum_{j=1}^n G_i^T G_j C_j \tilde{P}_{ji} + 2 \sum_{j=1}^n G_i^T G_j \hat{C}_j \hat{P}_{ji} \\ &\quad + 2 \sum_{j,k=1}^n K_{ji}^T \hat{B}_j^T \tilde{Q}_{kj}^T E_k \tilde{P}_{ki} + 2 \sum_{j,k=1}^n K_{ji}^T \hat{B}_j^T \hat{Q}_{kj}^T \hat{E}_k \hat{P}_{ki}. \quad \diamond \end{aligned}$$

*Proof.* We proceed as in the proof of Theorem 2.44. We find similarly that the Lagrange function is

$$\begin{aligned}\mathcal{L} &= \text{tr}\left(\mathcal{C}\mathcal{P}\mathcal{C}^T - 2\mathcal{C}\tilde{\mathcal{P}}\hat{\mathcal{C}}^T + \hat{\mathcal{C}}\tilde{\mathcal{P}}\hat{\mathcal{C}}^T\right) \\ &\quad + \text{tr}\left(2\tilde{\mathcal{Q}}^T\mathcal{A}\tilde{\mathcal{P}}\hat{\mathcal{E}}^T + 2\tilde{\mathcal{Q}}^T\mathcal{E}\tilde{\mathcal{P}}\hat{\mathcal{A}}^T + 2\tilde{\mathcal{Q}}^T\mathcal{B}\hat{\mathcal{B}}^T\right) \\ &\quad + \text{tr}\left(2\hat{\mathcal{Q}}\hat{\mathcal{A}}\hat{\mathcal{P}}\hat{\mathcal{E}}^T + \hat{\mathcal{Q}}\hat{\mathcal{B}}\hat{\mathcal{B}}^T\right),\end{aligned}$$

and after inserting (5.3), (5.10),

$$\begin{aligned}\mathcal{L} &= \text{tr}\left(G\mathcal{C}\mathcal{P}\mathcal{C}^T G^T - 2G\mathcal{C}\tilde{\mathcal{P}}\hat{\mathcal{C}}^T G^T + G\hat{\mathcal{C}}\tilde{\mathcal{P}}\hat{\mathcal{C}}^T G^T\right) \\ &\quad + \text{tr}\left(2\tilde{\mathcal{Q}}^T(A + BKC)\tilde{\mathcal{P}}\hat{\mathcal{E}}^T + 2\tilde{\mathcal{Q}}^T E\tilde{\mathcal{P}}\left(\hat{A} + \hat{B}K\hat{C}\right)^T + 2\tilde{\mathcal{Q}}^T BFF^T \hat{B}^T\right) \\ &\quad + \text{tr}\left(2\hat{\mathcal{Q}}\left(\hat{A} + \hat{B}K\hat{C}\right)\hat{\mathcal{P}}\hat{\mathcal{E}}^T + \hat{\mathcal{Q}}\hat{B}FF^T \hat{B}^T\right).\end{aligned}$$

Now we find the gradients

$$\begin{aligned}\nabla_{\hat{E}_i} \mathcal{L} &= 2 \sum_{j=1}^n \tilde{Q}_{ji}^T A_j \tilde{P}_{ji} + 2 \sum_{j,k=1}^n \tilde{Q}_{ji}^T B_j K_{jk} C_k \tilde{P}_{ki} \\ &\quad + 2 \sum_{j=1}^n \hat{Q}_{ij} \hat{A}_j \hat{P}_{ji} + 2 \sum_{j,k=1}^n \hat{Q}_{ij} \hat{B}_j K_{jk} \hat{C}_k \hat{P}_{ki}, \\ \nabla_{\hat{A}_i} \mathcal{L} &= 2 \sum_{j=1}^n \tilde{Q}_{ji}^T E_j \tilde{P}_{ji} + 2 \sum_{j=1}^n \hat{Q}_{ij} \hat{E}_j \hat{P}_{ji}, \\ \nabla_{\hat{B}_i} \mathcal{L} &= 2 \sum_{j,k=1}^n \tilde{Q}_{ji}^T E_j \tilde{P}_{jk} \hat{C}_k^T K_{ik}^T + 2 \sum_{j=1}^n \tilde{Q}_{ji}^T B_j F_j F_i^T \\ &\quad + 2 \sum_{j,k=1}^n \hat{Q}_{ij} \hat{E}_j \hat{P}_{jk} \hat{C}_k^T K_{ik}^T + 2 \sum_{j=1}^n \hat{Q}_{ji}^T \hat{B}_j F_j F_i^T, \\ \nabla_{\hat{C}_i} \mathcal{L} &= -2 \sum_{j=1}^n G_i^T G_j C_j \tilde{P}_{ji} + 2 \sum_{j=1}^n G_i^T G_j \hat{C}_j \hat{P}_{ji} \\ &\quad + 2 \sum_{j,k=1}^n K_{ji}^T \hat{B}_j^T \tilde{Q}_{kj}^T E_k \tilde{P}_{ki} + 2 \sum_{j,k=1}^n K_{ji}^T \hat{B}_j^T \hat{Q}_{kj}^T \hat{E}_k \hat{P}_{ki}.\end{aligned}$$

Using Lemma 5.6 completes the proof.  $\square$

The following corollary gives Wilson-type necessary optimality conditions for  $\mathcal{H}_2$ -optimal subsystem reduction for interconnected systems.

**Corollary 5.8:**

Let (5.9) be an  $\mathcal{H}_2$ -optimal reduced network systems for (5.2). Then

$$\begin{aligned}
 0 &= \sum_{j=1}^n \tilde{Q}_{ji}^T E_j \tilde{P}_{ji} + \sum_{j=1}^n \hat{Q}_{ij} \hat{E}_j \hat{P}_{ji}, \\
 0 &= \sum_{j=1}^n \tilde{Q}_{ji}^T A_j \tilde{P}_{ji} + \sum_{j=1}^n \hat{Q}_{ij} \hat{A}_j \hat{P}_{ji} \\
 &\quad + \sum_{j,k=1}^n \tilde{Q}_{ji}^T B_j K_{jk} C_k \tilde{P}_{ki} + \sum_{j,k=1}^n \hat{Q}_{ij} \hat{B}_j K_{jk} \hat{C}_k \hat{P}_{ki}, \\
 0 &= \sum_{j=1}^n \tilde{Q}_{ji}^T B_j F_j F_i^T + \sum_{j=1}^n \hat{Q}_{ji}^T \hat{B}_j F_j F_i^T \\
 &\quad + \sum_{j,k=1}^n \tilde{Q}_{ji}^T E_j \tilde{P}_{jk} \hat{C}_k^T K_{ik}^T + \sum_{j,k=1}^n \hat{Q}_{ij} \hat{E}_j \hat{P}_{jk} \hat{C}_k^T K_{ik}^T, \\
 0 &= - \sum_{j=1}^n G_i^T G_j C_j \tilde{P}_{ji} + \sum_{j=1}^n G_i^T G_j \hat{C}_j \hat{P}_{ji} \\
 &\quad + \sum_{j,k=1}^n K_{ji}^T \hat{B}_j^T \tilde{Q}_{kj}^T E_k \tilde{P}_{ki} + \sum_{j,k=1}^n K_{ji}^T \hat{B}_j^T \hat{Q}_{kj}^T \hat{E}_k \hat{P}_{ki}. \quad \diamond
 \end{aligned}$$

### 5.3.2 Multi-agent systems

Here, we consider multi-agent systems of the form

$$\begin{aligned}
 \mathcal{E} \dot{x}(t) &= \mathcal{A}x(t) + \mathcal{B}u(t), \\
 y(t) &= \mathcal{C}x(t),
 \end{aligned} \tag{5.13a}$$

with

$$\mathcal{E} = \mathfrak{M} \otimes E, \quad \mathcal{A} = \mathfrak{M} \otimes A - \mathfrak{L} \otimes BKC, \quad \mathcal{B} = \mathfrak{B} \otimes B, \quad \mathcal{C} = \mathfrak{C} \otimes C, \tag{5.13b}$$

where  $E, A \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times m}$ ,  $C \in \mathbb{R}^{p \times n}$  and  $K \in \mathbb{R}^{m \times p}$ . We want to find a ROM with reduced agents

$$\begin{aligned}
 \hat{\mathcal{E}} \dot{\hat{x}}(t) &= \hat{\mathcal{A}}\hat{x}(t) + \hat{\mathcal{B}}u(t), \\
 \hat{y}(t) &= \hat{\mathcal{C}}\hat{x}(t),
 \end{aligned} \tag{5.14a}$$

with

$$\hat{\mathcal{E}} = \mathfrak{M} \otimes \hat{E}, \quad \hat{\mathcal{A}} = \mathfrak{M} \otimes \hat{A} - \mathfrak{L} \otimes \hat{B}K\hat{C}, \quad \hat{\mathcal{B}} = \mathfrak{B} \otimes \hat{B}, \quad \hat{\mathcal{C}} = \mathfrak{C} \otimes \hat{C}, \tag{5.14b}$$

where  $\hat{E}, \hat{A} \in \mathbb{R}^{r \times r}$ ,  $\hat{B} \in \mathbb{R}^{r \times m}$ , and  $\hat{C} \in \mathbb{R}^{p \times r}$  for some  $r < n$ .

We proceed similarly to the previous section.

**Theorem 5.9:**

Let (5.13) and (5.14) be asymptotically stable systems. Then for the squared  $\mathcal{H}_2$ -error, we have

$$\begin{aligned}
 \nabla_{\widehat{E}}\mathcal{J} &= 2 \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widetilde{Q}_{ji}^T A \widetilde{P}_{ji} + 2 \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widehat{Q}_{ij} \widehat{A} \widehat{P}_{ji} \\
 &\quad - 2 \sum_{i,j,k=1}^n \mathbf{m}_i [\mathfrak{L}]_{jk} \widetilde{Q}_{ji}^T B K C \widetilde{P}_{ki} - 2 \sum_{i,j,k=1}^n \mathbf{m}_i [\mathfrak{L}]_{jk} \widehat{Q}_{ij} \widehat{B} K \widehat{C} \widehat{P}_{ki}, \\
 \nabla_{\widehat{A}}\mathcal{J} &= 2 \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widetilde{Q}_{ji}^T E \widetilde{P}_{ji} + 2 \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widehat{Q}_{ij} \widehat{E} \widehat{P}_{ji}, \\
 \nabla_{\widehat{B}}\mathcal{J} &= -2 \sum_{i,j,k=1}^n \mathbf{m}_j [\mathfrak{L}]_{ik} \widetilde{Q}_{ji}^T E \widetilde{P}_{jk} \widehat{C}^T K^T - 2 \sum_{i,j,k=1}^n \mathbf{m}_j [\mathfrak{L}]_{ik} \widehat{Q}_{ij} \widehat{E} \widehat{P}_{jk} \widehat{C}^T K^T \\
 &\quad + 2 \sum_{i,j=1}^n [\mathfrak{B} \mathfrak{B}^T]_{ji} \widetilde{Q}_{ji}^T B + 2 \sum_{i,j=1}^n [\mathfrak{B} \mathfrak{B}^T]_{ji} \widehat{Q}_{ij} \widehat{B}, \\
 \nabla_{\widehat{C}}\mathcal{J} &= -2K \widehat{B}^T \sum_{i,j,k=1}^n \mathbf{m}_j [\mathfrak{L}]_{ik} \widetilde{Q}_{ji}^T E \widetilde{P}_{jk} - 2K \widehat{B}^T \sum_{i,j,k=1}^n \mathbf{m}_j [\mathfrak{L}]_{ik} \widehat{Q}_{ij} \widehat{E} \widehat{P}_{jk} \\
 &\quad - 2C \sum_{j,k=1}^n [\mathfrak{e}^T \mathfrak{e}]_{jk} \widetilde{P}_{jk} + 2\widehat{C} \sum_{j,k=1}^n [\mathfrak{e}^T \mathfrak{e}]_{jk} \widehat{P}_{jk}. \quad \diamond
 \end{aligned}$$

*Proof.* As in the proof of 5.8, the Lagrange function is

$$\begin{aligned}
 \mathcal{L} &= \text{tr} \left( (\mathfrak{e} \otimes C) P (\mathfrak{e}^T \otimes C^T) - 2(\mathfrak{e} \otimes C) \widetilde{P} (\mathfrak{e}^T \otimes \widehat{C}^T) + (\mathfrak{e} \otimes \widehat{C}) \widehat{P} (\mathfrak{e}^T \otimes \widehat{C}^T) \right) \\
 &\quad + \text{tr} \left( 2\widetilde{Q}^T (\mathfrak{M} \otimes A - \mathfrak{L} \otimes B K C) \widetilde{P} (\mathfrak{M} \otimes \widehat{E}^T) \right) \\
 &\quad + \text{tr} \left( 2\widetilde{Q}^T (\mathfrak{M} \otimes E) \widetilde{P} (\mathfrak{M} \otimes \widehat{A}^T - \mathfrak{L}^T \otimes \widehat{C}^T K^T \widehat{B}^T) \right) \\
 &\quad + \text{tr} \left( 2\widetilde{Q}^T (\mathfrak{B} \otimes B) (\mathfrak{B}^T \otimes \widehat{B}^T) \right) \\
 &\quad + \text{tr} \left( 2\widehat{Q} (\mathfrak{M} \otimes \widehat{A} - \mathfrak{L} \otimes \widehat{B} K \widehat{C}) \widehat{P} (\mathfrak{M} \otimes \widehat{E}^T) + \widehat{Q} (\mathfrak{B} \otimes \widehat{B}) (\mathfrak{B}^T \otimes \widehat{B}^T) \right).
 \end{aligned}$$

The gradients are

$$\begin{aligned}
 \nabla_{\widehat{E}}\mathcal{L} &= 2 \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widetilde{Q}_{ji}^T A \widetilde{P}_{ji} + 2 \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widehat{Q}_{ij} \widehat{A} \widehat{P}_{ji} \\
 &\quad - 2 \sum_{i,j,k=1}^n \mathbf{m}_i [\mathfrak{L}]_{jk} \widetilde{Q}_{ji}^T B K C \widetilde{P}_{ki} - 2 \sum_{i,j,k=1}^n \mathbf{m}_i [\mathfrak{L}]_{jk} \widehat{Q}_{ij} \widehat{B} K \widehat{C} \widehat{P}_{ki},
 \end{aligned}$$

$$\begin{aligned}
 \nabla_{\widehat{A}} \mathcal{L} &= 2 \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widetilde{Q}_{ji}^T E \widetilde{P}_{ji} + 2 \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widehat{Q}_{ij} \widehat{E} \widehat{P}_{ji}, \\
 \nabla_{\widehat{B}} \mathcal{L} &= -2 \sum_{i,j,k=1}^n \mathbf{m}_j [\mathcal{L}]_{ik} \widetilde{Q}_{ji}^T E \widetilde{P}_{jk} \widehat{C}^T K^T - 2 \sum_{i,j,k=1}^n \mathbf{m}_j [\mathcal{L}]_{ik} \widehat{Q}_{ij} \widehat{E} \widehat{P}_{jk} \widehat{C}^T K^T \\
 &\quad + 2 \sum_{i,j=1}^n [\mathfrak{B} \mathfrak{B}^T]_{ji} \widetilde{Q}_{ji}^T B + 2 \sum_{i,j=1}^n [\mathfrak{B} \mathfrak{B}^T]_{ji} \widehat{Q}_{ij} \widehat{B}, \\
 \nabla_{\widehat{C}} \mathcal{L} &= -2K \widehat{B}^T \sum_{i,j,k=1}^n \mathbf{m}_j [\mathcal{L}]_{ik} \widetilde{Q}_{ji}^T E \widetilde{P}_{jk} - 2K \widehat{B}^T \sum_{i,j,k=1}^n \mathbf{m}_j [\mathcal{L}]_{ik} \widehat{Q}_{ij} \widehat{E} \widehat{P}_{jk} \\
 &\quad - 2C \sum_{j,k=1}^n [\mathbf{e}^T \mathbf{e}]_{jk} \widetilde{P}_{jk} + 2\widehat{C} \sum_{j,k=1}^n [\mathbf{e}^T \mathbf{e}]_{jk} \widehat{P}_{jk}. \quad \square
 \end{aligned}$$

As a direct consequence of the previous theorem, we get the Wilson-type necessary optimality conditions for  $\mathcal{H}_2$ -optimal MOR of multi-agent systems.

**Corollary 5.10:**

Let (5.14) be an  $\mathcal{H}_2$ -optimal reduced multi-agent system for (5.13). Then

$$\begin{aligned}
 0 &= \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widetilde{Q}_{ji}^T A \widetilde{P}_{ji} + \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widehat{Q}_{ij} \widehat{A} \widehat{P}_{ji} \\
 &\quad - \sum_{i,j,k=1}^n \mathbf{m}_i [\mathcal{L}]_{jk} \widetilde{Q}_{ji}^T B K C \widetilde{P}_{ki} - \sum_{i,j,k=1}^n \mathbf{m}_i [\mathcal{L}]_{jk} \widehat{Q}_{ij} \widehat{B} K \widehat{C} \widehat{P}_{ki}, \\
 0 &= \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widetilde{Q}_{ji}^T E \widetilde{P}_{ji} + \sum_{i,j=1}^n \mathbf{m}_i \mathbf{m}_j \widehat{Q}_{ij} \widehat{E} \widehat{P}_{ji}, \\
 0 &= - \sum_{i,j,k=1}^n \mathbf{m}_j [\mathcal{L}]_{ik} \widetilde{Q}_{ji}^T E \widetilde{P}_{jk} \widehat{C}^T K^T - \sum_{i,j,k=1}^n \mathbf{m}_j [\mathcal{L}]_{ik} \widehat{Q}_{ij} \widehat{E} \widehat{P}_{jk} \widehat{C}^T K^T \\
 &\quad + \sum_{i,j=1}^n [\mathfrak{B} \mathfrak{B}^T]_{ji} \widetilde{Q}_{ji}^T B + \sum_{i,j=1}^n [\mathfrak{B} \mathfrak{B}^T]_{ji} \widehat{Q}_{ij} \widehat{B}, \\
 0 &= -K \widehat{B}^T \sum_{i,j,k=1}^n \mathbf{m}_j [\mathcal{L}]_{ik} \widetilde{Q}_{ji}^T E \widetilde{P}_{jk} - K \widehat{B}^T \sum_{i,j,k=1}^n \mathbf{m}_j [\mathcal{L}]_{ik} \widehat{Q}_{ij} \widehat{E} \widehat{P}_{jk} \\
 &\quad - C \sum_{j,k=1}^n [\mathbf{e}^T \mathbf{e}]_{jk} \widetilde{P}_{jk} + \widehat{C} \sum_{j,k=1}^n [\mathbf{e}^T \mathbf{e}]_{jk} \widehat{P}_{jk}. \quad \diamond
 \end{aligned}$$

## 5.4 Conclusion

In Section 5.2, we developed a stability-preserving method consisting of applying BRBT to subsystems. The preservation of asymptotic stability is based on a sufficient condition. As such, the method can only be applied to asymptotically stable interconnected systems which additionally satisfy this sufficient condition. On the other hand, it could be possible in practical applications to design the interconnected systems to satisfy this condition. This would not only allow the use of the developed method, but also give an a priori guarantee for the asymptotic stability of the original interconnected system.

In Section 5.3, we derived Wilson-type necessary optimality conditions for  $\mathcal{H}_2$ -optimal subsystem reduction of interconnected and multi-agent systems. These conditions can be used with a gradient-based optimization, with backtracking line search to ensure asymptotic stability of the ROM, to find locally  $\mathcal{H}_2$ -optimal ROMs.

# CHAPTER 6

## $\mathcal{H}_2$ -OPTIMAL MODEL ORDER REDUCTION OF FURTHER STRUCTURED SYSTEMS

### Contents

---

6.1	Introduction . . . . .	113
6.2	Second-order systems . . . . .	114
6.2.1	Wilson-type conditions . . . . .	114
6.2.2	Interpolatory conditions . . . . .	118
6.3	Port-Hamiltonian systems . . . . .	123
6.3.1	Wilson-type conditions . . . . .	125
6.3.2	Interpolatory conditions . . . . .	126
6.4	Linear parametric systems . . . . .	128
6.4.1	$\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal model order reduction . . . . .	128
6.4.2	Interpolatory conditions . . . . .	131
6.5	Linear time-delay systems . . . . .	132
6.5.1	Wilson-type conditions . . . . .	134
6.6	Conclusion . . . . .	138

---

## 6.1 Introduction

In this chapter, we discuss  $\mathcal{H}_2$ -optimal MOR of structured systems which are not represented using networks as in the previous chapters. Just as before, we are interested in structure-preserving MOR, which can be beneficial for preserving the physical interpretation of the system. Additionally, simulation and optimization methods tailored for specific system structures can be used for the ROMs.

Second-order systems are considered in Section 6.2, port-Hamiltonian systems in Section 6.3, parametric systems in Section 6.4, and time-delay systems in Section 6.5.

