



HAL
open science

Contrôle multi-objectifs d'ordre réduit

Christian Fischer

► **To cite this version:**

Christian Fischer. Contrôle multi-objectifs d'ordre réduit. Autre [cs.OH]. École Nationale Supérieure des Mines de Paris, 2011. Français. NNT : 2011ENMP0033 . pastel-00644122

HAL Id: pastel-00644122

<https://pastel.hal.science/pastel-00644122>

Submitted on 23 Nov 2011

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

École doctorale n°84 :
Sciences et technologies de l'information et de la communication

Doctorat ParisTech

T H È S E

pour obtenir le grade de docteur délivré par

l'École Nationale Supérieure des Mines de Paris

Spécialité « Contrôle, Optimisation et Prospective »

présentée et soutenue publiquement par

Christian FISCHER

le 27 Juillet 2011

Contrôle multi-objectifs d'Ordre réduit

Directeur de thèse : **Nadia MAÏZI**
Co-encadrement de la thèse : **Marc BORDIER**

Jury

M. Michel ZASADZINSKI, Professeur, Université Henri Poincaré, Cosnes-et-Romain
M. Alban QUADRAT, Chargé de Recherche, INRIA Saclay
Mme Nadia MAÏZI, Professeur, Mines ParisTech
M. Marc BORDIER, Chargé de Recherche, Mines ParisTech
Mme Martine OLIVI, Chargé de Recherche, INRIA Sophia Antipolis
M. Rémi DRAI, Docteur-Ingénieur, ESA, Noordwijk, Pays-Bas

Président
Rapporteur
Examinateur
Examinateur
Membre invité
Membre invité

Contents

Remerciements	v
Résumé	vii
1 Introduction	1
1.1 Motivation	1
1.2 Overview	1
1.3 Notation	3
2 Controller Parametrizations	5
2.1 Ring of Transfer Functions	5
2.1.1 Transfer Functions	5
2.1.2 Ring Structure	6
2.2 Controller Parametrization	7
2.2.1 The Control Loop	7
2.2.2 Connectivity	9
2.2.3 Doubly Coprime Factorization	9
2.2.4 Loop Factorization	11
2.2.5 Loop Decomposition	11
2.3 Passage to State Space	13
2.3.1 Realization of Transfer Functions	13
2.3.2 Realization of the Control Loop	14
2.4 Fixed Order Controller Parametrization	14
2.4.1 Need for a new Parametrization	14
2.4.2 Fixed Order Loop Factorization	15
2.4.3 Existence of Solutions	19
2.4.4 Relation to Estimator Corrector Type Controllers	21
2.4.5 Parametrization of Fixed Order Controllers	25
2.4.6 Controllability and Observability	26
2.4.7 Necessary Condition for Loop Ring Membership	27
3 Stable Observable Pairs	29
3.1 Lossless Systems	29
3.1.1 Transfer Function Representation	29
3.1.2 State Space Representation	30
3.1.3 Discrete Lossless Completion	32
3.1.4 Link to Conjugate Realization	34
3.2 Nudelman Interpolation Problem	35

3.2.1	Interpolation Problem Definition	35
3.2.2	Interpolation Problem Solution	39
3.2.3	Link to Conjugate Realization	44
3.2.4	Parametrization of Stable Observable Pairs	46
4	Multiobjective Optimization	49
4.1	Problem Formulation	49
4.1.1	Extended Loops	49
4.1.2	System Norms	51
4.1.3	Optimization Problem	52
4.2	Parametrized Optimization	53
4.2.1	Passage between Controller and Parameter	54
4.2.2	Parametrized Extended Loops	55
4.2.3	Parametrized Optimization Problem	55
4.2.4	Evaluation Step	56
4.2.5	Actualization Step	57
4.2.6	Variation of the Constraint Function	57
4.2.7	Variation of the Extended Loops	59
4.2.8	Variation of the Parameter	61
4.2.9	Realness Constraint	64
4.3	Sufficient Conditions	68
4.3.1	Sufficient Condition for Parameter Construction	68
4.3.2	Sufficient Condition for Subgradient Calculation	68
4.3.3	Sufficient Condition for Controller Reconstruction	69
4.3.4	Sufficient Condition for L1 Parameter Construction	69
4.4	Complete Algorithm	70
4.4.1	Initialization	72
4.4.2	Optimization	73
4.4.3	Finalization	73
5	Conclusion	75
A	Algebra	77
A.1	Ring Structure	77
A.1.1	Rings and Units	77
A.1.2	Ideals and Equivalence	79
A.1.3	Division	81
A.1.4	Factorization	84
A.1.5	Extension to Matrices	85
A.1.6	Extension to Fractions	91
A.2	Matrix Calculus	94
A.2.1	Blocks	94
A.2.2	Spectra	95
A.2.3	Operations	97
A.3	Matrix Equations	101
A.3.1	Underdetermined Linear Equations	101
A.3.2	Sylvester Equations	101
A.3.3	Stein Equations	102

B	Analysis	103
B.1	Complex Analysis	103
B.1.1	Analytic Functions	103
B.1.2	Isolated Singularity	103
B.1.3	Residue	103
B.1.4	Identity Theorem	104
B.1.5	Open Mapping Theorem	104
B.1.6	Maximum Modulus Theorem	104
B.2	Norms	105
B.2.1	Norms for Constants	105
B.2.2	Norms for Functions with Real Variables	106
B.2.3	Norms for Functions with Complex Variables	106
B.3	Optimization	107
B.3.1	Problem Formulation	107
B.3.2	Lagrange Function	107
B.3.3	Primal and Dual Problem	108
B.3.4	First Order Perturbation Analysis	108
C	Systems	111
C.1	Linear Systems	111
C.1.1	System Representations	111
C.1.2	Parseval Theorem	112
C.1.3	L2/H2 System Norm	113
C.1.4	Star Norm	114
C.1.5	Realness	116
C.1.6	Stability	117
C.1.7	Controllability	117
C.1.8	Observability	118
C.1.9	Equivalence	118
C.1.10	Change of Basis	118
C.1.11	Kalman Decomposition	118
C.1.12	McMillan Degree	119
C.2	Connections	120
C.2.1	Flow Inversion	120
C.2.2	Parallel Connection	121
C.2.3	Serial Connection	121
C.2.4	Feedback Loop	122
C.2.5	Common Input	123
C.2.6	Common Output	124
C.2.7	Bridge	125
C.3	Factorizations	126
C.3.1	Right Factorization	126
C.3.2	Left Factorization	127

Remerciements

Tous mes remerciements vont aux membres de l'équipe du Centre de Mathématiques Appliquées pour leur accueil chaleureux et pour la patience avec laquelle ils m'ont répondu à toutes mes questions,

spécialement à Professeur Nadia Maïzi qui m'a donné l'opportunité d'effectuer ma thèse au sein de son laboratoire et qui m'a confié un travail exigeant mais particulièrement intéressant,

à mes tuteurs de thèse Monsieur Marc Bordier et Monsieur Jean-Paul Marmorat pour leur créativité exceptionnelle face aux problèmes scientifiques et pour la bonne atmosphère de travail qu'ils ont su créer,

à Madame Valérie Roy pour l'expérience édifiante comme tuteur dans le cadre des projets MIG de l'école que je ne veux pas m'en passer,

à Madame Josiane Bedos, à Madame Dominique Micollier, à Monsieur Jean-Charles Bonin et à Monsieur Gilles Guerassimoff pour leur soutien, leur aide directe et leurs conseils précieux,

et à mes camarades avec lesquels j'ai travaillé ensemble, qui n'ont jamais hésité de m'aider et qui, grâce à leur humeur, ont rendu mon séjour inoubliable.

Résumé

Le présent rapport décrit une méthode de résolution d'un problème de contrôle multi-objectifs d'un système linéaire multi-entrées multi-sorties. Étant donné un système, on cherche un régulateur tel que la dynamique du système rebouclé sur son régulateur vérifie certaines propriétés. Outre ces contraintes sur la dynamique, on exige que le régulateur optimal ait une complexité – mesurée par son degré de McMillan – donnée. Ce problème, réputé difficile, a fait l'objet, depuis le début des années 2000, de nombreuses recherches [3], [8] et il existe désormais des boîtes à outils Matlab telles que *Hinfstruct* et *Hifoo* qui, à ordre de régulateur fixé, en donnent une solution. Essentiellement, les algorithmes proposés reposent sur des techniques d'optimisation non différentiable du critère par rapport au régulateur lui-même. Ce rapport présente une autre façon de résoudre ce problème. À la différence des techniques mentionnées ci-dessus, on construit un calcul différentiel fondé sur un paramétrage minimal des régulateurs stabilisants. La démarche proposée connecte des idées de trois domaines différents présentés au cours des trois chapitres principaux. Le chapitre 2 développe, à partir du concept de factorisation sur un anneau initié par le travail de [13], une méthode générique de factorisation à ordre donné d'une boucle fermée en faisant apparaître un paramétrage nouveau du régulateur. Le chapitre 3 introduit la notion de système conservatif qui, selon [2], forment une variété différentiable. Cette propriété, appliquée au paramétrage développé à la fin du chapitre précédent, permet alors de bouger le correcteur d'une manière différentiable : pour un régulateur réel d'ordre n et de dimension $p \times m$, le nombre de paramètres réels nécessaires à sa représentation est $n(p + m)$ à comparer avec les $n^2 + n(m + p)$ utilisés par [3] et [8]. Le chapitre 4 introduit le problème d'optimisation multi-objectifs dans lequel les spécifications sont exprimées au moyen de quelques normes communes de système. En appliquant les résultats trouvés en [4] et en utilisant le calcul différentiel développé à la fin du chapitre précédent, il devient possible de calculer un sous gradient du critère qui, par une méthode de descente, permet d'améliorer le régulateur relativement aux spécifications données tout en gardant son ordre fixé.

Dans le chapitre 2, étant donné l'unité \mathbf{P} , on cherche un régulateur \mathbf{K} tel que les éléments formant la matrice de transfert \mathbf{H} de la boucle fermée construite à partir ces deux systèmes linéaires, soient des fonctions de transfert propres ayant toutes leurs pôles dans une zone Z prédéfinie par l'utilisateur. L'ensemble de ces fonctions de transfert a une structure d'anneau et les résultats de la théorie des anneaux présentés en A.1 peuvent être appliqués afin d'obtenir une description de tous les régulateurs qui satisfont cette exigence de conception. En particulier, on montre comment ils engendrent une factorisation de la boucle fermée sur l'anneau de référence et comment l'ensemble des régulateurs factorisants

peut être décrit au moyen du paramètre de Youla-Kucera Q . Le passage à la représentation d'état, qui est numériquement plus pratique, introduit une nouvelle variable à savoir la dimension de l'espace d'état de la représentation. Bien qu'une factorisation soit toujours possible, celle-ci se paye, en général, par des réalisations de grande taille pour chacun des facteurs et par des régulateurs complexes d'ordre élevé. Une factorisation de la boucle fermée à ordre fixé est alors introduite comme un cas spécial de la factorisation générale obtenu en fixant les dimensions des espaces d'états de l'unité et du régulateur. L'avantage de fixer ces dimensions consiste en ce que les régulateurs à ordre réduit, avec une dimension d'espace d'état égale ou inférieure à celle de l'unité, peuvent être paramétrés directement, ce qui est intéressant si la situation de contrôle demande des régulateurs à complexité réduite. Enfin, on montre pourquoi les paramétrages à ordre fixé Q_L, Q_R partagent avec le paramètre Q la propriété d'être chacun un élément de l'anneau pour que la boucle H le soit aussi, un fait qui permettra à Q_L, Q_R d'être eux-mêmes paramétrés au moyen des systèmes conservatifs au cours du chapitre 3.

Le chapitre 3 introduit la notion de système conservatif et fournit un paramétrage des paires stables observables en résolvant un problème d'interpolation du type Nudelman. Il a été montré en [2] que les systèmes conservatifs et les paires stables observables qui peuvent en être extraites forment une variété différentiable. Le régulateur K pouvant ne pas être stable, ce paramétrage ne peut pas être appliqué directement à K . Cependant il a été montré à la fin du chapitre 2 que le paramètre gauche Q_L devait être stable pour que la boucle H soit stable et observable pour que le régulateur K soit observable. Le paramétrage de la paire stable observable contenue dans Q_L permet alors de faire varier de façon infinitésimale Q_L tout entier et de construire un calcul différentiel sur Q_L ce qui sera important au cours l'étape d'optimisation du chapitre 4.

Le chapitre 4 définit un problème d'optimisation multi-objectifs. Pour une unité P donnée, un régulateur K doit être trouvé tel qu'il soit optimal relativement à une fonction de coût donnée et des contraintes données. Ces objectifs d'optimisation sont formulés en termes de normes L1, L2/H2, H_∞ de système. L'avantage d'utiliser les normes L1, L2/H2, H_∞ consiste en ce que leurs calculs peuvent eux-même être énoncés sous forme de problèmes d'optimisation semi-définis standards ce qui permet de calculer les sensibilités locales ainsi que le sous gradient du critère d'optimisation multi-objectifs. La formulation du problème d'optimisation en termes du paramétrage élaboré au cours des chapitres 2 et 3 permet alors de bouger de façon infinitésimale le paramètre Q_L à la place de K dans la direction descendante indiquée par le sous gradient. Dès lors que la tâche d'optimisation est complète, le régulateur K peut être récupéré du paramètre Q_L . Ceci permet une amélioration locale du régulateur initial K_0 donné vers un régulateur supérieur K^* relativement à la fonction de coût et dans le cadre des contraintes.

Pour conclure cette introduction, disons qu'une méthode d'amélioration d'un régulateur relativement à des spécifications multi-objectifs formulée en termes des normes de système L1, L2/H2, H_∞ , est présentée. Trois éléments clés distinguent l'approche proposée. Premièrement, un paramétrage des régulateurs à ordre fixé au moyen d'un paramétrage des systèmes stables et observables, basé sur le travail en [13]. Deuxièmement, un paramétrage de toutes les paires stables observables qui permet un ajustement infinitésimal du régulateur basé sur le tra-

vail en [2]. Et troisièmement, l'utilisation d'une information de sensibilité pour calculer le sous gradient local du problème d'optimisation pour que le régulateur puisse être ajusté de la bonne façon, basé sur le travail en [4]. Bien que beaucoup de mesures de performance pour quantifier les spécifications d'un problème d'optimisation multi-objectifs existent, les normes de système L1, L2/H2, H_∞ sont possiblement les plus importantes. Elles ont été choisies dans ce rapport à cause de leur pertinence pratique et aussi pour illustrer le champ d'application de l'approche proposée. Pour un système donné, si les normes H2 et H_∞ se traduisent par des inégalités matricielles linéaires, la norme L1 s'obtient comme inégalité matricielle bilinéaire. Malgré ces différences, la méthode que nous proposons traite ces contraintes d'une manière uniforme, comme une partie d'une contrainte semi-définie plus générale. La nécessité pour l'approche proposée d'accéder à l'information de sensibilité ne limite nullement les contraintes possibles à ces trois normes. De fait, la solution décrite dans cette thèse est basée sur le paramétrage gauche \mathbf{Q}_L du régulateur auquel correspond un paramétrage des paires stables observables. Cependant l'approche duale est aussi possible. Elle utilise le paramétrage droit \mathbf{Q}_R du régulateur et le paramétrage des paires stables contrôlables. Dans les applications pratiques on peut imaginer d'utiliser \mathbf{Q}_L et le paramétrage des paires stables observables dans le cas où \mathbf{K} a plus d'entrées que de sorties vu que les matrices de sortie $\mathbf{C}_K, \mathbf{C}_{Q_L}$ ont des dimensions inférieures à celles des matrices d'entrée $\mathbf{B}_K, \mathbf{B}_{Q_R}$ et que le paramétrage des paires stables observables de $(\mathbf{A}_{Q_L}, \mathbf{C}_{Q_L})$ devient alors plus simple que le paramétrage des paires stables contrôlables de $(\mathbf{A}_{Q_R}, \mathbf{B}_{Q_R})$ et d'utiliser \mathbf{Q}_R et le paramétrage des paires stables contrôlables dans l'autre cas. Les équations de Riccati (2.21), (2.22) peuvent posséder des solutions $\mathbf{R}_{KP}, \mathbf{R}_{PK}$ multiples pour une unité \mathbf{P} et un régulateur \mathbf{K} donnés. Les conséquences sur le comportement de l'algorithme de descente ne sont pas encore connues. Il peut être possible que ces comportements soient différents pour des $\mathbf{R}_{KP}, \mathbf{R}_{PK}$ différents de sorte que, démarrant d'un point selle ou d'un maximum de la fonction de coût par exemple, le résultat de l'optimisation locale soit totalement différent.

Chapter 1

Introduction

This report describes an attempt to solve a multiobjective optimization problem in the domain of multi input multi output linear time invariant control. This problem consists of a specified plant whose output signals are fed back to its input signals by an unspecified controller so that multiple specified objectives are met. As an additional constraint, the resulting optimal controller is required to have a specified complexity, measured by its order which is, simply put, the dimension of its associated state space.

1.1 Motivation

The idea for this approach started while trying to improve a working controller for a satellite with respect to a specified cost function and constrained by multiple specified design requirements. The environment in which the controller has to run on board of a satellite strongly affects the choice of suitable controllers. Objects in space operate in a very low density environment, almost vacuum. On the one hand, the low material density surrounding a satellite implies almost no fluid damping of oscillations of the main body and its most often long elastic appendages, like solar panels. Modeling a satellite therefore leads to a complex plant of high order with almost undamped elastic modes that can not be neglected. On the other hand, the almost lack of a material screen around the satellite has it exposed to strong electromagnetic radiation. The onboard computers therefore need to be radiation hardened which leads to a lower computation performance as a tradeoff. This explains why in this environment the controller is required to be of a specified order, almost always lower than the order of the plant, which is the satellite model. While common sense suggests that it must be easier to find an optimal controller of high order since there are more parameters and thus more freedom to play with, the converse argument indicates that it also might be more difficult to find a reduced order controller that matches the specifications. This challenge is the motivation for this work.

1.2 Overview

There exist methods *Hinfstruct* and *Hifoo* which have been cast in Matlab toolboxes of the same name, that allow the controller order to be fixed, as presented

in [3] and [8]. These algorithms apply nonsmooth optimization directly to the controller. Another attempt to solve this problem is described in this report. In contrast to the aforementioned techniques, a differentiable parametrization of the controller is developed so that smooth optimization can be applied to the controller parameter and thus indirectly to the controller. It connects ideas from three different areas that are presented in the three main chapters.

Chapter 2 introduces the concept of feedback loop factorizations on algebraic rings pioneered by the work in [13]. However, this approach, while providing a parametrization of all controllers, in general leads to high order controllers that more or less cancel out the plant dynamics and replace it with any desired dynamics. The controller is thus in general working against the nature of the plant. To avoid this, a new approach is presented that uses fixed orders from ground up, starting with fixed order loop factorizations and ending with a fixed order controller parametrization. The parametrization of fixed order controllers inherits an important property from the aforementioned parametrization of all controllers in that the parameter needs to be stable for the feedback loop to be stable. In addition to that the fixed order parameter needs to be observable for the controller to be observable. Hence the fixed order controller parameter is necessarily stable and observable, a fact that is important for the next chapter.

Chapter 3 introduces the notion of square lossless systems and shows that any stable observable pair of a state space realization can be completed to form a lossless system. A parametrization of stable observable pairs is then obtained by solving an interpolation problem for lossless systems. It has been shown in [2] that lossless systems and the stable observable pairs that can be extracted from them form a differentiable manifold. Thus, parametrizing the stable observable pair of the fixed order controller parameter found in the previous chapter allows to move the fixed order controller parameter smoothly, a fact that is important for the next chapter.

Chapter 4 introduces the multiobjective optimization setup with objectives formulated in terms of some common system norms. This problem turns out to be nonconvex and thus difficult to solve. However, fixing the controller parameter leads to a convex solvable problem. Solving the multiobjective optimization problem for a current controller parameter value then provides a local sensitivity information in form of the optimal Lagrange multiplier value. By applying the results in [4] this sensitivity information can be used to calculate the subgradient of the cost function for the current controller parameter value. Due to the smooth parametrization of the controller parameter found in the previous chapter, the information provided by the subgradient enables a local improvement of the controller parameter value with respect to the cost function. Then the multiobjective optimization problem is again solved for the updated controller parameter value and the process repeats until satisfactory results are achieved or no further progress is made.

To ease the reading, an annex is provided with some of the most frequently used formulas and theorems in algebra, analysis and linear time invariant system theory.

1.3 Notation

Within the scope of this report bold lowercase letters represent vectors, bold uppercase letters represent matrices and non bold letters represent scalars. The constants $\mathbf{0}$ and $\mathbf{1}$ represent the additive and multiplicative matrix identities respectively. The matrices \mathbf{A}^T and \mathbf{A}^H represent the transposed and transconjugated matrix \mathbf{A} respectively. The functions $\text{tr } \mathbf{A}$ and $\det \mathbf{A}$ denote the trace and determinant of matrix \mathbf{A} respectively. The functions $\text{rnk } \mathbf{A}$ and $\text{col } \mathbf{A}$ return the rank and number of columns of matrix \mathbf{A} respectively. In an abuse of language, a matrix \mathbf{A} is called *stable* if and only if $\text{Re } \lambda < 0$ for continuous systems or $|\lambda| < 1$ for discrete systems for all eigenvalues λ of matrix \mathbf{A} . The sets $\mathbb{C}, \mathbb{R}, \mathbb{Z}, \mathbb{N}$ denote the sets of complex, real, entire and natural numbers respectively, the set \mathbb{N}_n denotes the natural numbers up to and including the natural number n . A zero subscript denotes an initial value, a star superscript denotes an optimal value.

Chapter 2

Controller Parametrizations

In this chapter, a feedback loop \mathbf{H} is constructed from two linear systems, a given plant \mathbf{P} and a controller \mathbf{K} which is to be designed. The controller has to assure that the transfer matrix of the feedback loop becomes a matrix of proper transfer functions whose poles are confined to a user defined region Z . These transfer functions form the algebraic structure of a ring. Results of the ring theory provided in A.1 are then applied to obtain a description of all controllers that meet this design requirement. It is shown how these controllers admit a general factorization of the closed feedback loop on the given ring. A parametrization of all these factorizing controllers with a parameter \mathbf{Q} is then given in transfer function representation. Passing to the numerically more practical state space representation of linear systems introduces a new variable which is the state space dimension or order of a representation. Albeit a factorization is always possible, it is shown that this in general has to be paid by high state space dimensions of the factors which leads to complex controllers of high state dimension. A special fixed order factorization of the feedback loop is then presented in state space representation, which is the special case of the general factorization in that the state space dimensions of plant and controller are fixed. Fixing the dimension has the advantage that reduced order controllers, with state dimension equal or less than that of the plant, can be parametrized directly, which is of interest if controllers of low complexity are to be found. The fixed order parameters $\mathbf{Q}_L, \mathbf{Q}_R$ share the property of the parameter \mathbf{Q} in that the parameter needs to be a ring member for the loop \mathbf{H} to be a ring member, a fact that allows the parameter itself to be parametrized in chapter 3.

2.1 Ring of Transfer Functions

This section serves as a bridge between algebraic ring structures and the transfer functions that arise in the study of linear time invariant systems.

2.1.1 Transfer Functions

A linear time invariant system may be expressed either in time domain, where it is described by a system of linear differential equations over time $t \in \mathbb{R}$ or $k \in \mathbb{Z}$, or in frequency domain where it is described by a system of linear



Figure 2.1: A linear time invariant system P with input \mathbf{u} and output \mathbf{y} .

algebraic equations over frequency $s \in \mathbb{C}$ or $z \in \mathbb{C}$, for continuous or discrete systems respectively. The link between the two descriptions is the Laplace transformation.

Figure 2.1 shows a block diagram of a linear time invariant system. The input \mathbf{u} and output \mathbf{y} are in general vectors containing multiple scalar inputs and outputs respectively. Due to the linearity, the system description in frequency domain can be brought in the following form:

$$\mathbf{y}(x) = \mathbf{P}(x)\mathbf{u}(x) \quad (2.1)$$

with $x \in \{s, z\}$ for continuous and discrete systems respectively. For reason of clarity the arguments s, z are omitted in the following. The matrix function \mathbf{P} between input and output is called the *transfer matrix*. Its entries, the scalar *transfer functions*, are rational functions of frequency $x \in \{s, z\}$

$$y_j = [\mathbf{P}]_{ij}u_i = k \frac{\prod_{k=0}^m (x - a_k)}{\prod_{l=0}^n (x - b_l)} u_i$$

with zeros $a_k \in \mathbb{C}$, poles $b_l \in \mathbb{C}$ and static gain $k \in \mathbb{C}$. A physically realizable transfer function is characterized by $m \leq n$ and is called *proper*, so the numerator polynomial is never of higher degree than the denominator polynomial, see also section 2.3.1. Transfer functions are either *stable* or *unstable*. A continuous transfer function is stable if and only if all its poles b_l have negative real part $\operatorname{Re} b_l < 0, \forall b_l$. A discrete transfer function is stable if and only if all its poles b_l lie inside the unit circle $|b_l| < 1, \forall b_l$.

The transfer function may have common zeros and poles. In this case they can be mutually cancelled and the degree of the numerator and denominator polynomials can be decreased until they have no common zeros and poles. Then the transfer function is said to be *minimal* since the degree of the numerator and denominator polynomials can not be decreased further.

2.1.2 Ring Structure

Verification of (A.1) to (A.9) as well as (A.13) and (A.14) shows that proper transfer functions with all poles contained in a zone Z form a commutative integral ring R_Z . It has been shown in [6] that these transfer functions more precisely form a *principal* commutative integral ring, so that every ideal is principal. The multiplicative inverse of proper transfer functions with $m = n$ and all poles and zeros contained in Z are again proper transfer functions with all poles in Z , so they are the units in this ring.

Definition Let R_Z be the ring of transfer functions with all poles contained in a selected zone $Z \subset \mathbb{C}$ and let $F_Z = \frac{R_Z}{R_Z \setminus \{0\}}$ be the field of fractions on R_Z .

Proposition All physically realizable and thus proper transfer functions are a member of F_Z .

Proof Let an arbitrary physically realizable transfer function be given by

$$P = k \frac{\prod_{k=0}^m (x - c_k)}{\prod_{l=0}^n (x - d_l)}$$

with $x \in \{s, z\}$ for continuous and discrete systems respectively. This transfer function is physically realizable if it is proper so that $m \leq n$. Then let $c_1, \dots, c_n \in Z$ and define

$$A = k \frac{\prod_{k=0}^m (x - a_k)}{\prod_{l=0}^n (x - c_l)}, \quad B = \frac{\prod_{l=0}^n (x - b_l)}{\prod_{l=0}^n (x - c_l)}$$

then obviously $P = \frac{A}{B}$ with $A, B \in R_Z$ and therefore $P \in F_Z$.

Remark 1 This means that the ring theory of chapter A.1 can now be applied to proper transfer functions with all poles in Z . The extension to matrices of section A.1.5 allows the treatment of entire proper transfer matrices with all poles in Z . The extension to fractions of section A.1.6 allows the treatment of general transfer matrices with poles anywhere. That way, the transfer matrix of any linear time invariant system can be described by means of proper transfer functions with poles in Z .

Remark 2 If the zone Z coincides with the left half plane for continuous systems or the inner unit disk for discrete systems, then proper transfer functions with all poles in this zone Z are stable. Therefore stability, which is a main objective in linear control, arises naturally in this representation but such a ring is only a special case as the zone Z may also be defined so that a maximal settling time or a minimal amount of damping is respected.

2.2 Controller Parametrization

This section introduces the concept of loop factorization which provides a way of parametrizing the set of all linear controllers in transfer function representation that keep the loop formed with a linear plant in the selected underlying ring R_Z . Details can be found in [13].

2.2.1 The Control Loop

Figure 2.2 shows the control loop underlying this problem. The plant $P \in F_Z^{m,n}$ is controlled by the controller $K \in F_Z^{n,m}$. The input signals \mathbf{v} and \mathbf{w} might be the reference and measurement noise, respectively and the output signals \mathbf{y} and \mathbf{u} are the plant and controller outputs, respectively.

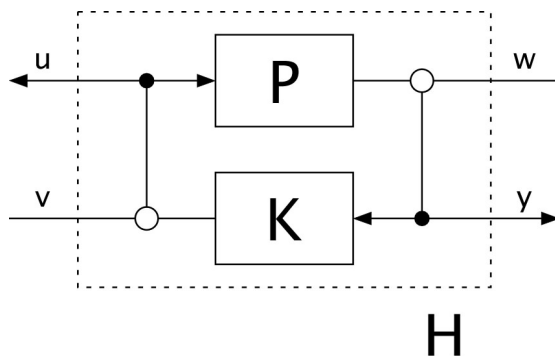


Figure 2.2: The plant P and the controller K form the control loop.