## 6.2 Second-order systems

Here, we consider LTI second-order systems

$$\begin{aligned} M\ddot{x}(t) + E\dot{x}(t) + Kx(t) &= Bu(t), \\ y(t) &= C_p x(t) + C_v \dot{x}(t), \end{aligned} \quad (6.1)$$

where  $M, E, K \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times m}$ ,  $C_p, C_v \in \mathbb{R}^{p \times n}$ ,  $x(t) \in \mathbb{R}^n$ ,  $u(t) \in \mathbb{R}^m$ , and  $y(t) \in \mathbb{R}^p$ . We assume that  $M$  is invertible and that the matrix pencil  $\lambda^2 M + \lambda E + K$  is asymptotically stable. These systems appear, e.g., when analyzing mechanical or electrical systems. We want to find a ROM of the same structure

$$\begin{aligned} \widehat{M}\ddot{\hat{x}}(t) + \widehat{E}\dot{\hat{x}}(t) + \widehat{K}\hat{x}(t) &= \widehat{B}u(t), \\ \hat{y}(t) &= \widehat{C}_p \hat{x}(t) + \widehat{C}_v \dot{\hat{x}}(t), \end{aligned} \quad (6.2)$$

where  $\widehat{M}, \widehat{E}, \widehat{K} \in \mathbb{R}^{r \times r}$ ,  $\widehat{B} \in \mathbb{R}^{r \times m}$ , and  $\widehat{C}_p, \widehat{C}_v \in \mathbb{R}^{p \times r}$ , with  $r \ll n$ , such that  $\|H - \widehat{H}\|_{\mathcal{H}_2}$  is minimized. In particular, the matrix pencil  $\lambda^2 \widehat{M} + \lambda \widehat{E} + \widehat{K}$  should be asymptotically stable.

There is some work towards  $\mathcal{H}_2$ -optimal MOR of second-order systems. Beattie and Gugercin [BG09] showed it is possible to find a structured ROM which interpolates the original model. Wyatt [Wya12] proposed several iterative methods based on IRKA using the result from [BG09], but  $\mathcal{H}_2$ -optimality was not investigated. Beattie and Benner [BB14] derived interpolation-based necessary  $\mathcal{H}_2$ -optimality conditions for second-order systems, but the ROM was restricted to be modally-damped, i.e., such that  $\widehat{M}^{-1}\widehat{E}$  and  $\widehat{M}^{-1}\widehat{K}$  are simultaneously diagonalizable. Additionally, finding an algorithm which would achieve the interpolation conditions remains an open problem. Sato [Sat17] developed a method based on Riemannian optimization for  $\mathcal{H}_2$ -optimal MOR of second-order systems with symmetric positive definite mass, damping, and stiffness matrices. We consider a more general setting here.

There are several balancing-based methods which preserve the second-order structure [MS96, CLVVD06, RS08a]. They can also find ROMs with small  $\mathcal{H}_2$ -errors, but there is no a priori error bound as for the standard BT.

We derive Wilson-type necessary  $\mathcal{H}_2$ -optimality conditions for second-order systems in Section 6.2.1 and corresponding interpolatory conditions in Section 6.2.2.

### 6.2.1 Wilson-type conditions

In [BB14, Section 5], the authors focus on interpolation-based necessary optimality conditions for  $\mathcal{H}_2$ -optimal MOR of second-order systems. Here, we use the Gramian-based approach to derive  $\mathcal{H}_2$ -optimality conditions similar to Wilson conditions for first-order systems. In the next section, we derive interpolatory conditions and compare them to the results from [BB14, Section 5].



As is well known, the second-order system (6.1) has an equivalent first-order representation

$$\underbrace{\begin{bmatrix} I & 0 \\ 0 & M \end{bmatrix}}_{\mathcal{E}} \dot{z}(t) = \underbrace{\begin{bmatrix} 0 & I \\ -K & -E \end{bmatrix}}_{\mathcal{A}} z(t) + \underbrace{\begin{bmatrix} 0 \\ B \end{bmatrix}}_{\mathcal{B}} u(t),$$

$$y(t) = \underbrace{\begin{bmatrix} C_p & C_v \end{bmatrix}}_{\mathcal{C}} z(t),$$

where  $z(t) = \text{col}(x(t), \dot{x}(t))$ ,  $\mathcal{E}, \mathcal{A} \in \mathbb{R}^{2n \times 2n}$ ,  $\mathcal{B} \in \mathbb{R}^{2n \times m}$ , and  $\mathcal{C} \in \mathbb{R}^{p \times 2n}$ . Let  $\mathcal{P} \in \mathbb{R}^{2n \times 2n}$  and  $\mathcal{Q} \in \mathbb{R}^{2n \times 2n}$  be the controllability and observability Gramians of this system, i.e., solutions to the following Lyapunov equations:

$$\mathcal{A}\mathcal{P}\mathcal{E}^T + \mathcal{E}\mathcal{P}\mathcal{A}^T + \mathcal{B}\mathcal{B}^T = 0,$$

$$\mathcal{A}^T\mathcal{Q}\mathcal{E} + \mathcal{E}^T\mathcal{Q}\mathcal{A} + \mathcal{C}^T\mathcal{C} = 0.$$

An equivalent first-order representation of (6.2) is

$$\underbrace{\begin{bmatrix} I & 0 \\ 0 & \widehat{M} \end{bmatrix}}_{\widehat{\mathcal{E}}} \dot{\widehat{z}}(t) = \underbrace{\begin{bmatrix} 0 & I \\ -\widehat{K} & -\widehat{E} \end{bmatrix}}_{\widehat{\mathcal{A}}} \widehat{z}(t) + \underbrace{\begin{bmatrix} 0 \\ \widehat{B} \end{bmatrix}}_{\widehat{\mathcal{B}}} u(t),$$

$$\widehat{y}(t) = \underbrace{\begin{bmatrix} \widehat{C}_p & \widehat{C}_v \end{bmatrix}}_{\widehat{\mathcal{C}}} \widehat{z}(t),$$

with  $\widehat{\mathcal{E}}, \widehat{\mathcal{A}} \in \mathbb{R}^{2r \times 2r}$ ,  $\widehat{\mathcal{B}} \in \mathbb{R}^{2r \times m}$ ,  $\widehat{\mathcal{C}} \in \mathbb{R}^{p \times 2r}$ , and Gramians  $\widehat{\mathcal{P}}, \widehat{\mathcal{Q}} \in \mathbb{R}^{2r \times 2r}$  satisfying

$$\widehat{\mathcal{A}}\widehat{\mathcal{P}}\widehat{\mathcal{E}}^T + \widehat{\mathcal{E}}\widehat{\mathcal{P}}\widehat{\mathcal{A}}^T + \widehat{\mathcal{B}}\widehat{\mathcal{B}}^T = 0, \quad (6.3a)$$

$$\widehat{\mathcal{A}}^T\widehat{\mathcal{Q}}\widehat{\mathcal{E}} + \widehat{\mathcal{E}}^T\widehat{\mathcal{Q}}\widehat{\mathcal{A}} + \widehat{\mathcal{C}}^T\widehat{\mathcal{C}} = 0. \quad (6.3b)$$

The error system is

$$\underbrace{\begin{bmatrix} \mathcal{E} & 0 \\ 0 & \widehat{\mathcal{E}} \end{bmatrix}}_{\mathcal{E}_{\text{err}}} \begin{bmatrix} \dot{z}(t) \\ \dot{\widehat{z}}(t) \end{bmatrix} = \underbrace{\begin{bmatrix} \mathcal{A} & 0 \\ 0 & \widehat{\mathcal{A}} \end{bmatrix}}_{\mathcal{A}_{\text{err}}} \begin{bmatrix} z(t) \\ \widehat{z}(t) \end{bmatrix} + \underbrace{\begin{bmatrix} \mathcal{B} \\ \widehat{\mathcal{B}} \end{bmatrix}}_{\mathcal{B}_{\text{err}}} u(t)$$

$$y(t) - \widehat{y}(t) = \underbrace{\begin{bmatrix} \mathcal{C} & -\widehat{\mathcal{C}} \end{bmatrix}}_{\mathcal{C}_{\text{err}}} \begin{bmatrix} z(t) \\ \widehat{z}(t) \end{bmatrix},$$

with Gramians

$$\mathcal{P}_{\text{err}} = \begin{bmatrix} \mathcal{P} & \widetilde{\mathcal{P}} \\ \widetilde{\mathcal{P}}^T & \widehat{\mathcal{P}} \end{bmatrix}, \quad \mathcal{Q}_{\text{err}} = \begin{bmatrix} \mathcal{Q} & \widetilde{\mathcal{Q}} \\ \widetilde{\mathcal{Q}}^T & \widehat{\mathcal{Q}} \end{bmatrix},$$

where  $\tilde{\mathcal{P}} \in \mathbb{R}^{2n \times 2r}$  and  $\tilde{\mathcal{Q}} \in \mathbb{R}^{2n \times 2r}$  satisfy Sylvester equations

$$\mathcal{A}\tilde{\mathcal{P}}\hat{\mathcal{E}}^T + \varepsilon\tilde{\mathcal{P}}\hat{\mathcal{A}}^T + \mathcal{B}\hat{\mathcal{B}}^T = 0, \quad (6.4a)$$

$$\mathcal{A}^T\tilde{\mathcal{Q}}\hat{\mathcal{E}} + \varepsilon^T\tilde{\mathcal{Q}}\hat{\mathcal{A}} - \mathcal{C}^T\hat{\mathcal{C}} = 0. \quad (6.4b)$$

We denote

$$\tilde{\mathcal{P}} = \begin{bmatrix} \tilde{P}_{pp} & \tilde{P}_{pv} \\ \tilde{P}_{vp} & \tilde{P}_{vv} \end{bmatrix}, \quad \hat{\mathcal{P}} = \begin{bmatrix} \hat{P}_{pp} & \hat{P}_{pv} \\ \hat{P}_{vp} & \hat{P}_{vv} \end{bmatrix}, \quad \tilde{\mathcal{Q}} = \begin{bmatrix} \tilde{Q}_{pp} & \tilde{Q}_{pv} \\ \tilde{Q}_{vp} & \tilde{Q}_{vv} \end{bmatrix}, \quad \hat{\mathcal{Q}} = \begin{bmatrix} \hat{Q}_{pp} & \hat{Q}_{pv} \\ \hat{Q}_{vp} & \hat{Q}_{vv} \end{bmatrix},$$

where  $\tilde{P}_{ij}, \tilde{Q}_{ij} \in \mathbb{R}^{n \times r}$  and  $\hat{P}_{ij}, \hat{Q}_{ij} \in \mathbb{R}^{r \times r}$  for  $i, j \in \{p, v\}$ .

We find the squared  $\mathcal{H}_2$ -error as in Theorem 5.7.

**Theorem 6.1:**

Consider the second-order system (6.1). Let (6.2) be asymptotically stable. Then for the squared  $\mathcal{H}_2$ -error  $\mathcal{J}$ , we have

$$\begin{aligned} \nabla_{\hat{M}}\mathcal{J} &= -2\tilde{Q}_{vv}^T K \tilde{P}_{pv} + 2\tilde{Q}_{pv}^T \tilde{P}_{vv} - 2\tilde{Q}_{vv}^T E \tilde{P}_{vv} - 2\hat{Q}_{vv} \hat{K} \hat{P}_{pv} + 2\hat{Q}_{pv}^T \hat{P}_{vv} - 2\hat{Q}_{vv} \hat{E} \hat{P}_{vv}, \\ \nabla_{\hat{E}}\mathcal{J} &= -2\tilde{Q}_{pv}^T \tilde{P}_{pv} - 2\tilde{Q}_{vv}^T M \tilde{P}_{vv} - 2\hat{Q}_{pv}^T \hat{P}_{pv} - 2\hat{Q}_{vv} \hat{M} \hat{P}_{vv}, \\ \nabla_{\hat{K}}\mathcal{J} &= -2\tilde{Q}_{pv}^T \tilde{P}_{pp} - 2\tilde{Q}_{vv}^T M \tilde{P}_{vp} - 2\hat{Q}_{pv}^T \hat{P}_{pp} - 2\hat{Q}_{vv} \hat{M} \hat{P}_{vp}, \\ \nabla_{\hat{B}}\mathcal{J} &= 2\tilde{Q}_{vv}^T B + 2\hat{Q}_{vv} \hat{B}, \\ \nabla_{\hat{C}_p}\mathcal{J} &= -2C_p \tilde{P}_{pp} - 2C_v \tilde{P}_{vp} + 2\hat{C}_p \hat{P}_{pp} + 2\hat{C}_v \hat{P}_{vp}, \\ \nabla_{\hat{C}_v}\mathcal{J} &= -2C_p \tilde{P}_{pv} - 2C_v \tilde{P}_{vv} + 2\hat{C}_p \hat{P}_{pv} + 2\hat{C}_v \hat{P}_{vv}. \end{aligned} \quad \diamond$$

*Proof.* Similar to the proof of Theorem 5.7, the Lagrange function is

$$\begin{aligned} \mathcal{L} &= \text{tr}(\mathcal{C}\mathcal{P}\mathcal{C}^T - 2\mathcal{C}\tilde{\mathcal{P}}\hat{\mathcal{C}}^T + \hat{\mathcal{C}}\hat{\mathcal{P}}\hat{\mathcal{C}}^T) \\ &\quad + \text{tr}(\tilde{\Lambda}^T \mathcal{A}\tilde{\mathcal{P}}\hat{\mathcal{E}}^T + \tilde{\Lambda}^T \varepsilon\tilde{\mathcal{P}}\hat{\mathcal{A}}^T + \tilde{\Lambda}^T \mathcal{B}\hat{\mathcal{B}}^T) \\ &\quad + \text{tr}(\hat{\Lambda}^T \hat{\mathcal{A}}\hat{\mathcal{P}}\hat{\mathcal{E}}^T + \hat{\Lambda}^T \varepsilon\hat{\mathcal{P}}\hat{\mathcal{A}}^T + \hat{\Lambda}^T \hat{\mathcal{B}}\hat{\mathcal{B}}^T), \end{aligned}$$

where  $\tilde{\Lambda} \in \mathbb{R}^{2n \times 2r}$  and  $\hat{\Lambda} \in \mathbb{R}^{2r \times 2r}$  are the Lagrange multipliers. From (6.4b), (6.3b), and

$$\begin{aligned} \nabla_{\tilde{\mathcal{P}}}\mathcal{L} &= -2\mathcal{C}^T\hat{\mathcal{C}} + \mathcal{A}^T\tilde{\Lambda}\hat{\mathcal{E}} + \varepsilon^T\tilde{\Lambda}\hat{\mathcal{A}}, \\ \nabla_{\hat{\mathcal{P}}}\mathcal{L} &= \hat{\mathcal{C}}^T\hat{\mathcal{C}} + \hat{\mathcal{A}}^T\hat{\Lambda}\hat{\mathcal{E}} + \hat{\mathcal{E}}^T\hat{\Lambda}\hat{\mathcal{A}}, \end{aligned}$$

it follows that  $\tilde{\Lambda} = 2\tilde{Q}$  and  $\hat{\Lambda} = \hat{Q}$ . Inserting this into the Lagrange function, we find

$$\begin{aligned}
 \mathcal{L} &= \text{tr} \left( \mathcal{C}\mathcal{P}\mathcal{C}^T - 2 [C_p \ C_v] \tilde{\mathcal{P}} \begin{bmatrix} \hat{C}_p & \hat{C}_v \end{bmatrix}^T + \begin{bmatrix} \hat{C}_p & \hat{C}_v \end{bmatrix} \hat{\mathcal{P}} \begin{bmatrix} \hat{C}_p & \hat{C}_v \end{bmatrix}^T \right) \\
 &+ \text{tr} \left( 2\tilde{Q}^T \begin{bmatrix} 0 & I \\ -K & -E \end{bmatrix} \tilde{\mathcal{P}} \begin{bmatrix} I & 0 \\ 0 & \hat{M} \end{bmatrix}^T + 2\tilde{Q}^T \begin{bmatrix} I & 0 \\ 0 & M \end{bmatrix} \tilde{\mathcal{P}} \begin{bmatrix} 0 & I \\ -\hat{K} & -\hat{E} \end{bmatrix}^T \right) \\
 &+ \text{tr} \left( 2\tilde{Q}^T \begin{bmatrix} 0 \\ B \end{bmatrix} \begin{bmatrix} 0 \\ \hat{B} \end{bmatrix}^T \right) \\
 &+ \text{tr} \left( 2\hat{Q} \begin{bmatrix} 0 & I \\ -\hat{K} & -\hat{E} \end{bmatrix} \hat{\mathcal{P}} \begin{bmatrix} I & 0 \\ 0 & \hat{M} \end{bmatrix}^T + \hat{Q} \begin{bmatrix} 0 \\ \hat{B} \end{bmatrix} \begin{bmatrix} 0 \\ \hat{B} \end{bmatrix}^T \right) \\
 &= \text{tr} \left( \mathcal{C}\mathcal{P}\mathcal{C}^T - 2C_p\tilde{P}_{pp}\hat{C}_p^T - 2C_p\tilde{P}_{pv}\hat{C}_v^T - 2C_v\tilde{P}_{vp}\hat{C}_p^T - 2C_v\tilde{P}_{vv}\hat{C}_v^T \right) \\
 &+ \text{tr} \left( \hat{C}_p\hat{P}_{pp}\hat{C}_p^T + 2\hat{C}_p\hat{P}_{pv}\hat{C}_v^T + \hat{C}_v\hat{P}_{vv}\hat{C}_v^T \right) \\
 &+ \text{tr} \left( -2\tilde{Q}_{vp}^T K \tilde{P}_{pp} + 2\tilde{Q}_{pp}^T \tilde{P}_{vp} - 2\tilde{Q}_{vp}^T E \tilde{P}_{vp} \right) \\
 &+ \text{tr} \left( -2\tilde{Q}_{vv}^T K \tilde{P}_{pv} \hat{M}^T + 2\tilde{Q}_{pv}^T \tilde{P}_{vv} \hat{M}^T - 2\tilde{Q}_{vv}^T E \tilde{P}_{vv} \hat{M}^T \right) \\
 &+ \text{tr} \left( 2\tilde{Q}_{pp}^T \tilde{P}_{pv} + 2\tilde{Q}_{vp}^T M \tilde{P}_{vv} \right) \\
 &+ \text{tr} \left( -2\tilde{Q}_{pv}^T \tilde{P}_{pp} \hat{K}^T - 2\tilde{Q}_{pv}^T \tilde{P}_{pv} \hat{E}^T - 2\tilde{Q}_{vv}^T M \tilde{P}_{vp} \hat{K}^T - 2\tilde{Q}_{vv}^T M \tilde{P}_{vv} \hat{E}^T \right) \\
 &+ \text{tr} \left( 2\tilde{Q}_{vv}^T B \hat{B}^T \right) \\
 &+ \text{tr} \left( -2\hat{Q}_{vp}^T \hat{K} \hat{P}_{pp} + 2\hat{Q}_{pp}^T \hat{P}_{vp} - 2\hat{Q}_{vp}^T \hat{E} \hat{P}_{vp} \right) \\
 &+ \text{tr} \left( -2\hat{Q}_{vv}^T \hat{K} \hat{P}_{pv} \hat{M}^T + 2\hat{Q}_{pv}^T \hat{P}_{vv} \hat{M}^T - 2\hat{Q}_{vv}^T \hat{E} \hat{P}_{vv} \hat{M}^T \right) \\
 &+ \text{tr} \left( \hat{Q}_{vv}^T \hat{B} \hat{B}^T \right).
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 \nabla_{\hat{M}} \mathcal{L} &= -2\tilde{Q}_{vv}^T K \tilde{P}_{pv} + 2\tilde{Q}_{pv}^T \tilde{P}_{vv} - 2\tilde{Q}_{vv}^T E \tilde{P}_{vv} - 2\hat{Q}_{vv} \hat{K} \hat{P}_{pv} + 2\hat{Q}_{pv}^T \hat{P}_{vv} - 2\hat{Q}_{vv} \hat{E} \hat{P}_{vv}, \\
 \nabla_{\hat{E}} \mathcal{L} &= -2\tilde{Q}_{pv}^T \tilde{P}_{pv} - 2\tilde{Q}_{vv}^T M \tilde{P}_{vv} - 2\tilde{Q}_{pv}^T \hat{P}_{pv} - 2\hat{Q}_{vv} \hat{M} \hat{P}_{vv}, \\
 \nabla_{\hat{K}} \mathcal{L} &= -2\tilde{Q}_{pv}^T \tilde{P}_{pp} - 2\tilde{Q}_{vv}^T M \tilde{P}_{vp} - 2\tilde{Q}_{pv}^T \hat{P}_{pp} - 2\hat{Q}_{vv} \hat{M} \hat{P}_{vp}, \\
 \nabla_{\hat{B}} \mathcal{L} &= 2\tilde{Q}_{vv}^T B + 2\hat{Q}_{vv} \hat{B}, \\
 \nabla_{\hat{C}_p} \mathcal{L} &= -2C_p \tilde{P}_{pp} - 2C_v \tilde{P}_{vp} + 2\hat{C}_p \hat{P}_{pp} + 2\hat{C}_v \hat{P}_{vp}, \\
 \nabla_{\hat{C}_v} \mathcal{L} &= -2C_p \tilde{P}_{pv} - 2C_v \tilde{P}_{vv} + 2\hat{C}_p \hat{P}_{pv} + 2\hat{C}_v \hat{P}_{vv}.
 \end{aligned}$$

The statement follows from Lemma 5.6.  $\square$

As a direct consequence, we obtain the Wilson-type necessary optimality conditions.

**Corollary 6.2:**

Consider the second-order system (6.1). Let (6.2) be its locally  $\mathcal{H}_2$ -optimal ROM. Then

$$\begin{aligned}
 & \begin{bmatrix} \tilde{Q}_{pv} \\ \tilde{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} I & 0 \\ 0 & M \end{bmatrix} \tilde{\mathcal{P}} + \begin{bmatrix} \hat{Q}_{pv} \\ \hat{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} I & 0 \\ 0 & \hat{M} \end{bmatrix} \hat{\mathcal{P}} = 0, \\
 & \begin{bmatrix} \tilde{Q}_{pv} \\ \tilde{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} 0 & I \\ -K & -E \end{bmatrix} \begin{bmatrix} \tilde{P}_{pv} \\ \tilde{P}_{vv} \end{bmatrix} + \begin{bmatrix} \hat{Q}_{pv} \\ \hat{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} 0 & I \\ -\hat{K} & -\hat{E} \end{bmatrix} \begin{bmatrix} \hat{P}_{pv} \\ \hat{P}_{vv} \end{bmatrix} = 0, \\
 & \begin{bmatrix} \tilde{Q}_{pv} \\ \tilde{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} 0 \\ B \end{bmatrix} + \begin{bmatrix} \hat{Q}_{pv} \\ \hat{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} 0 \\ \hat{B} \end{bmatrix} = 0, \\
 & [C_p \ C_v] \tilde{\mathcal{P}} - [\hat{C}_p \ \hat{C}_v] \hat{\mathcal{P}} = 0. \quad \diamond
 \end{aligned}$$

We can also consider second-order systems (6.1) with  $C_v = 0$ , where one would then want enforce  $\hat{C}_v = 0$  in (6.2).