Proposition The loop shown in figure 2.2 has a transfer matrix H given by

$$H = \begin{bmatrix} -K & 1 \\ 1 & -P \end{bmatrix}^{-1} \quad (2.2)$$

or equivalently by

$$\begin{aligned} H &= \begin{bmatrix} P(1 - KP)^{-1} & 1 + P(1 - KP)^{-1}K \\ (1 - KP)^{-1} & (1 - KP)^{-1}K \end{bmatrix} = \\ &= \begin{bmatrix} P \\ 1 \end{bmatrix} (1 - KP)^{-1} \begin{bmatrix} 1 & K \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \end{aligned} \quad (2.3)$$

or equivalently by

$$\begin{aligned} H &= \begin{bmatrix} (1 - PK)^{-1}P & (1 - PK)^{-1} \\ 1 + K(1 - PK)^{-1}P & K(1 - PK)^{-1} \end{bmatrix} = \\ &= \begin{bmatrix} 1 \\ K \end{bmatrix} (1 - PK)^{-1} \begin{bmatrix} P & 1 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} \end{aligned} \quad (2.4)$$

Proof From figure 2.2 it can be seen that

$$\left. \begin{array}{l} y = w + Pu \\ u = v + Ky \end{array} \right\} \Rightarrow \begin{bmatrix} 1 & -P \\ -K & 1 \end{bmatrix} \begin{bmatrix} y \\ u \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix}$$

The transfer matrix H may then be defined by

$$\begin{bmatrix} y \\ u \end{bmatrix} = H \begin{bmatrix} v \\ w \end{bmatrix}$$

so that (2.2) follows. The equivalent expressions (2.3), (2.4) then follow by application of (A.59) and use of the relations (A.56).

Remark The loop H is symmetric with respect to P, K as can be seen by (2.2). In the following, the plant P is assumed to be given and the controller K is the variable to assure that the requirements posed on H are met.

2.2.2 Connectivity

Proposition A controller \mathbf{K} can be connected to a plant \mathbf{P} to form a feedback loop \mathbf{H} as given in (2.2) if and only if

$$\exists(\mathbf{1} - \mathbf{PK})^{-1} \Leftrightarrow \exists(\mathbf{1} - \mathbf{KP})^{-1} \Leftrightarrow \exists\mathbf{H} = \begin{bmatrix} -\mathbf{K} & \mathbf{1} \\ \mathbf{1} & -\mathbf{P} \end{bmatrix}^{-1} \quad (2.5)$$

In that case the loop \mathbf{H} is said to be *well connected*.

Proof This follows from (2.3) or (2.4).

Remark From now on, all feedback loops \mathbf{H} are supposed to be well connected.

2.2.3 Doubly Coprime Factorization

Definition For a given plant \mathbf{P} , all controllers \mathbf{K} that lead to $\mathbf{H} \in R_Z^{m+n, m+n}$, with R_Z being the selected ring of proper transfer functions with all poles contained in a zone Z , form the set of *factorizing controllers*

$$\mathbf{K} \in K(\mathbf{P}) \Leftrightarrow \mathbf{H} \in R_Z^{m+n, m+n} \quad (2.6)$$

Proposition A controller \mathbf{K} is a factorizing controller if and only if there exist $\mathbf{M}_i \in R_Z^{m, n}$, $\mathbf{X}_i \in R_Z^{n, m}$, $\mathbf{N}_L, \mathbf{Y}_R \in R_Z^{m, m}$, $\mathbf{N}_R, \mathbf{Y}_L \in R_Z^{n, n}$ with $i \in \{R, L\}$ that satisfy the doubly coprime factorization

$$\mathbf{1} = \begin{bmatrix} \mathbf{M}_L & \mathbf{N}_L \\ \mathbf{Y}_L & \mathbf{X}_L \end{bmatrix} \begin{bmatrix} -\mathbf{X}_R & \mathbf{N}_R \\ \mathbf{Y}_R & -\mathbf{M}_R \end{bmatrix} = \begin{bmatrix} -\mathbf{X}_L & \mathbf{Y}_L \\ \mathbf{N}_L & -\mathbf{M}_L \end{bmatrix} \begin{bmatrix} \mathbf{M}_R & \mathbf{Y}_R \\ \mathbf{N}_R & \mathbf{X}_R \end{bmatrix} \quad (2.7)$$

with $\begin{cases} \mathbf{P} = \mathbf{M}_R \mathbf{N}_R^{-1} = \mathbf{N}_L^{-1} \mathbf{M}_L \\ \mathbf{K} = \mathbf{X}_R \mathbf{Y}_R^{-1} = \mathbf{Y}_L^{-1} \mathbf{X}_L \end{cases}$

Furthermore all valid matrices $\mathbf{X}_i, \mathbf{Y}_i$ and so \mathbf{K} are a function of a free parameter matrix $\mathbf{Q} \in R_Z^{n, m}$ and can be given by

$$\begin{aligned} \begin{bmatrix} \mathbf{X}_R \\ \mathbf{Y}_R \end{bmatrix} &= \begin{bmatrix} \mathbf{X}_{R0} & \mathbf{N}_R \\ \mathbf{Y}_{R0} & \mathbf{M}_R \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \mathbf{Q} \end{bmatrix} && \text{with } \mathbf{Y}_R \text{ invertible} \\ \begin{bmatrix} \mathbf{X}_L & \mathbf{Y}_L \end{bmatrix} &= \begin{bmatrix} \mathbf{1} & \mathbf{Q} \end{bmatrix} \begin{bmatrix} \mathbf{X}_{L0} & \mathbf{Y}_{L0} \\ \mathbf{N}_L & \mathbf{M}_L \end{bmatrix} && \text{with } \mathbf{Y}_L \text{ invertible} \end{aligned} \quad (2.8)$$

with $\mathbf{X}_{i0}, \mathbf{Y}_{i0}$ being a set of valid initial matrices that satisfy (2.7) and allow the construction of all other solutions $\mathbf{X}_i, \mathbf{Y}_i$ using (2.8).

Proof Let plant $\mathbf{P} \in F_Z^{m, n}$ and controller $\mathbf{K} \in F_Z^{n, m}$ be factorized in coprime transfer matrices that are members of the underlying ring R_Z

$$\begin{aligned} \mathbf{P} &= \mathbf{M}_R \mathbf{N}_R^{-1} = \mathbf{N}_L^{-1} \mathbf{M}_L \\ \mathbf{K} &= \mathbf{X}'_R \mathbf{Y}'_R{}^{-1} = \mathbf{Y}'_L{}^{-1} \mathbf{X}'_L \end{aligned}$$

using (A.50) with $M_i \in R_Z^{m,n}$; $X'_i \in R_Z^{n,m}$; $N_L, Y'_R \in R_Z^{m,m}$; $N_R, Y'_L \in R_Z^{n,n}$. With this factorization, H can be expressed as

$$\begin{aligned} H &= \begin{bmatrix} M_R N_R^{-1} \\ \mathbf{1} \end{bmatrix} (\mathbf{1} - Y'_L{}^{-1} X'_L M_R N_R^{-1})^{-1} \begin{bmatrix} \mathbf{1} & Y'_L{}^{-1} X'_L \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} = \\ &= \begin{bmatrix} M_R \\ N_R \end{bmatrix} (Y'_L N_R - X'_L M_R)^{-1} \begin{bmatrix} Y'_L & X'_L \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (2.9) \end{aligned}$$

according to (2.3), and as

$$\begin{aligned} H &= \begin{bmatrix} \mathbf{1} \\ X'_R Y'_R{}^{-1} \end{bmatrix} (\mathbf{1} - N_L^{-1} M_L X'_R Y'_R{}^{-1})^{-1} \begin{bmatrix} N_L^{-1} M_L & \mathbf{1} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} \end{bmatrix} = \\ &= \begin{bmatrix} Y'_R \\ X'_R \end{bmatrix} (N_L Y'_R - M_L X'_R)^{-1} \begin{bmatrix} M_L & N_L \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} \end{bmatrix} \quad (2.10) \end{aligned}$$

according to (2.4). The transfer matrix H has its poles in the predefined zone Z and is realizable if and only if $H \in R_Z^{m+n, m+n}$ that is if the elements of H are in R_Z and thus are all proper transfer functions with poles in Z . It can be seen from the two expressions for H above that this is true if and only if the two inverses $(Y'_L N_R - X'_L M_R)^{-1}$ and $(N_L Y'_R - M_L X'_R)^{-1}$ are matrices over R_Z which is equivalent of stating that $Y'_L N_R - X'_L M_R$ and $N_L Y'_R - M_L X'_R$ must be units over R_Z so that

$$\begin{aligned} Y'_L N_R - X'_L M_R &= U_L \\ N_L Y'_R - M_L X'_R &= U_R \end{aligned} \quad (2.11)$$

with $U_R \in R_Z^{m,m}$, $U_L \in R_Z^{n,n}$ being unit matrices. The set of matrices

$$\begin{aligned} \begin{bmatrix} X_R \\ Y_R \end{bmatrix} &= \begin{bmatrix} X'_R \\ Y'_R \end{bmatrix} U_R^{-1} && \text{right coprime} \\ \begin{bmatrix} X_L & Y_L \end{bmatrix} &= U_L^{-1} \begin{bmatrix} X'_L & Y'_L \end{bmatrix} && \text{left coprime} \end{aligned} \quad (2.12)$$

also form a coprime factorization of the controller since $K = X_R Y_R^{-1} = Y_L^{-1} X_L$. The constraints for loop ring membership (2.11) can now be written as

$$\begin{aligned} Y_L N_R - X_L M_R &= \mathbf{1} \\ N_L Y_R - M_L X_R &= \mathbf{1} \end{aligned}$$

This requirement together with the coprime factorizations $P = M_R N_R^{-1} = N_L^{-1} M_L$ and $K = X_R Y_R^{-1} = Y_L^{-1} X_L$ which lead to $M_L N_R - N_L M_R = \mathbf{0}$ and $X_L Y_R - Y_L X_R = \mathbf{0}$ can be represented by the matrix equation of (2.7). According to (A.51) the set of matrices X_i, Y_i forms a coprime completion of the plant's coprime factorization. With the additional requirement of Y_i to be invertible, the set of matrices X_i, X_i also form a coprime factorization of the controller. Hence, this represents a doubly coprime factorization of P, K . According to (A.52), all coprime factors X_i, Y_i of the controller K that solve this doubly coprime factorization for P can be parametrized by a parameter $Q \in R_Z^{n,m}$ and can be brought into the block matrix form of (2.8).

2.2.4 Loop Factorization

Proposition All controllers \mathbf{K} from the set $K(\mathbf{P})$ of solutions of the doubly coprime factorization (2.7) factorize the loop \mathbf{H} in the form of

$$\begin{aligned} \mathbf{H} &= \begin{bmatrix} \mathbf{M}_R \\ \mathbf{N}_R \end{bmatrix} \begin{bmatrix} \mathbf{Y}_L & \mathbf{X}_L \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \\ \mathbf{H} &= \begin{bmatrix} \mathbf{Y}_R \\ \mathbf{X}_R \end{bmatrix} \begin{bmatrix} \mathbf{M}_L & \mathbf{N}_L \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} \end{bmatrix} \end{aligned} \quad (2.13)$$

with $\mathbf{M}_i \in R_Z^{m,n}$, $\mathbf{X}_i \in R_Z^{n,m}$, $\mathbf{N}_L, \mathbf{Y}_R \in R_Z^{m,m}$, $\mathbf{N}_R, \mathbf{Y}_L \in R_Z^{n,n}$ and $i \in \{R, L\}$.

Proof This immediately follows from the above constraints for ring membership (2.11) of the loop given by (2.9), (2.10) using the factorization (2.12).

Remark An alternative form of \mathbf{H} in the form of a serial connection can be given by

$$\begin{aligned} \mathbf{H} &= \begin{bmatrix} \mathbf{Y}_R & \mathbf{0} \\ \mathbf{0} & \mathbf{N}_R \end{bmatrix} \begin{bmatrix} -\mathbf{X}_R & \mathbf{N}_R \\ \mathbf{Y}_R & -\mathbf{M}_R \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{Y}_R & \mathbf{0} \\ \mathbf{0} & \mathbf{N}_R \end{bmatrix} \begin{bmatrix} \mathbf{M}_L & \mathbf{N}_L \\ \mathbf{Y}_L & \mathbf{X}_L \end{bmatrix} \\ \mathbf{H} &= \begin{bmatrix} -\mathbf{X}_L & \mathbf{Y}_L \\ \mathbf{N}_L & -\mathbf{M}_L \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{Y}_L & \mathbf{0} \\ \mathbf{0} & \mathbf{N}_L \end{bmatrix} = \begin{bmatrix} \mathbf{M}_R & \mathbf{Y}_R \\ \mathbf{N}_R & \mathbf{X}_R \end{bmatrix} \begin{bmatrix} \mathbf{Y}_L & \mathbf{0} \\ \mathbf{0} & \mathbf{N}_L \end{bmatrix} \end{aligned} \quad (2.14)$$

which can be directly derived from (2.2). First using the coprime factorizations of \mathbf{P} and \mathbf{K} gives

$$\begin{aligned} \mathbf{H} &= \begin{bmatrix} -\mathbf{K} & \mathbf{1} \\ \mathbf{1} & -\mathbf{P} \end{bmatrix}^{-1} = \\ &= \begin{bmatrix} \mathbf{Y}_R & \mathbf{0} \\ \mathbf{0} & \mathbf{N}_R \end{bmatrix} \begin{bmatrix} -\mathbf{X}_R & \mathbf{N}_R \\ \mathbf{Y}_R & -\mathbf{M}_R \end{bmatrix}^{-1} = \begin{bmatrix} -\mathbf{X}_L & \mathbf{Y}_L \\ \mathbf{N}_L & -\mathbf{M}_L \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{Y}_L & \mathbf{0} \\ \mathbf{0} & \mathbf{N}_L \end{bmatrix} \end{aligned}$$

Then using the fact that the two block matrices in (2.7) are mutual inverses as constraint for ring membership this immediately leads to (2.14).

2.2.5 Loop Decomposition

Proposition All controllers \mathbf{K} from the set $K(\mathbf{P})$ provide a decomposition of the loop \mathbf{H} that is linear in the parameter \mathbf{Q}

$$\mathbf{H}(\mathbf{Q}) = \mathbf{H}(\mathbf{0}) + \begin{bmatrix} \mathbf{M}_R \\ \mathbf{N}_R \end{bmatrix} \mathbf{Q} \begin{bmatrix} \mathbf{M}_L & \mathbf{N}_L \end{bmatrix} \quad (2.15)$$

with $\mathbf{H}(\mathbf{0})$ being the loop of plant \mathbf{P} fed back by the initial controller $\mathbf{K}(\mathbf{0})$. The decomposition is illustrated by figure 2.3.

Proof The controller parametrization works in the way that if one valid controller $\mathbf{K}(\mathbf{0})$ is known, then all other valid controllers $\mathbf{K}(\mathbf{Q})$ that keep the loop formed with a given plant \mathbf{P} in the selected ring R_Z can be found by varying the free ring transfer matrix \mathbf{Q} according to (2.8). This variation will change

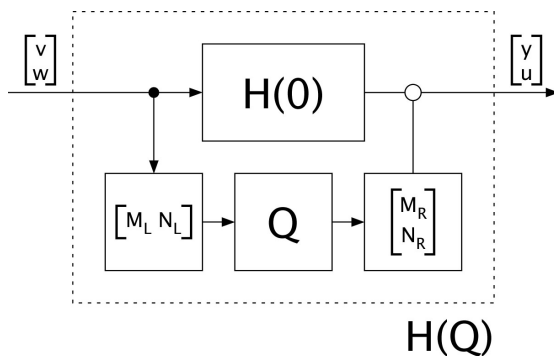


Figure 2.3: Decomposition of the feedback loop by extraction of the free parameter Q .

the closed loop transfer matrix \mathbf{H} which therefore is also a function of Q . This function can be derived by substituting (2.8) in the loop factorization (2.13). The first expression for \mathbf{H}

$$\mathbf{H} = \begin{bmatrix} \mathbf{M}_R \\ \mathbf{N}_R \end{bmatrix} [\mathbf{Y}_L \quad \mathbf{X}_L] + \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

together with

$$[\mathbf{Y}_L \quad \mathbf{X}_L] = [\mathbf{1} \quad \mathbf{Q}] \begin{bmatrix} \mathbf{Y}_{L0} & \mathbf{X}_{L0} \\ \mathbf{M}_L & \mathbf{N}_L \end{bmatrix}$$

from (2.8) leads to

$$\mathbf{H}(Q) = \underbrace{\begin{bmatrix} \mathbf{M}_R \\ \mathbf{N}_R \end{bmatrix} [\mathbf{Y}_{L0} \quad \mathbf{X}_{L0}] + \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}}_{\mathbf{H}(0)} + \begin{bmatrix} \mathbf{M}_R \\ \mathbf{N}_R \end{bmatrix} \mathbf{Q} [\mathbf{M}_L \quad \mathbf{N}_L]$$

which shows (2.15). The second expression for \mathbf{H}

$$\mathbf{H} = \begin{bmatrix} \mathbf{Y}_R \\ \mathbf{X}_R \end{bmatrix} [\mathbf{M}_L \quad \mathbf{N}_L] + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} \end{bmatrix}$$

together with

$$\begin{bmatrix} \mathbf{Y}_R \\ \mathbf{X}_R \end{bmatrix} = \begin{bmatrix} \mathbf{Y}_{R0} & \mathbf{M}_R \\ \mathbf{X}_{R0} & \mathbf{N}_R \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \mathbf{Q} \end{bmatrix}$$

from (2.8) leads to

$$\mathbf{H}(Q) = \underbrace{\begin{bmatrix} \mathbf{Y}_{R0} \\ \mathbf{X}_{R0} \end{bmatrix} [\mathbf{M}_L \quad \mathbf{N}_L] + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} \end{bmatrix}}_{\mathbf{H}(0)} + \begin{bmatrix} \mathbf{M}_R \\ \mathbf{N}_R \end{bmatrix} \mathbf{Q} [\mathbf{M}_L \quad \mathbf{N}_L]$$

which also shows (2.15).

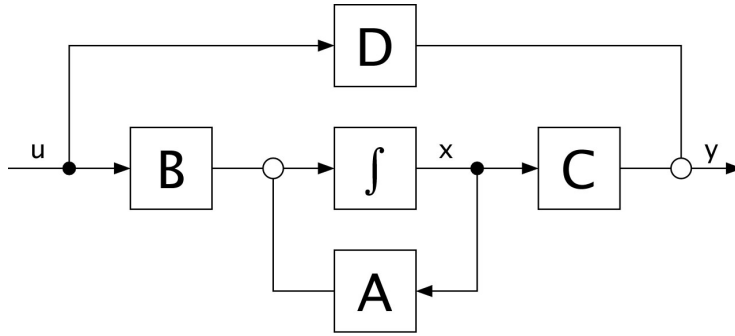


Figure 2.4: Realization of a linear time invariant system with input \mathbf{u} , output \mathbf{y} and system state \mathbf{x} .

2.3 Passage to State Space

This section shows that a linear time invariant systems in transfer function representation may also be described in state space representation. This has an advantage in computer aided design of controllers since constant matrices instead of polynomials are treated.

2.3.1 Realization of Transfer Functions

Proper transfer functions are physically realizable with the three common elements adder, gain and integrator, whereas improper transfer functions would require a differentiator that is not realizable. Any realizable system such as the plant in figure 2.1 may be written in state space representation which is therefore called a *realization*. In time domain it is given by

$$\begin{aligned}\dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) \\ \mathbf{y}(t) &= \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)\end{aligned}\tag{2.16}$$

for continuous systems and

$$\begin{aligned}\mathbf{x}[k+1] &= \mathbf{A}\mathbf{x}[k] + \mathbf{B}\mathbf{u}[k] \\ \mathbf{y}[k] &= \mathbf{C}\mathbf{x}[k] + \mathbf{D}\mathbf{u}[k]\end{aligned}\tag{2.17}$$

for discrete systems, with system state vectors $\mathbf{x}(t)$, $\mathbf{x}[k]$ and constant matrices $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$. For reason of clarity the argument t is omitted in the following. This representation is readily realizable with adders, gains and integrators or delays as shown in figure 2.4. The link between this state space representation in time domain and the transfer matrix \mathbf{P} of the plant in frequency domain is the Laplace transformation. By transforming (2.16), (2.17) and eliminating the state \mathbf{x} it can be compared to (2.1) which then leads to

$$\mathbf{P}(x) = \mathbf{C}(x\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D} =: \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \hline \mathbf{C} & \mathbf{D} \end{array} \right]$$

with $x \in \{s, z\}$ for continuous and discrete systems respectively. Since $(x\mathbf{1} - \mathbf{A})^{-1} = \text{adj}(x\mathbf{1} - \mathbf{A})(\det(x\mathbf{1} - \mathbf{A}))^{-1}$ the eigenvalues of \mathbf{A} are the poles of $\mathbf{P}(x)$. The static gain is $\mathbf{P}(\infty) = \mathbf{D}$.

The equivalent of a transfer function having common zeros and poles is a state space representation having uncontrollable or unobservable states. Like transfer functions with no common zeros and poles, a completely controllable and observable system is called *minimal*. However, whole transfer matrices describing systems with multiple inputs and outputs may have no common zeros and poles despite being non minimal. This is one of the reasons for the description of systems in state space representation. The other reason is that constant matrices instead of polynomials are treated which is a huge advantage in computer aided controller design.

2.3.2 Realization of the Control Loop

Now let the state space representations of plant \mathbf{P} , controller \mathbf{K} and loop \mathbf{H} as in figure 2.2 be given by

$$\mathbf{P} = \left[\begin{array}{c|c} \mathbf{A}_P & \mathbf{B}_P \\ \hline \mathbf{C}_P & \mathbf{D}_P \end{array} \right], \quad \mathbf{K} = \left[\begin{array}{c|c} \mathbf{A}_K & \mathbf{B}_K \\ \hline \mathbf{C}_K & \mathbf{D}_K \end{array} \right], \quad \mathbf{H} = \left[\begin{array}{c|c} \mathbf{A}_H & \mathbf{B}_H \\ \hline \mathbf{C}_H & \mathbf{D}_H \end{array} \right] \quad (2.18)$$

then with (C.48), the state space matrices of \mathbf{H} amount to

$$\begin{aligned} \mathbf{A}_H &= \begin{bmatrix} \mathbf{A}_P + \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{B}_P \mathbf{E}_{KP} \mathbf{C}_K \\ \mathbf{B}_K \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{A}_K + \mathbf{B}_K \mathbf{D}_P \mathbf{E}_{KP} \mathbf{C}_K \end{bmatrix} \\ \mathbf{B}_H &= \begin{bmatrix} \mathbf{B}_P \mathbf{E}_{KP} & \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \\ \mathbf{B}_K \mathbf{D}_P \mathbf{E}_{KP} & \mathbf{B}_K \mathbf{E}_{PK} \end{bmatrix} \\ \mathbf{C}_H &= \begin{bmatrix} \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{D}_P \mathbf{E}_{KP} \mathbf{C}_K \\ \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{E}_{KP} \mathbf{C}_K \end{bmatrix} \\ \mathbf{D}_H &= \begin{bmatrix} \mathbf{D}_P \mathbf{E}_{KP} & \mathbf{E}_{PK} \\ \mathbf{E}_{KP} & \mathbf{D}_K \mathbf{E}_{PK} \end{bmatrix} \end{aligned} \quad (2.19)$$

The inverses $\mathbf{E}_{ij} = (\mathbf{1} - \mathbf{D}_i \mathbf{D}_j)^{-1}$ with $i, j \in \{P, K\}$ must exist for \mathbf{P} and \mathbf{K} to be placed in a feedback loop due to (2.5) and $\mathbf{D}_P = \mathbf{P}(\infty)$, $\mathbf{D}_K = \mathbf{K}(\infty)$. A useful relation is $\mathbf{E}_{ij} \mathbf{D}_i \mathbf{D}_j = \mathbf{D}_i \mathbf{E}_{ji} \mathbf{D}_j = \mathbf{D}_i \mathbf{D}_j \mathbf{E}_{ij} = \mathbf{E}_{ij} - \mathbf{1}$ which will be applied throughout this text.

2.4 Fixed Order Controller Parametrization

This section shows why the controller parametrization in transfer function representation is not well adapted to problems formulated in state space representations and presents a new fixed order parametrization derived from the set of all linear controllers in state space representation that factorize the loop formed with a linear plant on a ring R_Z .

2.4.1 Need for a new Parametrization

Transferring the controller parametrization to state space description reveals a practical problem. The loop decomposition (2.15) shows that by starting with an initial controller $\mathbf{K}(\mathbf{0})$ that leads to an initial loop $\mathbf{H}(\mathbf{0})$, all possible loops are found by varying a parameter \mathbf{Q} . This means, in order to obtain a new dynamic in $\mathbf{H}(\mathbf{Q})$, the initial dynamic of $\mathbf{H}(\mathbf{0})$ has to be cancelled out by \mathbf{Q} .

Cancellations of poles and zeros in transfer function representation correspond to reductions of the state dimension in state space representation. If a system's state dimension can be reduced, then it is not minimal. Therefore $\mathbf{H}(\mathbf{Q})$ is in general not minimal. To give an idea of the state dimensions encountered, look at the way a new controller $\mathbf{K}(\mathbf{Q})$ is constructed from an initial $\mathbf{K}(\mathbf{0})$ as given in (2.7) and (2.8).

The new controller's factors $\mathbf{X}_i, \mathbf{Y}_i$ with $i \in \{L, R\}$ are obtained by (2.8) which is a serial connection between the plant's factors $\mathbf{M}_i, \mathbf{N}_i$ of state dimension p_P and the parameter \mathbf{Q} of state dimension p_Q and then a parallel connection of the result with the initial controller's factors $\mathbf{X}_{i0}, \mathbf{Y}_{i0}$ of state dimension p_{K0} . According to the formulas for serial (C.46) and parallel (C.44) connection, each such connection has the state of the connecting systems combined so that $p_{X_i, Y_i} = p_P + p_{K0} + p_Q$. The formulas for right (C.59) and left (C.61) factorization show that the state dimension of $\mathbf{K}(\mathbf{Q})$ is the same as its factors so $p_K = p_P + p_{K0} + p_Q$. With the formula for the feedback loop (C.48), the state dimension of $\mathbf{H}(\mathbf{Q})$ is $p_H = 2p_P + p_{K0} + p_Q$. Using this parametrization in state space description therefore lead to very large matrices that have to be treated by the multiobjective optimization stage later on, which is not desirable. Thus, in order to obtain reduced order controllers directly, the state space expansion to large non minimal systems has to be avoided if possible.

The use of the parameter \mathbf{Q} to parametrize the solutions of the doubly coprime factorization (2.7) is therefore adapted to the transfer function representation. To find the new parametrization adapted to the state space representation, the doubly coprime factorization has to be solved directly in state space form. Then it is shown in the following that among all possible solutions there exist fixed order solutions so that the loop \mathbf{H} has a factorization (2.13) with fixed order factors. Reduced order controllers can then be found by fixing the controller state dimension to a value lower or equal to the plant state dimension so that $p_K \leq p_P$. Since reduced controllers are a subset of all controllers, this comes at the cost of not parametrizing all possible controllers.

2.4.2 Fixed Order Loop Factorization

Definition For a given plant \mathbf{P} with realization $(\mathbf{A}_P, \mathbf{B}_P, \mathbf{C}_P, \mathbf{D}_P)$, all controllers $\mathbf{K} \in K(\mathbf{P})$ with a state space representation of fixed dimension p_K form the set of p_K -factorizing controllers

$$\mathbf{K} \in K_{p_K}(\mathbf{A}_P, \mathbf{B}_P, \mathbf{C}_P, \mathbf{D}_P) \Leftrightarrow \mathbf{K} \in K(\mathbf{P}) \text{ with state dimension } p_K \quad (2.20)$$

Proposition For a given plant \mathbf{P} with realization $(\mathbf{A}_P, \mathbf{B}_P, \mathbf{C}_P, \mathbf{D}_P)$, a controller $\mathbf{K} \in K(\mathbf{P})$ is a p_K -factorizing controller if and only if at least one of the following Riccati equations has a solution

$$\mathbf{0} = [-\mathbf{R}_{KP} \quad \mathbf{1}] \mathbf{A}_H \begin{bmatrix} \mathbf{1} \\ \mathbf{R}_{KP} \end{bmatrix} \quad (2.21)$$

$$\mathbf{0} = [\mathbf{1} \quad -\mathbf{R}_{PK}] \mathbf{A}_H \begin{bmatrix} \mathbf{R}_{PK} \\ \mathbf{1} \end{bmatrix} \quad (2.22)$$

If \mathbf{R}_{KP} exists, the control loop \mathbf{H} has a fixed order loop factorization in the form of

$$\begin{aligned} \mathbf{H} &= \begin{bmatrix} \mathbf{M}_R \\ \mathbf{N}_R \end{bmatrix} [\mathbf{Y}_L \ \mathbf{X}_L] + \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \text{ with} & (2.23) \\ \begin{bmatrix} \mathbf{M}_R \\ \mathbf{N}_R \end{bmatrix} &= \left[\begin{array}{c|c} \frac{\mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P}{\mathbf{C}_P + \mathbf{D}_P \mathbf{F}_P} & \frac{\mathbf{B}_P \mathbf{U}_P}{\mathbf{D}_P \mathbf{U}_P} \\ \mathbf{F}_P & \mathbf{U}_P \end{array} \right] \\ [\mathbf{Y}_L \ \mathbf{X}_L] &= \left[\begin{array}{c|c} \frac{\mathbf{A}_K + \mathbf{L}_K \mathbf{C}_K}{\mathbf{V}_K \mathbf{C}_K} & \frac{\mathbf{L}_K \ \mathbf{B}_K + \mathbf{L}_K \mathbf{D}_K}{\mathbf{V}_K \ \mathbf{V}_K \mathbf{D}_K} \end{array} \right] \\ \begin{bmatrix} \mathbf{0} & \mathbf{L}_K \\ \mathbf{F}_P & \mathbf{U}_P \mathbf{V}_K \end{bmatrix} &= \begin{bmatrix} [-\mathbf{R}_{KP} & \mathbf{1}] & \mathbf{0} \\ \mathbf{0} & [\mathbf{0} & \mathbf{1}] \end{bmatrix} \begin{bmatrix} \mathbf{A}_H & \mathbf{B}_H \\ \mathbf{C}_H & \mathbf{D}_H \end{bmatrix} \begin{bmatrix} \begin{bmatrix} \mathbf{1} \\ \mathbf{R}_{KP} \end{bmatrix} & \mathbf{0} \\ \mathbf{0} & \begin{bmatrix} \mathbf{1} \\ \mathbf{0} \end{bmatrix} \end{bmatrix} \end{aligned}$$

If \mathbf{R}_{PK} exists, the control loop \mathbf{H} has a fixed order loop factorization in the form of

$$\begin{aligned} \mathbf{H} &= \begin{bmatrix} \mathbf{Y}_R \\ \mathbf{X}_R \end{bmatrix} [\mathbf{M}_L \ \mathbf{N}_L] + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} \end{bmatrix} \text{ with} & (2.24) \\ [\mathbf{M}_L \ \mathbf{N}_L] &= \left[\begin{array}{c|c} \frac{\mathbf{A}_P + \mathbf{L}_P \mathbf{C}_P}{\mathbf{V}_P \mathbf{C}_P} & \frac{\mathbf{B}_P + \mathbf{L}_P \mathbf{D}_P \ \mathbf{L}_P}{\mathbf{V}_P \mathbf{D}_P \ \mathbf{V}_P} \end{array} \right] \\ \begin{bmatrix} \mathbf{Y}_R \\ \mathbf{X}_R \end{bmatrix} &= \left[\begin{array}{c|c} \frac{\mathbf{A}_K + \mathbf{B}_K \mathbf{F}_K}{\mathbf{F}_K} & \frac{\mathbf{B}_K \mathbf{U}_K}{\mathbf{U}_K} \\ \mathbf{C}_K + \mathbf{D}_K \mathbf{F}_K & \mathbf{D}_K \mathbf{U}_K \end{array} \right] \\ \begin{bmatrix} \mathbf{0} & \mathbf{L}_P \\ \mathbf{F}_K & \mathbf{U}_K \mathbf{V}_P \end{bmatrix} &= \begin{bmatrix} [\mathbf{1} \ -\mathbf{R}_{PK}] & \mathbf{0} \\ \mathbf{0} & [\mathbf{1} \ \mathbf{0}] \end{bmatrix} \begin{bmatrix} \mathbf{A}_H & \mathbf{B}_H \\ \mathbf{C}_H & \mathbf{D}_H \end{bmatrix} \begin{bmatrix} \begin{bmatrix} \mathbf{R}_{PK} \\ \mathbf{1} \end{bmatrix} & \mathbf{0} \\ \mathbf{0} & \begin{bmatrix} \mathbf{0} \\ \mathbf{1} \end{bmatrix} \end{bmatrix} \end{aligned}$$

The matrices $\mathbf{A}_H, \mathbf{B}_H, \mathbf{C}_H, \mathbf{D}_H$ are the state space matrices of the loop \mathbf{H} as given in (2.19).

Proof Condition (2.21) can be verified by using (C.59) and (C.61) to obtain the right and left factors of plant \mathbf{P} and controller \mathbf{K} of state dimension p_P and p_K respectively, see [7] for details. All other factors of that dimension are related to these factors by change of basis, the freedom which will be accounted for by using the basis transformation \mathbf{T} further below. Using these factors, the loop factorization can be given by

$$\begin{aligned} \mathbf{H} &= \begin{bmatrix} \mathbf{M}_R \\ \mathbf{N}_R \end{bmatrix} [\mathbf{Y}_L \ \mathbf{X}_L] + \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} = \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{array} \right] = \\ &= \left[\begin{array}{c|c} \frac{\mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P}{\mathbf{C}_P + \mathbf{D}_P \mathbf{F}_P} & \frac{\mathbf{B}_P \mathbf{U}_P \mathbf{V}_K \mathbf{C}_K}{\mathbf{D}_P \mathbf{U}_P \mathbf{V}_K \mathbf{C}_K} \mid \frac{\mathbf{B}_P \mathbf{U}_P \mathbf{V}_K}{\mathbf{U}_P \mathbf{V}_K} & \frac{\mathbf{B}_P \mathbf{U}_P \mathbf{V}_K \mathbf{D}_K}{\mathbf{D}_P \mathbf{U}_P \mathbf{V}_K \mathbf{D}_K + \mathbf{1}} \\ \mathbf{0} & \frac{\mathbf{A}_K + \mathbf{L}_K \mathbf{C}_K}{\mathbf{U}_P \mathbf{V}_K \mathbf{C}_K} & \mathbf{L}_K & \frac{\mathbf{B}_K + \mathbf{L}_K \mathbf{D}_K}{\mathbf{U}_P \mathbf{V}_K \mathbf{D}_K} \end{array} \right] \end{aligned}$$

using (C.46). The general non factorized feedback loop is given by (2.19) so that

$$\mathbf{H} = \left[\begin{array}{c|c} \mathbf{A}_H & \mathbf{B}_H \\ \mathbf{C}_H & \mathbf{D}_H \end{array} \right]$$

These two state space representations of \mathbf{H} have the same state dimension and must therefore be equivalent in the sense of (C.35) for the controller to be a member of $K(\mathbf{P})$. Then according to (C.36) there must exist an invertible basis transformation matrix \mathbf{T} so that $\mathbf{TA} = \mathbf{A}_H\mathbf{T}$, $\mathbf{TB} = \mathbf{B}_H$, $\mathbf{C} = \mathbf{C}_H\mathbf{T}$, $\mathbf{D} = \mathbf{D}_H$ which can now be calculated using the partitioning

$$\mathbf{T} = \begin{bmatrix} \mathbf{W} & \mathbf{X} \\ \mathbf{Y} & \mathbf{Z} \end{bmatrix}$$

The condition $\mathbf{D} = \mathbf{D}_H$ provides

$$\mathbf{U}_P\mathbf{V}_K = \mathbf{E}_{KP} \quad (2.25)$$

The condition $\mathbf{C} = \mathbf{C}_H\mathbf{T}$ provides

$$\begin{aligned} \mathbf{C}_P\mathbf{W} &= \mathbf{C}_P \\ \mathbf{C}_P\mathbf{X} &= \mathbf{0} \\ \mathbf{C}_K\mathbf{Z} &= \mathbf{C}_K \\ \mathbf{F}_P &= \mathbf{E}_{KP}(\mathbf{D}_K\mathbf{C}_P + \mathbf{C}_K\mathbf{Y}) \end{aligned} \quad (2.26)$$

The condition $\mathbf{TB} = \mathbf{B}_H$ provides

$$\begin{aligned} \mathbf{XB}_K &= \mathbf{0} \\ \mathbf{ZB}_K &= \mathbf{B}_K \\ \mathbf{XL}_K &= (\mathbf{1} - \mathbf{W})\mathbf{B}_P\mathbf{E}_{KP} \\ \mathbf{ZL}_K &= (\mathbf{B}_K\mathbf{D}_P - \mathbf{YB}_P)\mathbf{E}_{KP} \end{aligned} \quad (2.27)$$

The condition $\mathbf{TA} = \mathbf{A}_H\mathbf{T}$ provides

$$\begin{aligned} \mathbf{XA}_K &= \mathbf{A}_P\mathbf{X} \\ \mathbf{ZA}_K &= \mathbf{A}_K\mathbf{Z} \\ \mathbf{WA}_P &= \mathbf{A}_P\mathbf{W} + (\mathbf{1} - \mathbf{W})\mathbf{B}_P\mathbf{F}_P \end{aligned} \quad (2.28)$$

and

$$\begin{aligned} \mathbf{Y}(\mathbf{A}_P + \mathbf{B}_P\mathbf{D}_K\mathbf{E}_{PK}\mathbf{C}_P) + \mathbf{YB}_P\mathbf{E}_{KP}\mathbf{C}_K\mathbf{Y} &= \\ &= \mathbf{B}_K\mathbf{E}_{PK}\mathbf{C}_P + (\mathbf{A}_K + \mathbf{B}_K\mathbf{D}_P\mathbf{E}_{KP}\mathbf{C}_K)\mathbf{Y} \end{aligned} \quad (2.29)$$

using (2.26) and (2.27). If there exists a solution $\mathbf{R}_{KP} = \mathbf{Y}$ to this Riccati equation, then a solution to (2.25), (2.26), (2.27), (2.28) also exists in the form of

$$\begin{aligned} \mathbf{W} &= \mathbf{1} \\ \mathbf{X} &= \mathbf{0} \\ \mathbf{Y} &= \mathbf{R}_{KP} \\ \mathbf{Z} &= \mathbf{1} \\ \mathbf{U}_P\mathbf{V}_K &= \mathbf{E}_{KP} \\ \mathbf{F}_P &= \mathbf{E}_{KP}(\mathbf{D}_K\mathbf{C}_P + \mathbf{C}_K\mathbf{R}_{KP}) \\ \mathbf{L}_K &= (\mathbf{B}_K\mathbf{D}_P - \mathbf{R}_{KP}\mathbf{B}_P)\mathbf{E}_{KP} \end{aligned}$$

Hence \mathbf{T} is invertible so that the state space representations $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ and $(\mathbf{A}_H, \mathbf{B}_H, \mathbf{C}_H, \mathbf{D}_H)$ are equivalent and thus \mathbf{H} can be factorized in the form of (2.23) by a controller of state dimension p_K . This verifies that (2.21) is the condition for \mathbf{K} to be a member of $K_{p_K}(\mathbf{A}_P, \mathbf{B}_P, \mathbf{C}_P, \mathbf{D}_P) \subset K(\mathbf{P})$. The proof of condition (2.22) is similar.

Remark 1 The Riccati equations (2.21) and (2.22) can also be written in the Radon form

$$\begin{aligned} \mathbf{A}_H \begin{bmatrix} \mathbf{1} \\ \mathbf{R}_{KP} \end{bmatrix} &= \begin{bmatrix} \mathbf{1} \\ \mathbf{R}_{KP} \end{bmatrix} (\mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P) \\ [\mathbf{1} \quad -\mathbf{R}_{PK}] \mathbf{A}_H &= (\mathbf{A}_P + \mathbf{L}_P \mathbf{C}_P) [\mathbf{1} \quad -\mathbf{R}_{PK}] \end{aligned}$$

which follows by calculating

$$\begin{aligned} \mathbf{A}_H \begin{bmatrix} \mathbf{1} \\ \mathbf{R}_{KP} \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_P + \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P + \mathbf{B}_P \mathbf{E}_{KP} \mathbf{C}_K \mathbf{R}_{KP} \\ \mathbf{B}_K \mathbf{E}_{PK} \mathbf{C}_P + (\mathbf{A}_K + \mathbf{B}_K \mathbf{D}_P \mathbf{E}_{KP} \mathbf{C}_K) \mathbf{R}_{KP} \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{A}_P + \mathbf{B}_P \mathbf{E}_{KP} (\mathbf{D}_K \mathbf{C}_P + \mathbf{C}_K \mathbf{R}_{KP}) \\ \mathbf{R}_{KP} (\mathbf{A}_P + \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P) + \mathbf{R}_{KP} \mathbf{B}_P \mathbf{E}_{KP} \mathbf{C}_K \mathbf{R}_{KP} \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P \\ \mathbf{R}_{KP} (\mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P) \end{bmatrix} = \begin{bmatrix} \mathbf{1} \\ \mathbf{R}_{KP} \end{bmatrix} (\mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P) \end{aligned}$$

using $\mathbf{F}_P = \mathbf{E}_{KP} (\mathbf{D}_K \mathbf{C}_P + \mathbf{C}_K \mathbf{R}_{KP})$ from (2.23) and the Riccati equation (2.21). The proof for the second line is similar, this time using $\mathbf{L}_P = (\mathbf{B}_P \mathbf{D}_K - \mathbf{R}_{PK} \mathbf{B}_K) \mathbf{E}_{PK}$ from (2.24) and the Riccati equation (2.22).

Remark 2 The existence of solutions of (2.21) and (2.22) does not depend on the selected state space basis of plant or controller which can be seen by taking all equivalent state space representations of \mathbf{P}, \mathbf{K} in the form of

$$\mathbf{P} = \left[\begin{array}{c|c} \mathbf{T}_P \mathbf{A}_P \mathbf{T}_P^{-1} & \mathbf{T}_P \mathbf{B}_P \\ \hline \mathbf{C}_P \mathbf{T}_P^{-1} & \mathbf{D}_P \end{array} \right], \quad \mathbf{K} = \left[\begin{array}{c|c} \mathbf{T}_K \mathbf{A}_K \mathbf{T}_K^{-1} & \mathbf{T}_K \mathbf{B}_K \\ \hline \mathbf{C}_K \mathbf{T}_K^{-1} & \mathbf{D}_K \end{array} \right]$$

with $\mathbf{T}_P, \mathbf{T}_K$ invertible. This leads to the loop dynamic matrix

$$\mathbf{A}_H(\mathbf{T}_P, \mathbf{T}_K) = \begin{bmatrix} \mathbf{T}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_K \end{bmatrix} \mathbf{A}_H \begin{bmatrix} \mathbf{T}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_K \end{bmatrix}^{-1}$$

with \mathbf{A}_H given in (2.19). The Riccati equations (2.21), (2.22) can then be written as

$$\begin{aligned} \mathbf{0} &= [-\mathbf{R}_{KP} \quad \mathbf{1}] \begin{bmatrix} \mathbf{T}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_K \end{bmatrix}^{-1} \mathbf{A}_H(\mathbf{T}_P, \mathbf{T}_K) \begin{bmatrix} \mathbf{T}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_K \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \mathbf{R}_{KP} \end{bmatrix} \\ \mathbf{0} &= [\mathbf{1} \quad -\mathbf{R}_{PK}] \begin{bmatrix} \mathbf{T}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_K \end{bmatrix}^{-1} \mathbf{A}_H(\mathbf{T}_P, \mathbf{T}_K) \begin{bmatrix} \mathbf{T}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_K \end{bmatrix} \begin{bmatrix} \mathbf{R}_{PK} \\ \mathbf{1} \end{bmatrix} \end{aligned}$$

Since $\mathbf{T}_P, \mathbf{T}_K$ are invertible, this can be reshaped to

$$\begin{aligned} \mathbf{0} &= [-\mathbf{T}_K \mathbf{R}_{KP} \mathbf{T}_P^{-1} \quad \mathbf{1}] \mathbf{A}_H(\mathbf{T}_P, \mathbf{T}_K) \begin{bmatrix} \mathbf{1} \\ \mathbf{T}_K \mathbf{R}_{KP} \mathbf{T}_P^{-1} \end{bmatrix} \\ \mathbf{0} &= [\mathbf{1} \quad -\mathbf{T}_P \mathbf{R}_{PK} \mathbf{T}_K^{-1}] \mathbf{A}_H(\mathbf{T}_P, \mathbf{T}_K) \begin{bmatrix} \mathbf{T}_P \mathbf{R}_{PK} \mathbf{T}_K^{-1} \\ \mathbf{1} \end{bmatrix} \end{aligned}$$

Hence the solutions of the transformed Riccati equations are given by

$$\begin{aligned}\mathbf{R}_{KP}(\mathbf{T}_P, \mathbf{T}_K) &= \mathbf{T}_K \mathbf{R}_{KP} \mathbf{T}_P^{-1} \\ \mathbf{R}_{PK}(\mathbf{T}_P, \mathbf{T}_K) &= \mathbf{T}_P \mathbf{R}_{PK} \mathbf{T}_K^{-1}\end{aligned}$$

so that these solutions exist if and only if $\mathbf{R}_{KP}, \mathbf{R}_{PK}$ exist.

2.4.3 Existence of Solutions

Proposition A sufficient condition for the solvability of both Riccati equations (2.21) and (2.22) is given by

$$\mathbf{A}_H \text{ is diagonalizable} \Rightarrow \mathbf{R}_{PK} \text{ and } \mathbf{R}_{KP} \text{ exist} \quad (2.30)$$

with \mathbf{A}_H being the dynamic matrix of the loop \mathbf{H} as given in (2.19). In this case both Riccati equations (2.21) and (2.22) have solutions and \mathbf{K} admits both loop factorizations (2.23) and (2.24). Some examples are given after the proof to illustrate the situation if \mathbf{A}_H is not diagonalizable.

Proof Since \mathbf{A}_H is diagonalizable we can write

$$\mathbf{A}_H \mathbf{V} = \mathbf{V} \mathbf{\Lambda} \quad (2.31)$$

with invertible \mathbf{V} containing the eigenvectors of \mathbf{A}_H and diagonal $\mathbf{\Lambda}$ containing the eigenvalues of \mathbf{A}_H . Since \mathbf{V} is invertible, it can be arranged in a way

$$\mathbf{V} = \begin{bmatrix} \mathbf{V}_{PP} & \mathbf{V}_{PK} \\ \mathbf{V}_{KP} & \mathbf{V}_{KK} \end{bmatrix}$$

with \mathbf{V}_{PP} invertible. Since $\mathbf{\Lambda}$ is diagonal we can now write (2.31) in the form

$$\mathbf{A}_H \begin{bmatrix} \mathbf{V}_{PP} & \mathbf{V}_{PK} \\ \mathbf{V}_{KP} & \mathbf{V}_{KK} \end{bmatrix} = \begin{bmatrix} \mathbf{V}_{PP} & \mathbf{V}_{PK} \\ \mathbf{V}_{KP} & \mathbf{V}_{KK} \end{bmatrix} \begin{bmatrix} \mathbf{\Lambda}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_K \end{bmatrix}$$

Right projection with $[\mathbf{1} \ \mathbf{0}]^T$ leads to

$$\mathbf{A}_H \begin{bmatrix} \mathbf{V}_{PP} \\ \mathbf{V}_{KP} \end{bmatrix} = \begin{bmatrix} \mathbf{V}_{PP} & \mathbf{V}_{PK} \\ \mathbf{V}_{KP} & \mathbf{V}_{KK} \end{bmatrix} \begin{bmatrix} \mathbf{\Lambda}_P \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{V}_{PP} \\ \mathbf{V}_{KP} \end{bmatrix} \mathbf{\Lambda}_P$$

Right multiplication with \mathbf{V}_{PP}^{-1} gives

$$\mathbf{A}_H \begin{bmatrix} \mathbf{1} \\ \mathbf{V}_{KP} \mathbf{V}_{PP}^{-1} \end{bmatrix} = \begin{bmatrix} \mathbf{1} \\ \mathbf{V}_{KP} \mathbf{V}_{PP}^{-1} \end{bmatrix} \mathbf{V}_{PP} \mathbf{\Lambda}_P \mathbf{V}_{PP}^{-1}$$

Left projection with $[-\mathbf{V}_{KP} \mathbf{V}_{PP}^{-1} \ \mathbf{1}]$ leads to the Riccati equation

$$[-\mathbf{V}_{KP} \mathbf{V}_{PP}^{-1} \ \mathbf{1}] \mathbf{A}_H \begin{bmatrix} \mathbf{1} \\ \mathbf{V}_{KP} \mathbf{V}_{PP}^{-1} \end{bmatrix} = \mathbf{0}$$

which is identical to (2.21) so that a solution $\mathbf{R}_{KP} = \mathbf{V}_{KP} \mathbf{V}_{PP}^{-1}$ has been found. The proof for \mathbf{R}_{PK} is similar.

Example 1 Let a plant \mathbf{P} and a controller \mathbf{K} be given by

$$\mathbf{P} = \left[\begin{array}{cc|cc} \lambda+1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{array} \right], \quad \mathbf{K} = \left[\begin{array}{cc|c} \lambda-1 & 1 & 1 \\ 1 & -1 & 0 \end{array} \right]$$

The matrices

$$\mathbf{E}_{PK} = 1, \quad \mathbf{E}_{KP} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}$$

are invertible so the loop \mathbf{H} is well connected. Its dynamic matrix is then given by

$$\mathbf{A}_H = \begin{bmatrix} \lambda & 0 \\ 1 & \lambda \end{bmatrix}$$

Of the Riccati equations (2.21) and (2.22) that simplify to

$$\begin{aligned} 0 &= 1 \\ 0 &= \mathbf{R}_{PK}^2 \end{aligned}$$

only the second one has a solution

$$\mathbf{R}_{PK} = 0$$

This is therefore an example of a controller that is a member of $K_1(\mathbf{P}) \subset K(\mathbf{P})$ albeit one that admits only one loop factorization.

Example 2 Let a plant \mathbf{P} and a controller \mathbf{K} be given by

$$\mathbf{P} = \left[\begin{array}{cc|cc} -1 & 0 & 1 & 0 \\ 0 & -2 & 0 & 1 \\ 1 & -6 & 0 & 3 \\ 0 & 1 & 0 & 0 \end{array} \right], \quad \mathbf{K} = \left[\begin{array}{cc|c} -3 & 1 & 0 \\ -3 & 1 & 0 \\ 1 & 0 & 2 \end{array} \right]$$

The matrices

$$\mathbf{E}_{PK} = \begin{bmatrix} 1 & 6 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{E}_{KP} = \begin{bmatrix} 1 & 3 \\ 0 & 1 \end{bmatrix}$$

are invertible so the loop \mathbf{H} is well connected. Its dynamic matrix is then given by

$$\mathbf{A}_H = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

However the Riccati equations (2.21) and (2.22) that simplify to

$$\begin{aligned} \mathbf{0} &= \mathbf{R}_{KP} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mathbf{R}_{KP} - [1 \quad 0] \\ \mathbf{0} &= \mathbf{R}_{PK} [1 \quad 0] \mathbf{R}_{PK} - \begin{bmatrix} 0 \\ 1 \end{bmatrix} \end{aligned}$$

have both no solution. This is therefore an example of a controller that is not a member of $K_1(\mathbf{P})$ since it admits no loop factorization.

Example 3 Let a plant \mathbf{P} and a controller \mathbf{K} be given by

$$\mathbf{P} = \left[\begin{array}{cc|cc} -1 & 0 & 1 & 0 \\ 0 & -2 & 0 & 1 \\ \hline 1 & -6 & 0 & 3 \\ 0 & 1 & 0 & 0 \end{array} \right], \quad \mathbf{K} = \left[\begin{array}{cc|cc} -3 & 0 & 1 & 0 \\ 0 & \lambda & 0 & 1 \\ \hline -3 & 0 & 1 & 0 \\ 1 & 0 & 0 & 2 \end{array} \right]$$

The plant and the matrices $\mathbf{E}_{PK}, \mathbf{E}_{KP}$ are identical to the ones in example 2 and the controller has the same transfer function than the one in example 2, but an unobservable state has been added to the controller. The loop dynamic matrix is now given by

$$\mathbf{A}_H = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & \lambda \end{bmatrix}$$

Of the Riccati equations (2.21) and (2.22) that simplify to

$$\begin{aligned} \mathbf{0} &= \mathbf{R}_{KP} \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} \mathbf{R}_{KP} - \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} 0 & 0 \\ 0 & \lambda \end{bmatrix} \mathbf{R}_{KP} \\ \mathbf{0} &= \mathbf{R}_{PK}^2 - \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} + \mathbf{R}_{PK} \begin{bmatrix} 0 & 0 \\ 0 & \lambda \end{bmatrix} \end{aligned}$$

the second one has a solution if $\lambda \neq 0$

$$\mathbf{R}_{PK} = \begin{bmatrix} 0 & 0 \\ -\frac{1}{\lambda} & -\lambda \end{bmatrix}$$

This is therefore an example of a controller that is a member of $K_2(\mathbf{P}) \subset K(\mathbf{P})$. This stands in contrast to example 2 since the only difference is an unobservable state with a new eigenvalue $\lambda \neq 0$ that is added to the controller, while the transfer functions are all exactly as in example 2. Together with the result of example 2 this means that \mathbf{K} is a member of $K_2(\mathbf{P}) \setminus K_1(\mathbf{P})$ and therefore \mathbf{K} is a loop factorizing controller but the dimension of its state space representation given in example 2 was too low. Thus a controller in state space representation can become loop factorizing by increasing its state dimension.

2.4.4 Relation to Estimator Corrector Type Controllers

An *estimator corrector* type control is characterized by exhibiting the *separation principle* which states that the control task can be separated into a state estimator part and a state corrector part and both parts can be designed independently, a fact that is reflected by the eigenvalues of the loop dynamic matrix \mathbf{A}_H being the union of the eigenvalues independently determined by the estimator and corrector stage. In the following it is shown that the fixed order factorized loop resembles an estimator corrector type controlled loop in that an estimator mismatch $\boldsymbol{\epsilon}$ is formed between the controller part of the state vector \boldsymbol{x}_K and a projection of the plant part of the state vector \boldsymbol{x}_P in the form of $\boldsymbol{\epsilon} = \dot{\boldsymbol{x}}_K - \mathbf{R}_{KP}\dot{\boldsymbol{x}}_P$ in the case of factorization (2.23) or with indices exchanged in the case of factorization (2.24). Despite the similarities, it does however not exhibit the separation principle, so that the notation *general estimator corrector* seems more appropriate.

Proposition For a given plant P , a controller

$$\mathbf{K} \text{ is of the general estimator corrector type } \Leftrightarrow \mathbf{R}_{KP} \text{ or } \mathbf{R}_{PK} \text{ exist} \quad (2.32)$$

with \mathbf{R}_{KP} and \mathbf{R}_{PK} being the solutions of the Riccati equations (2.21) and (2.22).