**Corollary 6.3:**

Consider the second-order system (6.1) with  $\hat{C}_v = 0$ . Let (6.2) be its locally  $\mathcal{H}_2$ -optimal ROM with  $\hat{C}_v = 0$ . Then

$$\begin{aligned}
 & \begin{bmatrix} \tilde{Q}_{pv} \\ \tilde{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} I & 0 \\ 0 & M \end{bmatrix} \tilde{\mathcal{P}} + \begin{bmatrix} \hat{Q}_{pv} \\ \hat{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} I & 0 \\ 0 & \hat{M} \end{bmatrix} \hat{\mathcal{P}} = 0, \\
 & \begin{bmatrix} \tilde{Q}_{pv} \\ \tilde{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} 0 & I \\ -K & -E \end{bmatrix} \begin{bmatrix} \tilde{P}_{pv} \\ \tilde{P}_{vv} \end{bmatrix} + \begin{bmatrix} \hat{Q}_{pv} \\ \hat{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} 0 & I \\ -\hat{K} & -\hat{E} \end{bmatrix} \begin{bmatrix} \hat{P}_{pv} \\ \hat{P}_{vv} \end{bmatrix} = 0, \\
 & \tilde{Q}_{vv}^T B + \hat{Q}_{vv}^T \hat{B} = 0, \\
 & C_p \tilde{P}_{pp} - \hat{C}_p \hat{P}_{pp} = 0. \quad \diamond
 \end{aligned}$$

*Proof.* Following the proof of Theorem 6.1, we see that it is enough to ignore  $\nabla_{\hat{C}_v} \mathcal{L}$  and replace  $C_v$  and  $\hat{C}_v$  with zero.  $\square$

Similarly, we can consider the case of  $C_p = 0$  and  $\hat{C}_p = 0$  and the results can be obtained as in Corollary 6.3.

## 6.2.2 Interpolatory conditions

The following theorem shows how to derive interpolatory conditions using Wilson-type conditions, similar to Theorem 2.45. As in [BB14, Section 5], we assume that the ROM is such that  $\hat{M}^{-1}\hat{E}$  and  $\hat{M}^{-1}\hat{K}$  are simultaneously diagonalizable. This is the case for second-order, modally damped dynamical systems.

**Theorem 6.4:**

Consider the second-order system (6.1). Let (6.2) be a locally  $\mathcal{H}_2$ -optimal ROM, such that  $\widehat{M}^{-1}\widehat{E}$  and  $\widehat{M}^{-1}\widehat{K}$  are simultaneously diagonalizable, i.e., there exist invertible  $S$  and  $T$  such that  $S^T\widehat{M}T = I_r$ ,  $S^T\widehat{E}T = -(\Lambda_1 + \Lambda_2)$ ,  $S^T\widehat{K}T = \Lambda_1\Lambda_2$ , where  $\Lambda_1 = \text{diag}(\lambda_{1,i})$  and  $\Lambda_2 = \text{diag}(\lambda_{2,i})$ . Additionally, let all  $\lambda_{1,i}$  and  $\lambda_{2,j}$  be pairwise distinct. Denote  $t_i = Te_i$ ,  $c_{p,i} = \widehat{C}_p t_i$ ,  $c_{v,i} = \widehat{C}_v t_i$ ,  $s_i = Se_i$ , and  $b_i^T = s_i^T \widehat{B}$ , which gives us

$$\widehat{H}(s) = \sum_{i=1}^r \frac{(c_{p,i} + sc_{v,i})b_i^T}{(s - \lambda_{1,i})(s - \lambda_{2,i})}.$$

Then,

$$\begin{aligned} (H(-\lambda_{1,i}) - H(-\lambda_{2,i}))b_i &= \left(\widehat{H}(-\lambda_{1,i}) - \widehat{H}(-\lambda_{2,i})\right)b_i, \\ (\lambda_{1,i}H(-\lambda_{1,i}) - \lambda_{2,i}H(-\lambda_{2,i}))b_i &= \left(\lambda_{1,i}\widehat{H}(-\lambda_{1,i}) - \lambda_{2,i}\widehat{H}(-\lambda_{2,i})\right)b_i, \\ (c_{p,i} + \lambda_{1,i}c_{v,i})^T (H(-\lambda_{1,i}) - H(-\lambda_{2,i})) &= (c_{p,i} + \lambda_{1,i}c_{v,i})^T \left(\widehat{H}(-\lambda_{1,i}) - \widehat{H}(-\lambda_{2,i})\right), \\ (c_{p,i} + \lambda_{1,i}c_{v,i})^T H'(-\lambda_{1,i})b_i &= (c_{p,i} + \lambda_{1,i}c_{v,i})^T \widehat{H}'(-\lambda_{1,i})b_i, \\ (c_{p,i} + \lambda_{2,i}c_{v,i})^T H'(-\lambda_{2,i})b_i &= (c_{p,i} + \lambda_{2,i}c_{v,i})^T \widehat{H}'(-\lambda_{2,i})b_i, \end{aligned}$$

for  $i = 1, 2, \dots, r$ . ◇

*Proof.* From the assumptions, we have  $\widehat{E}t_i = -(\lambda_{1,i} + \lambda_{2,i})\widehat{M}t_i$ ,  $s_i^T \widehat{E} = -(\lambda_{1,i} + \lambda_{2,i})s_i^T \widehat{M}$ ,  $\widehat{K}t_i = \lambda_{1,i}\lambda_{2,i}\widehat{M}t_i$ , and  $s_i^T \widehat{K} = \lambda_{1,i}\lambda_{2,i}s_i^T \widehat{M}$ . Postmultiplying (6.4a) by

$$\begin{bmatrix} \widehat{M}^T s_i & 0 \\ 0 & s_i \end{bmatrix} \begin{bmatrix} -\lambda_{2,i} & -\lambda_{1,i} \\ 1 & 1 \end{bmatrix}$$

and using

$$\begin{bmatrix} 0 & -\lambda_{1,i}\lambda_{2,i} \\ 1 & \lambda_{1,i} + \lambda_{2,i} \end{bmatrix} \begin{bmatrix} -\lambda_{2,i} & -\lambda_{1,i} \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} -\lambda_{2,i} & -\lambda_{1,i} \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \lambda_{1,i} & 0 \\ 0 & \lambda_{2,i} \end{bmatrix},$$

we find

$$\begin{aligned} \mathcal{A}\widetilde{\mathcal{P}} \begin{bmatrix} -\lambda_{2,i}\widehat{M}^T s_i & -\lambda_{1,i}\widehat{M}^T s_i \\ \widehat{M}^T s_i & \widehat{M}^T s_i \end{bmatrix} + \mathcal{E}\widetilde{\mathcal{P}} \begin{bmatrix} -\lambda_{2,i}\widehat{M}^T s_i & -\lambda_{1,i}\widehat{M}^T s_i \\ \widehat{M}^T s_i & \widehat{M}^T s_i \end{bmatrix} \begin{bmatrix} \lambda_{1,i} & 0 \\ 0 & \lambda_{2,i} \end{bmatrix} \\ + \mathcal{B} \begin{bmatrix} b_i & b_i \end{bmatrix} = 0. \end{aligned}$$

This gives us

$$\begin{aligned} -\lambda_{2,i} \begin{bmatrix} \widetilde{P}_{pp} \\ \widetilde{P}_{vp} \end{bmatrix} \widehat{M}^T s_i + \begin{bmatrix} \widetilde{P}_{pv} \\ \widetilde{P}_{vv} \end{bmatrix} \widehat{M}^T s_i &= (-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i, \\ -\lambda_{1,i} \begin{bmatrix} \widetilde{P}_{pp} \\ \widetilde{P}_{vp} \end{bmatrix} \widehat{M}^T s_i + \begin{bmatrix} \widetilde{P}_{pv} \\ \widetilde{P}_{vv} \end{bmatrix} \widehat{M}^T s_i &= (-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i, \end{aligned}$$

which implies

$$\begin{aligned}
 (\lambda_{1,i} - \lambda_{2,i}) \begin{bmatrix} \tilde{P}_{pp} \\ \tilde{P}_{vp} \end{bmatrix} \widehat{M}^\top s_i &= (-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i - (-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i, \\
 (\lambda_{1,i} - \lambda_{2,i}) \begin{bmatrix} \tilde{P}_{pv} \\ \tilde{P}_{vv} \end{bmatrix} \widehat{M}^\top s_i &= \lambda_{1,i} (-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i - \lambda_{2,i} (-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i.
 \end{aligned} \tag{6.5}$$

Similarly, we find

$$\begin{aligned}
 (\lambda_{1,i} - \lambda_{2,i}) \begin{bmatrix} \widehat{P}_{pp} \\ \widehat{P}_{vp} \end{bmatrix} \widehat{M}^\top s_i &= \left(-\lambda_{1,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}}\right)^{-1} \widehat{\mathcal{B}}b_i - \left(-\lambda_{2,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}}\right)^{-1} \widehat{\mathcal{B}}b_i, \\
 (\lambda_{1,i} - \lambda_{2,i}) \begin{bmatrix} \widehat{P}_{pv} \\ \widehat{P}_{vv} \end{bmatrix} \widehat{M}^\top s_i &= \lambda_{1,i} \left(-\lambda_{1,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}}\right)^{-1} \widehat{\mathcal{B}}b_i - \lambda_{2,i} \left(-\lambda_{2,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}}\right)^{-1} \widehat{\mathcal{B}}b_i.
 \end{aligned} \tag{6.6}$$

Analogously, postmultiplying (6.4b) by

$$\begin{bmatrix} t_i & 0 \\ 0 & t_i \end{bmatrix} \begin{bmatrix} 1 & 1 \\ \lambda_{1,i} & \lambda_{2,i} \end{bmatrix}$$

and using

$$\begin{bmatrix} 0 & 1 \\ -\lambda_{1,i}\lambda_{2,i} & \lambda_{1,i} + \lambda_{2,i} \end{bmatrix} \begin{bmatrix} 1 & 1 \\ \lambda_{1,i} & \lambda_{2,i} \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ \lambda_{1,i} & \lambda_{2,i} \end{bmatrix} \begin{bmatrix} \lambda_{1,i} & 0 \\ 0 & \lambda_{2,i} \end{bmatrix},$$

we find

$$\begin{aligned}
 \mathcal{A}^\top \widetilde{Q} \begin{bmatrix} t_i & t_i \\ \lambda_{1,i} \widehat{M}t_i & \lambda_{2,i} \widehat{M}t_i \end{bmatrix} \begin{bmatrix} t_i & 0 \\ 0 & \widehat{M}t_i \end{bmatrix} + \mathcal{E}^\top \widetilde{Q} \begin{bmatrix} t_i & t_i \\ \lambda_{1,i} \widehat{M}t_i & \lambda_{2,i} \widehat{M}t_i \end{bmatrix} \begin{bmatrix} \lambda_{1,i} & 0 \\ 0 & \lambda_{2,i} \end{bmatrix} \\
 - \mathcal{C}^\top \begin{bmatrix} c_{p,i} + \lambda_{1,i}c_{v,i} & c_{p,i} + \lambda_{2,i}c_{v,i} \end{bmatrix} &= 0.
 \end{aligned}$$

This in turn gives us

$$\begin{aligned}
 \begin{bmatrix} \widetilde{Q}_{pp} \\ \widetilde{Q}_{vp} \end{bmatrix} t_i + \lambda_{1,i} \begin{bmatrix} \widetilde{Q}_{pv} \\ \widetilde{Q}_{vv} \end{bmatrix} \widehat{M}t_i &= -(-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-\top} \mathcal{C}^\top (c_{p,i} + \lambda_{1,i}c_{v,i}), \\
 \begin{bmatrix} \widetilde{Q}_{pp} \\ \widetilde{Q}_{vp} \end{bmatrix} t_i + \lambda_{2,i} \begin{bmatrix} \widetilde{Q}_{pv} \\ \widetilde{Q}_{vv} \end{bmatrix} \widehat{M}t_i &= -(-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-\top} \mathcal{C}^\top (c_{p,i} + \lambda_{2,i}c_{v,i}),
 \end{aligned}$$

which implies

$$\begin{aligned}
 (\lambda_{1,i} - \lambda_{2,i}) \begin{bmatrix} \widetilde{Q}_{pv} \\ \widetilde{Q}_{vv} \end{bmatrix} \widehat{M}t_i &= -(-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-\top} \mathcal{C}^\top (c_{p,i} + \lambda_{1,i}c_{v,i}) \\
 &\quad + (-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-\top} \mathcal{C}^\top (c_{p,i} + \lambda_{2,i}c_{v,i}).
 \end{aligned} \tag{6.7}$$

We similarly find

$$\begin{aligned}
 (\lambda_{1,i} - \lambda_{2,i}) \begin{bmatrix} \widehat{Q}_{pv} \\ \widehat{Q}_{vv} \end{bmatrix} \widehat{M}t_i &= \left( -\lambda_{1,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}} \right)^{-\text{T}} \widehat{\mathcal{C}}^{\text{T}}(c_{p,i} + \lambda_{1,i}c_{v,i}) \\
 &\quad - \left( -\lambda_{2,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}} \right)^{-\text{T}} \widehat{\mathcal{C}}^{\text{T}}(c_{p,i} + \lambda_{2,i}c_{v,i}).
 \end{aligned} \tag{6.8}$$

Now we use the Wilson-type conditions from Corollary 6.2 to derive the interpolatory conditions. Postmultiplying

$$[C_p \ C_v] \widetilde{\mathcal{P}} - \begin{bmatrix} \widehat{C}_p & \widehat{C}_v \end{bmatrix} \widehat{\mathcal{P}} = 0,$$

by  $\text{diag}(\widehat{M}^{\text{T}}s_i, \widehat{M}^{\text{T}}s_i)$  and using (6.5) and (6.6), we get

$$\begin{aligned}
 &\mathcal{C}(-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i - \mathcal{C}(-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i \\
 &\quad = \widehat{\mathcal{C}}(-\lambda_{1,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}})^{-1} \widehat{\mathcal{B}}b_i - \widehat{\mathcal{C}}(-\lambda_{2,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}})^{-1} \widehat{\mathcal{B}}b_i, \\
 &\lambda_{1,i}\mathcal{C}(-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i - \lambda_{2,i}\mathcal{C}(-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i \\
 &\quad = \lambda_{1,i}\widehat{\mathcal{C}}(-\lambda_{1,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}})^{-1} \widehat{\mathcal{B}}b_i - \lambda_{2,i}\widehat{\mathcal{C}}(-\lambda_{2,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}})^{-1} \widehat{\mathcal{B}}b_i,
 \end{aligned}$$

i.e.,

$$\begin{aligned}
 (H(-\lambda_{1,i}) - H(-\lambda_{2,i})) b_i &= \left( \widehat{H}(-\lambda_{1,i}) - \widehat{H}(-\lambda_{2,i}) \right) b_i, \\
 (\lambda_{1,i}H(-\lambda_{1,i}) - \lambda_{2,i}H(-\lambda_{2,i})) b_i &= \left( \lambda_{1,i}\widehat{H}(-\lambda_{1,i}) - \lambda_{2,i}\widehat{H}(-\lambda_{2,i}) \right) b_i,
 \end{aligned}$$

which are exactly the first two optimality conditions. Premultiplying

$$\begin{bmatrix} \widetilde{Q}_{pv} \\ \widetilde{Q}_{vv} \end{bmatrix}^{\text{T}} \begin{bmatrix} 0 \\ B \end{bmatrix} + \begin{bmatrix} \widehat{Q}_{pv} \\ \widehat{Q}_{vv} \end{bmatrix}^{\text{T}} \begin{bmatrix} 0 \\ \widehat{B} \end{bmatrix} = 0,$$

by  $t_i^{\text{T}} \widehat{M}^{\text{T}}$  and using (6.7) and (6.8), we get

$$\begin{aligned}
 &(c_{p,i} + \lambda_{1,i}c_{v,i})^{\text{T}} \mathcal{C}(-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B} - (c_{p,i} + \lambda_{2,i}c_{v,i})^{\text{T}} \mathcal{C}(-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B} \\
 &\quad = (c_{p,i} + \lambda_{1,i}c_{v,i})^{\text{T}} \widehat{\mathcal{C}}(-\lambda_{1,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}})^{-1} \widehat{\mathcal{B}} - (c_{p,i} + \lambda_{2,i}c_{v,i})^{\text{T}} \widehat{\mathcal{C}}(-\lambda_{2,i}\widehat{\mathcal{E}} - \widehat{\mathcal{A}})^{-1} \widehat{\mathcal{B}},
 \end{aligned}$$

i.e.,

$$(c_{p,i} + \lambda_{1,i}c_{v,i})^{\text{T}} (H(-\lambda_{1,i}) - H(-\lambda_{2,i})) = (c_{p,i} + \lambda_{1,i}c_{v,i})^{\text{T}} \left( \widehat{H}(-\lambda_{1,i}) - \widehat{H}(-\lambda_{2,i}) \right),$$

which is the third optimality condition. From

$$\begin{bmatrix} \tilde{Q}_{pv} \\ \tilde{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} I & 0 \\ 0 & M \end{bmatrix} \tilde{\mathcal{P}} + \begin{bmatrix} \hat{Q}_{pv} \\ \hat{Q}_{vv} \end{bmatrix}^T \begin{bmatrix} I & 0 \\ 0 & \hat{M} \end{bmatrix} \hat{\mathcal{P}} = 0,$$

we get

$$\begin{aligned} & ((c_{p,i} + \lambda_{1,i}c_{v,i})^T \mathcal{C} (-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} - (c_{p,i} + \lambda_{2,i}c_{v,i})^T \mathcal{C} (-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1}) \\ & \cdot \mathcal{E} ((-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i - (-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i) \\ & = \left( (c_{p,i} + \lambda_{1,i}c_{v,i})^T \hat{\mathcal{C}} \left( -\lambda_{1,i}\hat{\mathcal{E}} - \hat{\mathcal{A}} \right)^{-1} - (c_{p,i} + \lambda_{2,i}c_{v,i})^T \hat{\mathcal{C}} \left( -\lambda_{2,i}\hat{\mathcal{E}} - \hat{\mathcal{A}} \right)^{-1} \right) \\ & \cdot \hat{\mathcal{E}} \left( \left( -\lambda_{1,i}\hat{\mathcal{E}} - \hat{\mathcal{A}} \right)^{-1} \hat{\mathcal{B}}b_i - \left( -\lambda_{2,i}\hat{\mathcal{E}} - \hat{\mathcal{A}} \right)^{-1} \hat{\mathcal{B}}b_i \right), \\ & ((c_{p,i} + \lambda_{1,i}c_{v,i})^T \mathcal{C} (-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} - (c_{p,i} + \lambda_{2,i}c_{v,i})^T \mathcal{C} (-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1}) \\ & \cdot \mathcal{E} (\lambda_{1,i}(-\lambda_{1,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i - \lambda_{2,i}(-\lambda_{2,i}\mathcal{E} - \mathcal{A})^{-1} \mathcal{B}b_i) \\ & = \left( (c_{p,i} + \lambda_{1,i}c_{v,i})^T \hat{\mathcal{C}} \left( -\lambda_{1,i}\hat{\mathcal{E}} - \hat{\mathcal{A}} \right)^{-1} - (c_{p,i} + \lambda_{2,i}c_{v,i})^T \hat{\mathcal{C}} \left( -\lambda_{2,i}\hat{\mathcal{E}} - \hat{\mathcal{A}} \right)^{-1} \right) \\ & \cdot \hat{\mathcal{E}} \left( \lambda_{1,i} \left( -\lambda_{1,i}\hat{\mathcal{E}} - \hat{\mathcal{A}} \right)^{-1} \hat{\mathcal{B}}b_i - \lambda_{2,i} \left( -\lambda_{2,i}\hat{\mathcal{E}} - \hat{\mathcal{A}} \right)^{-1} \hat{\mathcal{B}}b_i \right), \end{aligned}$$

thus

$$\begin{aligned} & (c_{p,i} + \lambda_{1,i}c_{v,i})^T H'(-\lambda_{1,i})b_i + (c_{p,i} + \lambda_{2,i}c_{v,i})^T H'(-\lambda_{2,i})b_i \\ & = (c_{p,i} + \lambda_{1,i}c_{v,i})^T \hat{H}'(-\lambda_{1,i})b_i + (c_{p,i} + \lambda_{2,i}c_{v,i})^T \hat{H}'(-\lambda_{2,i})b_i, \\ & \lambda_{1,i}(c_{p,i} + \lambda_{1,i}c_{v,i})^T H'(-\lambda_{1,i})b_i + \lambda_{2,i}(c_{p,i} + \lambda_{2,i}c_{v,i})^T H'(-\lambda_{2,i})b_i \\ & = \lambda_{1,i}(c_{p,i} + \lambda_{1,i}c_{v,i})^T \hat{H}'(-\lambda_{1,i})b_i + \lambda_{2,i}(c_{p,i} + \lambda_{2,i}c_{v,i})^T \hat{H}'(-\lambda_{2,i})b_i, \end{aligned}$$

i.e., since  $\lambda_{1,i} \neq \lambda_{2,i}$ ,

$$\begin{aligned} & (c_{p,i} + \lambda_{1,i}c_{v,i})^T H'(-\lambda_{1,i})b_i = (c_{p,i} + \lambda_{1,i}c_{v,i})^T \hat{H}'(-\lambda_{1,i})b_i, \\ & (c_{p,i} + \lambda_{2,i}c_{v,i})^T H'(-\lambda_{2,i})b_i = (c_{p,i} + \lambda_{2,i}c_{v,i})^T \hat{H}'(-\lambda_{2,i})b_i, \end{aligned}$$

which are the final two interpolatory optimality conditions.  $\square$

Here are the interpolatory conditions with additional assumptions that  $C_v = 0$  and  $\hat{C}_v = 0$ , as derived in [BB14, Section 5] using a different approach.

**Corollary 6.5:**

Let the assumptions in Theorem 6.4 hold. Additionally, let  $C_v = 0$  and  $\hat{C}_v = 0$ . Then,

$$\begin{aligned} & (H(-\lambda_{1,i}) - H(-\lambda_{2,i}))b_i = \left( \hat{H}(-\lambda_{1,i}) - \hat{H}(-\lambda_{2,i}) \right)b_i, \\ & c_{p,i}^T (H(-\lambda_{1,i}) - H(-\lambda_{2,i})) = c_{p,i}^T \left( \hat{H}(-\lambda_{1,i}) - \hat{H}(-\lambda_{2,i}) \right), \\ & c_{p,i}^T H'(-\lambda_{1,i})b_i = c_{p,i}^T \hat{H}'(-\lambda_{1,i})b_i, \\ & c_{p,i}^T H'(-\lambda_{2,i})b_i = c_{p,i}^T \hat{H}'(-\lambda_{2,i})b_i, \end{aligned}$$



for  $i = 1, 2, \dots, r$ . ◇

*Proof.* Proceed as the proof of Theorem 6.4, using the Wilson-type conditions from Corollary 6.3. □

## 6.3 Port-Hamiltonian systems

We consider LTI, input-state-output port-Hamiltonian systems (see [vdSJ14]) without algebraic constraints of the form

$$\begin{aligned} \dot{x}(t) &= (J - R)Sx(t) + Bu(t), \\ y(t) &= B^T Sx(t), \end{aligned} \tag{6.9}$$

where,  $J, R, S \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times m}$ ,  $J^T = -J$ ,  $R \succcurlyeq 0$ , and  $S \succ 0$ . Matrix  $S$  is called the *energy matrix* and  $R$  the *dissipation matrix*.

### Remark 6.6:

The standard notation in literature for the energy matrix is  $Q$ . Since we use  $Q$  for the observability Gramian, we decided to use  $S$  to denote the energy matrix. ◇

The *Hamiltonian*  $\frac{1}{2}x(t)^T Sx(t)$  defines the total energy of the system. The state  $x(t) \in \mathbb{R}^n$  are called the *energy variables*, while  $u(t), y(t) \in \mathbb{R}^m$  the *power variables*. The inner product of the power variables  $u(t)^T y(t)$  is the power supplied to the system. The system is passive, i.e., the rate of increase of total energy is bounded by the power provided to the system, as seen from

$$\frac{d}{dt} \left( \frac{1}{2} x(t)^T Sx(t) \right) = u(t)^T y(t) - x(t)^T SRSx(t) \leq u(t)^T y(t).$$

We are looking for a ROM

$$\begin{aligned} \hat{x}(t) &= \left( \hat{J} - \hat{R} \right) \hat{S} \hat{x}(t) + \hat{B} u(t), \\ \hat{y}(t) &= \hat{B}^T \hat{S} \hat{x}(t), \end{aligned} \tag{6.10}$$

which preserves the port-Hamiltonian structure in (6.9) (in particular, such that  $\hat{J}^T = -\hat{J}$ ,  $\hat{R} \succcurlyeq 0$ ,  $\hat{S} \succ 0$ ) and minimizes the  $\mathcal{H}_2$ -error  $\|H - \hat{H}\|_{\mathcal{H}_2}$ , where  $H$  and  $\hat{H}$  are respectively the transfer function of (6.9) and (6.10). Additionally, we assume  $(J - R)S$  is Hurwitz, which is guaranteed when  $R \succ 0$ .

An equivalent form of the port-Hamiltonian system (6.9) is using the *co-energy variables*  $e(t) = Sx(t)$ , given by

$$\begin{aligned} S^{-1} \dot{e}(t) &= (J - R)e(t) + Bu(t), \\ y(t) &= B^T e(t). \end{aligned} \tag{6.11}$$

From this, we see that a system is port-Hamiltonian if it has an equivalent representation  $(\tilde{E}; \tilde{A}, \tilde{B}, \tilde{B}^T)$  with  $\tilde{E} \succ 0$  and  $\tilde{A} + \tilde{A}^T \preceq 0$ . Clearly, applying a Galerkin projection to (6.11) with some full-rank matrix  $V \in \mathbb{R}^{n \times r}$  will preserve the structure, since

$$(V^T J V)^T = -V^T J V, \quad V^T R V \succcurlyeq 0, \quad \text{and} \quad (V^T S^{-1} V)^{-1} \succ 0.$$

In [GPBvdS12], the authors propose projecting the matrices by

$$\hat{J} = W^T J W, \quad \hat{R} = W^T R W, \quad \hat{S} = V^T S V, \quad \text{and} \quad \hat{B} = W^T B,$$

where  $W = S V (V^T S V)^{-1}$ , giving the ROM

$$((W^T J W - W^T R W) V^T S V, W^T B, B^T W V^T S V). \quad (6.12)$$

It can be seen that system (6.12) is equivalent to projecting (6.11) by Galerkin projection  $V_G = S V$  and to projecting (6.9) by Petrov-Galerkin projection

$$\begin{aligned} V_{\text{PG}} &= V = V (V^T S V)^{-1} V^T S V, \\ W_{\text{PG}} &= S V (V^T S V)^{-1}, \end{aligned}$$

also showed in the proof of [GPBvdS12, Theorem 7].

If  $V$  is chosen such that

$$\text{im}(V) = \text{im}([( \sigma_1 I - (J - R) S )^{-1} B b_1 \quad \cdots \quad ( \sigma_r I - (J - R) S )^{-1} B b_r]),$$

then, according to Theorem 2.41, we have

$$H(\sigma_i) b_i = \hat{H}(\sigma_i) b_i, \quad i = 1, 2, \dots, r.$$

Based on this, [GPBvdS12] proposes an iterative algorithm similar to IRKA, called IRKA-PH, which finds a ROM with a transfer function

$$\hat{H}(s) = \sum_{i=1}^r \frac{c_i b_i^T}{s - \lambda_i}$$

which, upon convergence, satisfies

$$H(-\lambda_i) b_i = \hat{H}(-\lambda_i) b_i, \quad i = 1, 2, \dots, r.$$

This is one of the interpolatory necessary  $\mathcal{H}_2$ -optimality conditions for unstructured first-order systems (see Theorem 2.42). The idea is that IRKA-PH should give a ROM which is close to  $\mathcal{H}_2$ -optimal. But, it is not clear a priori how close it will be.

Theorem 1 in [BB14] shows the following interpolatory necessary  $\mathcal{H}_2$ -optimality conditions.

**Theorem 6.7 ([BB14, Theorem 1]):**

Suppose that  $\widehat{H}$  is an  $\mathcal{H}_2$ -optimal ROM with  $\widehat{R} \succ 0$ , has  $r$  distinct poles, and is represented as  $\widehat{H}(s) = \sum_{i=1}^r \frac{c_i b_i^T}{s - \lambda_i}$ . Then

$$\begin{aligned} c_i^T (H(-\lambda_i) - H(-\lambda_j)) b_j &= c_i^T \left( \widehat{H}(-\lambda_i) - \widehat{H}(-\lambda_j) \right) b_j, \\ c_i^T H'(-\lambda_i) b_i &= c_i^T \widehat{H}'(-\lambda_i) b_i, \end{aligned}$$

for  $i, j = 1, 2, \dots, r$ . ◇

Similar to the interpolatory necessary optimality conditions for second-order systems, an algorithm which would satisfy these conditions is not known.

In the following sections, we will find Wilson-type optimality conditions and derive interpolatory conditions as in Theorem 6.7.

### 6.3.1 Wilson-type conditions

Similar to Section 6.2.1, we first derive gradients of the squared  $\mathcal{H}_2$ -error.

**Theorem 6.8:**

Let (6.10) be asymptotically stable, with  $\widehat{S} = I$ ,  $\widehat{J} = \widehat{J}_2 - \widehat{J}_2^T$ , and  $\widehat{R} = \widehat{R}_2 \widehat{R}_2^T$ . Then for the squared  $\mathcal{H}_2$ -error  $\mathcal{J}$ , we have

$$\begin{aligned} \nabla_{\widehat{J}_2} \mathcal{J} &= 2\widetilde{Q}^T \widetilde{P} - 2\widetilde{P}^T \widetilde{Q} + 2\widehat{Q} \widehat{P} - 2\widehat{P} \widehat{Q}, \\ \nabla_{\widehat{R}_2} \mathcal{J} &= -2\widetilde{Q}^T \widetilde{P} \widehat{R}_2 - 2\widetilde{P}^T \widetilde{Q} \widehat{R}_2 - 2\widehat{Q} \widehat{P} \widehat{R}_2 - 2\widehat{P} \widehat{Q} \widehat{R}_2, \\ \nabla_{\widehat{B}} \mathcal{J} &= -2\widetilde{P}^T S B + 2\widehat{P} \widehat{B} + 2\widetilde{Q}^T B + 2\widehat{Q} \widehat{B}. \end{aligned} \quad \diamond$$

*Proof.* As in the proof of Theorem 2.44, we find the Lagrange function is

$$\begin{aligned} \mathcal{L} &= \text{tr} \left( B^T S P S B - 2B^T S \widetilde{P} \widehat{B} + \widehat{B}^T \widehat{P} \widehat{B} \right) \\ &\quad + \text{tr} \left( 2\widetilde{Q}^T (J - R) S \widetilde{P} + 2\widetilde{Q}^T \widetilde{P} \left( \widehat{J}_2^T - \widehat{J}_2 - \widehat{R}_2 \widehat{R}_2^T \right) + 2\widetilde{Q}^T B \widehat{B}^T \right) \\ &\quad + \text{tr} \left( 2\widehat{Q} \left( \widehat{J}_2 - \widehat{J}_2^T - \widehat{R}_2 \widehat{R}_2^T \right) \widehat{P} + \widehat{Q} \widehat{B} \widehat{B}^T \right). \end{aligned}$$

Gradients with respect to the reduced matrices are

$$\begin{aligned} \nabla_{\widehat{J}_2} \mathcal{L} &= 2\widetilde{Q}^T \widetilde{P} - 2\widetilde{P}^T \widetilde{Q} + 2\widehat{Q} \widehat{P} - 2\widehat{P} \widehat{Q}, \\ \nabla_{\widehat{R}_2} \mathcal{L} &= -2\widetilde{Q}^T \widetilde{P} \widehat{R}_2 - 2\widetilde{P}^T \widetilde{Q} \widehat{R}_2 - 2\widehat{Q} \widehat{P} \widehat{R}_2 - 2\widehat{P} \widehat{Q} \widehat{R}_2, \\ \nabla_{\widehat{B}} \mathcal{L} &= -2\widetilde{P}^T S B + 2\widehat{P} \widehat{B} + 2\widetilde{Q}^T B + 2\widehat{Q} \widehat{B}. \end{aligned}$$

The result follows from Lemma 5.6. □

As a consequence, we get the Wilson-type necessary optimality conditions for  $\mathcal{H}_2$ -optimal structure-preserving MOR of port-Hamiltonian systems.

**Theorem 6.9:**

Let (6.10), with  $\widehat{S} = I$ , be an  $\mathcal{H}_2$ -optimal ROM for (6.9) with  $\widehat{R} \succ 0$ . Then

$$\begin{aligned} \widetilde{Q}^T \widetilde{P} + \widehat{Q} \widehat{P} &= 0, \\ \text{sym} \left( \widetilde{Q}^T (J - R) S \widetilde{P} + \widehat{Q}^T (\widehat{J} - \widehat{R}) \widehat{P} \right) &= 0, \\ \left( \widetilde{Q}^T - \widetilde{P}^T S \right) B + \left( \widehat{Q} + \widehat{P} \right) \widehat{B} &= 0. \end{aligned} \quad \diamond$$

*Proof.* From Theorem 6.8, we have

$$\begin{aligned} 0 &= \widetilde{Q}^T \widetilde{P} - \widetilde{P}^T \widetilde{Q} + \widehat{Q} \widehat{P} - \widehat{P} \widehat{Q}, \\ 0 &= \widetilde{Q}^T \widetilde{P} \widehat{R}_2 + \widetilde{P}^T \widetilde{Q} \widehat{R}_2 + \widehat{Q} \widehat{P} \widehat{R}_2 + \widehat{P} \widehat{Q} \widehat{R}_2, \\ 0 &= -\widetilde{P}^T S B + \widehat{P} \widehat{B} + \widetilde{Q}^T B + \widehat{Q} \widehat{B}. \end{aligned}$$

Using that  $\widehat{R}_2$  is invertible, it follows that

$$\widetilde{Q}^T \widetilde{P} + \widehat{Q} \widehat{P} = 0, \quad (6.13)$$

$$\left( \widetilde{Q}^T - \widetilde{P}^T S \right) B + \left( \widehat{Q} + \widehat{P} \right) \widehat{B} = 0. \quad (6.14)$$

Summing

$$\begin{aligned} \widetilde{Q}^T \left( (J - R) S \widetilde{P} + \widetilde{P} \left( \widehat{J}_2^T - \widehat{J}_2 - \widehat{R}_2 \widehat{R}_2^T \right) + B \widehat{B}^T \right) &= 0, \\ \widehat{Q} \left( \left( \widehat{J}_2 - \widehat{J}_2^T - \widehat{R}_2 \widehat{R}_2^T \right) \widehat{P} + \widehat{P} \left( \widehat{J}_2^T - \widehat{J}_2 - \widehat{R}_2 \widehat{R}_2^T \right) + \widehat{B} \widehat{B}^T \right) &= 0, \\ \widetilde{P}^T \left( S(-J - R) \widetilde{Q} + \widetilde{Q} \left( \widehat{J}_2 - \widehat{J}_2^T - \widehat{R}_2 \widehat{R}_2^T \right) - S B \widehat{B}^T \right) &= 0, \\ \widehat{P} \left( \left( \widehat{J}_2^T - \widehat{J}_2 - \widehat{R}_2 \widehat{R}_2^T \right) \widehat{Q} + \widehat{Q} \left( \widehat{J}_2 - \widehat{J}_2^T - \widehat{R}_2 \widehat{R}_2^T \right) + \widehat{B} \widehat{B}^T \right) &= 0, \end{aligned}$$

and using (6.13) and (6.14) gives

$$\text{sym} \left( \widetilde{Q}^T (J - R) S \widetilde{P} + \widehat{Q} \left( \widehat{J}_2 - \widehat{J}_2^T - \widehat{R}_2 \widehat{R}_2^T \right) \widehat{P} \right) = 0. \quad \square$$

### 6.3.2 Interpolatory conditions

Here, we derive interpolatory conditions using Theorem 6.9. Assuming diagonalizability of  $\widehat{J} - \widehat{R}$ , we obtain results about bitangential interpolation from [BB14, Theorem 1] (see Theorem 6.7). Under an additional assumption, we derive a tangential interpolation condition.

**Theorem 6.10:**

Let the assumptions in Theorem 6.9 hold. Furthermore, let  $T$  be an invertible matrix such that  $T^{-1}(\widehat{J} - \widehat{R})T = \Lambda = \text{diag}(\lambda_i)$ , where  $\lambda_i$ -s are pairwise distinct. Denote  $t_i = Te_i$ ,  $s_i^T = e_i^T T^{-1}$ ,  $c_i = \widehat{B}^T t_i$ , and  $b_i^T = s_i^T \widehat{B}$  so that  $\widehat{H}(s) = \sum_{i=1}^r \frac{c_i b_i^T}{s - \lambda_i}$ . Then

$$\begin{aligned} c_i^T (H(-\lambda_i) - H(-\lambda_j)) b_j &= c_i^T \left( \widehat{H}(-\lambda_i) - \widehat{H}(-\lambda_j) \right) b_j, \\ c_i^T H'(-\lambda_i) b_i &= c_i^T H'(-\lambda_i) b_i, \end{aligned}$$

for  $i, j = 1, 2, \dots, r$ .

Additionally, if  $T$  can be chosen to be unitary, which is equivalent to  $\widehat{J} - \widehat{R}$  being normal and to  $\widehat{J}$  and  $\widehat{R}$  commuting, then  $\widehat{H}(s) = \sum_{i=1}^r \frac{\bar{b}_i b_i^T}{s - \lambda_i}$  and also

$$(H(-\lambda_i) + H(-\lambda_i)^*) b_i = \left( \widehat{H}(-\lambda_i) + \widehat{H}(-\lambda_i)^* \right) b_i,$$

for  $i = 1, 2, \dots, r$ . ◇

*Proof.* From the assumptions, we have  $(\widehat{J} - \widehat{R})t_i = \lambda_i t_i$  and  $s_i^T (\widehat{J} - \widehat{R}) = \lambda_i s_i^T$ . From

$$(J - R)S\tilde{P}s_i + \tilde{P} \left( -\widehat{J} - \widehat{R} \right) s_i + B\widehat{B}^T s_i = 0,$$

we get

$$(J - R)S\tilde{P}s_i + \lambda_i \tilde{P}s_i + B\widehat{B}^T s_i = 0,$$

and then

$$\begin{aligned} \tilde{P}s_i &= (-\lambda_i I - (J - R)S)^{-1} B b_i, \\ \widehat{P}s_i &= \left( -\lambda_i I - (\widehat{J} - \widehat{R}) \right)^{-1} \widehat{B} b_i, \\ \tilde{Q}t_i &= -(-\lambda_i I - (J - R)S)^{-T} S B c_i, \\ \widehat{Q}t_i &= \left( -\lambda_i I - (\widehat{J} - \widehat{R}) \right)^{-T} \widehat{B} c_i. \end{aligned}$$

From this and proceeding as in the proof of Theorem 2.45, we find

$$\begin{aligned} t_i^T \left( \tilde{Q}^T \tilde{P} + \widehat{Q} \widehat{P} \right) s_i &= c_i^T (H'(-\lambda_i) - H'(-\lambda_i)) b_i, \\ t_i^T \left( \tilde{Q}^T \tilde{P} + \widehat{Q} \widehat{P} \right) s_j &= c_i^T \left( \frac{H(-\lambda_i) - H(-\lambda_j)}{(-\lambda_j) - (-\lambda_j)} - \frac{\widehat{H}(-\lambda_i) - \widehat{H}(-\lambda_j)}{(-\lambda_j) - (-\lambda_j)} \right) b_j. \end{aligned}$$

Let additionally  $T^{-1} = T^*$ , which is possible if and only if  $\widehat{J} - \widehat{R}$  is normal. We can see that  $(\widehat{J} - \widehat{R})(\widehat{J} - \widehat{R})^T = (\widehat{J} - \widehat{R})^T(\widehat{J} - \widehat{R})$  is equivalent to  $\widehat{J}\widehat{R} = \widehat{R}\widehat{J}$ , i.e.,  $\widehat{J}$  and  $\widehat{R}$

commute. From  $T^{-1} = T^*$ , we have  $s_i = \bar{t}_i$  and  $c_i = \bar{b}_i$ . Then

$$\begin{aligned}
 & t_i^T \left( \tilde{Q}^T - \tilde{P}^T S \right) B + t_i^T \left( \hat{Q} + \hat{P} \right) \hat{B} \\
 &= t_i^T \tilde{Q}^T B - s_i^* \tilde{P}^T S B + t_i^T \hat{Q} \hat{B} + s_i^* \hat{P} \hat{B} \\
 &= -c_i^T B^T S (-\lambda_i I - (J - R)S)^{-1} B - b_i^* B^T (-\bar{\lambda}_i I - (J - R)S)^{-T} S B \\
 &\quad + c_i^T \hat{B}^T \left( -\lambda_i I - (\hat{J} - \hat{R}) \right)^{-1} \hat{B} + b_i^* \hat{B}^T \left( -\bar{\lambda}_i I - (\hat{J} - \hat{R}) \right)^{-T} \hat{B} \\
 &= -b_i^* (H(-\lambda_i) + H(-\lambda_i)^*) + b_i^* \left( \hat{H}(-\lambda_i) + \hat{H}(-\lambda_i)^* \right). \quad \square
 \end{aligned}$$

## 6.4 Linear parametric systems

Consider a parameterized LTI (PLTI) system

$$\begin{aligned}
 E(\mathbf{p})\dot{x}(t, \mathbf{p}) &= A(\mathbf{p})x(t, \mathbf{p}) + B(\mathbf{p})u(t), \\
 y(t, \mathbf{p}) &= C(\mathbf{p})x(t, \mathbf{p}),
 \end{aligned} \tag{6.15}$$

where  $\mathbf{p} \in \mathbf{P}$  is the parameter,  $\mathbf{P} \subset \mathbb{R}^d$  is a compact set, and  $E(\mathbf{p}), A(\mathbf{p}) \in \mathbb{R}^{n \times n}$ ,  $B(\mathbf{p}) \in \mathbb{R}^{n \times m}$ ,  $C(\mathbf{p}) \in \mathbb{R}^{p \times n}$  are continuous matrix-valued functions. We assume  $E(\mathbf{p})$  is invertible and  $\sigma(A(\mathbf{p}), E(\mathbf{p})) \subset \mathbb{C}_-$ , for all  $\mathbf{p} \in \mathbf{P}$ .