Proof Assume the existence of $\mathbf{R}_{KP}, \mathbf{R}_{PK}$. Then the basis of the loop's state space representation \mathbf{H} as given in (2.19) can be changed using matrices

$$\mathbf{T}_{KP} = \begin{bmatrix} \mathbf{1} & \mathbf{R}_{KP} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}, \quad \mathbf{T}_{PK} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{R}_{PK} & \mathbf{1} \end{bmatrix}$$

so that

$$\mathbf{H} = \left[\begin{array}{c|c} \mathbf{T}_{KP}^{-1} \mathbf{A}_H \mathbf{T}_{KP} & \mathbf{T}_{KP}^{-1} \mathbf{B}_H \\ \hline \mathbf{C}_H \mathbf{T}_{KP} & \mathbf{D}_H \end{array} \right] = \left[\begin{array}{c|c} \mathbf{T}_{PK}^{-1} \mathbf{A}_H \mathbf{T}_{PK} & \mathbf{T}_{PK}^{-1} \mathbf{B}_H \\ \hline \mathbf{C}_H \mathbf{T}_{PK} & \mathbf{D}_H \end{array} \right]$$

Using (2.23) and (2.24) this leads to the following equivalent state space representations of the closed loop system

$$\begin{aligned} \begin{bmatrix} \dot{\mathbf{x}}_K - \mathbf{R}_{KP} \dot{\mathbf{x}}_P \\ \dot{\mathbf{x}}_P \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_K + \mathbf{L}_K \mathbf{C}_K & \mathbf{0} \\ \mathbf{B}_P \mathbf{E}_{KP} \mathbf{C}_K & \mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P \end{bmatrix} \begin{bmatrix} \mathbf{x}_K - \mathbf{R}_{KP} \mathbf{x}_P \\ \mathbf{x}_P \end{bmatrix} + \\ &+ \begin{bmatrix} \mathbf{B}_K + \mathbf{L}_K \mathbf{D}_K & \mathbf{L}_K \\ \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} & \mathbf{B}_P \mathbf{E}_{KP} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} \\ \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix} &= \begin{bmatrix} \mathbf{E}_{KP} \mathbf{C}_K & \mathbf{F}_P \\ \mathbf{D}_P \mathbf{E}_{KP} \mathbf{C}_K & \mathbf{C}_P + \mathbf{D}_P \mathbf{F}_P \end{bmatrix} \begin{bmatrix} \mathbf{x}_K - \mathbf{R}_{KP} \mathbf{x}_P \\ \mathbf{x}_P \end{bmatrix} + \\ &+ \begin{bmatrix} \mathbf{D}_K \mathbf{E}_{PK} & \mathbf{E}_{KP} \\ \mathbf{E}_{PK} & \mathbf{D}_P \mathbf{E}_{KP} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} \end{aligned} \quad (2.33)$$

for \mathbf{R}_{KP} and

$$\begin{aligned} \begin{bmatrix} \dot{\mathbf{x}}_P - \mathbf{R}_{PK} \dot{\mathbf{x}}_K \\ \dot{\mathbf{x}}_K \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_P + \mathbf{L}_P \mathbf{C}_P & \mathbf{0} \\ \mathbf{B}_K \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{A}_K + \mathbf{B}_K \mathbf{F}_K \end{bmatrix} \begin{bmatrix} \mathbf{x}_P - \mathbf{R}_{PK} \mathbf{x}_K \\ \mathbf{x}_K \end{bmatrix} + \\ &+ \begin{bmatrix} \mathbf{B}_P + \mathbf{L}_P \mathbf{D}_P & \mathbf{L}_P \\ \mathbf{B}_K \mathbf{D}_P \mathbf{E}_{KP} & \mathbf{B}_K \mathbf{E}_{PK} \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \mathbf{w} \end{bmatrix} \\ \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix} &= \begin{bmatrix} \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{F}_K \\ \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{C}_K + \mathbf{D}_K \mathbf{F}_K \end{bmatrix} \begin{bmatrix} \mathbf{x}_P - \mathbf{R}_{PK} \mathbf{x}_K \\ \mathbf{x}_K \end{bmatrix} + \\ &+ \begin{bmatrix} \mathbf{D}_P \mathbf{E}_{KP} & \mathbf{E}_{PK} \\ \mathbf{E}_{KP} & \mathbf{D}_K \mathbf{E}_{PK} \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \mathbf{w} \end{bmatrix} \end{aligned} \quad (2.34)$$

for \mathbf{R}_{PK} . It can be seen by the structure of the transformed closed loop system (2.33), that the fixed order factorization describes a form of general estimator corrector type control in which the plant's state \mathbf{x}_P is projected onto the range of \mathbf{P}_{KP} and then compared to the controller's state \mathbf{x}_K . The fact that the estimator mismatch $\mathbf{x}_K - \mathbf{R}_{KP} \mathbf{x}_P$ becomes an independent state of \mathbf{H} due to the block triangular structure of the dynamic matrix \mathbf{A}_H already shows the separation property of estimator corrector type control. However, unlike with actual estimator corrector type control, the matrices $\mathbf{F}_P, \mathbf{L}_K$ are in general not independent design variables but are rather both a function of \mathbf{R}_{KP} according to

(2.23), so the separation principle that helps in the design of common estimator corrector type controllers is in general not valid. The same is true in the dual case (2.34) for $\mathbf{F}_K, \mathbf{L}_P$ that are a function of \mathbf{R}_{PK} according to (2.24).

Proposition For a given plant \mathbf{P} , a controller

$$\mathbf{K} \text{ is of the estimator corrector type } \Leftrightarrow \mathbf{R}_{KP} \text{ and } \mathbf{R}_{PK} \text{ invertible} \quad (2.35)$$

with \mathbf{R}_{KP} and \mathbf{R}_{PK} being the solutions of the Riccati equations (2.21) and (2.22).

Proof Assume the existence of an invertible \mathbf{R}_{PK} . Then \mathbf{R}_{PK} solves

$$\begin{bmatrix} \mathbf{1} & -\mathbf{R}_{PK} \end{bmatrix} \mathbf{A}_H \begin{bmatrix} \mathbf{R}_{PK} \\ \mathbf{1} \end{bmatrix} = \mathbf{0}$$

according to (2.24). Left multiplication by $-\mathbf{R}_{PK}^{-1}$ and right multiplication by \mathbf{R}_{PK}^{-1} leads to

$$\begin{bmatrix} -\mathbf{R}_{PK}^{-1} & \mathbf{1} \end{bmatrix} \mathbf{A}_H \begin{bmatrix} \mathbf{1} \\ \mathbf{R}_{PK}^{-1} \end{bmatrix} = \mathbf{0}$$

Comparing this to (2.23) shows that there exists a solution $\mathbf{R}_{KP} = \mathbf{R}_{PK}^{-1}$. This also works in the reverse direction so that \mathbf{R}_{PK}^{-1} exists if and only if \mathbf{R}_{KP}^{-1} exists. The state space basis transformation $\mathbf{x}'_K = \mathbf{R}_{PK} \mathbf{x}_K$ with $\mathbf{R}_{KP} = \mathbf{R}_{PK}^{-1}$ and accordingly

$$\begin{bmatrix} \mathbf{A}'_K & \mathbf{B}'_K \\ \mathbf{C}'_K & \mathbf{D}_K \end{bmatrix} = \begin{bmatrix} \mathbf{R}_{PK} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix} \begin{bmatrix} \mathbf{R}_{KP} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}$$

together with $\mathbf{F}'_K = \mathbf{F}_K \mathbf{R}_{KP}$ and $\mathbf{L}'_K = \mathbf{R}_{PK} \mathbf{L}_K$ leads to

$$\begin{aligned} \mathbf{F}_P &= \mathbf{E}_{KP}(\mathbf{D}_K \mathbf{C}_P + \mathbf{C}'_K) \\ \mathbf{L}_P &= (\mathbf{B}_P \mathbf{D}_K - \mathbf{B}'_K) \mathbf{E}_{PK} \\ \mathbf{F}'_K &= \mathbf{E}_{PK}(\mathbf{D}_P \mathbf{C}'_K + \mathbf{C}_P) \\ \mathbf{L}'_K &= (\mathbf{B}'_K \mathbf{D}_P - \mathbf{B}_P) \mathbf{E}_{KP} \end{aligned}$$

using (2.23), (2.24). This further gives

$$\begin{aligned} \mathbf{A}'_K + \mathbf{B}'_K \mathbf{F}'_K &= \mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P \\ \mathbf{A}'_K + \mathbf{L}'_K \mathbf{C}'_K &= \mathbf{A}_P + \mathbf{L}_P \mathbf{C}_P \\ \mathbf{B}'_K + \mathbf{L}'_K \mathbf{D}_K &= -\mathbf{L}_P \\ &\quad -\mathbf{L}'_K = \mathbf{B}_P + \mathbf{L}_P \mathbf{D}_P \\ \mathbf{C}'_K + \mathbf{D}_K \mathbf{F}'_K &= \mathbf{F}_P \\ \mathbf{F}'_K &= \mathbf{C}_P + \mathbf{D}_P \mathbf{F}_P \end{aligned}$$

Since $\mathbf{R}_{KP}, \mathbf{R}_{PK}$ are assumed to exist, the state space representations of the closed loop system can be written in the forms (2.33), (2.34) which can now be

given by

$$\begin{aligned}
\begin{bmatrix} \dot{\mathbf{x}}'_K - \dot{\mathbf{x}}_P \\ \dot{\mathbf{x}}_P \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_P + \mathbf{L}_P \mathbf{C}_P & \mathbf{0} \\ \mathbf{B}_P (\mathbf{F}_P - \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P) & \mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P \end{bmatrix} \begin{bmatrix} \mathbf{x}'_K - \mathbf{x}_P \\ \mathbf{x}_P \end{bmatrix} + \\
&+ \begin{bmatrix} -\mathbf{L}_P & -(\mathbf{B}_P + \mathbf{L}_P \mathbf{D}_P) \\ \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} & \mathbf{B}_P \mathbf{E}_{KP} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix} \\
\begin{bmatrix} \mathbf{u} \\ \mathbf{y} \end{bmatrix} &= \begin{bmatrix} \mathbf{F}_P - \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{F}_P \\ \mathbf{D}_P (\mathbf{F}_P - \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P) & \mathbf{C}_P + \mathbf{D}_P \mathbf{F}_P \end{bmatrix} \begin{bmatrix} \mathbf{x}'_K - \mathbf{x}_P \\ \mathbf{x}_P \end{bmatrix} + \\
&+ \begin{bmatrix} \mathbf{D}_K \mathbf{E}_{PK} & \mathbf{E}_{KP} \\ \mathbf{E}_{PK} & \mathbf{D}_P \mathbf{E}_{KP} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix}
\end{aligned} \tag{2.36}$$

or alternatively by

$$\begin{aligned}
\begin{bmatrix} \dot{\mathbf{x}}_P - \dot{\mathbf{x}}'_K \\ \dot{\mathbf{x}}'_K \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_P + \mathbf{L}_P \mathbf{C}_P & \mathbf{0} \\ (\mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} - \mathbf{L}_P) \mathbf{C}_P & \mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P \end{bmatrix} \begin{bmatrix} \mathbf{x}_P - \mathbf{x}'_K \\ \mathbf{x}'_K \end{bmatrix} + \\
&+ \begin{bmatrix} \mathbf{B}_P + \mathbf{L}_P \mathbf{D}_P & \mathbf{L}_P \\ (\mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} - \mathbf{L}_P) \mathbf{D}_P & \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} - \mathbf{L}_P \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \mathbf{w} \end{bmatrix} \\
\begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix} &= \begin{bmatrix} \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{C}_P + \mathbf{D}_P \mathbf{F}_P \\ \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{F}_P \end{bmatrix} \begin{bmatrix} \mathbf{x}_P - \mathbf{x}'_K \\ \mathbf{x}'_K \end{bmatrix} + \\
&+ \begin{bmatrix} \mathbf{D}_P \mathbf{E}_{KP} & \mathbf{E}_{PK} \\ \mathbf{E}_{KP} & \mathbf{D}_K \mathbf{E}_{PK} \end{bmatrix} \begin{bmatrix} \mathbf{v} \\ \mathbf{w} \end{bmatrix}
\end{aligned} \tag{2.37}$$

In contrast to (2.33), (2.34), the matrices $\mathbf{L}_P, \mathbf{F}_P$ are now independent parameters so that the separation principle is now valid and thus the controller is of the estimator corrector type. The controller can then be given explicitly

$$\begin{aligned}
\mathbf{A}'_K &= \mathbf{A}_P - \mathbf{B}_P \mathbf{D}_K \mathbf{C}_P + \mathbf{B}_P \mathbf{E}_{KP}^{-1} \mathbf{F}_P + \mathbf{L}_P \mathbf{E}_{PK}^{-1} \mathbf{C}_P + \mathbf{L}_P \mathbf{D}_P \mathbf{E}_{KP}^{-1} \mathbf{F}_P \\
\mathbf{B}'_K &= \mathbf{B}_P \mathbf{D}_K - \mathbf{L}_P \mathbf{E}_{PK}^{-1} \\
\mathbf{C}'_K &= \mathbf{E}_{KP}^{-1} \mathbf{F}_P - \mathbf{D}_K \mathbf{C}_P
\end{aligned} \tag{2.38}$$

as a function of the design parameters $\mathbf{L}_P, \mathbf{F}_P$. The dynamics of the controller can then be given explicitly as well

$$\begin{aligned}
\text{estimator} &\begin{cases} \dot{\mathbf{x}}_K = \mathbf{A}_P \mathbf{x}_K + \mathbf{B}_P (\mathbf{u} - \mathbf{v}) + \mathbf{L} (\mathbf{y} - \mathbf{y}_K) \\ \mathbf{y}_K = \mathbf{C}_P \mathbf{x}_K + \mathbf{D}_P \mathbf{u}_K \end{cases} \\
\text{corrector} &\begin{cases} \mathbf{u} - \mathbf{v} = \mathbf{u}_K + \mathbf{D}_K (\mathbf{y} - \mathbf{y}_K) \\ \mathbf{u}_K = \mathbf{F} \mathbf{x}_K \end{cases}
\end{aligned}$$

with $\mathbf{F} = \mathbf{F}_P, \mathbf{L} = -\mathbf{L}_P \mathbf{E}_{PK}^{-1}$ and it can be seen that \mathbf{K} is the generic estimator corrector type controller, consisting of a Luenberger state estimator with arbitrary parameter \mathbf{L} and a static state corrector with arbitrary parameter \mathbf{F} .

Now suppose \mathbf{K} is of this generic estimator corrector type

$$\mathbf{K} = \left[\begin{array}{c|c} \mathbf{A}'_K & \mathbf{B}'_K \\ \hline \mathbf{C}'_K & \mathbf{D}_K \end{array} \right]$$

with $\mathbf{A}'_K, \mathbf{B}'_K, \mathbf{C}'_K$ as given in (2.38) for arbitrary $\mathbf{F}_P, \mathbf{L}_P, \mathbf{D}_K$. Then according to (2.19) and using this generic estimator corrector type controller \mathbf{K} the

dynamic matrix of the loop \mathbf{H} gives

$$\mathbf{A}_H = \begin{bmatrix} \mathbf{A}_P + \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P & \mathbf{B}_P \mathbf{F}_P - \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P \\ \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P - \mathbf{L}_P \mathbf{C}_P & \mathbf{A}_P + \mathbf{B}_P (\mathbf{F}_P - \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P) + \mathbf{L}_P \mathbf{C}_P \end{bmatrix}$$

Then clearly

$$\begin{bmatrix} -\mathbf{1} & \mathbf{1} \end{bmatrix} \mathbf{A}_H \begin{bmatrix} \mathbf{1} \\ \mathbf{1} \end{bmatrix} = \mathbf{0} = \begin{bmatrix} \mathbf{1} & -\mathbf{1} \end{bmatrix} \mathbf{A}_H \begin{bmatrix} \mathbf{1} \\ \mathbf{1} \end{bmatrix}$$

so that the Riccati equations (2.21) and (2.22) thus have the trivial solution $\mathbf{R}_{KP} = \mathbf{1} = \mathbf{R}_{PK}$. This means that the generic estimator corrector type controller, which is given explicitly in (2.38) and which can therefore be constructed for any given plant \mathbf{P} , is a factorizing controller of state dimension $p_K = p_P$, hence (2.35) follows.

Remark If the state space representation of a given plant \mathbf{P} of dimension p_P is minimal, which is equivalent to it having full McMillan degree $\delta(\mathbf{P}) = p_P$, then there always exists a factorizing controller of state space dimension $p_K = p_P$ independent of the selected underlying ring R_Z so that the Riccati equations (2.21) and (2.22) have a solution and

$$\mathbf{P} \text{ minimal} \Rightarrow K_{p_P}(\mathbf{A}_P, \mathbf{B}_P, \mathbf{C}_P, \mathbf{D}_P) \text{ is nonempty} \quad (2.39)$$

which due to (2.20) implies that $K(\mathbf{P})$ is nonempty and the doubly coprime factorization (2.7) has a solution. This stems from the fact that if the state space representation of \mathbf{P} is minimal, it is controllable so that the eigenvalues of $\mathbf{A}_P + \mathbf{B}_P \mathbf{F}_P$ can be assured to lie in a zone $Z \subset \mathbb{C}$ and it is observable so that the eigenvalues of $\mathbf{A}_P + \mathbf{L}_P \mathbf{C}_P$ can be assured to lie in the same zone Z . Hence, due to the block triangular structure of the loop dynamic matrix in (2.36) and (2.37), $\mathbf{H} \in R_Z^{m+n, m+n}$, with R_Z being the ring of proper transfer functions with all poles contained in a predefined zone Z so that (2.39) follows.

2.4.5 Parametrization of Fixed Order Controllers

There is a reversible relation between the controller \mathbf{K} and the factors $\mathbf{X}_L, \mathbf{X}_R$ in the form of $\mathbf{K} = \mathbf{X}_R \mathbf{Y}_R^{-1} = \mathbf{Y}_L^{-1} \mathbf{X}_L$ which leads to the idea of the following change of variables.

Definition The factors $\mathbf{X}_L, \mathbf{X}_R$ of the controller \mathbf{K} which appear in the fixed order loop factorizations (2.23), (2.24) can be made independent of $\mathbf{U}_K, \mathbf{V}_K$ by defining

$$\mathbf{Q}_L = \mathbf{U}_P \mathbf{X}_L = \left[\begin{array}{c|c} \mathbf{A}_{QL} & \mathbf{B}_{QL} \\ \mathbf{C}_{QL} & \mathbf{D}_{QL} \end{array} \right] \text{ with} \quad (2.40)$$

$$\begin{bmatrix} \mathbf{A}_{QL} & \mathbf{B}_{QL} \\ \mathbf{C}_{QL} & \mathbf{D}_{QL} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{L}_K \\ \mathbf{0} & \mathbf{E}_{KP} \end{bmatrix} \begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix}$$

$$\mathbf{Q}_R = \mathbf{X}_R \mathbf{V}_P = \left[\begin{array}{c|c} \mathbf{A}_{QR} & \mathbf{B}_{QR} \\ \mathbf{C}_{QR} & \mathbf{D}_{QR} \end{array} \right] \text{ with} \quad (2.41)$$

$$\begin{bmatrix} \mathbf{A}_{QR} & \mathbf{B}_{QR} \\ \mathbf{C}_{QR} & \mathbf{D}_{QR} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{F}_K & \mathbf{E}_{PK} \end{bmatrix}$$

Since they are based on the controller's factors $\mathbf{X}_L, \mathbf{X}_R$ there is a relation between these variables and the controller \mathbf{K} in the form of

$$\mathbf{K} = \mathbf{Q}_R(\mathbf{Y}_R\mathbf{V}_P)^{-1} = (\mathbf{U}_P\mathbf{Y}_L)^{-1}\mathbf{Q}_L \quad (2.42)$$

Proposition The new variables \mathbf{Q}_L and \mathbf{Q}_R that parametrize the controller \mathbf{K} as defined in (2.40) and (2.41) respectively, lead to the loop's state space matrices $\mathbf{A}_H, \mathbf{B}_H, \mathbf{C}_H, \mathbf{D}_H$ becoming linear in $\mathbf{A}_{Q_i}, \mathbf{B}_{Q_i}, \mathbf{C}_{Q_i}, \mathbf{D}_{Q_i}$ with $i \in \{L, R\}$. In addition to that the quadratic Riccati equations (2.21) and (2.22) turn into linear Sylvester equations for \mathbf{R}_{KP} and \mathbf{R}_{PK} respectively.

The following tables translate recurring terms appearing in the state space matrices of the loop (2.19) and the Riccati conditions (2.21), (2.22) in the new parametrizations.

Original Term	Term using \mathbf{Q}_L
Riccati (2.21)	$\mathbf{R}_{KP}\mathbf{A}_P = \mathbf{A}_{QL}\mathbf{R}_{KP} + \mathbf{B}_{QL}\mathbf{C}_P$
$\mathbf{A}_K + \mathbf{B}_K\mathbf{D}_P\mathbf{E}_{KP}\mathbf{C}_K$	$\mathbf{A}_{QL} + \mathbf{R}_{KP}\mathbf{B}_P\mathbf{C}_{QL}$
$\mathbf{B}_K\mathbf{E}_{PK}$	$\mathbf{B}_{QL} + \mathbf{R}_{KP}\mathbf{B}_P\mathbf{D}_{QL}$
$\mathbf{E}_{KP}\mathbf{C}_K$	\mathbf{C}_{QL}
$\mathbf{E}_{KP}\mathbf{D}_K = \mathbf{D}_K\mathbf{E}_{PK}$	\mathbf{D}_{QL}
$\mathbf{E}_{KP} = \mathbf{U}_P\mathbf{V}_K$	$1 + \mathbf{D}_{QL}\mathbf{D}_P$
$\mathbf{E}_{PK} = \mathbf{U}_K\mathbf{V}_P$	$1 + \mathbf{D}_P\mathbf{D}_{QL}$
\mathbf{L}_K	$\mathbf{B}_{QL}\mathbf{D}_P - \mathbf{R}_{KP}\mathbf{B}_P$
\mathbf{F}_P	$\mathbf{D}_{QL}\mathbf{C}_P + \mathbf{C}_{QL}\mathbf{R}_{KP}$
$\mathbf{A}_K + \mathbf{L}_K\mathbf{C}_K$	\mathbf{A}_{QL}
$\mathbf{B}_K + \mathbf{L}_K\mathbf{D}_K$	\mathbf{B}_{QL}

(2.43)

Original Term	Term using \mathbf{Q}_R
Riccati (2.22)	$\mathbf{A}_P\mathbf{R}_{PK} = \mathbf{R}_{PK}\mathbf{A}_{QR} - \mathbf{B}_P\mathbf{C}_{QR}$
$\mathbf{A}_K + \mathbf{B}_K\mathbf{D}_P\mathbf{E}_{KP}\mathbf{C}_K$	$\mathbf{A}_{QR} - \mathbf{B}_{QR}\mathbf{C}_P\mathbf{R}_{PK}$
$\mathbf{E}_{KP}\mathbf{C}_K$	$\mathbf{C}_{QR} - \mathbf{D}_{QR}\mathbf{C}_P\mathbf{R}_{PK}$
$\mathbf{B}_K\mathbf{E}_{PK}$	\mathbf{B}_{QR}
$\mathbf{D}_K\mathbf{E}_{PK} = \mathbf{E}_{KP}\mathbf{D}_K$	\mathbf{D}_{QR}
$\mathbf{E}_{PK} = \mathbf{U}_K\mathbf{V}_P$	$1 + \mathbf{D}_P\mathbf{D}_{QR}$
$\mathbf{E}_{KP} = \mathbf{U}_P\mathbf{V}_K$	$1 + \mathbf{D}_{QR}\mathbf{D}_P$
\mathbf{F}_K	$\mathbf{D}_P\mathbf{C}_{QR} + \mathbf{C}_P\mathbf{R}_{PK}$
\mathbf{L}_P	$\mathbf{B}_P\mathbf{D}_{QR} - \mathbf{R}_{PK}\mathbf{B}_{QR}$
$\mathbf{A}_K + \mathbf{B}_K\mathbf{F}_K$	\mathbf{A}_{QR}
$\mathbf{C}_K + \mathbf{D}_K\mathbf{F}_K$	\mathbf{C}_{QR}

(2.44)

Proof This follows if the definitions (2.40) and (2.41) are applied on the loop's state space matrices (2.19) and the Riccati conditions (2.21), (2.22).

2.4.6 Controllability and Observability

Proposition The right and left parametrizations respectively preserve controllability and observability of the controller

$$\begin{aligned} \mathbf{K} \text{ controllable} &\Leftrightarrow \mathbf{Q}_R \text{ controllable} \\ \mathbf{K} \text{ observable} &\Leftrightarrow \mathbf{Q}_L \text{ observable} \end{aligned} \quad (2.45)$$

Proof This can be verified by application of (C.33) on \mathbf{K} , so that \mathbf{K} is unobservable if and only if there exists an eigenvector \mathbf{x} with $\mathbf{A}_K \mathbf{x} = \lambda \mathbf{x}$ and $\mathbf{C}_K \mathbf{x} = \mathbf{0}$ which is equivalent to $\mathbf{A}_K \mathbf{x} + \mathbf{L}_K (\mathbf{C}_K \mathbf{x}) = \lambda \mathbf{x}$ and $\mathbf{E}_{KP} \mathbf{C}_K \mathbf{x} = \mathbf{0}$ since \mathbf{E}_{KP} is invertible. This again is equivalent to $\mathbf{A}_{QL} \mathbf{x} = \lambda \mathbf{x}$ and $\mathbf{C}_{QL} \mathbf{x} = \mathbf{0}$ which is true if and only if \mathbf{Q}_L is unobservable and thus verifies the second line in (2.45). The proof for \mathbf{Q}_R is similar.

2.4.7 Necessary Condition for Loop Ring Membership

Proposition Ring membership of the parameters is a necessary condition for loop ring membership

$$\mathbf{K} \in K(\mathbf{P}) \Rightarrow \mathbf{Q}_i \in R_Z^{n,m} \quad (2.46)$$

for $i \in \{L, R\}$.

Proof This is an immediate consequence of (2.7) which for $\mathbf{K} \in K(\mathbf{P})$ requires $\mathbf{X}_i \in R_Z^{n,m}$. Since $\mathbf{Q}_L = \mathbf{U}_P \mathbf{X}_L$ and $\mathbf{Q}_R = \mathbf{X}_R \mathbf{V}_P$ with $\mathbf{U}_P, \mathbf{V}_P$ being constant matrices. This can also be deduced from the triangular structure of the base transformed dynamic matrix \mathbf{A}_H of the loop \mathbf{H} in (2.33) and (2.34) where the dynamic matrices $\mathbf{A}_{QL} = \mathbf{A}_K + \mathbf{L}_K \mathbf{C}_K$ and $\mathbf{A}_{QR} = \mathbf{A}_K + \mathbf{B}_K \mathbf{F}_K$ of the parameters \mathbf{Q}_L and \mathbf{Q}_R are blocks on the diagonal.

Chapter 3

Stable Observable Pairs

This chapter introduces the notion of square lossless systems and provides a parametrization of stable observable pairs by solving a Nudelman interpolation problem. It has been shown in [2] that lossless systems and the stable observable pairs that can be extracted from them form a differentiable manifold. The controller \mathbf{K} is in general not stable, so this parametrization can not be applied to \mathbf{K} directly. However, it has been shown at the end of chapter 2 that the left parameter \mathbf{Q}_L needs to be stable for the loop \mathbf{H} to be stable and it needs to be observable for the controller \mathbf{K} to be observable. The parametrization of the stable observable pair contained in \mathbf{Q}_L then allows to vary the entire parameter \mathbf{Q}_L infinitesimally. This enables differential calculus to be carried out with \mathbf{Q}_L which will be important in the optimization stage in chapter 4.

3.1 Lossless Systems

This section develops representations of discrete square lossless systems from an energy conservation law. It further shows that any stable observable pair can be completed to a discrete square lossless system.

3.1.1 Transfer Function Representation

If the output energy of a system \mathbf{L} equals its input energy then it is energy conserving and \mathbf{L} is called a *lossless system*.

Proposition A transfer function \mathbf{L} represents a discrete stable lossless system if and only if

$$\mathbf{L}^H(e^{j\mathbb{R}})\mathbf{L}(e^{j\mathbb{R}}) = \mathbf{1} \quad (3.1)$$

with its analytic continuation

$$\mathbf{L}^H(\bar{z}^{-1})\mathbf{L}(z) = \mathbf{1} \text{ almost everywhere} \quad (3.2)$$

Proof To verify this, the energy conserving property has to be exploited. The H2 norm of a signal provides a measure of the energy contained in the signal.

A linear time invariant system $\mathbf{y} = \mathbf{L}\mathbf{u}$ is lossless if the energy of the output signal equals the energy of the input signal

$$\|\mathbf{y}\|_{H_2} = \|\mathbf{u}\|_{H_2}$$

According to (B.16) and (B.13) the H2 norm of a signal \mathbf{y} analytic in $D = \{z \in \mathbb{C} : |z| \geq 1\}$ is given by

$$\|\mathbf{y}(D)\|_{H_2} = \|\mathbf{y}(e^{j2\pi\mathbb{R}})\|_{L_2} = \sqrt{\int_{\mathbb{R}} \|\mathbf{y}(e^{j2\pi t})\|^2 dt}$$

The simplest result will be obtained when using the matrix 2 norm (B.7) which can be given by

$$\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^H \mathbf{x}}$$

for any vector \mathbf{x} . The H2 norm of the signal \mathbf{y} is then given by

$$\|\mathbf{y}(D)\|_{H_2} = \sqrt{\int_{\mathbb{R}} \mathbf{y}^H(e^{j2\pi t}) \mathbf{y}(e^{j2\pi t}) dt}$$

Applied to the energy conservation condition of a lossless system provides

$$\begin{aligned} \|\mathbf{y}(D)\|_{H_2} &= \sqrt{\int_{\mathbb{R}} \mathbf{y}^H(e^{j2\pi t}) \mathbf{y}(e^{j2\pi t}) dt} = \\ &= \sqrt{\int_{\mathbb{R}} \mathbf{u}^H(e^{j2\pi t}) \mathbf{L}^H(e^{j2\pi t}) \mathbf{L}(e^{j2\pi t}) \mathbf{u}(e^{j2\pi t}) dt} \stackrel{!}{=} \\ &\stackrel{!}{=} \sqrt{\int_{\mathbb{R}} \mathbf{u}^H(e^{j2\pi t}) \mathbf{u}(e^{j2\pi t}) dt} = \|\mathbf{u}(D)\|_{H_2} \end{aligned}$$

The general solution for all inputs \mathbf{u} is

$$\mathbf{L}^H(e^{j2\pi t}) \mathbf{L}(e^{j2\pi t}) = \mathbf{1}; \quad \forall t \in \mathbb{R}$$

so that (3.1) follows. In order to show the analytic continuation, define

$$\mathbf{F}(z) = \mathbf{L}^H(\bar{z}^{-1}) \mathbf{L}(z) - \mathbf{1}$$

which is analytic except at its poles and by (3.1) satisfies $\mathbf{F}(e^{j\mathbb{R}}) = \mathbf{0}$ on the unit circle. Since $\mathbf{L}(z)$ is stable, it has all its poles strictly inside the unit circle and consequently $\mathbf{L}^H(\bar{z}^{-1})$ has all its poles strictly outside the unit circle so that $\mathbf{F}(z)$ is analytic on a small open neighborhood of the unit circle. Then (B.4) applies and states that if $\mathbf{F}(z) = \mathbf{0}$ is true on the open neighborhood of the unit circle, then it is also true almost everywhere, that is everywhere except at its poles. This then verifies (3.2).