### 6.4.1 $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal model order reduction

Following [BBBG11], we define the  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -norm of (6.15) with

$$\|H\|_{\mathcal{H}_2 \otimes \mathcal{L}_2(\mathcal{P})}^2 := \frac{1}{2\pi} \int_{\mathbf{P}} \int_{-\infty}^{\infty} \|H(\mathbf{z}\omega, \mathbf{p})\|_{\mathbb{F}}^2 d\omega d\mathbf{p} = \int_{\mathbf{P}} \|H(\cdot, \mathbf{p})\|_{\mathcal{H}_2}^2 d\mathbf{p},$$

where  $H$  is the parameterized transfer function of (6.15):

$$H(s, \mathbf{p}) = C(\mathbf{p})(sE(\mathbf{p}) - A(\mathbf{p}))^{-1}B(\mathbf{p}).$$

The norm is well-defined since  $\|H(\cdot, \mathbf{p})\|_{\mathcal{H}_2}^2$  is continuous with respect to  $\mathbf{p}$  and  $\mathbf{P}$  is compact.

We are interested in finding an  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal ROM

$$\begin{aligned}
 \hat{E}(\mathbf{p})\hat{x}(t, \mathbf{p}) &= \hat{A}(\mathbf{p})\hat{x}(t, \mathbf{p}) + \hat{B}(\mathbf{p})u(t), \\
 \hat{y}(t, \mathbf{p}) &= \hat{C}(\mathbf{p})\hat{x}(t, \mathbf{p}),
 \end{aligned} \tag{6.16}$$

with a parameterized transfer function  $\hat{H}$ , where  $\hat{E}(\mathbf{p}), \hat{A}(\mathbf{p}) \in \mathbb{R}^{r \times r}$ ,  $\hat{B}(\mathbf{p}) \in \mathbb{R}^{r \times m}$ ,  $\hat{C}(\mathbf{p}) \in \mathbb{R}^{p \times r}$ . But, clearly,  $\|H - \hat{H}\|_{\mathcal{H}_2 \otimes \mathcal{L}_2(\mathcal{P})}$  is minimized if and only if  $\|H(\cdot, \mathbf{p}) - \hat{H}(\cdot, \mathbf{p})\|_{\mathcal{H}_2}$  is minimized for all  $\mathbf{p} \in \mathbf{P}$ . This would mean that, in this setting,  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal

MOR would consist of performing  $\mathcal{H}_2$ -optimal MOR for every parameter value  $\mathbf{p} \in \mathbf{P}$ , which is a bottleneck if the ROM (6.16) needs to be computed for many parameter values.

To overcome this, we need to restrict the structure of the ROM (6.16). One possible approach is to take that the parametric matrices are separable (see [Haa17, Definition 2.6])

$$\begin{aligned}\widehat{E}(\mathbf{p}) &= \sum_{i=1}^{q_{\widehat{E}}} \theta_i^{\widehat{E}}(\mathbf{p}) \widehat{E}_i, & \widehat{A}(\mathbf{p}) &= \sum_{i=1}^{q_{\widehat{A}}} \theta_i^{\widehat{A}}(\mathbf{p}) \widehat{A}_i, \\ \widehat{B}(\mathbf{p}) &= \sum_{i=1}^{q_{\widehat{B}}} \theta_i^{\widehat{B}}(\mathbf{p}) \widehat{B}_i, & \widehat{C}(\mathbf{p}) &= \sum_{i=1}^{q_{\widehat{C}}} \theta_i^{\widehat{C}}(\mathbf{p}) \widehat{C}_i.\end{aligned}\tag{6.17}$$

where  $\theta_i^{\widehat{E}}, \theta_i^{\widehat{A}}, \theta_i^{\widehat{B}}, \theta_i^{\widehat{C}}: \mathbf{P} \rightarrow \mathbb{R}$  are given continuous functions, and then to optimize the non-parametric matrices  $\widehat{E}_i, \widehat{A}_i, \widehat{B}_i, \widehat{C}_i$ . Clearly, with the form (6.17), if additionally  $q_{\widehat{E}}, q_{\widehat{A}}, q_{\widehat{B}}, q_{\widehat{C}}$  are small and functions  $\theta_i^{\widehat{E}}, \theta_i^{\widehat{A}}, \theta_i^{\widehat{B}}, \theta_i^{\widehat{C}}$  are easy to compute, the ROM (6.16) can be computed for many parameter values in an efficient manner.

Parameter separability is a beneficial property for full-order matrices in projection-based methods. To see this, let  $A(\mathbf{p}) = \sum_{i=1}^{q_A} \theta_i^A(\mathbf{p}) A_i$ . Then, when using a Petrov-Galerkin projection to find reduced matrices, we have  $W^T A(\mathbf{p}) V = \sum_{i=1}^{q_A} \theta_i^A(\mathbf{p}) W^T A_i V$ . Therefore, after precomputing matrices  $W^T A_i V$ , the matrix  $W^T A(\mathbf{p}) V$  can be assembled for any parameter value  $\mathbf{p} \in \mathbf{P}$  with time complexity which is independent of the order of the original model. If the assumption of parameter separability is not satisfied, an approximation is done, e.g., using the empirical interpolation method [BMNP04] (see also [Haa17, Section 2.3.7]). Here, we will only assume that the ROM has separable parametric matrices. This enables structure preservation when full-order matrices are separable and otherwise avoids an intermediate approximation step.

The following theorem gives the gradients of the squared  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -error.

**Theorem 6.11:**

Let (6.16) be of the form (6.17) and asymptotically stable for all  $\mathbf{p} \in \mathbf{P}$ . Then for the the squared  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -error, we have

$$\begin{aligned}\nabla_{\widehat{E}_i} \mathcal{J} &= 2 \int_{\mathbf{P}} \theta_i^{\widehat{E}}(\mathbf{p}) \left( \widetilde{Q}(\mathbf{p})^T A(\mathbf{p}) \widetilde{P}(\mathbf{p}) + \widehat{Q}(\mathbf{p}) \widehat{A}(\mathbf{p}) \widehat{P}(\mathbf{p}) \right) d\mathbf{p}, \\ \nabla_{\widehat{A}_i} \mathcal{J} &= 2 \int_{\mathbf{P}} \theta_i^{\widehat{A}}(\mathbf{p}) \left( \widetilde{Q}(\mathbf{p})^T E(\mathbf{p}) \widetilde{P}(\mathbf{p}) + \widehat{Q}(\mathbf{p}) \widehat{E}(\mathbf{p}) \widehat{P}(\mathbf{p}) \right) d\mathbf{p}, \\ \nabla_{\widehat{B}_i} \mathcal{J} &= 2 \int_{\mathbf{P}} \theta_i^{\widehat{B}}(\mathbf{p}) \left( \widetilde{Q}(\mathbf{p})^T B(\mathbf{p}) + \widehat{Q}(\mathbf{p}) \widehat{B}(\mathbf{p}) \right) d\mathbf{p}, \\ \nabla_{\widehat{C}_i} \mathcal{J} &= 2 \int_{\mathbf{P}} \theta_i^{\widehat{C}}(\mathbf{p}) \left( -C(\mathbf{p}) \widetilde{P}(\mathbf{p}) + \widehat{C}(\mathbf{p}) \widehat{P}(\mathbf{p}) \right) d\mathbf{p}. \quad \diamond\end{aligned}$$

*Proof.* Since  $\tilde{P}$  and  $\hat{P}$  are continuous functions over a compact set  $\mathbf{P}$ , we have  $\tilde{P} \in \mathcal{L}_2(\mathbf{P}; \mathbb{R}^{n \times r})$  and  $\hat{P} \in \mathcal{L}_2(\mathbf{P}; \mathbb{R}^{r \times r})$ . Similar to Theorem 2.44, we find the Lagrange function is

$$\begin{aligned} \mathcal{L} = & \int_{\mathbf{P}} \text{tr} \left( C(\mathbf{p})P(\mathbf{p})C(\mathbf{p})^T - 2C(\mathbf{p})\tilde{P}(\mathbf{p})\hat{C}(\mathbf{p})^T + \hat{C}(\mathbf{p})\hat{P}(\mathbf{p})\hat{C}(\mathbf{p})^T \right) d\mathbf{p} \\ & + \int_{\mathbf{P}} \text{tr} \left( \tilde{\Lambda}(\mathbf{p})^T A(\mathbf{p})\tilde{P}(\mathbf{p})\hat{E}(\mathbf{p})^T + \tilde{\Lambda}(\mathbf{p})^T E(\mathbf{p})\tilde{P}(\mathbf{p})\hat{A}(\mathbf{p})^T + \tilde{\Lambda}(\mathbf{p})^T B(\mathbf{p})\hat{B}(\mathbf{p})^T \right) d\mathbf{p} \\ & + \int_{\mathbf{P}} \text{tr} \left( \hat{\Lambda}(\mathbf{p})^T \hat{A}(\mathbf{p})\hat{P}(\mathbf{p})\hat{E}(\mathbf{p})^T + \hat{\Lambda}(\mathbf{p})^T \hat{E}(\mathbf{p})\hat{P}(\mathbf{p})\hat{A}(\mathbf{p})^T + \hat{\Lambda}(\mathbf{p})^T \hat{B}(\mathbf{p})\hat{B}(\mathbf{p})^T \right) d\mathbf{p}. \end{aligned}$$

We see that the gradients with respect to  $\tilde{P}$  and  $\hat{P}$  are

$$\begin{aligned} \nabla_{\tilde{P}} \mathcal{L} &= -2C^T \hat{C} + A^T \tilde{\Lambda} \hat{E} + E^T \tilde{\Lambda} \hat{A} \in \mathcal{L}_2(\mathbf{P}; \mathbb{R}^{n \times r}), \\ \nabla_{\hat{P}} \mathcal{L} &= \hat{C}^T \hat{C} + \hat{A}^T \hat{\Lambda} \hat{E} + \hat{E}^T \hat{\Lambda} \hat{A} \in \mathcal{L}_2(\mathbf{P}; \mathbb{R}^{r \times r}). \end{aligned}$$

Equating with zero, it follows that  $\tilde{\Lambda} = 2\tilde{Q}$  and  $\hat{\Lambda} = \hat{Q}$ . Therefore, the Lagrange function simplifies to

$$\begin{aligned} \mathcal{L} = & \int_{\mathbf{P}} \text{tr} \left( C(\mathbf{p})P(\mathbf{p})C(\mathbf{p})^T - 2C(\mathbf{p})\tilde{P}(\mathbf{p})\hat{C}(\mathbf{p})^T + \hat{C}(\mathbf{p})\hat{P}(\mathbf{p})\hat{C}(\mathbf{p})^T \right) d\mathbf{p} \\ & + \int_{\mathbf{P}} \text{tr} \left( 2\tilde{Q}(\mathbf{p})^T A(\mathbf{p})\tilde{P}(\mathbf{p})\hat{E}(\mathbf{p})^T + 2\tilde{Q}(\mathbf{p})^T E(\mathbf{p})\tilde{P}(\mathbf{p})\hat{A}(\mathbf{p})^T + 2\tilde{Q}(\mathbf{p})^T B(\mathbf{p})\hat{B}(\mathbf{p})^T \right) d\mathbf{p} \\ & + \int_{\mathbf{P}} \text{tr} \left( 2\hat{Q}(\mathbf{p})\hat{A}(\mathbf{p})\hat{P}(\mathbf{p})\hat{E}(\mathbf{p})^T + \hat{Q}(\mathbf{p})\hat{B}(\mathbf{p})\hat{B}(\mathbf{p})^T \right) d\mathbf{p}. \end{aligned}$$

Finally, the gradients with respect to reduced matrices are

$$\begin{aligned} \nabla_{\hat{E}_i} \mathcal{L} &= 2 \int_{\mathbf{P}} \theta_i^{\hat{E}}(\mathbf{p}) \left( \tilde{Q}(\mathbf{p})^T A(\mathbf{p})\tilde{P}(\mathbf{p}) + \hat{Q}(\mathbf{p})\hat{A}(\mathbf{p})\hat{P}(\mathbf{p}) \right) d\mathbf{p}, \\ \nabla_{\hat{A}_i} \mathcal{L} &= 2 \int_{\mathbf{P}} \theta_i^{\hat{A}}(\mathbf{p}) \left( \tilde{Q}(\mathbf{p})^T E(\mathbf{p})\tilde{P}(\mathbf{p}) + \hat{Q}(\mathbf{p})\hat{E}(\mathbf{p})\hat{P}(\mathbf{p}) \right) d\mathbf{p}, \\ \nabla_{\hat{B}_i} \mathcal{L} &= 2 \int_{\mathbf{P}} \theta_i^{\hat{B}}(\mathbf{p}) \left( \tilde{Q}(\mathbf{p})^T B(\mathbf{p}) + \hat{Q}(\mathbf{p})\hat{B}(\mathbf{p}) \right) d\mathbf{p}, \\ \nabla_{\hat{C}_i} \mathcal{L} &= 2 \int_{\mathbf{P}} \theta_i^{\hat{C}}(\mathbf{p}) \left( -C(\mathbf{p})\tilde{P}(\mathbf{p}) + \hat{C}(\mathbf{p})\hat{P}(\mathbf{p}) \right) d\mathbf{p}. \end{aligned}$$

The claim follows from Lemma 5.6.  $\square$

Directly, we obtain the Wilson-type necessary optimality conditions for  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal MOR.



**Corollary 6.12:**

Let (6.16) be an  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal ROM with the form (6.17) for (6.15). Then

$$\begin{aligned} \int_{\mathbf{P}} \theta_i^{\hat{A}}(\mathbf{p}) \left( \tilde{Q}(\mathbf{p})^T E(\mathbf{p}) \tilde{P}(\mathbf{p}) + \hat{Q}(\mathbf{p}) \hat{E}(\mathbf{p}) \hat{P}(\mathbf{p}) \right) d\mathbf{p} &= 0, & i = 1, 2, \dots, q_{\hat{A}}, \\ \int_{\mathbf{P}} \theta_i^{\hat{E}}(\mathbf{p}) \left( \tilde{Q}(\mathbf{p})^T A(\mathbf{p}) \tilde{P}(\mathbf{p}) + \hat{Q}(\mathbf{p}) \hat{A}(\mathbf{p}) \hat{P}(\mathbf{p}) \right) d\mathbf{p} &= 0, & i = 1, 2, \dots, q_{\hat{E}}, \\ \int_{\mathbf{P}} \theta_i^{\hat{B}}(\mathbf{p}) \left( \tilde{Q}(\mathbf{p})^T B(\mathbf{p}) + \hat{Q}(\mathbf{p}) \hat{B}(\mathbf{p}) \right) d\mathbf{p} &= 0, & i = 1, 2, \dots, q_{\hat{B}}, \\ \int_{\mathbf{P}} \theta_i^{\hat{C}}(\mathbf{p}) \left( C(\mathbf{p}) \tilde{P}(\mathbf{p}) - \hat{C}(\mathbf{p}) \hat{P}(\mathbf{p}) \right) d\mathbf{p} &= 0, & i = 1, 2, \dots, q_{\hat{C}}. \quad \diamond \end{aligned}$$

### 6.4.2 Interpolatory conditions

Using Corollary 6.12, we derive interpolatory necessary optimality conditions. To have the pole-residue form of the ROM, we assume  $\hat{A}(\mathbf{p})$  and  $\hat{E}(\mathbf{p})$  can be diagonalized using parameter-independent transformation matrices. Note that we also assumed diagonalizability for second-order systems in Theorem 6.4 and port-Hamiltonian systems in Theorem 6.10.

**Theorem 6.13:**

Let (6.16) be an  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal ROM with the form (6.17) for (6.15). Assume that  $\{\hat{E}_1^{-1} \hat{E}_2, \dots, \hat{E}_1^{-1} \hat{E}_{q_{\hat{E}}}, \hat{E}_1^{-1} \hat{A}_1, \dots, \hat{E}_1^{-1} \hat{A}_{q_{\hat{A}}}\}$  is a simultaneously diagonalizable family of matrices. Let  $S$  and  $T$  be invertible matrices such that  $S^T \hat{E}_i T = \Lambda^{\hat{E}_i} = \text{diag}(\lambda_j^{\hat{E}_i})$  and  $S^T \hat{A}_i T = \Lambda^{\hat{A}_i} = \text{diag}(\lambda_j^{\hat{A}_i})$ , with  $\Lambda^{\hat{E}_i} = I_r$ . Furthermore, let

$$\begin{aligned} t_j &= T e_j, \quad s_j^T = e_j^T S, \\ \lambda_j^{\hat{A}}(\mathbf{p}) &= \sum_{i=1}^{q_{\hat{A}}} \theta_i^{\hat{A}}(\mathbf{p}) \lambda_j^{\hat{A}_i}, \quad \lambda_j^{\hat{E}}(\mathbf{p}) = \sum_{i=1}^{q_{\hat{E}}} \theta_i^{\hat{E}}(\mathbf{p}) \lambda_j^{\hat{E}_i}, \quad \lambda_j(\mathbf{p}) = \frac{\lambda_j^{\hat{A}}(\mathbf{p})}{\lambda_j^{\hat{E}}(\mathbf{p})}, \\ b_j(\mathbf{p}) &= \frac{1}{\lambda_j^{\hat{E}}(\mathbf{p})} \hat{B}(\mathbf{p})^T s_j, \quad \text{and } c_j(\mathbf{p}) = \frac{1}{\lambda_j^{\hat{E}}(\mathbf{p})} \hat{C}(\mathbf{p}) t_j, \end{aligned}$$

where  $\lambda_j(\mathbf{p})$  are pairwise distinct for almost all  $\mathbf{p} \in \mathbf{P}$ . Then  $\hat{H}(s, \mathbf{p}) = \sum_{j=1}^r \lambda_j^{\hat{E}}(\mathbf{p}) \frac{c_j(\mathbf{p}) b_j(\mathbf{p})^T}{s - \lambda_j(\mathbf{p})}$  and

$$\begin{aligned} \int_{\mathbf{P}} \theta_i^{\hat{C}}(\mathbf{p}) \cdot H(-\lambda_j(\mathbf{p}), \mathbf{p}) b_j(\mathbf{p}) d\mathbf{p} &= \int_{\mathbf{P}} \theta_i^{\hat{C}}(\mathbf{p}) \cdot \hat{H}(-\lambda_j(\mathbf{p}), \mathbf{p}) b_j(\mathbf{p}) d\mathbf{p}, & i \in [q_{\hat{C}}], \\ \int_{\mathbf{P}} \theta_i^{\hat{B}}(\mathbf{p}) \cdot c_j(\mathbf{p})^T H(-\lambda_j(\mathbf{p}), \mathbf{p}) d\mathbf{p} &= \int_{\mathbf{P}} \theta_i^{\hat{B}}(\mathbf{p}) \cdot c_j(\mathbf{p})^T \hat{H}(-\lambda_j(\mathbf{p}), \mathbf{p}) d\mathbf{p}, & i \in [q_{\hat{B}}], \\ \int_{\mathbf{P}} \theta_i^{\hat{A}}(\mathbf{p}) \cdot c_j(\mathbf{p})^T H'(-\lambda_j(\mathbf{p}), \mathbf{p}) b_j(\mathbf{p}) d\mathbf{p} &= \int_{\mathbf{P}} \theta_i^{\hat{A}}(\mathbf{p}) \cdot c_j(\mathbf{p})^T \hat{H}'(-\lambda_j(\mathbf{p}), \mathbf{p}) b_j(\mathbf{p}) d\mathbf{p}, & i \in [q_{\hat{A}}], \end{aligned}$$

for  $j = 1, 2, \dots, r$ , where  $[k] := \{1, 2, \dots, k\}$ .  $\diamond$

*Proof.* We have  $\widehat{E}_i t_j = \lambda_j^{\widehat{E}_i} \widehat{E}_1 t_j$ ,  $\widehat{A}_i t_j = \lambda_j^{\widehat{A}_i} \widehat{E}_1 t_j$ ,  $s_j^T \widehat{E}_i = \lambda_j^{\widehat{E}_i} s_j^T \widehat{E}_1$ ,  $s_j^T \widehat{A}_i = \lambda_j^{\widehat{A}_i} s_j^T \widehat{E}_1$ . Therefore, we also have  $\widehat{E}(\mathbf{p}) t_j = \lambda_j^{\widehat{E}}(\mathbf{p}) \widehat{E}_1 t_j$ ,  $\widehat{A}(\mathbf{p}) t_j = \lambda_j^{\widehat{A}}(\mathbf{p}) \widehat{E}_1 t_j$ ,  $s_j^T \widehat{E}_i(\mathbf{p}) = \lambda_j^{\widehat{E}_i}(\mathbf{p}) s_j^T \widehat{E}_1$ ,  $s_j^T \widehat{A}_i(\mathbf{p}) = \lambda_j^{\widehat{A}_i}(\mathbf{p}) s_j^T \widehat{E}_1$ . From

$$A(\mathbf{p}) \widetilde{P}(\mathbf{p}) \widehat{E}(\mathbf{p})^T s_j + E(\mathbf{p}) \widetilde{P}(\mathbf{p}) \widehat{A}(\mathbf{p})^T s_j + B(\mathbf{p}) \widehat{B}(\mathbf{p})^T s_j = 0,$$

it follows that

$$\left( \lambda_j^{\widehat{E}}(\mathbf{p}) A(\mathbf{p}) + \lambda_j^{\widehat{A}}(\mathbf{p}) E(\mathbf{p}) \right) \widetilde{P}(\mathbf{p}) \widehat{E}_1^T s_j + B(\mathbf{p}) \widehat{B}(\mathbf{p})^T s_j = 0.$$

Then

$$\widetilde{P}(\mathbf{p}) \widehat{E}_1^T s_j = (-\lambda_j(\mathbf{p}) E(\mathbf{p}) - A(\mathbf{p}))^{-1} B(\mathbf{p}) b_j(\mathbf{p}),$$

and similarly

$$\begin{aligned} \widehat{P}(\mathbf{p}) \widehat{E}_1^T s_j &= \left( -\lambda_j(\mathbf{p}) \widehat{E}(\mathbf{p}) - \widehat{A}(\mathbf{p}) \right)^{-1} \widehat{B}(\mathbf{p}) b_j(\mathbf{p}), \\ \widetilde{Q}(\mathbf{p}) \widehat{E}_1 t_j &= (-\lambda_j(\mathbf{p}) E(\mathbf{p}) - A(\mathbf{p}))^{-T} C(\mathbf{p})^T c_j(\mathbf{p}), \\ \widehat{Q}(\mathbf{p}) \widehat{E}_1 t_j &= \left( -\lambda_j(\mathbf{p}) \widehat{E}(\mathbf{p}) - \widehat{A}(\mathbf{p}) \right)^{-T} \widehat{C}(\mathbf{p})^T c_j(\mathbf{p}). \end{aligned}$$

From

$$\begin{aligned} \int_{\mathbf{P}} \theta_i^{\widehat{C}}(\mathbf{p}) \left( C(\mathbf{p}) \widetilde{P}(\mathbf{p}) - \widehat{C}(\mathbf{p}) \widehat{P}(\mathbf{p}) \right) \widehat{E}_1^T s_j \, \mathrm{d}\mathbf{p} &= 0, \\ \int_{\mathbf{P}} \theta_i^{\widehat{B}}(\mathbf{p}) t_j^T \widehat{E}_1^T \left( \widetilde{Q}(\mathbf{p})^T B(\mathbf{p}) + \widehat{Q}(\mathbf{p}) \widehat{B}(\mathbf{p}) \right) \, \mathrm{d}\mathbf{p} &= 0, \\ \int_{\mathbf{P}} \theta_i^{\widehat{A}}(\mathbf{p}) t_j^T \widehat{E}_1^T \left( \widetilde{Q}(\mathbf{p})^T E(\mathbf{p}) \widetilde{P}(\mathbf{p}) + \widehat{Q}(\mathbf{p}) \widehat{E}(\mathbf{p}) \widehat{P}(\mathbf{p}) \right) \widehat{E}_1^T s_j \, \mathrm{d}\mathbf{p} &= 0, \end{aligned}$$

we find the statement of the theorem.  $\square$

## 6.5 Linear time-delay systems

We consider a linear time-delay (LTD) system with a single delay

$$\begin{aligned} E\dot{x}(t) &= A_0 x(t) + A_\tau x(t - \tau) + Bu(t), \\ y(t) &= Cx(t), \end{aligned} \tag{6.18}$$

where  $E, A_0, A_\tau \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times m}$ , and  $C \in \mathbb{R}^{p \times n}$ , with  $E$  invertible. Additionally, we assume the system to be exponentially stable.