3.1.2 State Space Representation

A discrete lossless system $\mathbf{L}(s)$ is called *square* if its transfer matrix has an equal number of columns (inputs) and lines (outputs).

Proposition A minimal realization $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ represents a discrete stable minimal lossless system if and only if

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^H \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \text{ with } \mathbf{G} > 0 \quad (3.3)$$

where $\mathbf{G}, \mathbf{G}^{-1}$ are the system's Gram observability and controllability matrices respectively.

Proof This can be verified by applying the lossless condition (3.1) on a transfer function $\mathbf{L}(z)$ given in terms of the state space matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$ as in (C.7):

$$\mathbf{1} = \mathbf{L}^H(z)\mathbf{L}(z) = (\mathbf{C}(z\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D})^H (\mathbf{C}(z\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D})$$

for $|z| = 1 \Leftrightarrow \bar{z} = z^{-1}$. This develops into

$$\begin{aligned} \mathbf{1} = \mathbf{L}^H(z)\mathbf{L}(z) &= \mathbf{B}^H(z^{-1}\mathbf{1} - \mathbf{A}^H)^{-1}\mathbf{C}^H\mathbf{C}(z\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + \\ &+ \mathbf{B}^H(z^{-1}\mathbf{1} - \mathbf{A}^H)^{-1}\mathbf{C}^H\mathbf{D} + \mathbf{D}^H\mathbf{C}(z\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}^H\mathbf{D} \end{aligned} \quad (3.4)$$

Since $\mathbf{L}(z)$ is required to be stable, the inverse $(z\mathbf{1} - \mathbf{A})^{-1} = z^{-1}(\mathbf{1} - z^{-1}\mathbf{A})^{-1}$ and thus also $(\bar{z}\mathbf{1} - \mathbf{A}^H)^{-1} = (z^{-1}\mathbf{1} - \mathbf{A}^H)^{-1} = z(\mathbf{1} - z\mathbf{A}^H)^{-1}$ exist for $|z| = 1 \Leftrightarrow \bar{z} = z^{-1}$. And since the realization of $\mathbf{L}(z)$ is furthermore required to be minimal and thus observable, according to (C.34) there is a unique Gram observability matrix $\mathbf{G} > 0$ that solves

$$\mathbf{G} - \mathbf{A}^H\mathbf{G}\mathbf{A} = \mathbf{C}^H\mathbf{C} \quad (3.5)$$

In order to simplify the lossless condition (3.4), the expansion of the following product is helpful:

$$\begin{aligned} (z^{-1}\mathbf{1} - \mathbf{A}^H)\mathbf{G}(z\mathbf{1} - \mathbf{A}) &= (\mathbf{1} - z\mathbf{A}^H)\mathbf{G}(\mathbf{1} - z^{-1}\mathbf{A}) = \\ &= \mathbf{G}(\mathbf{1} - z^{-1}\mathbf{A}) + (\mathbf{1} - z\mathbf{A}^H)\mathbf{G} + \mathbf{A}^H\mathbf{G}\mathbf{A} - \mathbf{G} = \\ &= \mathbf{G}(\mathbf{1} - z^{-1}\mathbf{A}) + (\mathbf{1} - z\mathbf{A}^H)\mathbf{G} - \mathbf{C}^H\mathbf{C} \end{aligned}$$

using (3.5). Left multiplication by $(\mathbf{1} - z\mathbf{A}^H)^{-1}$ and right multiplication by $(\mathbf{1} - z^{-1}\mathbf{A})^{-1}$ leads to

$$\begin{aligned} (\mathbf{1} - z\mathbf{A}^H)^{-1}\mathbf{C}^H\mathbf{C}(\mathbf{1} - z^{-1}\mathbf{A})^{-1} &= (z^{-1}\mathbf{1} - \mathbf{A}^H)^{-1}\mathbf{C}^H\mathbf{C}(z\mathbf{1} - \mathbf{A})^{-1} = \\ &= \mathbf{G}(\mathbf{1} - z^{-1}\mathbf{A})^{-1} + (\mathbf{1} - z\mathbf{A}^H)^{-1}\mathbf{G} - \mathbf{G} \end{aligned}$$

The matrix relation (A.56) can be used to obtain

$$\begin{aligned} (\mathbf{1} - z^{-1}\mathbf{A})^{-1} &= \mathbf{1} + z^{-1}\mathbf{A}(\mathbf{1} - z^{-1}\mathbf{A})^{-1} = \mathbf{1} + \mathbf{A}(z\mathbf{1} - \mathbf{A})^{-1} \\ (\mathbf{1} - z\mathbf{A}^H)^{-1} &= \mathbf{1} + (\mathbf{1} - z\mathbf{A}^H)^{-1}z\mathbf{A}^H = \mathbf{1} + (z^{-1}\mathbf{1} - \mathbf{A}^H)^{-1}\mathbf{A}^H \end{aligned}$$

and thus

$$(z^{-1}\mathbf{1} - \mathbf{A}^H)^{-1}\mathbf{C}^H\mathbf{C}(z\mathbf{1} - \mathbf{A})^{-1} = \mathbf{G}\mathbf{A}(z\mathbf{1} - \mathbf{A})^{-1} + (z^{-1}\mathbf{1} - \mathbf{A}^H)^{-1}\mathbf{A}^H\mathbf{G} + \mathbf{G}$$

Substituting this in (3.4) gives

$$\mathbf{M}^H (z\mathbf{1} - \mathbf{A})^{-1} \mathbf{B} + \mathbf{B}^H (z^{-1}\mathbf{1} - \mathbf{A}^H)^{-1} \mathbf{M} = \mathbf{N}$$

with $\mathbf{M} = \mathbf{A}^H \mathbf{G} \mathbf{B} + \mathbf{C}^H \mathbf{D}$ and $\mathbf{N} = \mathbf{1} - \mathbf{B}^H \mathbf{G} \mathbf{B} - \mathbf{D}^H \mathbf{D}$. According to (A.67), the inverses can now be developed in geometric series:

$$\begin{aligned} (z\mathbf{1} - \mathbf{A})^{-1} &= z^{-1}(\mathbf{1} - z^{-1}\mathbf{A})^{-1} = z^{-1} \sum_{k=0}^{\infty} z^{-k} \mathbf{A}^k = \sum_{k=-\infty}^{-1} z^k \mathbf{A}^{-k-1} \\ (z^{-1}\mathbf{1} - \mathbf{A}^H)^{-1} &= z(\mathbf{1} - z\mathbf{A}^H)^{-1} = z \sum_{k=0}^{\infty} z^k \mathbf{A}^{Hk} = \sum_{k=1}^{\infty} z^k \mathbf{A}^{H-k-1} \end{aligned}$$

Now the lossless condition writes as

$$\mathbf{M}^H \left(\sum_{k=-\infty}^{-1} z^k \mathbf{A}^{-k-1} \right) \mathbf{B} + \mathbf{B}^H \left(\sum_{k=1}^{\infty} z^k \mathbf{A}^{Hk-1} \right) \mathbf{M} = \mathbf{N}$$

which is equal to

$$\sum_{k=-\infty}^{\infty} z^k \mathbf{F}_k = \mathbf{N} \text{ with } \mathbf{F}_k = \begin{cases} \mathbf{M}^H \mathbf{A}^{k-1} \mathbf{B}; & k \leq -1 \\ \mathbf{0}; & k = 0 \\ \mathbf{B}^H \mathbf{A}^{Hk-1} \mathbf{M}; & k \geq 1 \end{cases}$$

Comparing the coefficients provides

$$\begin{aligned} \mathbf{M}^H [\mathbf{B} \quad \mathbf{A}\mathbf{B} \quad \mathbf{A}^2\mathbf{B} \dots] &= \mathbf{0} \text{ for } k \neq 0 \\ \mathbf{N} &= \mathbf{0} \text{ for } k = 0 \end{aligned}$$

But since the realization of $\mathbf{L}(z)$ is required to be minimal and thus controllable, according to (C.31) the matrix $[\mathbf{B} \quad \mathbf{A}\mathbf{B} \quad \mathbf{A}^2\mathbf{B} \dots]$ has full line rank and thus a discrete square stable minimal lossless system is characterized by

$$\begin{aligned} \mathbf{M} &= \mathbf{A}^H \mathbf{G} \mathbf{B} + \mathbf{C}^H \mathbf{D} = \mathbf{0} \\ \mathbf{N} &= \mathbf{1} - \mathbf{B}^H \mathbf{G} \mathbf{B} - \mathbf{D}^H \mathbf{D} = \mathbf{0} \end{aligned}$$

in addition to (3.5). Assembling these three equations leads to the block matrix form of (3.3).

3.1.3 Discrete Lossless Completion

Proposition A stable observable pair (\mathbf{A}, \mathbf{C}) with Gram observability matrix \mathbf{G} can be completed by a pair (\mathbf{B}, \mathbf{D}) so that the realization $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ represents a discrete stable minimal lossless system

$$\begin{aligned} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^H \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} &= \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \\ \text{parametrized by } \begin{bmatrix} \mathbf{B} \\ \mathbf{D} \end{bmatrix} &= \begin{bmatrix} \mathbf{B}_0 \\ \mathbf{D}_0 \end{bmatrix} \mathbf{U} \end{aligned} \tag{3.6}$$

with \mathbf{U} unitary.

Proof This is due to the fact that observability of the pair $(\mathbf{A} \in \mathbb{C}^{n,n}, \mathbf{C} \in \mathbb{C}^{m,n})$ implies the existence of a Gram observability matrix $\mathbf{G} > 0$ so that

$$\mathbf{A}^H \mathbf{G} \mathbf{A} + \mathbf{C}^H \mathbf{C} = \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \mathbf{A} \\ \mathbf{C} \end{bmatrix}^H \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \mathbf{A} \\ \mathbf{C} \end{bmatrix} = \mathbf{G} \geq 0$$

according to (C.34) and (A.77). This shows that

$$\text{rk} \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \mathbf{A} \\ \mathbf{C} \end{bmatrix} = n$$

so that according to (A.80) there exist orthogonal completions

$$\begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \mathbf{A} \\ \mathbf{C} \end{bmatrix}_{\perp} = \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'_0 \\ \mathbf{Y}'_0 \end{bmatrix} \mathbf{T}$$

parametrized by invertible \mathbf{T} so that

$$\begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \mathbf{A} \\ \mathbf{C} \end{bmatrix} \underbrace{\left(\begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \mathbf{A} \\ \mathbf{C} \end{bmatrix}^H \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \mathbf{A} \\ \mathbf{C} \end{bmatrix} \right)^{-1}}_{\mathbf{G}^{-1}} \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \mathbf{A} \\ \mathbf{C} \end{bmatrix} + \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}^H \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}^H = \mathbf{1} \quad (3.7)$$

The restriction

$$\mathbf{T} = \left(\begin{bmatrix} \mathbf{X}'_0 \\ \mathbf{Y}'_0 \end{bmatrix}^H \begin{bmatrix} \mathbf{X}'_0 \\ \mathbf{Y}'_0 \end{bmatrix} \right)^{-\frac{1}{2}} \mathbf{U}$$

with \mathbf{U} unitary, reduces the set of orthogonal completions to

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \begin{bmatrix} \mathbf{X}_0 \\ \mathbf{Y}_0 \end{bmatrix} \mathbf{U} \quad (3.8)$$

with

$$\begin{bmatrix} \mathbf{X}_0 \\ \mathbf{Y}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{X}'_0 \\ \mathbf{Y}'_0 \end{bmatrix} \left(\begin{bmatrix} \mathbf{X}'_0 \\ \mathbf{Y}'_0 \end{bmatrix}^H \begin{bmatrix} \mathbf{X}'_0 \\ \mathbf{Y}'_0 \end{bmatrix} \right)^{-\frac{1}{2}}$$

so that the constraint

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}^H \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \mathbf{1}$$

is fulfilled. This together with (3.7) gives

$$\begin{bmatrix} \mathbf{G}^{\frac{1}{2}} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A} \mathbf{G}^{-1} \mathbf{A}^H & \mathbf{A} \mathbf{G}^{-1} \mathbf{C}^H \\ \mathbf{C} \mathbf{G}^{-1} \mathbf{A}^H & \mathbf{C} \mathbf{G}^{-1} \mathbf{C}^H \end{bmatrix} \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} + \begin{bmatrix} \mathbf{X} \mathbf{X}^H & \mathbf{X} \mathbf{Y}^H \\ \mathbf{Y} \mathbf{X}^H & \mathbf{Y} \mathbf{Y}^H \end{bmatrix} = \mathbf{1}$$

which can be written as

$$\begin{bmatrix} \mathbf{A} & \mathbf{G}^{-\frac{1}{2}} \mathbf{X} \\ \mathbf{C} & \mathbf{Y} \end{bmatrix} \begin{bmatrix} \mathbf{G}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{G}^{-\frac{1}{2}} \mathbf{X} \\ \mathbf{C} & \mathbf{Y} \end{bmatrix}^H = \begin{bmatrix} \mathbf{G}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \quad (3.9)$$

With the definition of $\mathbf{B} = \mathbf{G}^{-\frac{1}{2}} \mathbf{X}$, $\mathbf{D} = \mathbf{Y}$ and $\mathbf{B}_0 = \mathbf{G}^{-\frac{1}{2}} \mathbf{X}_0$, $\mathbf{D}_0 = \mathbf{Y}_0$ the orthogonal completions (3.8) are given by

$$\begin{bmatrix} \mathbf{B} \\ \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{G}^{-\frac{1}{2}} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \begin{bmatrix} \mathbf{G}^{-\frac{1}{2}} \mathbf{X}_0 \\ \mathbf{Y}_0 \end{bmatrix} \mathbf{U} = \begin{bmatrix} \mathbf{B}_0 \\ \mathbf{D}_0 \end{bmatrix} \mathbf{U}$$

and equation (3.9) can be written in the form of (3.3) so that (3.6) follows.

3.1.4 Link to Conjugate Realization

Proposition If an observable pair (\mathbf{A}, \mathbf{C}) has a matrix \mathbf{I} that solves

$$\begin{bmatrix} \mathbf{A} \\ \mathbf{C} \end{bmatrix} \mathbf{I} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A} \\ \mathbf{C} \end{bmatrix}} \quad (3.10)$$

then \mathbf{I} is unique and invertible.

Proof Assume there exists a matrix \mathbf{I} that solves (3.10). Since the pair (\mathbf{A}, \mathbf{C}) is observable the observability matrix

$$\mathbf{O} = \begin{bmatrix} \mathbf{C} \\ \mathbf{CA} \\ \vdots \end{bmatrix}$$

has full column rank and provides

$$\mathbf{O}\mathbf{I} = \overline{\mathbf{O}}$$

using (3.10). Suppose \mathbf{I} not uniquely determined by (3.10) then there should exist another $\mathbf{I}' \neq \mathbf{I}$ with $\mathbf{O}\mathbf{I}' = \overline{\mathbf{O}}$ hence $\mathbf{O}(\mathbf{I} - \mathbf{I}') = \mathbf{0}$ which contradicts \mathbf{O} having full column rank so that \mathbf{I} must be unique. Suppose \mathbf{I} not invertible then there should exist a vector $\mathbf{z} \neq \mathbf{0}$ with $\mathbf{I}\mathbf{z} = \mathbf{0}$ hence $\mathbf{O}\mathbf{I}\mathbf{z} = \mathbf{0}$ which contradicts \mathbf{O} having full column rank so that \mathbf{I} must be invertible.

Proposition If the observable pair (\mathbf{A}, \mathbf{C}) of a discrete stable minimal lossless system $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ has a matrix \mathbf{I} that solves (3.10) then it also solves

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}} \begin{bmatrix} \mathbf{I}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{U} \end{bmatrix} \quad (3.11)$$

with unitary \mathbf{U} given by

$$\mathbf{U} = \overline{\mathbf{B}}^H \mathbf{I}^H \mathbf{G} \mathbf{B} + \overline{\mathbf{D}}^H \mathbf{D} \quad (3.12)$$

and the Gram observability matrix $\overline{\mathbf{G}}$ of the conjugate pair $(\overline{\mathbf{A}}, \overline{\mathbf{C}})$ satisfying

$$\overline{\mathbf{G}} = \mathbf{I}^H \mathbf{G} \mathbf{I} \quad (3.13)$$

Hence the lossless system is real in the sense of (C.27) if and only if $\mathbf{U} = \mathbf{1}$ due to (C.28).

Proof Assume there exists a matrix \mathbf{I} that solves (3.10) for the observable pair (\mathbf{A}, \mathbf{C}) of the discrete stable minimal lossless system $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$. The Gram observability matrix $\mathbf{G} > \mathbf{0}$ is uniquely defined by

$$\mathbf{A}^H \mathbf{G} \mathbf{A} - \mathbf{G} + \mathbf{C}^H \mathbf{C} = \mathbf{0}$$

according to (C.34). With (3.10) the conjugate of this equation provides

$$\mathbf{A}^H \mathbf{I}^{-H} \overline{\mathbf{G}} \mathbf{I}^{-1} \mathbf{A} - \mathbf{I}^{-H} \mathbf{G} \mathbf{I}^{-1} + \mathbf{C}^H \mathbf{C} = \mathbf{0}$$

Subtracting this from (3.10) gives

$$\mathbf{A}^H(\mathbf{G} - \mathbf{I}^{-H}\overline{\mathbf{G}}\mathbf{I}^{-1})\mathbf{A} - (\mathbf{G} - \mathbf{I}^{-H}\mathbf{G}\mathbf{I}^{-1}) = \mathbf{0}$$

This is a Stein equation and according to (A.88) has a unique solution (3.13) since \mathbf{A} is stable. With (3.10) and (3.13) the conjugate of the lossless condition (3.3) provides

$$\begin{bmatrix} \mathbf{A} & \overline{\mathbf{I}\mathbf{B}} \\ \mathbf{C} & \overline{\mathbf{D}} \end{bmatrix}^H \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \overline{\mathbf{I}\mathbf{B}} \\ \mathbf{C} & \overline{\mathbf{D}} \end{bmatrix} = \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}$$

Thus (\mathbf{B}, \mathbf{D}) and $(\overline{\mathbf{I}\mathbf{B}}, \overline{\mathbf{D}})$ are two lossless completions to the pair (\mathbf{A}, \mathbf{C}) . Then according to (3.6) there exists a unitary \mathbf{U} so that

$$\begin{bmatrix} \mathbf{B} \\ \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{B} \\ \mathbf{D} \end{bmatrix}} \mathbf{U}$$

Together with (3.10) this can be written in the form (3.11). Due to (3.3) we can rewrite (3.11) in the form of

$$\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}^H \end{bmatrix} = \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^H \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}}$$

so that \mathbf{U} is uniquely defined and (3.12) follows.

3.2 Nudelman Interpolation Problem

This section introduces the Nudelman interpolation problem and provides a solution to it in the form of a constraint on general discrete square lossless systems. Extensive work on this subject has been done in [11] wherein it is also shown that the Nudelman interpolation problem is a matrix generalization of the Nevanlinna-Pick interpolation problem. A parametrization of the solutions of the Nudelman interpolation problem then leads to a parametrization of stable observable pairs.

3.2.1 Interpolation Problem Definition

The *Nudelman interpolation problem* consists of finding a discrete lossless system $\mathbf{L}(z)$ which for a given stable and observable pair $(\mathbf{A}_0, \mathbf{C}_0)$ and interpolation value \mathbf{W} solves

$$\mathbf{W} = \frac{1}{2\pi j} \oint_C \mathbf{L}^H(\bar{z}^{-1})\mathbf{C}_0(z\mathbf{1} - \mathbf{A}_0)^{-1}dz \quad (3.14)$$

with $C = \{z : |z| = 1\}$ being the unit circle and $\delta(\mathbf{L}) = \text{col } \mathbf{A}_0$, which means the McMillan degree of \mathbf{L} equals the state space dimension of \mathbf{A}_0 . Note that $\mathbf{L}(\bar{z}^{-1}) = \mathbf{L}^{-1}(z)$ due to (3.2).

Proposition The Nudelman interpolation problem (3.14) has an explicit solution $\mathbf{L}(z)$ if the existence condition

$$\exists \mathbf{G} > \mathbf{0} : \mathbf{A}_0^H \mathbf{G} \mathbf{A}_0 - \mathbf{G} + \mathbf{C}_0^H \mathbf{C}_0 = \mathbf{W}^H \mathbf{W} \quad (3.15)$$

is satisfied. The solution $\mathbf{L}(z)$ then has a stable minimal state space representation $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$, with \mathbf{A}_0, \mathbf{A} having same dimension, that is constrained by

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \mathbf{W} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{C}_0 \end{bmatrix} \quad (3.16)$$

with $\mathbf{L}(z) = \mathbf{C}(z\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$ being a square stable minimal lossless system with Gram observability matrix \mathbf{G} and

$$0 < \mathbf{G} \leq \mathbf{G}_0 \quad (3.17)$$

wherein \mathbf{G}_0 is the Gram observability matrix of the pair $(\mathbf{A}_0, \mathbf{C}_0)$.

Proof In order to verify (3.16), let $(\mathbf{A}', \mathbf{B}', \mathbf{C}', \mathbf{D}')$ be another minimal state space representation of \mathbf{L} with the dynamic matrices $\mathbf{A}_0, \mathbf{A}'$ having same dimension, then with $\mathbf{L}^H(\bar{z}^{-1}) = \mathbf{B}'^H(z^{-1}\mathbf{1} - \mathbf{A}'^H)^{-1}\mathbf{C}'^H + \mathbf{D}'^H$, the interpolation problem (3.14) writes as

$$\begin{aligned} \mathbf{W} = \mathbf{B}'^H \frac{1}{2\pi j} \oint_C (z^{-1}\mathbf{1} - \mathbf{A}'^H)^{-1} \mathbf{C}'^H \mathbf{C}_0 (z\mathbf{1} - \mathbf{A}_0)^{-1} dz + \\ + \mathbf{D}'^H \mathbf{C}_0 \frac{1}{2\pi j} \oint_C (z\mathbf{1} - \mathbf{A}_0)^{-1} dz \end{aligned}$$

For the first integral, using (A.56), we get

$$\begin{aligned} (z^{-1}\mathbf{1} - \mathbf{A}'^H)^{-1} &= z(\mathbf{1} - z\mathbf{A}'^H)^{-1} \\ (z\mathbf{1} - \mathbf{A}_0)^{-1} &= z^{-1}(\mathbf{1} - z^{-1}\mathbf{A}_0)^{-1} = z^{-1}((z\mathbf{1} - \mathbf{A}_0)^{-1}\mathbf{A}_0 + \mathbf{1}) \end{aligned}$$

For the second integral, using (A.67) on $(z\mathbf{1} - \mathbf{A}_0)^{-1}$ we get

$$(z\mathbf{1} - \mathbf{A}_0)^{-1} = z^{-1}(\mathbf{1} - z^{-1}\mathbf{A}_0)^{-1} = z^{-1} \sum_{k=0}^{\infty} z^{-k} \mathbf{A}_0^k = \sum_{k=0}^{\infty} z^{-(k+1)} \mathbf{A}_0^k$$

which converges on C since \mathbf{A}_0 is stable. Using (B.3) the residue of this function is $\mathbf{1}$ and the second integral can be evaluated as

$$\mathbf{D}'^H \mathbf{C}_0 \frac{1}{2\pi j} \oint_C (z\mathbf{1} - \mathbf{A}_0)^{-1} dz = \mathbf{D}'^H \mathbf{C}_0$$

so that the interpolation problem can now be written as

$$\begin{aligned} \mathbf{W} = \mathbf{B}'^H \frac{1}{2\pi j} \oint_C (\mathbf{1} - z\mathbf{A}'^H)^{-1} \mathbf{C}'^H \mathbf{C}_0 (z\mathbf{1} - \mathbf{A}_0)^{-1} dz \mathbf{A}_0 + \\ + \mathbf{B}'^H \frac{1}{2\pi j} \oint_C (\mathbf{1} - z\mathbf{A}'^H)^{-1} dz \mathbf{C}'^H \mathbf{C}_0 + \mathbf{D}'^H \mathbf{C}_0 \end{aligned}$$

For the first integral, define

$$\mathbf{T} = \frac{1}{2\pi j} \oint_C (\mathbf{1} - z\mathbf{A}'^H)^{-1} \mathbf{C}'^H \mathbf{C}_0 (z\mathbf{1} - \mathbf{A}_0)^{-1} dz \quad (3.18)$$

which is a square matrix since $\mathbf{A}_0, \mathbf{A}'$ are required to have same dimension. For the second integral, using (A.67) on $(\mathbf{1} - z\mathbf{A}'^H)^{-1}$ we get

$$(\mathbf{1} - z\mathbf{A}'^H)^{-1} = \sum_{k=0}^{\infty} z^k \mathbf{A}'^{Hk}$$

which converges on C since \mathbf{A}' is stable. Using (B.3) the residue of this function is $\mathbf{0}$ and the second integral can be evaluated as

$$\mathbf{B}'^H \frac{1}{2\pi j} \oint_C (\mathbf{1} - z\mathbf{A}'^H)^{-1} dz \mathbf{C}'^H \mathbf{C}_0 = \mathbf{0}$$

so that the interpolation problem can now be written as

$$\mathbf{W} = \mathbf{B}'^H \mathbf{T} \mathbf{A}_0 + \mathbf{D}'^H \mathbf{C}_0 \quad (3.19)$$

A matrix equation for \mathbf{T} as defined in (3.18) can be constructed by calculating

$$\mathbf{A}'^H \mathbf{T} \mathbf{A}_0 = \frac{1}{2\pi j} \oint_C \mathbf{A}'^H (\mathbf{1} - z\mathbf{A}'^H)^{-1} \mathbf{C}'^H \mathbf{C}_0 (z\mathbf{1} - \mathbf{A}_0)^{-1} \mathbf{A}_0 dz$$

Again using (A.56), we get

$$\begin{aligned} \mathbf{A}'^H (\mathbf{1} - z\mathbf{A}'^H)^{-1} &= z^{-1} z \mathbf{A}'^H (\mathbf{1} - z\mathbf{A}'^H)^{-1} = z^{-1} ((\mathbf{1} - z\mathbf{A}'^H)^{-1} - \mathbf{1}) \\ (z\mathbf{1} - \mathbf{A}_0)^{-1} \mathbf{A}_0 &= (\mathbf{1} - z^{-1} \mathbf{A}_0)^{-1} z^{-1} \mathbf{A}_0 = (\mathbf{1} - z^{-1} \mathbf{A}_0)^{-1} - \mathbf{1} \end{aligned}$$

so that

$$\begin{aligned} \mathbf{A}'^H \mathbf{T} \mathbf{A}_0 &= \mathbf{T} - \frac{1}{2\pi j} \oint_C z^{-1} (\mathbf{1} - z\mathbf{A}'^H)^{-1} dz \mathbf{C}'^H \mathbf{C}_0 - \\ &\quad - \mathbf{C}'^H \mathbf{C}_0 \frac{1}{2\pi j} \oint_C (z\mathbf{1} - \mathbf{A}_0)^{-1} dz + \frac{1}{2\pi j} \oint_C z^{-1} dz \mathbf{C}'^H \mathbf{C}_0 \end{aligned}$$

For the first integral, using (A.67) on $z^{-1}(\mathbf{1} - z\mathbf{A}'^H)^{-1}$ we get

$$z^{-1}(\mathbf{1} - z\mathbf{A}'^H)^{-1} = z^{-1} \sum_{k=0}^{\infty} z^k \mathbf{A}'^{Hk} = \sum_{k=0}^{\infty} z^{k-1} \mathbf{A}'^{Hk}$$

which converges on C since \mathbf{A}' is stable. Using (B.3) the residue of this function is $\mathbf{1}$ and the first integral can be evaluated as

$$\frac{1}{2\pi j} \oint_C z^{-1} (\mathbf{1} - z\mathbf{A}'^H)^{-1} dz \mathbf{C}'^H \mathbf{C}_0 = \mathbf{C}'^H \mathbf{C}_0$$

For the second integral, using (A.67) on $(z\mathbf{1} - \mathbf{A}_0)^{-1}$ we get

$$(z\mathbf{1} - \mathbf{A}_0)^{-1} = z^{-1} (\mathbf{1} - z^{-1} \mathbf{A}_0)^{-1} = z^{-1} \sum_{k=0}^{\infty} z^{-k} \mathbf{A}_0^k = \sum_{k=0}^{\infty} z^{-(k+1)} \mathbf{A}_0^k$$

which converges on C since \mathbf{A}_0 is stable. Using (B.3) the residue of this function is $\mathbf{1}$ and the second integral can be evaluated as

$$\mathbf{C}'^H \mathbf{C}_0 \frac{1}{2\pi j} \oint_C (z\mathbf{1} - \mathbf{A}_0)^{-1} dz = \mathbf{C}'^H \mathbf{C}_0$$

For the third integral, using (B.3) the residue of z^{-1} is 1 and the third integral can be evaluated as

$$\frac{1}{2\pi j} \oint_C z^{-1} dz C'^H C_0 = C'^H C_0$$

Assembling the expressions leads to a Stein equation for \mathbf{T}

$$\mathbf{A}'^H \mathbf{T} \mathbf{A}_0 - \mathbf{T} + \mathbf{C}'^H \mathbf{C}_0 = \mathbf{0} \quad (3.20)$$

which according to (A.88) uniquely determines \mathbf{T} since $\mathbf{A}_0, \mathbf{A}'$ are stable and thus have eigenvalues $|\lambda_A| < 1, |\lambda_{A'}| < 1$. Together, matrix equations (3.20) and (3.19) can be written in block matrix form

$$\begin{bmatrix} \mathbf{A}' & \mathbf{B}' \\ \mathbf{C}' & \mathbf{D}' \end{bmatrix}^H \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{C}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{T} \\ \mathbf{W} \end{bmatrix} \quad (3.21)$$

Since $(\mathbf{A}', \mathbf{B}', \mathbf{C}', \mathbf{D}')$ is a minimal state space representation of square lossless \mathbf{L} , the adjoint Gram observability matrix $\mathbf{G}' > 0$ exists and (3.3) applies, so that (3.21) provides

$$\begin{aligned} \mathbf{T}^H \mathbf{G}'^{-1} \mathbf{T} + \mathbf{W}^H \mathbf{W} &= \begin{bmatrix} \mathbf{T} \\ \mathbf{W} \end{bmatrix}^H \begin{bmatrix} \mathbf{G}' & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{T} \\ \mathbf{W} \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{C}_0 \end{bmatrix}^H \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^H \begin{bmatrix} \mathbf{A}' & \mathbf{B}' \\ \mathbf{C}' & \mathbf{D}' \end{bmatrix} \begin{bmatrix} \mathbf{G}' & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A}' & \mathbf{B}' \\ \mathbf{C}' & \mathbf{D}' \end{bmatrix}^H \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{C}_0 \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{C}_0 \end{bmatrix}^H \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^H \begin{bmatrix} \mathbf{G}' & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{C}_0 \end{bmatrix} = \\ &= \mathbf{A}_0^H \mathbf{T}^H \mathbf{G}'^{-1} \mathbf{T} \mathbf{A}_0 + \mathbf{C}_0^H \mathbf{C}_0 \end{aligned}$$

Comparing this to the given existence condition (3.15) provides

$$\mathbf{T}^H \mathbf{G}'^{-1} \mathbf{T} = \mathbf{G} > 0 \quad (3.22)$$

This means, if the existence condition (3.15) is satisfied, the square matrix \mathbf{T} has full rank which implies \mathbf{T} to be invertible. Hence (3.21) can be written as

$$\left(\begin{bmatrix} \mathbf{T}^H & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}' & \mathbf{B}' \\ \mathbf{C}' & \mathbf{D}' \end{bmatrix} \begin{bmatrix} \mathbf{T}^H & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \right)^H \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{C}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{1} \\ \mathbf{W} \end{bmatrix}$$

which makes it clear that \mathbf{T} is a state space basis transformation matrix as in (C.36) that naturally leads to a new state space representation $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ with

$$\begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{T}^H & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}' & \mathbf{B}' \\ \mathbf{C}' & \mathbf{D}' \end{bmatrix} \begin{bmatrix} \mathbf{T}^H & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1}$$

so that (3.21) leads to (3.16). The old state space representation $(\mathbf{A}', \mathbf{B}', \mathbf{C}', \mathbf{D}')$ being stable minimal lossless with Gram observability matrix \mathbf{G}' , (3.3) provides

$$\begin{bmatrix} \mathbf{A}' & \mathbf{B}' \\ \mathbf{C}' & \mathbf{D}' \end{bmatrix}^H \begin{bmatrix} \mathbf{G}' & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}' & \mathbf{B}' \\ \mathbf{C}' & \mathbf{D}' \end{bmatrix} = \begin{bmatrix} \mathbf{G}' & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}$$

Transformed to the new state space representation and using (3.22) gives

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^H \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}$$

from which can be seen that \mathbf{G} is the Gram observability matrix of the state space representation $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ of the discrete lossless system \mathbf{L} .