The transfer function of this system is

$$H(s) = C \left( sE - A_0 - e^{-\tau s} A_\tau \right)^{-1} B.$$

As shown in [JVM11], the  $\mathcal{H}_2$ -norm  $\|H\|_{\mathcal{H}_2}$  can be computed using the Gramians similarly as for LTI systems without delay. The Gramians of the system (6.18) are solutions to a boundary value problem involving a delay differential Lyapunov equation:

$$\dot{P}(t)E^T = P(t)A_0^T + P(t-\tau)A_\tau^T, \quad t \geq 0, \quad (6.19a)$$

$$P(-t) = P(t)^T, \quad (6.19b)$$

$$-BB^T = A_0P(0)E^T + EP(0)A_0^T + A_\tau P(\tau)E^T + EP(-\tau)A_\tau^T, \quad (6.19c)$$

and

$$\dot{Q}(t)E = Q(t)A_0 + Q(t-\tau)A_\tau, \quad t \geq 0, \quad (6.20a)$$

$$Q(-t) = Q(t)^T, \quad (6.20b)$$

$$-C^TC = A_0^TQ(0)E + E^TQ(0)A_0 + A_\tau^TQ(\tau)E + E^TQ(-\tau)A_\tau. \quad (6.20c)$$

The  $\mathcal{H}_2$ -norm can then be computed as (see [JVM11, Theorem 1])

$$\|H\|_{\mathcal{H}_2}^2 = \text{tr}(CP(0)C^T) = \text{tr}(B^TQ(0)B).$$

The system (6.19) can be solved analytically, as done in [JVM11, Section III.A]. First, notice that differentiating (6.19b) gives us  $-\dot{P}(-t) = \dot{P}(t)^T$ . Now, transposing (6.19a), substituting  $t$  with  $-t + \tau$  and using the previous expression, we find

$$E\dot{P}(t-\tau) = -A_0P(t-\tau) - A_\tau P(t), \quad t \leq \tau.$$

Defining  $z(t) = \text{col}(\text{vec}(P(t)), \text{vec}(P(t-\tau)))$  gives us a system of ODEs with a boundary value condition

$$\begin{aligned} \begin{bmatrix} E \otimes I & 0 \\ 0 & I \otimes E \end{bmatrix} \dot{z}(t) &= \begin{bmatrix} A_0 \otimes I & A_\tau \otimes I \\ -I \otimes A_\tau & -I \otimes A_0 \end{bmatrix} z(t), \quad t \in [0, \tau], \\ \begin{bmatrix} -\text{vec}(BB^T) \\ 0 \end{bmatrix} &= \begin{bmatrix} A_0 \otimes E & A_\tau \otimes E \\ I & 0 \end{bmatrix} z(0) + \begin{bmatrix} E \otimes A_\tau & E \otimes A_0 \\ 0 & -I \end{bmatrix} z(\tau), \end{aligned}$$

which can be used to find the analytic solution. Therefore, an equivalent form of (6.19) is

$$\dot{P}(t)E^T = P(t)A_0^T + P(t-\tau)A_\tau^T, \quad t \in [0, \tau], \quad (6.21a)$$

$$E\dot{P}(t-\tau) = -A_0P(t-\tau) - A_\tau P(t), \quad t \in [0, \tau], \quad (6.21b)$$

$$-BB^T = A_0P(0)E^T + EP(0)A_0^T + A_\tau P(\tau)E^T + EP(-\tau)A_\tau^T, \quad (6.21c)$$

and similarly, system (6.20) is equivalent to

$$\dot{Q}(t)E = Q(t)A_0 + Q(t-\tau)A_\tau, \quad t \in [0, \tau], \quad (6.22a)$$

$$E^T\dot{Q}(t-\tau) = -A_0^TQ(t-\tau) - A_\tau^TQ(t), \quad t \in [0, \tau], \quad (6.22b)$$

$$-C^TC = A_0^TQ(0)E + E^TQ(0)A_0 + A_\tau^TQ(\tau)E + E^TQ(-\tau)A_\tau. \quad (6.22c)$$

We will prefer this form in the following.

### 6.5.1 Wilson-type conditions

We want to find  $\mathcal{H}_2$ -optimality conditions for the ROM

$$\begin{aligned}\widehat{E}\dot{\widehat{x}}(t) &= \widehat{A}_0\widehat{x}(t) + \widehat{A}_\tau\widehat{x}(t - \tau) + \widehat{B}u(t), \\ \widehat{y}(t) &= \widehat{C}\widehat{x}(t).\end{aligned}\tag{6.23}$$

with  $\widehat{E}, \widehat{A}_0, \widehat{A}_\tau \in \mathbb{R}^{r \times r}$ ,  $\widehat{B} \in \mathbb{R}^{r \times m}$ , and  $\widehat{C} \in \mathbb{R}^{p \times r}$ , where  $\widehat{E}$  is invertible. Just as for the full-order model (FOM), the Gramians of the ROM (6.23) satisfy

$$\begin{aligned}\dot{\widehat{P}}(t)\widehat{E}^\top &= \widehat{P}(t)\widehat{A}_0^\top + \widehat{P}(t - \tau)\widehat{A}_\tau^\top, \quad t \in [0, \tau], \\ \widehat{E}^\top\dot{\widehat{P}}(t - \tau) &= -\widehat{A}_0^\top\widehat{P}(t - \tau) - \widehat{A}_\tau^\top\widehat{P}(t), \quad t \in [0, \tau], \\ -\widehat{B}\widehat{B}^\top &= \widehat{A}_0\widehat{P}(0)\widehat{E}^\top + \widehat{E}\widehat{P}(0)\widehat{A}_0^\top + \widehat{A}_\tau\widehat{P}(\tau)\widehat{E}^\top + \widehat{E}\widehat{P}(-\tau)\widehat{A}_\tau^\top,\end{aligned}$$

and

$$\begin{aligned}\dot{\widehat{Q}}(t)\widehat{E} &= \widehat{Q}(t)\widehat{A}_0 + \widehat{Q}(t - \tau)\widehat{A}_\tau, \quad t \in [0, \tau], \\ \widehat{E}^\top\dot{\widehat{Q}}(t - \tau) &= -\widehat{A}_0^\top\widehat{Q}(t - \tau) - \widehat{A}_\tau^\top\widehat{Q}(t), \quad t \in [0, \tau], \\ -\widehat{C}^\top\widehat{C} &= \widehat{A}_0^\top\widehat{Q}(0)\widehat{E} + \widehat{E}^\top\widehat{Q}(0)\widehat{A}_0 + \widehat{A}_\tau^\top\widehat{Q}(\tau)\widehat{E} + \widehat{E}^\top\widehat{Q}(-\tau)\widehat{A}_\tau.\end{aligned}$$

We can define the error system as

$$\begin{aligned}\begin{bmatrix} E & 0 \\ 0 & \widehat{E} \end{bmatrix} \begin{bmatrix} \dot{x}(t) \\ \dot{\widehat{x}}(t) \end{bmatrix} &= \begin{bmatrix} A_0 & 0 \\ 0 & \widehat{A}_0 \end{bmatrix} \begin{bmatrix} x(t) \\ \widehat{x}(t) \end{bmatrix} + \begin{bmatrix} A_\tau & 0 \\ 0 & \widehat{A}_\tau \end{bmatrix} \begin{bmatrix} x(t - \tau) \\ \widehat{x}(t - \tau) \end{bmatrix} + \begin{bmatrix} B \\ \widehat{B} \end{bmatrix} u(t), \\ y(t) - \widehat{y}(t) &= \begin{bmatrix} C & -\widehat{C} \end{bmatrix} \begin{bmatrix} x(t) \\ \widehat{x}(t) \end{bmatrix}.\end{aligned}$$

We find the Gramians of the error system have the form

$$\begin{bmatrix} P(t) & \widetilde{P}(t) \\ \widetilde{P}(-t)^\top & \widehat{P}(t) \end{bmatrix}, \quad \begin{bmatrix} Q(t) & \widetilde{Q}(t) \\ \widetilde{Q}(-t)^\top & \widehat{Q}(t) \end{bmatrix},$$

where

$$\begin{aligned}\dot{\widetilde{P}}(t)\widehat{E}^\top &= \widetilde{P}(t)\widehat{A}_0^\top + \widetilde{P}(t - \tau)\widehat{A}_\tau^\top, \\ \widehat{E}^\top\dot{\widetilde{P}}(t - \tau) &= -\widehat{A}_0^\top\widetilde{P}(t - \tau) - \widehat{A}_\tau^\top\widetilde{P}(t), \\ -\widehat{B}\widehat{B}^\top &= \widehat{A}_0\widetilde{P}(0)\widehat{E}^\top + \widehat{E}\widetilde{P}(0)\widehat{A}_0^\top + \widehat{A}_\tau\widetilde{P}(\tau)\widehat{E}^\top + \widehat{E}\widetilde{P}(-\tau)\widehat{A}_\tau^\top,\end{aligned}$$

and

$$\begin{aligned}\dot{\widetilde{Q}}(t)\widehat{E} &= \widetilde{Q}(t)\widehat{A}_0 + \widetilde{Q}(t - \tau)\widehat{A}_\tau, \\ \widehat{E}^\top\dot{\widetilde{Q}}(t - \tau) &= -\widehat{A}_0^\top\widetilde{Q}(t - \tau) - \widehat{A}_\tau^\top\widetilde{Q}(t), \\ C^\top\widehat{C} &= \widehat{A}_0^\top\widetilde{Q}(0)\widehat{E} + \widehat{E}^\top\widetilde{Q}(0)\widehat{A}_0 + \widehat{A}_\tau^\top\widetilde{Q}(\tau)\widehat{E} + \widehat{E}^\top\widetilde{Q}(-\tau)\widehat{A}_\tau.\end{aligned}$$

**Theorem 6.14:**

Let (6.23) be asymptotically stable. Then for the squared  $\mathcal{H}_2$ -error  $\mathcal{J}$ , we have

$$\begin{aligned}
 \nabla_{\hat{E}} \mathcal{J} &= -2 \int_0^\tau \tilde{Q}(\tau-t)^\top A_\tau \dot{\tilde{P}}(t) dt + 2\tilde{Q}(0)^\top A_0 \tilde{P}(0) + 2\tilde{Q}(0)^\top A_\tau \tilde{P}(\tau) \\
 &\quad - 2 \int_0^\tau \hat{Q}(\tau-t)^\top \hat{A}_\tau \dot{\hat{P}}(t) dt + 2\hat{Q}(0)^\top \hat{A}_0 \hat{P}(0) + 2\hat{Q}(0)^\top \hat{A}_\tau \hat{P}(\tau), \\
 \nabla_{\hat{A}_0} \mathcal{J} &= 2 \int_0^\tau \tilde{Q}(\tau-t)^\top A_\tau \tilde{P}(t) dt + 2\tilde{Q}(0)^\top E \tilde{P}(0) \\
 &\quad + 2 \int_0^\tau \hat{Q}(\tau-t)^\top \hat{A}_\tau \hat{P}(t) dt + 2\hat{Q}(0)^\top \hat{E} \hat{P}(0), \\
 \nabla_{\hat{A}_\tau} \mathcal{J} &= 2 \int_0^\tau \tilde{Q}(\tau-t)^\top A_\tau \tilde{P}(t-\tau) dt + 2\tilde{Q}(0)^\top E \tilde{P}(-\tau) \\
 &\quad + 2 \int_0^\tau \hat{Q}(\tau-t)^\top \hat{A}_\tau \hat{P}(t-\tau) dt + 2\hat{Q}(0)^\top \hat{E} \hat{P}(-\tau), \\
 \nabla_{\hat{B}} \mathcal{J} &= 2\tilde{Q}(0)^\top B + 2\hat{Q}(0)^\top \hat{B}, \\
 \nabla_{\hat{C}} \mathcal{J} &= -2C\tilde{P}(0) + 2\hat{C}\hat{P}(0). \quad \diamond
 \end{aligned}$$

*Proof.* We proceed similar to the proof of Theorem 5.7. The  $\mathcal{H}_2$ -error is given by

$$\begin{aligned}
 \|H - \hat{H}\|_{\mathcal{H}_2}^2 &= \text{tr} \left( \begin{bmatrix} C & -\hat{C} \end{bmatrix} \begin{bmatrix} P(0) & \tilde{P}(0) \\ \tilde{P}(0)^\top & \hat{P}(0) \end{bmatrix} \begin{bmatrix} C^\top \\ -\hat{C}^\top \end{bmatrix} \right) \\
 &= \text{tr} \left( CP(0)C^\top - 2C\tilde{P}(0)\hat{C}^\top + \hat{C}\hat{P}(0)\hat{C}^\top \right).
 \end{aligned}$$

Therefore, we consider the optimization problem

$$\begin{aligned}
 &\underset{\hat{E}, \hat{A}_0, \hat{A}_\tau, \hat{B}, \hat{C}, \tilde{P}, \hat{P}}{\text{minimize}} && \text{tr} \left( CP(0)C^\top - 2C\tilde{P}(0)\hat{C}^\top + \hat{C}\hat{P}(0)\hat{C}^\top \right), \\
 &\text{subject to} && \dot{\tilde{P}}(t)\hat{E}^\top = \tilde{P}(t)\hat{A}_0^\top + \tilde{P}(t-\tau)\hat{A}_\tau^\top, \quad t \in [0, \tau], \\
 & && E\dot{\tilde{P}}(t-\tau) = -A_0\tilde{P}(t-\tau) - A_\tau\tilde{P}(t), \quad t \in [0, \tau], \\
 & && -B\hat{B}^\top = A_0\tilde{P}(0)\hat{E}^\top + E\tilde{P}(0)\hat{A}_0^\top + A_\tau\tilde{P}(\tau)\hat{E}^\top + E\tilde{P}(-\tau)\hat{A}_\tau^\top, \\
 & && \dot{\hat{P}}(t)\hat{E}^\top = \hat{P}(t)\hat{A}_0^\top + \hat{P}(t-\tau)\hat{A}_\tau^\top, \quad t \in [0, \tau], \\
 & && \hat{E}\dot{\hat{P}}(t-\tau) = -\hat{A}_0\hat{P}(t-\tau) - \hat{A}_\tau\hat{P}(t), \quad t \in [0, \tau], \\
 & && -\hat{B}\hat{B}^\top = \hat{A}_0\hat{P}(0)\hat{E}^\top + \hat{E}\hat{P}(0)\hat{A}_0^\top + \hat{A}_\tau\hat{P}(\tau)\hat{E}^\top + \hat{E}\hat{P}(-\tau)\hat{A}_\tau^\top,
 \end{aligned}$$

where we have  $\tilde{P} \in H^1([-\tau, \tau]; \mathbb{R}^{n \times r})$  and  $\hat{P} \in H^1([-\tau, \tau]; \mathbb{R}^{r \times r})$ . The Lagrange func-

tion is

$$\begin{aligned}
 \mathcal{L} = & \operatorname{tr}\left(CP(0)C^T - 2C\tilde{P}(0)\hat{C}^T + \hat{C}\hat{P}(0)\hat{C}^T\right) \\
 & + \int_0^\tau \operatorname{tr}\left(\tilde{\Lambda}_1(t)^T \left(\dot{\tilde{P}}(t)\hat{E}^T - \tilde{P}(t)\hat{A}_0^T - \tilde{P}(t-\tau)\hat{A}_\tau^T\right)\right) dt \\
 & + \int_0^\tau \operatorname{tr}\left(\tilde{\Lambda}_2(t)^T \left(E\dot{\tilde{P}}(t-\tau) + A_0\tilde{P}(t-\tau) + A_\tau\tilde{P}(t)\right)\right) dt \\
 & + \operatorname{tr}\left(\tilde{\Lambda}_3^T \left(A_0\tilde{P}(0)\hat{E}^T + E\tilde{P}(0)\hat{A}_0^T + A_\tau\tilde{P}(\tau)\hat{E}^T + E\tilde{P}(-\tau)\hat{A}_\tau^T + B\hat{B}^T\right)\right) \\
 & + \int_0^\tau \operatorname{tr}\left(\hat{\Lambda}_1(t)^T \left(\dot{\hat{P}}(t)\hat{E}^T - \hat{P}(t)\hat{A}_0^T - \hat{P}(t-\tau)\hat{A}_\tau^T\right)\right) dt \\
 & + \int_0^\tau \operatorname{tr}\left(\hat{\Lambda}_2(t)^T \left(\hat{E}\dot{\hat{P}}(t-\tau) + \hat{A}_0\hat{P}(t-\tau) + \hat{A}_\tau\hat{P}(t)\right)\right) dt \\
 & + \operatorname{tr}\left(\hat{\Lambda}_3^T \left(\hat{A}_0\hat{P}(0)\hat{E}^T + \hat{E}\hat{P}(0)\hat{A}_0^T + \hat{A}_\tau\hat{P}(\tau)\hat{E}^T + \hat{E}\hat{P}(-\tau)\hat{A}_\tau^T + \hat{B}\hat{B}^T\right)\right),
 \end{aligned}$$

where  $\tilde{\Lambda}_1, \tilde{\Lambda}_2 \in H^1([0, \tau]; \mathbb{R}^{n \times r})$ ,  $\tilde{\Lambda}_3 \in \mathbb{R}^{n \times r}$ ,  $\hat{\Lambda}_1, \hat{\Lambda}_2 \in H^1([0, \tau]; \mathbb{R}^{r \times r})$ , and  $\hat{\Lambda}_3 \in \mathbb{R}^{r \times r}$  are Lagrange multipliers.