Condition (3.17) can be readily verified by subtracting the solution existence condition (3.15) from the Stein equation for the Gram observability matrix \mathbf{G}_0 of the pair $(\mathbf{A}_0, \mathbf{C}_0)$ in (C.34) which gives

$$\mathbf{A}_0^H (\mathbf{G}_0 - \mathbf{G}) \mathbf{A}_0 - (\mathbf{G}_0 - \mathbf{G}) + \mathbf{W}^H \mathbf{W} = \mathbf{0} \quad (3.23)$$

According to (A.89), this new symmetric Stein equation has a unique solution $\mathbf{G}_0 - \mathbf{G} \geq 0 \Leftrightarrow \mathbf{G} \leq \mathbf{G}_0$ since $\mathbf{W}^H \mathbf{W} \geq 0$ and \mathbf{A}_0 is stable and thus has eigenvalues $|\lambda_A| < 1$.

3.2.2 Interpolation Problem Solution

Proposition If there exists a matrix $\mathbf{G} > 0$ that verifies (3.15), the interpolation problem (3.14) has a lossless solution \mathbf{L} with state space representation $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ constrained by (3.16). Then the general solution is given by

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_0 & \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix} \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{G} & \mathbf{W}^H \\ \mathbf{W} & -\mathbf{1} \end{bmatrix} \quad (3.24)$$

with

$$\begin{aligned} \mathbf{S} &= (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0)^{-\frac{1}{2}} (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}} \\ \mathbf{T} &= (\mathbf{G} + \mathbf{W}^H \mathbf{W})^{-1} \end{aligned} \quad (3.25)$$

and $(\mathbf{B}_0, \mathbf{D}_0)$ being any lossless completion to the observable pair $(\mathbf{A}_0, \mathbf{C}_0)$ with Gram observability matrix \mathbf{G}_0 , as in (3.6).

Proof To verify this, first note that according to (3.6), an observable pair $(\mathbf{A}_0, \mathbf{C}_0)$ with Gram observability matrix $\mathbf{G}_0 > 0$ can always be completed by $(\mathbf{B}_0, \mathbf{D}_0)$ so that

$$\begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix}^H \begin{bmatrix} \mathbf{G}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{G}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \quad (3.26)$$

Now the constraint (3.16) can be written in the form

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{X} \\ \mathbf{W} & \mathbf{Y} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix} \quad (3.27)$$

so that the state space matrices $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ of the interpolation problem solution \mathbf{L} are a function of yet unknown matrices \mathbf{X}, \mathbf{Y} that have to be determined in the following.

Since \mathbf{L} must be lossless and thus satisfy (3.3), this condition can be used on (3.27) to obtain a constraint for \mathbf{X}, \mathbf{Y} :

$$\begin{aligned} \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} &= \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-H} = \\ &= \begin{bmatrix} \mathbf{1} & \mathbf{X} \\ \mathbf{W} & \mathbf{Y} \end{bmatrix} \mathbf{H}^{-1} \begin{bmatrix} \mathbf{1} & \mathbf{X} \\ \mathbf{W} & \mathbf{Y} \end{bmatrix}^H \end{aligned} \quad (3.28)$$

with

$$\mathbf{H} = \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix}^H \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{H}_1 & \mathbf{H}_2 \\ \mathbf{H}_2^H & \mathbf{H}_3 \end{bmatrix} \quad (3.29)$$

The matrix \mathbf{H}_1 is invertible which can be seen by

$$\mathbf{H}_1 = \mathbf{A}_0^H \mathbf{G} \mathbf{A}_0 + \mathbf{C}_0^H \mathbf{C}_0 = \mathbf{G} + \mathbf{W}^H \mathbf{W} > 0 \quad (3.30)$$

using (3.15). Thus, according to (A.57), \mathbf{H} can be factorized into

$$\mathbf{H} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{E}^H & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{H}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{F} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{E} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \quad (3.31)$$

with $\mathbf{E} = \mathbf{H}_1^{-1} \mathbf{H}_2$ and $\mathbf{F} = \mathbf{H}_3 - \mathbf{H}_2^H \mathbf{H}_1^{-1} \mathbf{H}_2 = \mathbf{F}^H$, which is invertible since \mathbf{H} is invertible, so that (3.28) can now be written as

$$\begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{1} & \mathbf{X} \\ \mathbf{W} & \mathbf{Y} \end{bmatrix} \begin{bmatrix} \mathbf{1} & -\mathbf{E} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{H}_1^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{F}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ -\mathbf{E}^H & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{X} \\ \mathbf{W} & \mathbf{Y} \end{bmatrix}^H$$

With the definition of

$$\begin{aligned} \mathbf{U} &= \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \mathbf{W} \end{bmatrix} = \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \\ \mathbf{W} \end{bmatrix} \\ \mathbf{V} &= \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \end{aligned} \quad (3.32)$$

this gives

$$\begin{bmatrix} \mathbf{U} & \mathbf{V} \end{bmatrix} \begin{bmatrix} \mathbf{1} & -\mathbf{E} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} (\mathbf{U}^H \mathbf{U})^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{F}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ -\mathbf{E}^H & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{U}^H \\ \mathbf{V}^H \end{bmatrix} = \mathbf{1}$$

and thus the constraint for \mathbf{X}, \mathbf{Y} appearing in \mathbf{V} can be written in the form of

$$\mathbf{1} - \mathbf{U}(\mathbf{U}^H \mathbf{U})^{-1} \mathbf{U}^H = (\mathbf{V} - \mathbf{U} \mathbf{E}) \mathbf{F}^{-1} (\mathbf{V} - \mathbf{U} \mathbf{E})^H \quad (3.33)$$

According to (A.80), there exists a \mathbf{U}_\perp with

$$\mathbf{U}_\perp (\mathbf{U}_\perp^H \mathbf{U}_\perp)^{-1} \mathbf{U}_\perp^H = \mathbf{1} - \mathbf{U}(\mathbf{U}^H \mathbf{U})^{-1} \mathbf{U}^H$$

since \mathbf{U} as defined in (3.32) has full rank. In order to find \mathbf{U}_\perp , calculate

$$\mathbf{1} - \mathbf{U}(\mathbf{U}^H \mathbf{U})^{-1} \mathbf{U}^H = \begin{bmatrix} \mathbf{1} - \mathbf{G}^{\frac{1}{2}} (\mathbf{U}^H \mathbf{U})^{-1} \mathbf{G}^{\frac{1}{2}} & -\mathbf{G}^{\frac{1}{2}} (\mathbf{U}^H \mathbf{U})^{-1} \mathbf{W}^H \\ -\mathbf{W} (\mathbf{U}^H \mathbf{U})^{-1} \mathbf{G}^{\frac{1}{2}} & \mathbf{1} - \mathbf{W} (\mathbf{U}^H \mathbf{U})^{-1} \mathbf{W}^H \end{bmatrix}$$

With

$$(\mathbf{U}^H \mathbf{U})^{-1} = (\mathbf{G} + \mathbf{W}^H \mathbf{W})^{-1} = \mathbf{G}^{-1} - \mathbf{G}^{-1} \mathbf{W}^H \mathbf{\Delta} \mathbf{W} \mathbf{G}^{-1}$$

and $\mathbf{\Delta} = (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-1}$ and further using (A.56) this can be written as

$$\begin{aligned} \mathbf{1} - \mathbf{U}(\mathbf{U}^H \mathbf{U})^{-1} \mathbf{U}^H &= \begin{bmatrix} \mathbf{G}^{-\frac{1}{2}} \mathbf{W}^H \mathbf{\Delta} \mathbf{W} \mathbf{G}^{-\frac{1}{2}} & -\mathbf{G}^{-\frac{1}{2}} \mathbf{W}^H \mathbf{\Delta} \\ -\mathbf{\Delta} \mathbf{W} \mathbf{G}^{-\frac{1}{2}} & \mathbf{\Delta} \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{G}^{-\frac{1}{2}} \mathbf{W}^H \\ -\mathbf{1} \end{bmatrix} \left(\begin{bmatrix} \mathbf{G}^{-\frac{1}{2}} \mathbf{W}^H \\ -\mathbf{1} \end{bmatrix}^H \begin{bmatrix} \mathbf{G}^{-\frac{1}{2}} \mathbf{W}^H \\ -\mathbf{1} \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{G}^{-\frac{1}{2}} \mathbf{W}^H \\ -\mathbf{1} \end{bmatrix}^H = \\ &= \mathbf{U}_\perp (\mathbf{U}_\perp^H \mathbf{U}_\perp)^{-1} \mathbf{U}_\perp^H \end{aligned}$$

so that

$$\mathbf{U}_\perp = \begin{bmatrix} \mathbf{G}^{-\frac{1}{2}} \mathbf{W}^H \\ -\mathbf{1} \end{bmatrix} \text{ with } \mathbf{U}^H \mathbf{U}_\perp = \begin{bmatrix} \mathbf{G}^{\frac{1}{2}} \\ \mathbf{W} \end{bmatrix}^H \begin{bmatrix} \mathbf{G}^{-\frac{1}{2}} \mathbf{W}^H \\ -\mathbf{1} \end{bmatrix} = \mathbf{0} \quad (3.34)$$

is an orthogonal completion of \mathbf{U} with full rank. Now (3.33) can be written as

$$\mathbf{U}_\perp (\mathbf{U}_\perp^H \mathbf{U}_\perp)^{-1} \mathbf{U}_\perp^H = (\mathbf{V} - \mathbf{U} \mathbf{E}) \mathbf{F}^{-1} (\mathbf{V} - \mathbf{U} \mathbf{E})^H$$

which according to (A.77) can be uniquely expanded to

$$\mathbf{U}_\perp (\mathbf{U}_\perp^H \mathbf{U}_\perp)^{-\frac{1}{2}} \left(\mathbf{U}_\perp (\mathbf{U}_\perp^H \mathbf{U}_\perp)^{-\frac{1}{2}} \right)^H = (\mathbf{V} - \mathbf{U} \mathbf{E}) \mathbf{F}^{-\frac{1}{2}} \left((\mathbf{V} - \mathbf{U} \mathbf{E}) \mathbf{F}^{-\frac{1}{2}} \right)^H$$

so that according to (A.78), \mathbf{V} containing the unknown \mathbf{X}, \mathbf{Y} can be extracted in the form of

$$\mathbf{V} = \mathbf{U} \mathbf{E} + \mathbf{U}_\perp (\mathbf{U}_\perp^H \mathbf{U}_\perp)^{-\frac{1}{2}} \mathbf{Z} \mathbf{F}^{\frac{1}{2}}$$

with unitary parameter \mathbf{Z} . With (3.32) and (3.34) this leads to

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \begin{bmatrix} \mathbf{1} \\ \mathbf{W} \end{bmatrix} \mathbf{E} + \begin{bmatrix} \mathbf{G}^{-1} \mathbf{W}^H \\ -\mathbf{1} \end{bmatrix} \mathbf{S}^H$$

with

$$\mathbf{S}^H = (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}} \mathbf{Z} \mathbf{F}^{\frac{1}{2}} \quad (3.35)$$

and further to

$$\begin{bmatrix} \mathbf{1} & \mathbf{X} \\ \mathbf{W} & \mathbf{Y} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{G}^{-1} \mathbf{W}^H \\ \mathbf{W} & -\mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{E} \\ \mathbf{0} & \mathbf{S}^H \end{bmatrix}$$

The solution (3.27) now becomes independent of \mathbf{X}, \mathbf{Y} and using (3.26) can be written as

$$\begin{aligned} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} &= \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-H} \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix}^{-H} \begin{bmatrix} \mathbf{1} & \mathbf{X} \\ \mathbf{W} & \mathbf{Y} \end{bmatrix}^H \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{G}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix} \begin{bmatrix} \mathbf{G}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1} & \mathbf{X} \\ \mathbf{W} & \mathbf{Y} \end{bmatrix}^H \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{A}_0 \mathbf{G}_0^{-1} & \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 \\ \mathbf{C}_0 \mathbf{G}_0^{-1} & \mathbf{D}_0 \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{E}^H & \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{G} & \mathbf{W}^H \\ \mathbf{W} & -\mathbf{1} \end{bmatrix} \end{aligned}$$

thus

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{G}^{-1}\mathbf{G}_0\mathbf{A}_0\mathbf{G}_0^{-1} + \mathbf{G}^{-1}\mathbf{G}_0\mathbf{B}_0\mathbf{E}^H & \mathbf{G}^{-1}\mathbf{G}_0\mathbf{B}_0\mathbf{S} \\ \mathbf{C}_0\mathbf{G}_0^{-1} + \mathbf{D}_0\mathbf{E}^H & \mathbf{D}_0\mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{G} & \mathbf{W}^H \\ \mathbf{W} & -\mathbf{1} \end{bmatrix} \quad (3.36)$$

The matrix \mathbf{E} can be calculated using (3.29) and (3.30) and gives

$$\mathbf{E} = \mathbf{H}_1^{-1}\mathbf{H}_2 = (\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1}(\mathbf{A}_0^H\mathbf{G}\mathbf{B}_0 + \mathbf{C}_0^H\mathbf{D}_0)$$

which can also be written as

$$\mathbf{E} = (\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1}\mathbf{A}_0^H(\mathbf{G} - \mathbf{G}_0)\mathbf{B}_0 \quad (3.37)$$

by taking the expression $\mathbf{C}_0^H\mathbf{D}_0$ from (3.26). Rewriting (3.26) in the form of

$$\begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix} \begin{bmatrix} \mathbf{G}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix}^H = \begin{bmatrix} \mathbf{G}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1}$$

provides the following equations

$$\begin{aligned} \mathbf{B}_0\mathbf{B}_0^H &= \mathbf{G}_0^{-1} - \mathbf{A}_0\mathbf{G}_0^{-1}\mathbf{A}_0^H \\ \mathbf{D}_0\mathbf{B}_0^H &= -\mathbf{C}_0\mathbf{G}_0^{-1}\mathbf{A}_0^H \end{aligned} \quad (3.38)$$

Due to (3.37) and (3.38), the matrix $\mathbf{G}_0\mathbf{B}_0\mathbf{E}^H$ appearing in the solution (3.36) can be expressed as

$$\begin{aligned} \mathbf{G}_0\mathbf{B}_0\mathbf{E}^H &= \mathbf{G}_0\mathbf{B}_0\mathbf{B}_0^H(\mathbf{G} - \mathbf{G}_0)\mathbf{A}_0(\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1} = \\ &= ((\mathbf{G} - \mathbf{G}_0)\mathbf{A}_0 - \mathbf{G}_0\mathbf{A}_0\mathbf{G}_0^{-1}\mathbf{A}_0^H(\mathbf{G} - \mathbf{G}_0)\mathbf{A}_0)(\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1} \end{aligned}$$

Using (3.23) this leads to

$$\mathbf{G}_0\mathbf{B}_0\mathbf{E}^H = \mathbf{G}\mathbf{A}_0(\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1} - \mathbf{G}_0\mathbf{A}_0\mathbf{G}_0^{-1} \quad (3.39)$$

Due to (3.37) and (3.38), the matrix $\mathbf{D}_0\mathbf{E}^H$ appearing in the solution (3.36) can be expressed as

$$\begin{aligned} \mathbf{D}_0\mathbf{E}^H &= \mathbf{D}_0\mathbf{B}_0^H(\mathbf{G} - \mathbf{G}_0)\mathbf{A}_0(\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1} = \\ &= \mathbf{C}_0\mathbf{G}_0^{-1}\mathbf{A}_0^H(\mathbf{G}_0 - \mathbf{G})\mathbf{A}_0(\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1} \end{aligned}$$

Using (3.23) this leads to

$$\mathbf{D}_0\mathbf{E}^H = \mathbf{C}_0(\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1} - \mathbf{C}_0\mathbf{G}_0^{-1} \quad (3.40)$$

Substituting the expressions (3.39), (3.40) in (3.36) leads to

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_0(\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1} & \mathbf{G}^{-1}\mathbf{G}_0\mathbf{B}_0\mathbf{S} \\ \mathbf{C}_0(\mathbf{G} + \mathbf{W}^H\mathbf{W})^{-1} & \mathbf{D}_0\mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{G} & \mathbf{W}^H \\ \mathbf{W} & -\mathbf{1} \end{bmatrix}$$

so that (3.24) follows. The matrix \mathbf{F} can be calculated by comparing the expressions for \mathbf{H}^{-1} obtained by inverting (3.29) and (3.31). The first one gives

$$\begin{aligned} \mathbf{H}^{-1} &= \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix}^{-H} = \\ &= \begin{bmatrix} \mathbf{G}_0^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix}^H \begin{bmatrix} \mathbf{G}_0\mathbf{G}^{-1}\mathbf{G}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 & \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix} \begin{bmatrix} \mathbf{G}_0^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \end{aligned}$$

using (3.26). The second one gives

$$\mathbf{H}^{-1} = \begin{bmatrix} \mathbf{1} & -\mathbf{E} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{H}_1^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{F}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ -\mathbf{E}^H & \mathbf{1} \end{bmatrix}$$

Comparing the lower right block provides

$$\mathbf{F} = (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0)^{-1}$$

which can be substituted in (3.35) so that

$$\mathbf{S}(\mathbf{B}_0, \mathbf{D}_0, \mathbf{W}, \mathbf{Z}) = (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0)^{-\frac{1}{2}} \mathbf{Z} (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}} \quad (3.41)$$

with unitary parameter \mathbf{Z} . This can also be written as

$$\begin{aligned} \mathbf{S}(\mathbf{B}_0, \mathbf{D}_0, \mathbf{W}, \mathbf{Z}) &= \\ &= \mathbf{Z} \mathbf{Z}^H (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0)^{-\frac{1}{2}} \mathbf{Z} (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}} = \\ &= \mathbf{Z} (\mathbf{Z}^H (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0) \mathbf{Z})^{-\frac{1}{2}} (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}} = \\ &= \mathbf{Z} \mathbf{S}(\mathbf{B}_0 \mathbf{Z}, \mathbf{D}_0 \mathbf{Z}, \mathbf{W}, \mathbf{1}) \end{aligned}$$

using (A.76). The solution (3.24) then gives

$$\begin{aligned} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} &= f(\mathbf{A}_0, \mathbf{B}_0, \mathbf{C}_0, \mathbf{D}_0, \mathbf{W}, \mathbf{Z}) = \\ &= \begin{bmatrix} \mathbf{A}_0 & \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix} \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{Z} \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{G} & \mathbf{W}^H \\ \mathbf{W} & -\mathbf{1} \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{A}_0 & \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 \mathbf{Z} \\ \mathbf{C}_0 & \mathbf{D}_0 \mathbf{Z} \end{bmatrix} \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{G} & \mathbf{W}^H \\ \mathbf{W} & -\mathbf{1} \end{bmatrix} = \\ &= f(\mathbf{A}_0, \mathbf{B}_0 \mathbf{Z}, \mathbf{C}_0, \mathbf{D}_0 \mathbf{Z}, \mathbf{W}, \mathbf{1}) \end{aligned}$$

with $\mathbf{S} = \mathbf{S}(\mathbf{B}_0 \mathbf{Z}, \mathbf{D}_0 \mathbf{Z}, \mathbf{W}, \mathbf{1})$ given by (3.41). According to (3.6), this shows that the unitary \mathbf{Z} parametrizes all lossless completion pairs $(\mathbf{B}_0 \mathbf{Z}, \mathbf{D}_0 \mathbf{Z})$ and thus the solution can be written in the form of (3.24) with \mathbf{S} given by (3.25) and $(\mathbf{B}_0, \mathbf{D}_0)$ being any lossless completion pair.

Remark Note that for a unitary matrix \mathbf{U} and according to (A.76) we can write

$$\begin{aligned} &(\mathbf{U}^H (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0) \mathbf{U})^{-\frac{1}{2}} (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}} = \\ &= \mathbf{U}^H (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0)^{-\frac{1}{2}} \mathbf{U} (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}} \mathbf{U}^H \mathbf{U} = \\ &= \mathbf{U}^H (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0)^{-\frac{1}{2}} (\mathbf{U} (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H) \mathbf{U}^H)^{-\frac{1}{2}} \mathbf{U} \end{aligned}$$

so that the functions $\mathbf{S}(\mathbf{B}_0, \mathbf{D}_0, \mathbf{W}), \mathbf{T}(\mathbf{W})$ given in (3.25) satisfy

$$\begin{aligned} \mathbf{S}(\mathbf{B}_0 \mathbf{U}, \mathbf{D}_0 \mathbf{U}, \mathbf{W}) &= \mathbf{U}^H \mathbf{S}(\mathbf{B}_0, \mathbf{D}_0, \mathbf{U} \mathbf{W}) \mathbf{U} \\ \mathbf{T}(\mathbf{W}) &= \mathbf{T}(\mathbf{U} \mathbf{W}) \end{aligned}$$

hence (3.24) satisfies

$$\begin{aligned} \begin{bmatrix} A(B_0U, D_0U, W) & B(B_0U, D_0U, W) \\ C(B_0U, D_0U, W) & D(B_0U, D_0U, W) \end{bmatrix} &= \\ &= \begin{bmatrix} A(B_0, D_0, WU) & B(B_0, D_0, WU)U \\ C(B_0, D_0, WU) & D(B_0, D_0, WU)U \end{bmatrix} \end{aligned} \quad (3.42)$$

3.2.3 Link to Conjugate Realization

Proposition If the observable pair (A_0, C_0) of the discrete stable minimal lossless system (A_0, B_0, C_0, D_0) has a matrix I that solves

$$\begin{bmatrix} A_0 \\ C_0 \end{bmatrix} I = \begin{bmatrix} I & 0 \\ 0 & 1 \end{bmatrix} \overline{\begin{bmatrix} A_0 \\ C_0 \end{bmatrix}} \quad (3.43)$$

then I is unique and invertible and the observable pair (A, C) of the discrete stable minimal lossless system (A, B, C, D) defined by (3.24) satisfies

$$\begin{bmatrix} A \\ C \end{bmatrix} I = \begin{bmatrix} I & 0 \\ 0 & 1 \end{bmatrix} \overline{\begin{bmatrix} A \\ C \end{bmatrix}} \Leftrightarrow \overline{W} = UWI \quad (3.44)$$

with unitary U defined by

$$U = \overline{B}^H I^H G B + \overline{D}^H D = \overline{B}_0^H I^H G_0 B_0 + \overline{D}_0^H D_0 \quad (3.45)$$

Proof First suppose there exists a matrix I solving (3.43) then it is unique and invertible due to (3.10). According to (3.11) the matrix I also solves

$$\begin{bmatrix} A_0 & B_0 \\ C_0 & D_0 \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & 1 \end{bmatrix} \overline{\begin{bmatrix} A_0 & B_0 \\ C_0 & D_0 \end{bmatrix}} \begin{bmatrix} I^{-1} & 0 \\ 0 & U \end{bmatrix} \quad (3.46)$$

with unitary U given by

$$U = \overline{B}_0^H I^H G_0 B_0 + \overline{D}_0^H D_0 \quad (3.47)$$

according to (3.12) and Gram observability matrix G_0 satisfying

$$\overline{G}_0 = I^H G_0 I \quad (3.48)$$

according to (3.13).