Directional derivatives with respect to  $\tilde{P}$  and  $\hat{P}$ , in directions  $\tilde{D}$  and  $\hat{D}$ , are

$$\begin{aligned}
 d_{\tilde{P}}\mathcal{L}(\tilde{D}) = & \int_0^\tau \operatorname{tr}\left(\left(-\hat{E}^T\dot{\tilde{\Lambda}}_1(t)^T - \hat{A}_0^T\tilde{\Lambda}_1(t)^T + \tilde{\Lambda}_2(t)^T A_\tau\right) \tilde{D}(t)\right) dt \\
 & + \int_0^\tau \operatorname{tr}\left(\left(-\hat{A}_\tau^T\tilde{\Lambda}_1(t)^T - \dot{\tilde{\Lambda}}_2(t)^T E + \tilde{\Lambda}_2(t)^T A_0\right) \tilde{D}(t-\tau)\right) dt, \\
 & + \operatorname{tr}\left(\left(-2\hat{C}^T C - \hat{E}^T\tilde{\Lambda}_1(0)^T + \tilde{\Lambda}_2(\tau)^T E + \hat{E}^T\tilde{\Lambda}_3^T A_0 + \hat{A}_0^T\tilde{\Lambda}_3^T E\right) \tilde{D}(0)\right) \\
 & + \operatorname{tr}\left(\left(\hat{E}^T\tilde{\Lambda}_1(\tau)^T + \hat{E}^T\tilde{\Lambda}_3^T A_\tau\right) \tilde{D}(\tau)\right) \\
 & + \operatorname{tr}\left(\left(-\tilde{\Lambda}_2(0)^T E + \hat{A}_\tau^T\tilde{\Lambda}_3^T E\right) \tilde{D}(-\tau)\right) \\
 d_{\hat{P}}\mathcal{L}(\hat{D}) = & \int_0^\tau \operatorname{tr}\left(\left(-\hat{E}^T\dot{\hat{\Lambda}}_1(t)^T - \hat{A}_0^T\hat{\Lambda}_1(t)^T + \hat{\Lambda}_2(t)^T \hat{A}_\tau\right) \hat{D}(t)\right) dt \\
 & + \int_0^\tau \operatorname{tr}\left(\left(-\hat{A}_\tau^T\hat{\Lambda}_1(t)^T - \dot{\hat{\Lambda}}_2(t)^T \hat{E} + \hat{\Lambda}_2(t)^T \hat{A}_0\right) \hat{D}(t-\tau)\right) dt \\
 & + \operatorname{tr}\left(\left(\hat{C}^T \hat{C} - \hat{E}^T\hat{\Lambda}_1(0)^T + \hat{\Lambda}_2(\tau)^T \hat{E} + \hat{E}^T\hat{\Lambda}_3^T \hat{A}_0 + \hat{A}_0^T\hat{\Lambda}_3^T \hat{E}\right) \hat{D}(0)\right) \\
 & + \operatorname{tr}\left(\left(\hat{E}^T\hat{\Lambda}_1(\tau)^T + \hat{E}^T\hat{\Lambda}_3^T \hat{A}_\tau\right) \hat{D}(\tau)\right) \\
 & + \operatorname{tr}\left(\left(-\hat{\Lambda}_2(0)^T \hat{E} + \hat{A}_\tau^T\hat{\Lambda}_3^T \hat{E}\right) \hat{D}(-\tau)\right).
 \end{aligned}$$

Equating with zero, it follows that

$$\begin{aligned}\dot{\tilde{\Lambda}}_1(t)\hat{E} &= -\tilde{\Lambda}_1(t)\hat{A}_0 + A_\tau^\top \tilde{\Lambda}_2(t), \quad \text{for a.e. } t \in [0, \tau], \\ E^\top \dot{\tilde{\Lambda}}_2(t) &= A_0^\top \tilde{\Lambda}_2(t) - \tilde{\Lambda}_1(t)\hat{A}_\tau, \quad \text{for a.e. } t \in [0, \tau], \\ 2C^\top \hat{C} &= A_0^\top \tilde{\Lambda}_3 \hat{E} + E^\top \tilde{\Lambda}_3 \hat{A}_0 - \tilde{\Lambda}_1(0)\hat{E} + E^\top \tilde{\Lambda}_2(\tau), \\ \tilde{\Lambda}_1(\tau) &= -A_\tau^\top \tilde{\Lambda}_3, \\ \tilde{\Lambda}_2(0) &= \tilde{\Lambda}_3 \hat{A}_\tau,\end{aligned}$$

and

$$\begin{aligned}\dot{\hat{\Lambda}}_1(t)\hat{E} &= -\hat{\Lambda}_1(t)\hat{A}_0 + \hat{A}_\tau^\top \hat{\Lambda}_2(t), \quad \text{for a.e. } t \in [0, \tau], \\ \hat{E}^\top \dot{\hat{\Lambda}}_2(t) &= \hat{A}_0^\top \hat{\Lambda}_2(t) - \hat{\Lambda}_1(t)\hat{A}_\tau, \quad \text{for a.e. } t \in [0, \tau], \\ -\hat{C}^\top \hat{C} &= \hat{A}_0^\top \hat{\Lambda}_3 \hat{E} + \hat{E}^\top \hat{\Lambda}_3 \hat{A}_0 - \hat{\Lambda}_1(0)\hat{E} + \hat{E}^\top \hat{\Lambda}_2(\tau), \\ \hat{\Lambda}_1(\tau) &= -\hat{A}_\tau^\top \hat{\Lambda}_3, \\ \hat{\Lambda}_2(0) &= \hat{\Lambda}_3 \hat{A}_\tau.\end{aligned}$$

The solution is

$$\begin{aligned}\tilde{\Lambda}_1(t) &= -2A_\tau^\top \tilde{Q}(\tau - t), \quad \tilde{\Lambda}_2(t) = 2\tilde{Q}(-t)\hat{A}_\tau, \quad \tilde{\Lambda}_3 = 2\tilde{Q}(0), \\ \hat{\Lambda}_1(t) &= -\hat{A}_\tau^\top \hat{Q}(\tau - t), \quad \hat{\Lambda}_2(t) = \hat{Q}(-t)\hat{A}_\tau, \quad \hat{\Lambda}_3 = \hat{Q}(0).\end{aligned}$$

Gradients with respect to reduced matrices are

$$\begin{aligned}\nabla_{\hat{E}} \mathcal{L} &= \int_0^\tau \tilde{\Lambda}_1(t)^\top \dot{\tilde{P}}(t) dt + \tilde{\Lambda}_3^\top (A_0 \tilde{P}(0) + A_\tau \tilde{P}(\tau)) \\ &\quad + \int_0^\tau \hat{\Lambda}_1(t)^\top \dot{\hat{P}}(t) dt + \int_0^\tau \hat{\Lambda}_2(t) \dot{\hat{P}}(t - \tau)^\top dt \\ &\quad + \hat{\Lambda}_3^\top \hat{A}_0 \hat{P}(0) + \hat{\Lambda}_3 \hat{A}_0 \hat{P}(0)^\top + \hat{\Lambda}_3^\top \hat{A}_\tau \hat{P}(\tau) + \hat{\Lambda}_3 \hat{A}_\tau \hat{P}(-\tau)^\top \\ &= -2 \int_0^\tau \tilde{Q}(\tau - t)^\top A_\tau \dot{\tilde{P}}(t) dt + 2\tilde{Q}(0)^\top A_0 \tilde{P}(0) + 2\tilde{Q}(0)^\top A_\tau \tilde{P}(\tau) \\ &\quad - 2 \int_0^\tau \hat{Q}(\tau - t)^\top \hat{A}_\tau \dot{\hat{P}}(t) dt + 2\hat{Q}(0) \hat{A}_0 \hat{P}(0) + 2\hat{Q}(0) \hat{A}_\tau \hat{P}(\tau), \\ \nabla_{\hat{A}_0} \mathcal{L} &= - \int_0^\tau \tilde{\Lambda}_1(t)^\top \tilde{P}(t) dt + \tilde{\Lambda}_3^\top E \tilde{P}(0) \\ &\quad - \int_0^\tau \hat{\Lambda}_1(t)^\top \hat{P}(t) dt + \int_0^\tau \hat{\Lambda}_2(t) \hat{P}(t - \tau)^\top dt \\ &\quad + \hat{\Lambda}_3 \hat{E} \hat{P}(0)^\top + \hat{\Lambda}_3^\top \hat{E} \hat{P}(0) \\ &= 2 \int_0^\tau \tilde{Q}(\tau - t)^\top A_\tau \tilde{P}(t) dt + 2\tilde{Q}(0)^\top E \tilde{P}(0) \\ &\quad + 2 \int_0^\tau \hat{Q}(\tau - t)^\top \hat{A}_\tau \hat{P}(t) dt + 2\hat{Q}(0) \hat{E} \hat{P}(0),\end{aligned}$$

$$\begin{aligned}
 \nabla_{\hat{A}_\tau} \mathcal{L} &= - \int_0^\tau \tilde{\Lambda}_1(t)^\top \tilde{P}(t - \tau) dt + \tilde{\Lambda}_3^\top E \tilde{P}(-\tau) \\
 &\quad - \int_0^\tau \hat{\Lambda}_1(t)^\top \hat{P}(t - \tau) dt + \int_0^\tau \hat{\Lambda}_2(t) \hat{P}(t)^\top dt \\
 &\quad + \hat{\Lambda}_3 \hat{E} \hat{P}(\tau)^\top + \hat{\Lambda}_3^\top \hat{E} \hat{P}(-\tau) \\
 &= 2 \int_0^\tau \tilde{Q}(\tau - t)^\top A_\tau \tilde{P}(t - \tau) dt + 2 \tilde{Q}(0)^\top E \tilde{P}(-\tau) \\
 &\quad + 2 \int_0^\tau \hat{Q}(\tau - t)^\top \hat{A}_\tau \hat{P}(t - \tau) dt + 2 \hat{Q}(0) \hat{E} \hat{P}(-\tau), \\
 \nabla_{\hat{B}} \mathcal{L} &= \tilde{\Lambda}_3^\top B + \hat{\Lambda}_3 \hat{B} + \hat{\Lambda}_3^\top \hat{B} \\
 &= 2 \tilde{Q}(0)^\top B + 2 \hat{Q}(0) \hat{B}, \\
 \nabla_{\hat{C}} \mathcal{L} &= -2C \tilde{P}(0) + 2\hat{C} \hat{P}(0). \quad \square
 \end{aligned}$$

From the previous theorem, we directly find the Wilson-type necessary optimality conditions for  $\mathcal{H}_2$ -optimal MOR of LTD systems.

**Corollary 6.15:**

Let (6.23) be an  $\mathcal{H}_2$ -optimal ROM for (6.18). Then

$$\begin{aligned}
 0 &= \tilde{Q}(0)^\top E \tilde{P}(0) + \int_0^\tau \tilde{Q}(\tau - t)^\top A_\tau \tilde{P}(t) dt \\
 &\quad + \hat{Q}(0) \hat{E} \hat{P}(0) + \int_0^\tau \hat{Q}(\tau - t)^\top \hat{A}_\tau \hat{P}(t) dt, \\
 0 &= \tilde{Q}(0)^\top E \tilde{P}(-\tau) + \int_0^\tau \tilde{Q}(\tau - t)^\top A_\tau \tilde{P}(t - \tau) dt \\
 &\quad + \hat{Q}(0) \hat{E} \hat{P}(-\tau) + \int_0^\tau \hat{Q}(\tau - t)^\top \hat{A}_\tau \hat{P}(t - \tau) dt, \\
 0 &= \tilde{Q}(0)^\top A_0 \tilde{P}(0) + \tilde{Q}(0)^\top A_\tau \tilde{P}(\tau) - \int_0^\tau \tilde{Q}(\tau - t)^\top A_\tau \dot{\tilde{P}}(t) dt \\
 &\quad + \hat{Q}(0) \hat{A}_0 \hat{P}(0) + \hat{Q}(0) \hat{A}_\tau \hat{P}(\tau) - \int_0^\tau \hat{Q}(\tau - t)^\top \hat{A}_\tau \dot{\hat{P}}(t) dt, \\
 0 &= \tilde{Q}(0)^\top B + \hat{Q}(0) \hat{B}, \\
 0 &= C \tilde{P}(0) - \hat{C} \hat{P}(0). \quad \diamond
 \end{aligned}$$

## 6.6 Conclusion

We derived Wilson-type necessary optimality conditions for  $\mathcal{H}_2$ -optimal MOR of second-order systems, port-Hamiltonian systems, LTD systems with a single delay, and  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal MOR of PLTI systems. Similar approach can be taken for, e.g.,  $\mathcal{H}_2$ -optimal MOR of linear time-varying (LTV) systems and  $\mathcal{L}_2$ -optimal MOR of stationary parametric systems.



**Contents**


---

7.1 Summary . . . . .	139
7.2 Future research perspectives . . . . .	140

---

**7.1 Summary**

In this thesis, we have investigated structure-preserving MOR problems for different types of structured systems, mainly network systems.

In Chapter 3, we studied clustering-based MOR of multi-agent systems. In the first part, we have focused on linear multi-agent systems and considered the problem of  $\mathcal{H}_2$ -optimal clustering-based MOR. Based on the relaxation of the discrete optimization problem, we have proposed combining IRKA with a clustering algorithm. Next, we have generalized this to a framework of combining a projection-based method with a clustering algorithm and applied it to nonlinear multi-agent systems.

In Chapter 4, we have considered more theoretically the error due to clustering. First, we have looked at linear multi-agent systems and derived  $\mathcal{H}_2$  and  $\mathcal{H}_\infty$  error bounds when using an AEP. We have also proposed an extension to arbitrary partitions using the distance to a graph for which the partition becomes almost equitable. Next, we have considered nonlinear power systems, a type of nonlinear multi-agent systems. There, we derive equivalence conditions for clustering with zero error based on graph symmetries and equitable partitions.

In Chapter 5, we have studied the problem of subsystem reduction for linear network systems. In the first part, we extend a balancing-based MOR method which preserves stability for network systems satisfying a certain small-gain condition. Using the known a priori  $\mathcal{H}_\infty$  error bound, it allows automatic choice of the order of the reduced subsystems. In the second part, we considered  $\mathcal{H}_2$ -optimal subsystem reduction. Using the Gramian-based formulation of the  $\mathcal{H}_2$ -error, we have derived gradients with respect to matrices defining the ROM. Thereby, we have also obtained Wilson-type necessary

optimality conditions.

In Chapter 6, we have used the ideas from Chapter 5 to other structure-preserving  $\mathcal{H}_2$ -optimal MOR problems. In particular, we have considered structure-preserving MOR for second-order systems, port-Hamiltonian systems, and time-delay systems. Additionally, we have also considered  $\mathcal{H}_2 \otimes \mathcal{L}_2$ -optimal MOR for parametric systems. We have derived Wilson-type necessary optimality conditions and for some systems also the interpolatory optimality conditions.

## 7.2 Future research perspectives

We considered different aspects of structure-preserving MOR, which motivate future research in new methods for network systems and other structured systems.

In Chapter 3, we have found on a small-scale example that combining a projection-based method and a clustering algorithms gives a partition close to the optimal one. However, it is not clear whether this is true in general. Deriving an error bound for this approach would be an interesting goal for future research. Furthermore, comparing different clustering algorithms and with a theoretical analysis would be a worthwhile. Additionally, extending the framework to multi-agent system evolving over directed graphs should be possible.

In Chapter 4, we have derived error bounds for clustering-based MOR of certain linear multi-agent systems. Extending the results to more general multi-agent systems, e.g., with different output functions or with directed underlying graph, would be an interesting problem. Finding easy to compute error bounds for non-almost equitable partitions remains an open problem. For power systems, we have focused on reduction with zero error. Deriving general error bounds remains an open problem.

In the first part of Chapter 5, we have considered subsystem reduction for network systems satisfying a particular small-gain condition. Relaxing this condition, e.g., by dissipativity, would extend the applicability of this approach. In the second part, we have derived Wilson-type necessary  $\mathcal{H}_2$ -optimality conditions. Finding interpolatory conditions and developing an interpolatory method could be an interesting problem.

In Chapter 6, as in the previous, finding interpolatory conditions for  $\mathcal{H}_2$ -optimal MOR of time-delay systems would be an interesting research question. Additionally, efficient implementation for large-scale systems is an open problem.

## BIBLIOGRAPHY

- [AA13] A. M. Annaswamy and M. Amin. *IEEE Vision for Smart Grid Controls: 2030 and Beyond*. IEEE, June 2013. doi:10.1109/IEEEESTD.2013.6577608. 84
- [ABG10] A. C. Antoulas, C. A. Beattie, and S. Gugercin. Interpolatory model reduction of large-scale dynamical systems. In Javad Mohammadpour and Karolos M. Grigoriadis, editors, *Efficient Modeling and Control of Large-Scale Systems*, pages 3–58. Springer US, 2010. doi:10.1007/978-1-4419-5757-3\_1. 19, 20, 22, 23
- [ABGA13] B. Anić, C. Beattie, S. Gugercin, and A. C. Antoulas. Interpolatory weighted- $\mathcal{H}_2$  model reduction. *Automatica J. IFAC*, 49(5):1275–1280, 2013. doi:10.1016/j.automatica.2013.01.040. 105
- [Ant05] 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. 1, 13, 14, 15, 16, 17, 18, 19, 20
- [BB14] C. A. Beattie and P. Benner.  $\mathcal{H}_2$ -optimality conditions for structured dynamical systems. Preprint MPIMD/14-18, Max Planck Institute Magdeburg, 2014. Available from <http://www.mpi-magdeburg.mpg.de/preprints/>. 114, 118, 122, 124, 125, 126
- [BBBG11] U. Baur, C. A. Beattie, P. Benner, and S. Gugercin. Interpolatory projection methods for parameterized model reduction. *SIAM J. Sci. Comput.*, 33(5):2489–2518, 2011. doi:10.1137/090776925. 128
- [BBG15] T. Breiten, C. Beattie, and S. Gugercin. Near-optimal frequency-weighted interpolatory model reduction. *Syst. Control Lett.*, 78:8–18, 2015. doi:10.1016/j.sysconle.2015.01.005. 105
- [BBG19] T. Breiten, C. A. Beattie, and S. Gugercin.  $\mathcal{H}_2$ -gap model reduction for stabilizable and detectable systems. e-print 1909.13764, arXiv, 2019. math.NA. URL: <https://arxiv.org/abs/1909.13764>. 30
- [BFF<sup>+</sup>14] P. Benner, R. Findeisen, D. Flockerzi, U. Reichl, and K. Sundmacher, editors. *Large-Scale Networks in Engineering and Life Sciences*. Modeling and Simulation in Science, Engineering and Technology. Birkhäuser, Basel, CH, 2014. doi:10.1007/978-3-319-08437-4. 1

- [BG09] C. A. Beattie and S. Gugercin. Interpolatory projection methods for structure-preserving model reduction. *Syst. Control Lett.*, 58(3):225–232, 2009. doi:10.1016/j.sysconle.2008.10.016. 114
- [BG12] C. A. Beattie and S. Gugercin. Realization-independent  $\mathcal{H}_2$ -approximation. In *51st IEEE Conference on Decision and Control (CDC)*, pages 4953–4958, 2012. doi:10.1109/CDC.2012.6426344. 29
- [BGM17] C. Beattie, S. Gugercin, and V. Mehrmann. Model reduction for systems with inhomogeneous initial conditions. *Syst. Control Lett.*, 99:99–106, 2017. doi:10.1016/j.sysconle.2016.11.007. 19
- [BGW12] C. A. Beattie, S. Gugercin, and S. Wyatt. Inexact solves in interpolatory model reduction. *Linear Algebra Appl.*, 436(8):2916–2943, 2012. Special Issue dedicated to Danny Sorensen’s 65th birthday. doi:10.1016/j.laa.2011.07.015. 43
- [BKS11] P. Benner, M. Köhler, and J. Saak. Sparse-dense Sylvester equations in  $\mathcal{H}_2$ -model order reduction. Preprint MPIMD/11-11, Max Planck Institute Magdeburg, December 2011. Available from <http://www.mpi-magdeburg.mpg.de/preprints/>. 22, 29
- [BMNP04] M. Barrault, Y. Maday, N. C. Nguyen, and A. T. Patera. An ‘empirical interpolation’ method: application to efficient reduced-basis discretization of partial differential equations. *C. R. Math. Acad. Sci. Paris*, 339(9):667–672, 2004. doi:10.1016/j.crma.2004.08.006. 129
- [BMS05] P. Benner, V. Mehrmann, and D. C. Sorensen. *Dimension Reduction of Large-Scale Systems*, volume 45 of *Lect. Notes Comput. Sci. Eng.* Springer-Verlag, Berlin/Heidelberg, Germany, 2005. doi:10.1007/3-540-27909-1. 1
- [BNBG10] C. Boess, N. K. Nichols, and A. Bunse-Gerstner. Model order reduction for discrete unstable control systems using a balanced truncation approach. Preprint MPS\_2010\_06, University of Reading, 2010. Available from <https://www.reading.ac.uk/math-and-stats/research/math-preprints.aspx>. 29
- [BOCW17] P. Benner, M. Ohlberger, A. Cohen, and K. Willcox. *Model Reduction and Approximation*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017. doi:10.1137/1.9781611974829. 1
- [BP94] A. Berman and R. J. Plemmons. *Nonnegative Matrices in the Mathematical Sciences*. Society for Industrial and Applied Mathematics, 1994. doi:10.1137/1.9781611971262. 87

- 
- [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-Mitteilungen*, 36(1):32–52, 2013. doi:10.1002/gamm.201310003. 18
- [BSJ16] B. Besselink, H. Sandberg, and K. H. Johansson. Clustering-based model reduction of networked passive systems. *IEEE Trans. Autom. Control*, 61(10):2958–2973, 2016. doi:10.1109/TAC.2015.2505418. 2, 34, 40, 82
- [CDR07] D. M. Cardoso, C. Delorme, and P. Rama. Laplacian eigenvectors and eigenvalues and almost equitable partitions. *European J. Combin.*, 28(3):665–673, 2007. doi:10.1016/j.ejc.2005.03.006. 34
- [Chu87] E. K. Chu. The solution of the matrix equations  $AXB - CXD = E$  and  $(YA - DZ, YC - BZ) = (E, F)$ . *Linear Algebra Appl.*, 93(0):93–105, 1987. doi:10.1016/S0024-3795(87)90314-4. 8
- [CKS16] X. Cheng, Y. Kawano, and J. M. A. Scherpen. Graph structure-preserving model reduction of linear network systems. In *European Control Conference (ECC)*, pages 1970–1975, 2016. doi:10.1109/ECC.2016.7810580. 1, 34, 36, 82, 84
- [CKS17] X. Cheng, Y. Kawano, and J. M. A. Scherpen. Reduction of second-order network systems with structure preservation. *IEEE Trans. Autom. Control*, 62(10):5026–5038, 2017. doi:10.1109/TAC.2017.2679479. 1, 84
- [CKS18] X. Cheng, Y. Kawano, and J. M. A. Scherpen. Model reduction of multi-agent systems using dissimilarity-based clustering. *IEEE Trans. Autom. Control*, 2018. doi:10.1109/TAC.2018.2853578. 34
- [CLVVD06] Y. Chahlaoui, D. Lemonnier, A. Vandendorpe, and P. Van Dooren. Second-order balanced truncation. *Linear Algebra Appl.*, 415(2–3):373–384, 2006. doi:10.1016/j.laa.2004.03.032. 114
- [CM11] A. Chapman and M. Mesbahi. UAV flocking with wind gusts: Adaptive topology and model reduction. In *Proceedings of the American Control Conference*, pages 1045–1050, June 2011. doi:10.1109/ACC.2011.5990799. 39
- [Col12] R. Coleman. *Calculus on Normed Vector Spaces*. Universitext. Springer New York, 2012. doi:10.1007/978-1-4614-3894-6. 8, 9, 10, 11, 12
- [DB14] F. Dörfler and F. Bullo. Synchronization in complex networks of phase oscillators: A survey. *Automatica J. IFAC*, 50(6):1539–1564, 2014. doi:10.1016/j.automatica.2014.04.012. 84