To verify sufficiency suppose there exists a matrix I' that solves

$$\begin{bmatrix} A \\ C \end{bmatrix} I' = \begin{bmatrix} I' & 0 \\ 0 & 1 \end{bmatrix} \overline{\begin{bmatrix} A \\ C \end{bmatrix}}$$

then it is unique and invertible due to (3.10). According to (3.11) the matrix I' also solves

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} I' & 0 \\ 0 & 1 \end{bmatrix} \overline{\begin{bmatrix} A & B \\ C & D \end{bmatrix}} \begin{bmatrix} I'^{-1} & 0 \\ 0 & U' \end{bmatrix} \quad (3.49)$$

with unitary U' given by

$$U' = \overline{B}^H I'^H G B + \overline{D}^H D \quad (3.50)$$

according to (3.12) and Gram observability matrix \mathbf{G} satisfying

$$\overline{\mathbf{G}} = \mathbf{I}'^H \mathbf{G} \mathbf{I}' \quad (3.51)$$

according to (3.13). The constraint (3.16) can be rearranged in the form of

$$\begin{bmatrix} \mathbf{1} \\ \mathbf{W} \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{C}_0 \end{bmatrix} = \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^H \begin{bmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_0 \\ \mathbf{C}_0 \end{bmatrix}$$

using (3.3) since $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ is discrete stable minimal lossless. This provides

$$\begin{aligned} \mathbf{G} &= \mathbf{A}^H \mathbf{G} \mathbf{A}_0 + \mathbf{C}^H \mathbf{C}_0 \\ \mathbf{W} &= \mathbf{B}^H \mathbf{G} \mathbf{A}_0 + \mathbf{D}^H \mathbf{C}_0 \end{aligned} \quad (3.52)$$

With (3.46), (3.48) and (3.49), (3.51) the conjugate of (3.52) provides

$$\begin{aligned} \mathbf{G} \mathbf{I}' \mathbf{I}^{-1} &= \mathbf{A}^H \mathbf{G} \mathbf{I}' \mathbf{I}^{-1} \mathbf{A}_0 + \mathbf{C}^H \mathbf{C}_0 \\ \mathbf{U}'^H \overline{\mathbf{W}} \mathbf{I}^{-1} &= \mathbf{B}^H \mathbf{G} \mathbf{I}' \mathbf{I}^{-1} \mathbf{A}_0 + \mathbf{D}^H \mathbf{C}_0 \end{aligned}$$

Subtracting this from (3.52) gives

$$\begin{aligned} \mathbf{A}^H (\mathbf{G} - \mathbf{G} \mathbf{I}' \mathbf{I}^{-1}) \mathbf{A}_0 - (\mathbf{G} - \mathbf{G} \mathbf{I}' \mathbf{I}^{-1}) &= \mathbf{0} \\ \mathbf{B}^H (\mathbf{G} - \mathbf{G} \mathbf{I}' \mathbf{I}^{-1}) \mathbf{A}_0 - (\mathbf{W} - \mathbf{U}'^H \overline{\mathbf{W}} \mathbf{I}^{-1}) &= \mathbf{0} \end{aligned}$$

The first equation is a Stein equation and according to (A.88) has a unique solution $\mathbf{G} = \mathbf{G} \mathbf{I}' \mathbf{I}^{-1}$ since \mathbf{A}, \mathbf{A}_0 are stable. Since $\mathbf{G} > 0$ this provides

$$\mathbf{I}' = \mathbf{I} \quad (3.53)$$

The second equation then provides

$$\overline{\mathbf{W}} = \mathbf{U}' \mathbf{W} \mathbf{I} \quad (3.54)$$

With (3.46), (3.48), (3.49), (3.51) and (3.53), (3.54) it follows from (3.25) that $\mathbf{T} = \overline{\mathbf{I} \mathbf{T} \mathbf{I}^H}$ and the solution (3.24) can be written as

$$\begin{aligned} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}} \begin{bmatrix} \mathbf{I}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}' \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{A}_0 & \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix} \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}'^H \overline{\mathbf{S}} \mathbf{U}' \end{bmatrix} \begin{bmatrix} \mathbf{G} & \mathbf{W}^H \\ \mathbf{W} & -\mathbf{1} \end{bmatrix} \end{aligned}$$

This can be rewritten in the form

$$\begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{U}'^H \overline{\mathbf{S}} \mathbf{U}' \end{bmatrix} = \begin{bmatrix} \mathbf{A}_0 & \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 \\ \mathbf{C}_0 & \mathbf{D}_0 \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{G} & \mathbf{W}^H \\ \mathbf{W} & -\mathbf{1} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix}$$

hence

$$\overline{\mathbf{S}} = \mathbf{U} \mathbf{S} \mathbf{U}'^H \quad (3.55)$$

With (3.46), (3.48), (3.49), (3.51) and (3.53), (3.54) the conjugate of \mathbf{S} in (3.25) gives

$$\overline{\mathbf{S}} = \mathbf{U} (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0)^{-\frac{1}{2}} \mathbf{U}'^H \mathbf{U}' (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}} \mathbf{U}'^H$$

using (A.76). With (3.55) and the fact that the matrices $(\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0)^{-\frac{1}{2}}, (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}}$ are invertible this provides $\mathbf{U}^H \mathbf{U}' = \mathbf{1}$ hence

$$\mathbf{U}' = \mathbf{U} \quad (3.56)$$

Thus (3.45) follows from (3.46) and (3.49) using (3.53), (3.56) and (3.44) follows from (3.54) using (3.56), which verifies sufficiency.

To verify necessity suppose that parameter \mathbf{W} satisfies

$$\overline{\mathbf{W}} = \mathbf{U} \mathbf{W} \mathbf{I} \quad (3.57)$$

with unitary \mathbf{U} given by (3.47). With (3.43) and (3.57) the conjugate of (3.15) provides

$$\mathbf{A}_0^H \mathbf{I}^{-H} \overline{\mathbf{G}} \mathbf{I}^{-1} \mathbf{A}_0 - \mathbf{I}^{-H} \overline{\mathbf{G}} \mathbf{I}^{-1} + \mathbf{C}_0^H \mathbf{C}_0 = \mathbf{W}^H \mathbf{W}$$

Subtracting this from (3.15) gives

$$\mathbf{A}_0^H (\mathbf{G} - \mathbf{I}^{-H} \overline{\mathbf{G}} \mathbf{I}^{-1}) \mathbf{A}_0 - (\mathbf{G} - \mathbf{I}^{-H} \overline{\mathbf{G}} \mathbf{I}^{-1}) = \mathbf{0}$$

This is a Stein equation and according to (A.88) has a unique solution

$$\overline{\mathbf{G}} = \mathbf{I}^H \mathbf{G} \mathbf{I} \quad (3.58)$$

since \mathbf{A}_0 is stable. With (3.46), (3.48) and (3.58) the definitions (3.25) provide

$$\begin{aligned} \overline{\mathbf{S}} &= \mathbf{U} \mathbf{S} \mathbf{U}^H \\ \overline{\mathbf{T}} &= \mathbf{I}^{-1} \mathbf{T} \mathbf{I}^{-H} \end{aligned} \quad (3.59)$$

using (A.76). With (3.46), (3.48) and (3.57), (3.58), (3.59) it can be seen that the solution (3.24) satisfies

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}} \begin{bmatrix} \mathbf{I}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{U} \end{bmatrix}$$

so that (3.44) follows and according to (3.12) and with (3.47) matrix \mathbf{U} is given by (3.45), which verifies necessity.

3.2.4 Parametrization of Stable Observable Pairs

Proposition The set of all stable observable pairs (\mathbf{A}, \mathbf{C}) with Gram observability matrix $\mathbf{G} > 0$ is a differentiable manifold and can be parametrized by a matrix \mathbf{W} of same dimension than \mathbf{C} in the way

$$\begin{aligned} \mathbf{A} &= \mathbf{A}_0 \mathbf{T} \mathbf{G} + \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 \mathbf{S} \mathbf{W} \\ \mathbf{C} &= \mathbf{C}_0 \mathbf{T} \mathbf{G} + \mathbf{D}_0 \mathbf{S} \mathbf{W} \end{aligned} \quad (3.60)$$

with

$$\begin{aligned} \mathbf{S} &= (\mathbf{B}_0^H \mathbf{G}_0 \mathbf{G}^{-1} \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0)^{-\frac{1}{2}} (\mathbf{1} + \mathbf{W} \mathbf{G}^{-1} \mathbf{W}^H)^{-\frac{1}{2}} \\ \mathbf{T} &= (\mathbf{G} + \mathbf{W}^H \mathbf{W})^{-1} \end{aligned} \quad (3.61)$$

while $(\mathbf{A}_0, \mathbf{C}_0)$ is a stable observable pair with Gram observability matrix \mathbf{G}_0 and lossless completion $(\mathbf{B}_0, \mathbf{D}_0)$ as in (3.6). The Gram observability matrices are uniquely determined by

$$\begin{aligned} \mathbf{A}_0^H \mathbf{G}_0 \mathbf{A}_0 - \mathbf{G}_0 + \mathbf{C}_0^H \mathbf{C}_0 &= \mathbf{0} \\ \mathbf{A}_0^H \mathbf{G} \mathbf{A}_0 - \mathbf{G} + \mathbf{C}_0^H \mathbf{C}_0 &= \mathbf{W}^H \mathbf{W} \end{aligned} \quad (3.62)$$

and $\mathbf{G} \leq \mathbf{G}_0$ must obey the condition

$$\mathbf{G} > \mathbf{0} \quad (3.63)$$

Proof This parametrization is a particular consequence of the solution to the Nudelman interpolation problem (3.24), (3.25). It further follows from (3.42) that any lossless completion pair $(\mathbf{B}_0, \mathbf{D}_0)$ can be chosen since they are all parametrized by $(\mathbf{B}_0 \mathbf{U}, \mathbf{D}_0 \mathbf{U})$ according to (3.6). The effect of a different completion pair is only that in order to reach the same pair (\mathbf{A}, \mathbf{C}) a different parameter \mathbf{W} has to be used.

Remark Moreover it has been shown in [2] that the set of lossless systems and the set of stable observable pairs are *differentiable manifolds*. A differentiable manifold is a smooth curved surface of a dimension d that locally resembles the well known straight euclidean vector space \mathbb{R}^d . This means that locally, calculations can be carried out in the usual straight euclidean space instead of the curved manifold. The local link between the euclidean space and the manifold is provided by a local *chart*. A complete collection of local charts is called an *atlas* and globally describes the curved manifold like a complete collection of planar road maps describe the curved surface of earth. Topologically, these charts are an *homeomorphism* between an open subset of the curved manifold and an open subset of the straight \mathbb{R}^d whereby an homeomorphism is a bicontinuous function, that is a continuous function which has a continuous inverse function. Such an homeomorphism allows for a point on the manifold to find the corresponding point on the local chart and the other way round.

In this context, the stable observable pair $(\mathbf{A}_0, \mathbf{C}_0)$ with Gram observability matrix \mathbf{G}_0 is the center of a chart since for $\mathbf{W} = \mathbf{0}$ we get $\mathbf{G} = \mathbf{G}_0$ from (3.62) and further $\mathbf{S} = \mathbf{1}, \mathbf{T} = \mathbf{G}^{-1}$ from (3.61) and (3.3) which states that $\mathbf{B}_0^H \mathbf{G}_0 \mathbf{B}_0 + \mathbf{D}_0^H \mathbf{D}_0 = \mathbf{1}$ so that by (3.60) we get $\mathbf{A} = \mathbf{A}_0, \mathbf{C} = \mathbf{C}_0$ for $\mathbf{W} = \mathbf{0}$. The stable observable pair (\mathbf{A}, \mathbf{C}) is displaced from the center as a continuous function of \mathbf{W} while observing the condition (3.63). If this condition is close to be violated, another chart of the atlas has to be selected, which is done by centering on the current pair $(\mathbf{A}'_0 = \mathbf{A}, \mathbf{C}'_0 = \mathbf{C})$ balancing this realization $(\mathbf{A}_0 = \mathbf{J} \mathbf{A}'_0 \mathbf{J}^{-1}, \mathbf{C}_0 = \mathbf{C}'_0 \mathbf{J}^{-1})$ with \mathbf{J} defined by $\mathbf{G} = \mathbf{J}^H \mathbf{J}$ so that $\mathbf{G}_0 = \mathbf{1}$ and finally resetting $\mathbf{W} = \mathbf{0}$. Balancing is necessary because otherwise the Gram observability matrix of the new center pair remains $\mathbf{G}_0 = \mathbf{G}$ which is a poor choice since this value triggered the changing of the chart by almost violating condition (3.63). Since condition (3.63) defines an open set, it is fulfilled on the center at $\mathbf{W} = \mathbf{0}$ by $\mathbf{G} = \mathbf{G}_0 > \mathbf{0}$ but it will also be fulfilled in some neighborhood of the center, so that this condition does not prevent this parametrization from becoming differentiable.

Chapter 4

Multiobjective Optimization

This chapter sets up a multiobjective optimization problem. For a given plant P , a controller K has to be found which is optimal with respect to a given cost function and within given constraints. These optimization objectives are formulated in terms of L1, L2/H2, H_∞ system norms. The advantage of using the L1, L2/H2, H_∞ norms is that the calculation of these norms can themselves be set up as standard semidefinite optimization problems which allows the local sensitivities and eventually the subgradient of the multiobjective optimization problem to be calculated. Formulating the optimization problem in terms of the parametrization found in chapters 2 and 3 then allows to move the parameter Q_L instead of K infinitesimally in the descending direction indicated by the subgradient. Once the optimization task is complete, the controller K can be regained from the parameter Q_L . This enables a local improvement of a provided initial controller K_0 to a better one K^* with respect to the cost function and within the constraints.

4.1 Problem Formulation

In this section, an extended version of the plant is constructed that incorporates various filters and the plant itself. Together with the controller to be designed, this leads to an extended feedback loop that provides more freedom in the formulation of the optimization objectives. The objectives can then be defined on certain channels and over certain frequencies of this system.

4.1.1 Extended Loops

The plant P that is to be controlled by the controller K is extended by filter systems P_{zw_i} , P_{yw_i} , P_{zu_i} to form the extended plant P_i , which together with the controller K forms the extended loop H_i as shown in figure 4.1. These filters can be used to select performance channels and to emphasize certain frequencies. The fact that these filters become a part of the system to be optimized makes the formulation of objectives easier. A bridge structure (C.56) has been chosen for the extended plant P_i to cover all possible filtering situations. Many different

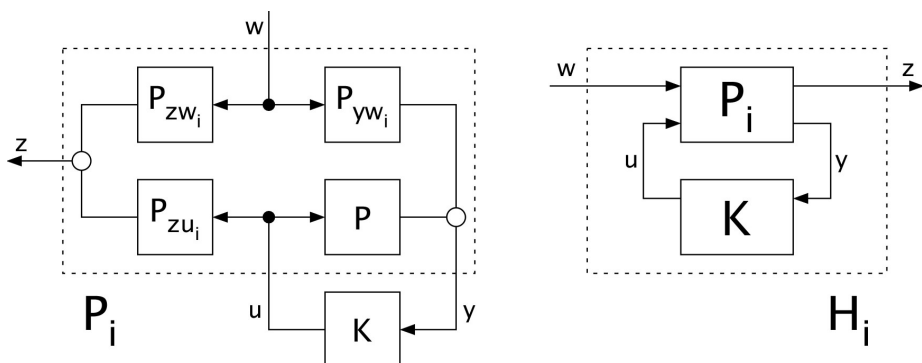


Figure 4.1: The plant P is extended by filters P_{zw_i} , P_{yw_i} , P_{zu_i} to form the extended plant P_i , which together with the controller K forms the extended loop H_i .

extended plants P_i and loops H_i can be defined and the index i is used to distinguish them.

Proposition The extended loop shown in figure 4.1 has a transfer matrix H_i given by

$$H_i = P_{zw_i} + P_{zu_i}(1 - KP)^{-1}KP_{yw_i} \quad (4.1)$$

A state space representation of H_i is given by

$$H_i = \left[\begin{array}{c|c} \mathbf{A}_i & \mathbf{B}_i \\ \hline \mathbf{C}_i & \mathbf{D}_i \end{array} \right] = \begin{array}{c} \left[\begin{array}{cccc} \mathbf{A}_{zw_i} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{yw_i} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{B}_{zu_i}\mathbf{D}_K\mathbf{E}_{PK}\mathbf{C}_{yw_i} & \mathbf{A}_{zu_i} & \mathbf{B}_{zu_i}\mathbf{D}_K\mathbf{E}_{PK}\mathbf{C}_P \\ \mathbf{0} & \mathbf{B}_P\mathbf{D}_K\mathbf{E}_{PK}\mathbf{C}_{yw_i} & \mathbf{0} & \mathbf{A}_P + \mathbf{B}_P\mathbf{D}_K\mathbf{E}_{PK}\mathbf{C}_P \end{array} \right] \dots \\ \left[\begin{array}{cc} \mathbf{0} & \mathbf{B}_{zw_i} \\ \mathbf{0} & \mathbf{B}_{yw_i} \\ \dots & \mathbf{B}_{zu_i}\mathbf{E}_{KP}\mathbf{C}_K \\ & \mathbf{B}_P\mathbf{E}_{KP}\mathbf{C}_K \\ \mathbf{A}_K + \mathbf{B}_K\mathbf{D}_P\mathbf{E}_{KP}\mathbf{C}_K & \mathbf{B}_K\mathbf{E}_{PK}\mathbf{D}_{yw_i} \\ \hline \mathbf{D}_{zu_i}\mathbf{E}_{KP}\mathbf{C}_K & \mathbf{D}_{zw_i} + \mathbf{D}_{zu_i}\mathbf{D}_K\mathbf{E}_{PK}\mathbf{D}_{yw_i} \end{array} \right] \end{array} \quad (4.2)$$

In the important special case of the filter system matrices satisfying

$$\begin{aligned} \mathbf{A}_{zw_i} &= \mathbf{A}_{yw_i} = \mathbf{A}_{zu_i} = \mathbf{A}_P \\ \mathbf{B}_{zw_i} &= \mathbf{B}_{yw_i} = \mathbf{B}_{w_i} \\ \mathbf{B}_{zu_i} &= \mathbf{B}_P \\ \mathbf{C}_{zw_i} &= \mathbf{C}_{zu_i} = \mathbf{C}_{z_i} \\ \mathbf{C}_{yw_i} &= \mathbf{C}_P \end{aligned} \quad (4.3)$$

this state space representation reduces to

$$\mathbf{H}_i = \left[\begin{array}{c|c} \mathbf{A}_i & \mathbf{B}_i \\ \hline \mathbf{C}_i & \mathbf{D}_i \end{array} \right] = \left[\begin{array}{c|c} \mathbf{A}_P + \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P & \dots \\ \mathbf{B}_K \mathbf{E}_{PK} \mathbf{C}_P & \dots \\ \hline \mathbf{C}_{z_i} + \mathbf{D}_{zu_i} \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P & \dots \\ \dots & \dots \\ \mathbf{B}_P \mathbf{E}_{KP} \mathbf{C}_K & \mathbf{B}_{w_i} + \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \mathbf{D}_{yw_i} \\ \mathbf{A}_K + \mathbf{B}_K \mathbf{D}_P \mathbf{E}_{KP} \mathbf{C}_K & \mathbf{B}_K \mathbf{E}_{PK} \mathbf{D}_{yw_i} \\ \hline \mathbf{D}_{zu_i} \mathbf{E}_{KP} \mathbf{C}_K & \mathbf{D}_{zw_i} + \mathbf{D}_{zu_i} \mathbf{D}_K \mathbf{E}_{PK} \mathbf{D}_{yw_i} \end{array} \right] \quad (4.4)$$

Proof The transfer matrix (4.1) follows from figure 4.1. Using \mathbf{H} as defined in (2.2) it can be seen that the transfer matrix can also be given by

$$\mathbf{H}_i = \mathbf{P}_{zw_i} + \mathbf{P}_{zu_i} (\mathbf{1} - \mathbf{K}\mathbf{P})^{-1} \mathbf{K}\mathbf{P}_{yw_i} = \mathbf{P}_{zw_i} + \mathbf{P}_{zu_i} \begin{bmatrix} \mathbf{0} & \mathbf{1} \end{bmatrix} \mathbf{H} \begin{bmatrix} \mathbf{0} \\ \mathbf{1} \end{bmatrix} \mathbf{P}_{yw_i}$$

Now let the state space representations of plant \mathbf{P} , controller \mathbf{K} and loop \mathbf{H} be given by (2.18) and (2.19) and the filters \mathbf{P}_{zw_i} , \mathbf{P}_{yw_i} , \mathbf{P}_{zu_i} be given by

$$\mathbf{P}_{kl_i} = \left[\begin{array}{c|c} \mathbf{A}_{kl_i} & \mathbf{B}_{kl_i} \\ \hline \mathbf{C}_{kl_i} & \mathbf{D}_{kl_i} \end{array} \right]$$

with $k \in \{y, z\}$ and $l \in \{u, w\}$. Then the state space representation (4.2) follows from this transfer matrix using (C.46), (C.44) and (2.19). In the special case given by (4.3) the state space representation (4.2) then further reduces to (4.4) by applying simple state transformations that allow to cut off uncontrollable and unobservable parts.

4.1.2 System Norms

In order to find an optimal controller, a method of comparing the effects of different controllers is required. One such method is to define norms on the input and output signals of the extended loop \mathbf{H}_i to measure its performance. The two most common signal norms are the L2 norm that measures the energy, hence a sort of average of the signal, and the L ∞ norm that measure the peak value of the signal. The following inequalities taken from [14] relate the norms of input and output

$$\begin{aligned} \|\mathbf{z}(t)\|_{L_2} &\leq \|\mathbf{H}_i(s)\|_{H_\infty} \|\mathbf{w}(t)\|_{L_2} \\ \|\mathbf{z}(t)\|_{L_\infty} &\leq \|\mathbf{H}_i(t)\|_{L_1} \|\mathbf{w}(t)\|_{L_\infty} \\ \|\mathbf{z}(t)\|_{L_\infty} &\leq \|\mathbf{H}_i(t)\|_{L_2} \|\mathbf{w}(t)\|_{L_2} = \|\mathbf{H}_i(s)\|_{H_2} \|\mathbf{w}(t)\|_{L_2} \end{aligned}$$

The correspondence between the L2 and H2 norm in the last inequality is due to the Parseval theorem (C.9). These inequalities show that \mathbf{H}_i needs to be stable for the L1, L2/H2, H ∞ norms to remain bounded. If the filters of the extended plant \mathbf{P}_i as in figure 4.1 are selected so that the input \mathbf{w}_i consists of external disturbances and the output \mathbf{z}_i contains error signals that are to be kept small, then the system norms L1, L2/H2, H ∞ can be used on \mathbf{H}_i to minimize or limit the effect of the disturbances on the error signals. This explains the interest to use the system norms L1, L2/H2, H ∞ as objectives in system specifications. These norms or best upper bounds on them are presented in the following in

the form of matrix inequalities for a discrete system \mathbf{H}_i . To our best knowledge there exists yet no matrix inequality form for the L1 norm. However, an upper bound, called the star norm, has been found by replacing the reachable set of the system state trajectories by the tightest ellipsoid that contains the reachable set, see [1], [9], [12] and section C.1.4. The intermediate variables $\alpha_i, \beta_i, \gamma_i, \mathbf{Z}_i$ are called decision variables. The variables α_i, β_i follow from the application of the so called S-lemma on quadratic constraints as demonstrated in [1] for continuous systems. The variables γ_i represent an upper bound on the norm and the hermitian and positive definite matrices \mathbf{Z}_i define Lyapunov functions or are Gram controllability or observability matrices.

L1 Norm An upper bound is the star norm derived in [1], [9] and leads to

$$\|\mathbf{H}_i[k]\|_{L1} \leq \|\mathbf{H}_i[k]\|_* \leq \gamma_i \quad (4.5)$$

if and only if $\exists \mathbf{Z}_i > 0, \alpha_i \in [0, 1], \beta_i \in [0, \gamma_i^2]$ so that

$$\begin{bmatrix} \mathbf{A}_i^H \mathbf{Z}_i \mathbf{A}_i + (\alpha_i - 1) \mathbf{Z}_i & \mathbf{A}_i^H \mathbf{Z}_i \mathbf{B}_i \\ \mathbf{B}_i^H \mathbf{Z}_i \mathbf{A}_i & \mathbf{B}_i^H \mathbf{Z}_i \mathbf{B}_i - \alpha_i \mathbf{1} \end{bmatrix} \leq 0 \quad (4.6)$$

$$\begin{bmatrix} \mathbf{C}_i^H \mathbf{C}_i - \beta_i \mathbf{Z}_i & \mathbf{C}_i^H \mathbf{D}_i \\ \mathbf{D}_i^H \mathbf{C}_i & \mathbf{D}_i^H \mathbf{D}_i + (\beta_i - \gamma_i^2) \mathbf{1} \end{bmatrix} \leq 0$$

L2/H2 Norm According to (C.12), (C.13) direct calculation leads to

$$\|\mathbf{H}_i[k]\|_{L2} = \|\mathbf{H}_i(z)\|_{H2} < \gamma_i \quad (4.7)$$

if and only if $\exists \mathbf{Z}_i > 0$ so that

$$\begin{aligned} \mathbf{A}_i^H \mathbf{Z}_i \mathbf{A}_i + \mathbf{C}_i^H \mathbf{C}_i - \mathbf{Z}_i &< 0 \\ \text{tr}(\mathbf{B}_i^H \mathbf{Z}_i \mathbf{B}_i + \mathbf{D}_i^H \mathbf{D}_i) - \gamma_i^2 &< 0 \end{aligned} \quad (4.8)$$

H ∞ Norm It follows from the strict bounded real lemma that

$$\|\mathbf{H}_i(z)\|_{H\infty} < \gamma_i \quad (4.9)$$

if and only if $\exists \mathbf{Z}_i > 0$ so that

$$\begin{bmatrix} \mathbf{A}_i^H \mathbf{Z}_i \mathbf{A}_i + \mathbf{C}_i^H \mathbf{C}_i - \mathbf{Z}_i & \mathbf{A}_i^H \mathbf{Z}_i \mathbf{B}_i + \mathbf{C}_i^H \mathbf{D}_i \\ \mathbf{B}_i^H \mathbf{Z}_i \mathbf{A}_i + \mathbf{D}_i^H \mathbf{C}_i & \mathbf{B}_i^H \mathbf{Z}_i \mathbf{B}_i + \mathbf{D}_i^H \mathbf{D}_i - \gamma_i^2 \mathbf{1} \end{bmatrix} < 0 \quad (4.10)$$

4.1.3 Optimization Problem

The typical optimization problem is formulated in terms of optimization objectives that have to be met. These objectives consist of a cost function and constraints. The constraints limit the set of all possible controllers to a subset of admissible controllers. Among these admissible controllers, the set of best controllers can then be found by minimizing the scalar cost function. Here, each objective is defined in terms of one of the norms L1, L2/H2, H ∞ presented in the previous section, acting on one extended loop \mathbf{H}_i . There can be as many different extended loops as objectives. Each objective comes with a set of decision variables $\alpha_i \geq 0, \beta_i \geq 0, \gamma_i > 0, \mathbf{Z}_i > 0$.

Proposition The multiobjective optimization problem of finding the controller \mathbf{K} that minimizes a cost function under constraints formulated in terms of the L1, L2/H2, H_∞ norms can be given in the general form

$$\mathbf{K}^* = \arg \inf_{\substack{\boldsymbol{\alpha} \geq 0 \\ \boldsymbol{\beta} \geq 0 \\ \gamma > 0 \\ \boldsymbol{\epsilon} > 0 \\ \mathbf{Z} > 0 \\ \mathbf{K}}} f(\boldsymbol{\gamma}) \text{ subject to } \mathbf{G}'(\boldsymbol{\alpha}, \boldsymbol{\beta}, \gamma, \boldsymbol{\epsilon}, \mathbf{Z}, \mathbf{K}) \leq 0 \quad (4.11)$$

The vectors $\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}$ contain the decision variables $\alpha_i, \beta_i, \gamma_i$ and the block diagonal matrix \mathbf{Z} contains the decision variables \mathbf{Z}_i of the different objectives. In addition to these, the vector $\boldsymbol{\epsilon}$ contains variables ϵ_i that are introduced to formally transform definite constraints of the form $\mathbf{M} < 0$ to semidefinite constraints by $\mathbf{M} + \epsilon \mathbf{1} \leq 0$ with $\epsilon > 0$.

Proof The decision variables γ_i in (4.5), (4.7), (4.9) represent upper bounds of the system norms, so the cost f will be formulated as a function of these. The constraints are then given by a combination of (4.5), (4.7), (4.9) together with user defined norm limits such as $\gamma \leq \gamma_{max}$. Aside from the decision variables $\alpha_i, \beta_i, \mathbf{Z}_i$, these constraints depend on the realizations $(\mathbf{A}_i, \mathbf{B}_i, \mathbf{C}_i, \mathbf{D}_i)$ of the extended loops \mathbf{H}_i which according to (4.1) are a function of \mathbf{K} . The ϵ_i are used to formally transform definite constraints $\mathbf{M}_i < 0$ into semidefinite constraints by $\mathbf{M}_i + \epsilon_i \mathbf{1} \leq 0$ with $\epsilon_i > 0$. All constraints can then be combined in the single unified semidefinite form $\mathbf{G}' \leq 0$ used in (4.11) by putting the separate semidefinite constraints as diagonal entries \mathbf{G}'_i into a block diagonal \mathbf{G}' . Then $\mathbf{G}' \leq 0$ if and only if all $\mathbf{G}'_i = [\mathbf{G}']_{ii} \leq 0$.

4.2 Parametrized Optimization

In theory, the solution of (4.11) yields the optimal controller \mathbf{K}^* but in practice, the problem is difficult to solve since the constraints are bilinear. The solution proposed in this section is therefore based on the idea of a local optimization. The fixed order parametrization developed in chapter 2 parametrizes all stable loops \mathbf{H} formed by the controller with a given plant \mathbf{P} with stable parameters $\mathbf{Q}_L, \mathbf{Q}_R$. The parameter \mathbf{Q}_L is observable if and only if \mathbf{K} is observable. It is therefore sensible to concentrate on stable observable \mathbf{Q}_L , so that the parametrization of stable observable pairs developed in chapter 3 can be applied. Since stable observable $\mathbf{Q} = \mathbf{Q}_L$ form a differentiable manifold, a subgradient can be calculated from a local sensitivity information so that a local optimization can be carried out. However, this requires the constraints to be formulated as linear matrix inequalities in order to get access to the sensitivity information. Therefore, a two step optimization is proposed. In an evaluation step the controller parameter \mathbf{Q} and the L1 parameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$ are kept constant so that the L1, L2/H2, H_∞ norm constraints become linear matrix inequalities in the remaining variables and the sensitivity information can be accessed. In an actualization step the local subgradient is calculated from the sensitivity matrix

so that the parameter values $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ can be changed in order to improve the cost function.