- [DP00] G. E. Dullerud and F. Paganini. *A Course in Robust Control Theory*. Springer New York, 2000. doi:[10.1007/978-1-4757-3290-0](https://doi.org/10.1007/978-1-4757-3290-0). 15, 16, 17
- [EFHO10] E. Estrada, M. Fox, D. J. Higham, and G.-L. Oppo, editors. *Network Science: Complexity in Nature and Technology*. Springer-Verlag, London, UK, 2010. doi:[10.1007/978-1-84996-396-1](https://doi.org/10.1007/978-1-84996-396-1). 1
- [Enn85] D. F. Enns. *Model Reduction for Control System Design*. PhD thesis, Stanford Univ., March 1985. URL: <https://ntrs.nasa.gov/search.jsp?R=19850014087>. 29
- [GAB08] S. Gugercin, A. C. Antoulas, and C. Beattie.  $\mathcal{H}_2$  model reduction for large-scale linear dynamical systems. *SIAM J. Matrix Anal. Appl.*, 30(2):609–638, 2008. doi:[10.1137/060666123](https://doi.org/10.1137/060666123). 22, 23
- [GPBvdS12] S. Gugercin, R. V. Polyuga, C. Beattie, and A. van der Schaft. Structure-preserving tangential interpolation for model reduction of port-Hamiltonian systems. *Automatica*, 48(9):1963–1974, 2012. doi:[10.1016/j.automatica.2012.05.052](https://doi.org/10.1016/j.automatica.2012.05.052). 124
- [GR01] C. Godsil and G. Royle. *Algebraic graph theory*, volume 207 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 2001. doi:[10.1007/978-1-4613-0163-9](https://doi.org/10.1007/978-1-4613-0163-9). 30, 31
- [GV13] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, fourth edition, 2013. 6, 7
- [Haa17] B. Haasdonk. Reduced basis methods for parametrized PDEs—a tutorial introduction for stationary and instationary problems. In P. Benner, A. Cohen, M. Ohlberger, and K. Willcox, editors, *Model Reduction and Approximation: Theory and Algorithms*, pages 65–136. SIAM, 2017. doi:[10.1137/1.9781611974829.ch2](https://doi.org/10.1137/1.9781611974829.ch2). 129
- [HJ85] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, Cambridge, 1985. 6, 7, 102, 103
- [HRA11] M. Heinkenschloss, T. Reis, and A. C. Antoulas. Balanced truncation model reduction for systems with inhomogeneous initial conditions. *Automatica J. IFAC*, 47(3):559–564, 2011. doi:[10.1016/j.automatica.2010.12.002](https://doi.org/10.1016/j.automatica.2010.12.002). 19
- [HV05] M. Hinze and S. Volkwein. Proper orthogonal decomposition surrogate models for nonlinear dynamical systems: Error estimates and suboptimal control. In P. Benner, V. Mehrmann, and D.C. Sorensen, editors,

- 
- Dimension Reduction of Large-Scale Systems*, volume 45 of *Lect. Notes Comput. Sci. Eng.*, pages 261–306. Springer-Verlag, Berlin/Heidelberg, Germany, 2005. [51](#)
- [HW79] J. A. Hartigan and M. A. Wong. Algorithm AS 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1):100–108, 1979. [doi:10.2307/2346830](#). [42](#)
- [III15] T. Ishizaki and J. Imura. Clustered model reduction of interconnected second-order systems. *Nonlinear Theory and Its Applications, IEICE*, 6(1):26–37, 2015. [doi:10.1587/nolta.6.26](#). [1](#), [83](#)
- [IKG<sup>+</sup>12] T. Ishizaki, K. Kashima, A. Girard, J. Imura, L. Chen, and K. Aihara. Clustering-based  $\mathcal{H}_2$ -state aggregation of positive networks and its application to reduction of chemical master equations. In *51st IEEE Conference on Decision and Control (CDC)*, pages 4175–4180, December 2012. [doi:10.1109/CDC.2012.6426793](#). [40](#)
- [IKG<sup>+</sup>15] T. Ishizaki, K. Kashima, A. Girard, J. Imura, L. Chen, and K. Aihara. Clustered model reduction of positive directed networks. *Automatica J. IFAC*, 59:238–247, 2015. [doi:10.1016/j.automatica.2015.06.027](#). [xi](#), [1](#), [34](#), [36](#), [82](#), [83](#), [84](#)
- [IKI16a] T. Ishizaki, R. Ku, and J. Imura. Clustered model reduction of networked dissipative systems. In *American Control Conference (ACC)*, pages 3662–3667, 2016. [doi:10.1109/ACC.2016.7525482](#). [xi](#), [2](#), [34](#), [35](#), [82](#), [83](#)
- [IKI16b] T. Ishizaki, R. Ku, and J. Imura. Eigenstructure analysis from symmetrical graph motives with application to aggregated controller design. In *55th IEEE Conference on Decision and Control (CDC)*, pages 5744–5749, 2016. [doi:10.1109/CDC.2016.7799152](#). [84](#)
- [IKIA14] T. Ishizaki, K. Kashima, J. Imura, and K. Aihara. Model reduction and clusterization of large-scale bidirectional networks. *IEEE Trans. Autom. Control*, 59(1):48–63, January 2014. [doi:10.1109/TAC.2013.2275891](#). [xi](#), [1](#), [34](#), [36](#), [39](#), [40](#), [56](#), [82](#), [83](#), [84](#)
- [JVM11] E. Jarlebring, J. Vanbiervliet, and W. Michiels. Characterizing and computing the  $\mathcal{H}_2$  norm of time-delay systems by solving the delay Lyapunov equation. *IEEE Trans. Autom. Control*, 56(4):814–825, April 2011. [doi:10.1109/TAC.2010.2067510](#). [133](#)
- [Kro39] G. Kron. *Tensor analysis of networks*. Wiley, New York, 1939. [87](#)

- [Kun94] P Kundur. *Power System Stability and Control*. McGraw-Hill, 1994. 84, 87
- [LDCH10] Z. Li, Z. Duan, G. Chen, and L. Huang. Consensus of multiagent systems and synchronization of complex networks: A unified viewpoint. *IEEE Trans. Circuits Syst. I, Regular Papers*, 57(1):213–224, January 2010. doi:10.1109/TCSI.2009.2023937. 1, 36
- [MBG10] C. Magruder, C. Beattie, and S. Gugercin. Rational krylov methods for optimal  $\mathcal{L}_2$  model reduction. In *49th IEEE Conference on Decision and Control (CDC)*, pages 6797–6802, December 2010. doi:10.1109/CDC.2010.5717454. 30
- [ME10] M. Mesbahi and M. Egerstedt. *Graph Theoretic Methods in Multiagent Networks*. Princeton Series in Applied Mathematics. Princeton University Press, Princeton, NJ, 2010. doi:10.1515/9781400835355. 30, 31, 33, 34, 59
- [MEB08] S. Martini, M. Egerstedt, and A. Bicchi. Controllability decompositions of networked systems through quotient graphs. In *47th IEEE Conference on Decision and Control (CDC)*, pages 5244–5249, December 2008. doi:10.1109/CDC.2008.4739213. 39
- [MEB10] S. Martini, M. Egerstedt, and A. Bicchi. Controllability analysis of multi-agent systems using relaxed equitable partitions. *Int. J. Syst., Control Commun.*, 2(1/2/3):100–121, January 2010. doi:10.1504/IJSCC.2010.031160. 39
- [MGB15] P. Mlinarić, S. Grundel, and P. Benner. Efficient model order reduction for multi-agent systems using QR decomposition-based clustering. In *54th IEEE Conference on Decision and Control (CDC)*, pages 4794–4799, December 2015. doi:10.1109/CDC.2015.7402967. 84
- [ML67] L. Meier and D. G. Luenberger. Approximation of linear constant systems. *IEEE Trans. Autom. Control*, 12(5):585–588, 1967. doi:10.1109/TAC.1967.1098680. 22
- [MN99] J. R. Magnus and H. Neudecker. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. Probability and Statistics. Wiley, 1999. 7
- [Mor05] L. Moreau. Stability of multiagent systems with time-dependent communication links. *IEEE Trans. Autom. Control*, 50(2):169–182, February 2005. doi:10.1109/TAC.2004.841888. 1



- 
- [MS96] D. G. Meyer and S. Srinivasan. Balancing and model reduction for second-order form linear systems. *IEEE Trans. Autom. Control*, 41(11):1632–1644, 1996. doi:10.1109/9.544000. 114
- [MS05] V. Mehrmann and T. Stykel. Balanced truncation model reduction for large-scale systems in descriptor form. In P. Benner, V. Mehrmann, and D. C. Sorensen, editors, *Dimension Reduction of Large-Scale Systems*, volume 45 of *Lect. Notes Comput. Sci. Eng.*, pages 83–115. Springer-Verlag, Berlin/Heidelberg, Germany, 2005. doi:10.1007/3-540-27909-1\\_3. 17
- [MTC13] N. Monshizadeh, H. L. Trentelman, and M. K. Camlibel. Stability and synchronization preserving model reduction of multi-agent systems. *Syst. Control Lett.*, 62(1):1–10, 2013. doi:10.1016/j.sysconle.2012.10.011. 2, 34, 36, 100, 102
- [MTC14] N. Monshizadeh, H. L. Trentelman, and M. K. Camlibel. Projection-based model reduction of multi-agent systems using graph partitions. *IEEE Trans. Control Netw. Syst.*, 1(2):145–154, June 2014. doi:10.1109/TCNS.2014.2311883. xi, 31, 34, 36, 37, 40, 44, 45, 46, 47, 54, 58, 61, 68, 80, 81, 97, 151
- [MZ10] C.-Q. Ma and J.-F. Zhang. Necessary and sufficient conditions for consensusability of linear multi-agent systems. *IEEE Trans. Autom. Control*, 55(5):1263–1268, May 2010. doi:10.1109/TAC.2010.2042764. 1
- [New10] M. E. J. Newman. *Networks: An Introduction*. Oxford University Press, Inc., New York, NY, USA, March 2010. doi:10.1093/acprof:oso/9780199206650.001.0001. 1
- [OJ88] P. C. Opdenacker and E. A. Jonckheere. A contraction mapping preserving balanced reduction scheme and its infinity norm error bounds. *IEEE Trans. Circuits Syst.*, 35(2):184–189, 1988. doi:10.1109/31.1720. xv, 101
- [OP05] R.C.L.F. Oliveira and P.L.D. Peres. Stability of polytopes of matrices via affine parameter-dependent Lyapunov functions: Asymptotically exact LMI conditions. *Linear Algebra Appl.*, 405:209–228, 2005. doi:10.1016/j.laa.2005.03.019. 80
- [ORW13] M. Opmeer, T. Reis, and W. Wollner. Finite-rank ADI iteration for operator Lyapunov equations. *SIAM J. Control Optim.*, 51(5):4084–4117, 2013. doi:10.1137/120885310. 29
- [OSM03] R. Olfati-Saber and R. M. Murray. Consensus protocols for networks of dynamic agents. In *Proceedings of the American Controls Conference*,

- volume 2, pages 951–956, June 2003. doi:10.1109/ACC.2003.1239709. 1
- [Ran96] A. Rantzer. On the Kalman-Yakubovich-Popov lemma. *Syst. Cont. Lett.*, 28(1):7–10, 1996. doi:10.1016/0167-6911(95)00063-1. 56
- [RJME09] A. Rahmani, M. Ji, M. Mesbahi, and M. Egerstedt. Controllability of multi-agent systems from a graph-theoretic perspective. *SIAM J. Control Optim.*, 48(1):162–186, 2009. doi:10.1137/060674909. 39, 54, 84
- [RS07] T. Reis and T. Stykel. Stability analysis and model order reduction of coupled systems. *Math. Comput. Model. Dyn. Syst.*, 13(5):413–436, 2007. doi:10.1080/13873950701189071. xii, 2, 99, 100, 102, 103, 104, 105
- [RS08a] T. Reis and T. Stykel. Balanced truncation model reduction of second-order systems. *Math. Comput. Model. Dyn. Syst.*, 14(5):391–406, 2008. doi:10.1080/13873950701844170. 114
- [RS08b] T. Reis and T. Stykel. A survey on model reduction of coupled systems. In W. H. A. Schilders, H. A. van der Vorst, and J. Rommes, editors, *Model Order Reduction: Theory, Research Aspects and Applications*, pages 133–155. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008. doi:10.1007/978-3-540-78841-6\_7. 2, 99
- [Sat17] K. Sato. Riemannian optimal model reduction of linear second-order systems. *IEEE Contr. Syst. Lett.*, 1(1):2–7, July 2017. doi:10.1109/LCSYS.2017.2698178. 114
- [Sch07] S. E. Schaeffer. Graph clustering. *Comput. Sci. Rev.*, 1(1):27–64, 2007. doi:10.1016/j.cosrev.2007.05.001. 39, 41
- [Sim16] V. Simoncini. Computational methods for linear matrix equations. *SIAM Review*, 58(3):377–441, 2016. doi:10.1137/130912839. 18
- [SM09] H. Sandberg and R. M. Murray. Model reduction of interconnected linear systems. *Optim. Control Appl. Methods*, 30(3):225–245, 2009. doi:10.1002/oca.854. 2, 99
- [TTM13] H. L. Trentelman, K. Takaba, and N. Monshizadeh. Robust synchronization of uncertain linear multi-agent systems. *IEEE Trans. Autom. Control*, 58(6):1511–1523, 2013. doi:10.1109/TAC.2013.2239011. 59
- [VDGA08] P. Van Dooren, K. Gallivan, and P.-A. Absil.  $\mathcal{H}_2$ -optimal model reduction of MIMO systems. *Appl. Math. Lett.*, 21:1267–1273, 2008. doi:10.1016/j.aml.2007.09.015. 22, 24, 26

- 
- [VDGA10] P. Van Dooren, K. A. Gallivan, and P.-A. Absil.  $\mathcal{H}_2$ -optimal model reduction with higher-order poles. *SIAM J. Matrix Anal. Appl.*, 31(5):2738–2753, 2010. doi:[10.1137/080731591](https://doi.org/10.1137/080731591). 22, 29
- [vdSJ14] A. van der Schaft and D. Jeltsema. *Port-Hamiltonian Systems Theory: An Introductory Overview*. now, 2014. doi:[10.1561/2600000002](https://doi.org/10.1561/2600000002). 123
- [VVD08] A. Vandendorpe and P. Van Dooren. Model reduction of interconnected systems. In W. H. A. Schilders, H. A. van der Vorst, and J. Rommes, editors, *Model Order Reduction: Theory, Research Aspects and Applications*, volume 13 of *Mathematics in Industry*, pages 305–321. Springer, Berlin, Heidelberg, 2008. doi:[10.1007/978-3-540-78841-6\\_14](https://doi.org/10.1007/978-3-540-78841-6_14). 2, 99, 105
- [Wil70] D. A. Wilson. Optimum solution of model-reduction problem. *Proceedings of the Institution of Electrical Engineers*, 117(6):1161–1165, 1970. doi:[10.1049/piee.1970.0227](https://doi.org/10.1049/piee.1970.0227). 22, 24
- [Wya12] S. Wyatt. *Issues in Interpolatory Model Reduction: Inexact Solves, Second-order Systems and DAEs*. PhD thesis, Virginia Polytechnic Institute and State University, Blacksburg, Virginia, USA, May 2012. URL: <http://hdl.handle.net/10919/27668>. 114
- [XC16] N. Xue and A. Chakraborty. Optimal control of large-scale networks using clustering based projections. e-print 1609.05265, arXiv, 2016. cs.SY. URL: <https://arxiv.org/abs/1609.05265>. 84
- [XZ11] Y. Xu and T. Zeng. Optimal  $\mathcal{H}_2$  model reduction for large scale MIMO systems via tangential interpolation. *Int. J. Numer. Anal. Model.*, 8(1):174–188, 2011. URL: <http://www.math.ualberta.ca/ijnam/Volume-8-2011/No-1-11/2011-01-10.pdf>. 22, 26
- [YCDAGX93] J. Yang, C. S. Chen, J. A. De Abreu-Garcia, and Y. Xu. Model reduction of unstable systems. *International Journal of Systems Science*, 24(12):2407–2414, 1993. doi:[10.1080/00207729308949638](https://doi.org/10.1080/00207729308949638). 29
- [ZCC14] S. Zhang, M. Cao, and M. K. Camlibel. Upper and lower bounds for controllable subspaces of networks of diffusively coupled agents. *IEEE Trans. Autom. Control*, 59(3):745–750, March 2014. doi:[10.1109/TAC.2013.2275666](https://doi.org/10.1109/TAC.2013.2275666). 34
- [ZDG96] K. Zhou, J. C. Doyle, and K. Glover. *Robust and Optimal Control*. Prentice-Hall, Upper Saddle River, NJ, 1996. 18
- [Zei85a] E. Zeidler. *Nonlinear Functional Analysis and its Applications I: Fixed Point Theorems*. Springer-Verlag, 1985. 8, 9, 10, 11, 12

- [Zei85b] E. Zeidler. *Nonlinear Functional Analysis and its Applications III: Variational Methods and Optimization*. Springer-Verlag, 1985. doi:[10.1007/978-1-4612-5020-3](https://doi.org/10.1007/978-1-4612-5020-3). 12
- [ZHD<sup>+</sup>01] H. Zha, X. He, C. Ding, H. Simon, and M. Gu. Spectral relaxation for k-means clustering. In *Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic*, pages 1057–1064, 2001. URL: <https://papers.nips.cc/paper/1992-spectral-relaxation-for-k-means-clustering.pdf>. xv, 41, 42
- [Zil91] A. Zilouchian. Balanced structures and model reduction of unstable systems. In *IEEE Proceedings of the SOUTHEASTCON '91*, pages 1198–1201 vol.2, April 1991. doi:[10.1109/SECON.1991.147956](https://doi.org/10.1109/SECON.1991.147956). 29
- [ZSW99] K. Zhou, G. Salomon, and E. Wu. Balanced realization and model reduction for unstable systems. *Internat. J. Robust Nonlinear Control*, 9(3):183–198, 1999. doi:[10.1002/\(SICI\)1099-1239\(199903\)9:3<183::AID-RNC399>3.0.CO;2-E](https://doi.org/10.1002/(SICI)1099-1239(199903)9:3<183::AID-RNC399>3.0.CO;2-E). 30

## STATEMENT OF SCIENTIFIC COOPERATIONS

This work is based on articles and reports (published and unpublished) that have been obtained in cooperation with various coauthors. To guarantee a fair assessment of this thesis, this statement clarifies the contributions that each individual coauthor has made. The following people contributed to the content of this work:

- Peter Benner (PB), Max Planck Institute for Dynamics of Complex Technical Systems;
- Aranya Chakraborty (AC), North Carolina State University;
- Sara Grundel (SG), Max Planck Institute for Dynamics of Complex Technical Systems;
- Takayuki Ishizaki (TI), Tokyo Institute of Technology;
- Hidde-Jan Jongsma (HJ), Netherlands Organisation for Applied Scientific Research;
- Harry L. Trentelman (HT), University of Groningen.

### Chapter 3

The work on this chapter started from PB pointing me to [MTC14]. All theoretical and computational results were obtained by myself, but proofread and improved by PB and SG.

### Chapter 4

HJ and HT extended the results about  $\mathcal{H}_2$ -error expressions for multi-agent systems with single-integrator agents from [MTC14] to upper bounds for  $\mathcal{H}_2$ -error for more general multi-agent systems (Theorem 4.5). I derived the expression for  $\mathcal{H}_\infty$ -error for special LTI systems in Lemma 4.2. Jointly with HJ and HT, we derived similar upper bounds for  $\mathcal{H}_\infty$ -error for symmetric multi-agent systems (Theorem 4.12 and Theorem 4.14). Based on an idea from HJ and HT, I derived an upper bound for non-almost equitable partitions in Section 4.2.7.

Section 4.3 is work I extended with help from TI. The problem setting was motivated by AC.

## **Chapter 5**

PB had the idea to extend to results from multi-agent systems to general network systems. All theoretical and computational results were obtained by myself, but proofread and improved by PB and SG.

## **Chapter 6**

The idea and results of this chapter were completely obtained by me.

## EHRENERKLÄRUNG

Ich versichere hiermit, dass ich die vorliegende Arbeit ohne unzulässige Hilfe Dritter und ohne Benutzung anderer als der angegebenen Hilfsmittel angefertigt habe; verwendete fremde und eigene Quellen sind als solche kenntlich gemacht.

Ich habe insbesondere nicht wissentlich:

- Ergebnisse erfunden oder widersprüchliche Ergebnisse verschwiegen,
- statistische Verfahren absichtlich missbraucht, um Daten in ungerechtfertigter Weise zu interpretieren,
- fremde Ergebnisse oder Veröffentlichungen plagiiert oder verzerrt wiedergegeben.

Mir ist bekannt, dass Verstöße gegen das Urheberrecht Unterlassungs- und Schadenersatzansprüche des Urhebers sowie eine strafrechtliche Ahndung durch die Strafverfolgungsbehörden begründen kann.

Die Arbeit wurde bisher weder im Inland noch im Ausland in gleicher oder ähnlicher Form als Dissertation eingereicht und ist als Ganzes auch noch nicht veröffentlicht.

Magdeburg, 30.10.2019

---

Petar Mlinarić