4.2.1 Passage between Controller and Parameter

Proposition If the Riccati equation

$$\begin{aligned} \mathbf{R}(\mathbf{A}_P + \mathbf{B}_P \mathbf{D}_K \mathbf{E}_{PK} \mathbf{C}_P) + \mathbf{R} \mathbf{B}_P \mathbf{E}_{KP} \mathbf{C}_K \mathbf{R} = \\ = \mathbf{B}_K \mathbf{E}_{PK} \mathbf{C}_P + (\mathbf{A}_K + \mathbf{B}_K \mathbf{D}_P \mathbf{E}_{KP} \mathbf{C}_K) \mathbf{R} \end{aligned} \quad (4.12)$$

has a solution $\mathbf{R} = \mathbf{R}_{KP}$ then the left parameter $\mathbf{Q} = \mathbf{Q}_L$ can be calculated from the controller \mathbf{K} by

$$\begin{aligned} \begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix} &= \begin{bmatrix} \mathbf{1} & \mathbf{L}_K \\ \mathbf{0} & \mathbf{E}_{KP} \end{bmatrix} \begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix} \\ \mathbf{L}_K &= (\mathbf{B}_K \mathbf{D}_P - \mathbf{R} \mathbf{B}_P) \mathbf{E}_{KP} \\ \mathbf{E}_{KP} &= (\mathbf{1} - \mathbf{D}_K \mathbf{D}_P)^{-1} \\ \mathbf{E}_{PK} &= (\mathbf{1} - \mathbf{D}_P \mathbf{D}_K)^{-1} \end{aligned} \quad (4.13)$$

If the Sylvester equation

$$\mathbf{R} \mathbf{A}_P = \mathbf{A}_Q \mathbf{R} + \mathbf{B}_Q \mathbf{C}_P \quad (4.14)$$

has a solution $\mathbf{R} = \mathbf{R}_{KP}$ then the controller \mathbf{K} can be calculated from the left parameter $\mathbf{Q} = \mathbf{Q}_L$ by

$$\begin{aligned} \begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix} &= \begin{bmatrix} \mathbf{1} & \mathbf{L}_K \\ \mathbf{0} & \mathbf{E}_{KP} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix} \\ \mathbf{L}_K &= \mathbf{B}_Q \mathbf{D}_P - \mathbf{R} \mathbf{B}_P \\ \mathbf{E}_{KP} &= \mathbf{1} + \mathbf{D}_Q \mathbf{D}_P \end{aligned} \quad (4.15)$$

Proof This is just a summary of (2.21), (2.23), (2.40) and (2.19) with $\mathbf{R} = \mathbf{R}_{KP}$.

Remark Note that the quadratic Riccati equation (4.12) becomes a linear Sylvester equation (4.14) if it is expressed in terms of the parametrization. The Riccati equation has to be solved at the start of the optimization when an initial \mathbf{Q}_0 has to be calculated from the initial \mathbf{K}_0 . The Sylvester equation has to be solved at the end of the optimization when a final \mathbf{K}^* has to be calculated from the final \mathbf{Q}^* .

4.2.2 Parametrized Extended Loops

Proposition A state space representation of the extended loop \mathbf{H}_i in terms of parameter $\mathbf{Q} = \mathbf{Q}_L$ is given by

$$\mathbf{H}_i = \left[\begin{array}{c|c} \mathbf{A}_i & \mathbf{B}_i \\ \hline \mathbf{C}_i & \mathbf{D}_i \end{array} \right] = \begin{array}{c} \left[\begin{array}{cccc} \mathbf{A}_{zw_i} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{yw_i} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{B}_{zu_i} \mathbf{D}_Q \mathbf{C}_{yw_i} & \mathbf{A}_{zu_i} & \mathbf{B}_{zu_i} \mathbf{D}_Q \mathbf{C}_P \\ \mathbf{0} & \mathbf{B}_P \mathbf{D}_Q \mathbf{C}_{yw_i} & \mathbf{0} & \mathbf{A}_P + \mathbf{B}_P \mathbf{D}_Q \mathbf{C}_P \\ \mathbf{0} & (\mathbf{B}_Q + \mathbf{R} \mathbf{B}_P \mathbf{D}_Q) \mathbf{C}_{yw_i} & \mathbf{0} & (\mathbf{B}_Q + \mathbf{R} \mathbf{B}_P \mathbf{D}_Q) \mathbf{C}_P \end{array} \right] \cdots \\ \hline \left[\begin{array}{cccc} \mathbf{C}_{zw_i} & \mathbf{D}_{zu_i} \mathbf{D}_Q \mathbf{C}_{yw_i} & \mathbf{C}_{zu_i} & \mathbf{D}_{zu_i} \mathbf{D}_Q \mathbf{C}_P \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \cdots & \mathbf{B}_{zu_i} \mathbf{C}_Q & \mathbf{B}_{zw_i} & \mathbf{B}_{yw_i} \\ & \mathbf{B}_P \mathbf{C}_Q & \mathbf{B}_{zu_i} \mathbf{D}_Q \mathbf{D}_{yw_i} & \mathbf{B}_P \mathbf{D}_Q \mathbf{D}_{yw_i} \\ & \mathbf{A}_Q + \mathbf{R} \mathbf{B}_P \mathbf{C}_Q & (\mathbf{B}_Q + \mathbf{R} \mathbf{B}_P \mathbf{D}_Q) \mathbf{D}_{yw_i} & (\mathbf{B}_Q + \mathbf{R} \mathbf{B}_P \mathbf{D}_Q) \mathbf{D}_{yw_i} \\ & \mathbf{D}_{zu_i} \mathbf{C}_Q & \mathbf{D}_{zw_i} + \mathbf{D}_{zu_i} \mathbf{D}_Q \mathbf{D}_{yw_i} & \mathbf{D}_{zw_i} + \mathbf{D}_{zu_i} \mathbf{D}_Q \mathbf{D}_{yw_i} \end{array} \right] \end{array} \quad (4.16)$$

In the special case (4.3) this state space representation reduces to

$$\mathbf{H}_i = \left[\begin{array}{c|c} \mathbf{A}_i & \mathbf{B}_i \\ \hline \mathbf{C}_i & \mathbf{D}_i \end{array} \right] = \left[\begin{array}{c|c} \mathbf{A}_P + \mathbf{B}_P \mathbf{D}_Q \mathbf{C}_P & \mathbf{0} \\ \hline (\mathbf{B}_Q + \mathbf{R} \mathbf{B}_P \mathbf{D}_Q) \mathbf{C}_P & \mathbf{0} \\ \mathbf{C}_{z_i} + \mathbf{D}_{zu_i} \mathbf{D}_Q \mathbf{C}_P & \mathbf{0} \\ \cdots & \mathbf{0} \\ \mathbf{B}_P \mathbf{C}_Q & \mathbf{B}_{w_i} + \mathbf{B}_P \mathbf{D}_Q \mathbf{D}_{yw_i} \\ \cdots & \mathbf{A}_Q + \mathbf{R} \mathbf{B}_P \mathbf{C}_Q & (\mathbf{B}_Q + \mathbf{R} \mathbf{B}_P \mathbf{D}_Q) \mathbf{D}_{yw_i} \\ \mathbf{D}_{zu_i} \mathbf{C}_Q & \mathbf{D}_{zw_i} + \mathbf{D}_{zu_i} \mathbf{D}_Q \mathbf{D}_{yw_i} \end{array} \right] \quad (4.17)$$

Proof The parametrized state space representations (4.16) and (4.17) respectively follow from (4.2) and (4.4) by replacing all terms containing state matrices of the controller \mathbf{K} according to (2.43) with $\mathbf{R} = \mathbf{R}_{KP}$.

4.2.3 Parametrized Optimization Problem

Proposition If condition (4.12) or (4.14) have a solution \mathbf{R} at all times during the optimization process, then the multiobjective optimization problem (4.11) is equivalent to the formulation in terms of the parameter \mathbf{Q} which can be given in the general form

$$\mathbf{Q}^* = \arg \inf_{\substack{\alpha \geq 0 \\ \beta \geq 0 \\ \gamma > 0 \\ \epsilon > 0 \\ \mathbf{Z} > 0 \\ \mathbf{Q}}} f(\gamma) \text{ subject to } \mathbf{G}(\alpha, \beta, \gamma, \epsilon, \mathbf{Z}, \mathbf{Q}) \leq 0 \quad (4.18)$$

Proof Problem (4.11) can be formulated in terms of \mathbf{Q} by formulating the extended loops \mathbf{H}_i as given in (4.2) in terms of \mathbf{Q} as given in (4.16). This leads to (4.18). If the conditions (4.12) or (4.14) have a solution \mathbf{R} at all times, then there is always a reversible relation between the controller \mathbf{K} and the parameter \mathbf{Q} expressed by (4.13) or (4.15) so that the problems (4.11) and (4.18) become equivalent.

4.2.4 Evaluation Step

However, despite being formulated in terms of the parameter \mathbf{Q} , problem (4.18) can not be solved directly with the currently known convex optimization techniques since these require the constraints formulated as linear matrix inequalities in all variables which is not the case as can be seen by (4.6), (4.8), (4.10). To circumvent this problem we propose an optimization consisting of two alternating steps, an evaluation and an actualization. The evaluation step calculates the local sensitivity matrix which is the optimal Lagrange multiplier associated to the constraint function for constant parameter values $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$.

Proposition For a given constant controller parameter value \mathbf{Q} and given constant L1 parameter values $\boldsymbol{\alpha}, \boldsymbol{\beta}$, the problem (4.18) reduces to an optimization problem for given constant extended loops \mathbf{H}_i . This means a search for the best upper norm bounds γ^* for a given system and may therefore be rather called a norm evaluation problem. If f is a cost function of the form $f(\boldsymbol{\gamma}) = \boldsymbol{\gamma}^T \mathbf{J} \boldsymbol{\gamma}$ with diagonal $\mathbf{J} > 0$ and if all objectives are formulated in terms of the L1, L2/H2, H_∞ norms this evaluation problem is convex with linear matrix inequality constraints. If this convex optimization problem exhibits strong duality then the local sensitivity matrix $\boldsymbol{\Omega}^*$ for the given parameter values $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ can be found and there is a relation in the form of

$$f(\boldsymbol{\gamma}^*) = g(\boldsymbol{\Omega}^*) = \inf_{\substack{\boldsymbol{\gamma} > 0 \\ \boldsymbol{\epsilon} > 0 \\ \mathbf{Z} > 0}} f(\boldsymbol{\gamma}) \text{ subject to } \mathbf{G}_{\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta}}(\boldsymbol{\gamma}, \boldsymbol{\epsilon}, \mathbf{Z}) \leq 0 \quad (4.19)$$

The function g is the Lagrange dual function to the cost function f as defined in (B.20). The sensitivity matrix $\boldsymbol{\Omega}$ is the Lagrange dual to the variable $\boldsymbol{\gamma}$ and is known as the Lagrange multiplier. The sensitivity matrix is hermitian. The evaluation constraint function $\mathbf{G}_{\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta}}$ equals the function \mathbf{G} from (4.18) but for $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ kept constant.

Proof For fixed parameters $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ and thus fixed $\mathbf{A}_i, \mathbf{B}_i, \mathbf{C}_i, \mathbf{D}_i$ the bilinear matrix inequalities in the L1, L2/H2, H_∞ norm objectives given in section 4.1.2 reduce to linear matrix inequalities in the variables γ_i^2, \mathbf{Z}_i . If the cost function takes the form $f(\boldsymbol{\gamma}) = \boldsymbol{\gamma}^T \mathbf{J} \boldsymbol{\gamma}$ with diagonal $\mathbf{J} > 0$ then $f(\boldsymbol{\gamma}) = \sum_i [\mathbf{J}]_{ii} \gamma_i^2$ becomes linear in γ_i^2 . Under these circumstances the optimization problem (4.18) consists of a linear cost function f in the variables γ_i^2 with linear matrix inequality constraints in the variables γ_i^2, \mathbf{Z}_i and therefore becomes a convex optimization problem. Most often, convex optimization problems exhibit strong duality, see [5] for extended information. The equality in (4.19) then follows if the optimization problem exhibits strong duality (B.22).

4.2.5 Actualization Step

The actualization step then calculates the local subgradient from the sensitivity matrix so that the parameter values $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ can be changed in order to improve the cost function.

Proposition If the optimization problem (4.19) exhibits strong duality and the local sensitivity matrix $\boldsymbol{\Omega}^*$ is known then the local subgradient of the cost function f can be calculated. This information can then be used to find the direction in which the parameter values $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ have to be actualized so that the cost function f decreases. The subgradient can be given by

$$df \leq \boldsymbol{\Omega}^* \bullet \delta \mathbf{G} \quad (4.20)$$

The scalar product \bullet of two matrices is as defined in (A.68).

Proof This is a result of the perturbation analysis of optimization problems that can be found in [4]. It is the generalization of a similar result in linear programming applied to a semidefinite programming. A simplified proof is given in section B.3. The key to understand these proofs is to note that the optimization problem (4.19) can be seen as an optimization problem with $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ as optimization parameters such as in (B.23). Then if the problem is convex and exhibits strong duality, the effects of a parameter perturbation on this optimization problem can be studied, which leads to (B.24) with Lagrange function defined as (B.25). In this context, the Lagrange function is given by

$$L(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\epsilon}, \mathbf{Z}, \boldsymbol{\Omega}, \mathbf{Q}) = f(\boldsymbol{\gamma}) + \boldsymbol{\Omega} \bullet \mathbf{G}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\epsilon}, \mathbf{Z}, \mathbf{Q})$$

with $\boldsymbol{\gamma}, \boldsymbol{\epsilon}, \mathbf{Z}$ forming the optimization variables, $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ the optimization parameters and $\boldsymbol{\Omega}$ the Lagrange multiplier. Then due to the linearity of the scalar product, the effect of parameter changes on the cost function takes the variational form

$$df \leq \delta L = \boldsymbol{\Omega}^* \bullet \delta \mathbf{G}$$

for a variation of the parameters $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$, which verifies (4.20).

Remark The total variation df of the cost function is defined as

$$df = f(\boldsymbol{\gamma}_+^*) - f(\boldsymbol{\gamma}^*)$$

and the partial variation $\delta \mathbf{G}$ of the constraint function is defined as

$$\delta \mathbf{G} = \mathbf{G}(\boldsymbol{\alpha}_+, \boldsymbol{\beta}_+, \boldsymbol{\gamma}^*, \boldsymbol{\epsilon}^*, \mathbf{Z}^*, \mathbf{Q}_+) - \mathbf{G}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}^*, \boldsymbol{\epsilon}^*, \mathbf{Z}^*, \mathbf{Q})$$

The values $\boldsymbol{\gamma}^*, \boldsymbol{\epsilon}^*, \mathbf{Z}^*$ are the optimal values of the decision variables obtained after the last evaluation step (4.19) for the current parameter values $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$. The value $\boldsymbol{\gamma}_+^*$ is an optimal value which would be obtained after the next evaluation step if the parameters were actualized to $\mathbf{Q}_+, \boldsymbol{\alpha}_+, \boldsymbol{\beta}_+$.

4.2.6 Variation of the Constraint Function

In order to obtain the subgradient in (4.20) the variation $\delta \mathbf{G}$ of the constraint function or at least the directional variation $\boldsymbol{\Omega}^* \bullet \delta \mathbf{G}$ has to be calculated.

Proposition If all objectives are formulated in terms of the L1, L2/H2, H_∞ norms then equation (4.20) can be written in the form

$$df \leq \operatorname{Re} \sum_i (\boldsymbol{\Omega}_{A_i} \bullet \delta \mathbf{A}_i + \boldsymbol{\Omega}_{B_i} \bullet \delta \mathbf{B}_i + \boldsymbol{\Omega}_{C_i} \bullet \delta \mathbf{C}_i + \boldsymbol{\Omega}_{D_i} \bullet \delta \mathbf{D}_i) + \boldsymbol{\Omega}_E \bullet \delta \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} \quad (4.21)$$

with $(\mathbf{A}_i, \mathbf{B}_i, \mathbf{C}_i, \mathbf{D}_i)$ being a realization of the extended loop \mathbf{H}_i and the sensibility matrices $\boldsymbol{\Omega}_{A_i}, \boldsymbol{\Omega}_{B_i}, \boldsymbol{\Omega}_{C_i}, \boldsymbol{\Omega}_{D_i}, \boldsymbol{\Omega}_E$ depending on the actual constraint function \mathbf{G} .

Proof Variables ϵ are used to formally transform definite constraints of the form $\mathbf{M} < 0$ into semidefinite constraints by $\mathbf{M} + \epsilon \mathbf{1} \leq 0$ with $\epsilon > 0$ and $\delta \epsilon = 0$ since ϵ does not depend on the optimization parameters $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$. Hence the directional variation of the definite constraint $\mathbf{M} < 0$ in semidefinite form is given by $\boldsymbol{\Omega} \bullet \delta(\mathbf{M} + \epsilon \mathbf{1}) = \boldsymbol{\Omega} \bullet \delta \mathbf{M}$. Then using (A.69), (A.70), (A.71), (A.72) and (A.81) the first order variations of the constraint function given in unified semidefinite form $\mathbf{G} \leq 0$ formulated in terms of the L1, L2/H2, H_∞ norm constraints can be calculated.

The first order directional variations of the L1 norm constraint (4.6) can then be given by

$$\begin{aligned} \boldsymbol{\Omega} \bullet \delta \mathbf{G}_{L1_i} &= \boldsymbol{\Omega}_1 \bullet \delta \begin{bmatrix} \mathbf{A}_i^H \mathbf{Z}_i \mathbf{A}_i + (\alpha_i - 1) \mathbf{Z}_i & \mathbf{A}_i^H \mathbf{Z}_i \mathbf{B}_i \\ \mathbf{B}_i^H \mathbf{Z}_i \mathbf{A}_i & \mathbf{B}_i^H \mathbf{Z}_i \mathbf{B}_i - \alpha_i \mathbf{1} \end{bmatrix} + \\ &+ \boldsymbol{\Omega}_2 \bullet \delta \begin{bmatrix} \mathbf{C}_i^H \mathbf{C}_i - \beta_i \mathbf{Z}_i & \mathbf{C}_i^H \mathbf{D}_i \\ \mathbf{D}_i^H \mathbf{C}_i & \mathbf{D}_i^H \mathbf{D}_i + (\beta_i - \gamma_i^2) \mathbf{1} \end{bmatrix} = \\ &= 2 \operatorname{Re}(\boldsymbol{\Omega}_1 \bullet ([\mathbf{A}_i \quad \mathbf{B}_i]^H \mathbf{Z}_i [\delta \mathbf{A}_i \quad \delta \mathbf{B}_i])) + \\ &+ 2 \operatorname{Re}(\boldsymbol{\Omega}_2 \bullet ([\mathbf{C}_i \quad \mathbf{D}_i]^H [\delta \mathbf{C}_i \quad \delta \mathbf{D}_i])) + \\ &+ \boldsymbol{\Omega}_1 \bullet \begin{bmatrix} \mathbf{Z}_i & \mathbf{0} \\ \mathbf{0} & -\mathbf{1} \end{bmatrix} \delta \alpha_i - \boldsymbol{\Omega}_2 \bullet \begin{bmatrix} \mathbf{Z}_i & \mathbf{0} \\ \mathbf{0} & -\mathbf{1} \end{bmatrix} \delta \beta_i \end{aligned}$$

which can be written in the form

$$\begin{aligned} \boldsymbol{\Omega} \bullet \delta \mathbf{G}_{L1_i} &= \\ &= \operatorname{Re}((2 \mathbf{Z}_i (\mathbf{A}_i \boldsymbol{\Omega}_{111} + \mathbf{B}_i \boldsymbol{\Omega}_{121})) \bullet \delta \mathbf{A}_i + (2 \mathbf{Z}_i (\mathbf{A}_i \boldsymbol{\Omega}_{112} + \mathbf{B}_i \boldsymbol{\Omega}_{122})) \bullet \delta \mathbf{B}_i + \\ &+ 2(\mathbf{C}_i \boldsymbol{\Omega}_{211} + \mathbf{D}_i \boldsymbol{\Omega}_{221}) \bullet \delta \mathbf{C}_i + 2(\mathbf{C}_i \boldsymbol{\Omega}_{212} + \mathbf{D}_i \boldsymbol{\Omega}_{222}) \bullet \delta \mathbf{D}_i) + \\ &+ \begin{bmatrix} \mathbf{Z}_i & \mathbf{0} \\ \mathbf{0} & -\mathbf{1} \\ -\mathbf{Z}_i & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \bullet \begin{bmatrix} \boldsymbol{\Omega}_1 \\ \boldsymbol{\Omega}_2 \end{bmatrix} \bullet \begin{bmatrix} \delta \alpha_i \\ \delta \beta_i \end{bmatrix} \quad (4.22) \end{aligned}$$

The first order directional variations of the L2/H2 norm constraint (4.8) can be given by

$$\begin{aligned} \boldsymbol{\Omega} \bullet \delta \mathbf{G}_{L2/H2_i} &= \boldsymbol{\Omega}_3 \bullet \delta(\mathbf{A}_i^H \mathbf{Z}_i \mathbf{A}_i + \mathbf{C}_i^H \mathbf{C}_i - \mathbf{Z}_i) + \\ &+ \boldsymbol{\Omega}_4 \bullet \delta(\operatorname{tr}(\mathbf{B}_i^H \mathbf{Z}_i \mathbf{B}_i + \mathbf{D}_i^H \mathbf{D}_i) - \gamma_i^2) = \\ &= 2 \operatorname{Re}(\boldsymbol{\Omega}_3 \bullet (\mathbf{A}_i^H \mathbf{Z}_i \delta \mathbf{A}_i + \mathbf{C}_i^H \delta \mathbf{C}_i)) + \\ &+ 2 \operatorname{Re}((\boldsymbol{\Omega}_4 \mathbf{1}) \bullet (\mathbf{B}_i^H \mathbf{Z}_i \delta \mathbf{B}_i + \mathbf{D}_i^H \delta \mathbf{D}_i)) \end{aligned}$$

which can be written in the form

$$\begin{aligned} \Omega \bullet \delta \mathbf{G}_{L2/H2_i} = & \operatorname{Re}((2\mathbf{Z}_i \mathbf{A}_i \Omega_3) \bullet \delta \mathbf{A}_i + (2\mathbf{Z}_i \mathbf{B}_i \Omega_4) \bullet \delta \mathbf{B}_i + \\ & + (2\mathbf{C}_i \Omega_3) \bullet \delta \mathbf{C}_i + (2\mathbf{D}_i \Omega_4) \bullet \delta \mathbf{D}_i) \end{aligned} \quad (4.23)$$

The first order directional variations of the H_∞ norm constraint (4.10) can be given by

$$\begin{aligned} \Omega \bullet \delta \mathbf{G}_{H_\infty_i} = & \Omega_5 \bullet \delta \begin{bmatrix} \mathbf{A}_i^H \mathbf{Z}_i \mathbf{A}_i + \mathbf{C}_i^H \mathbf{C}_i - \mathbf{Z}_i & \mathbf{A}_i^H \mathbf{Z}_i \mathbf{B}_i + \mathbf{C}_i^H \mathbf{D}_i \\ \mathbf{B}_i^H \mathbf{Z}_i \mathbf{A}_i + \mathbf{D}_i^H \mathbf{C}_i & \mathbf{B}_i^H \mathbf{Z}_i \mathbf{B}_i + \mathbf{D}_i^H \mathbf{D}_i - \gamma_i^2 \mathbf{1} \end{bmatrix} = \\ = & 2 \operatorname{Re}(\Omega_5 \bullet ([\mathbf{A}_i \ \mathbf{B}_i]^H \mathbf{Z}_i [\delta \mathbf{A}_i \ \delta \mathbf{B}_i] + [\mathbf{C}_i \ \mathbf{D}_i]^H [\delta \mathbf{C}_i \ \delta \mathbf{D}_i])) \end{aligned}$$

which can be written in the form

$$\begin{aligned} \Omega \bullet \delta \mathbf{G}_{H_\infty_i} = & \\ = & \operatorname{Re}((2\mathbf{Z}_i(\mathbf{A}_i \Omega_{511} + \mathbf{B}_i \Omega_{521})) \bullet \delta \mathbf{A}_i + (2\mathbf{Z}_i(\mathbf{A}_i \Omega_{512} + \mathbf{B}_i \Omega_{522})) \bullet \delta \mathbf{B}_i + \\ & + 2(\mathbf{C}_i \Omega_{511} + \mathbf{D}_i \Omega_{521}) \bullet \delta \mathbf{C}_i + 2(\mathbf{C}_i \Omega_{512} + \mathbf{D}_i \Omega_{522}) \bullet \delta \mathbf{D}_i) \end{aligned} \quad (4.24)$$

Due to the linearity of the scalar product, the directional variations of all L1, L2/H2, H_∞ norm constraints given in (4.22), (4.23), (4.24) and thus the directional variation of any constraint function \mathbf{G} that is formulated in terms of the L1, L2/H2, H_∞ norms can be written in the form (4.21).

4.2.7 Variation of the Extended Loops

In order to obtain the subgradient in (4.21) the directional variations $\Omega_{A_i} \bullet \delta \mathbf{A}_i$, $\Omega_{B_i} \bullet \delta \mathbf{B}_i$, $\Omega_{C_i} \bullet \delta \mathbf{C}_i$, $\Omega_{D_i} \bullet \delta \mathbf{D}_i$ of the extended loops have to be calculated.

Proposition If the Sylvester equations

$$\mathbf{F}_i \mathbf{A}_P^H - \mathbf{A}_Q^H \mathbf{F}_i = \mathbf{E}^H (\Omega_{A_i} \mathbf{E}_{R_i} + \Omega_{B_i} \mathbf{D}_{yw_i}^H \mathbf{D}_Q^H \mathbf{B}_P^H) \quad (4.25)$$

have solutions \mathbf{F}_i then equation (4.21) can be written in the form

$$df \leq \operatorname{Re}(\Omega_{A_Q} \bullet \delta \mathbf{A}_Q + \Omega_{B_Q} \bullet \delta \mathbf{B}_Q + \Omega_{C_Q} \bullet \delta \mathbf{C}_Q + \Omega_{D_Q} \bullet \delta \mathbf{D}_Q) + \Omega_E \bullet \delta \begin{bmatrix} \alpha \\ \beta \end{bmatrix} \quad (4.26)$$

with

$$\begin{aligned} \Omega_{A_Q} &= \sum_i (\mathbf{E}^H \Omega_{A_i} \mathbf{E} + \mathbf{F}_i \mathbf{R}^H) \\ \Omega_{B_Q} &= \sum_i (\mathbf{E}^H (\Omega_{A_i} \mathbf{E}_{B_i} + \Omega_{B_i} \mathbf{D}_{yw_i}^H) + \mathbf{F}_i \mathbf{C}_P^H) \\ \Omega_{C_Q} &= \sum_i (\mathbf{E}^H \Omega_{C_i} + \mathbf{D}_{zu_i}^H \Omega_{C_i}) \mathbf{E} \\ \Omega_{D_Q} &= \sum_i (\mathbf{E}^H (\Omega_{C_i} \mathbf{E}_{B_i} + \Omega_{B_i} \mathbf{D}_{yw_i}^H) + \mathbf{D}_{zu_i}^H (\Omega_{C_i} \mathbf{E}_{B_i} + \Omega_{D_i} \mathbf{D}_{yw_i}^H)) \end{aligned} \quad (4.27)$$

and

$$\mathbf{E} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{1} \end{bmatrix}, \mathbf{E}_{B_i} = \begin{bmatrix} \mathbf{0} \\ \mathbf{C}_{yw_i}^H \\ \mathbf{0} \\ \mathbf{C}_P^H \\ \mathbf{0} \end{bmatrix}, \mathbf{E}_{C_i} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{B}_{zu_i} \\ \mathbf{B}_P \\ \mathbf{RB}_P \end{bmatrix}, \mathbf{E}_{R_i} = (\mathbf{E}\mathbf{C}_Q^H + \mathbf{E}_{B_i}\mathbf{D}_Q^H)\mathbf{B}_P^H \quad (4.28)$$

which in the special case (4.3) reduce to $\mathbf{E}_{B_i} = \mathbf{E}_B, \mathbf{E}_{C_i} = \mathbf{E}_C, \mathbf{E}_{R_i} = \mathbf{E}_R$ with

$$\mathbf{E} = \begin{bmatrix} \mathbf{0} \\ \mathbf{1} \end{bmatrix}, \mathbf{E}_B = \begin{bmatrix} \mathbf{C}_P^H \\ \mathbf{0} \end{bmatrix}, \mathbf{E}_C = \begin{bmatrix} \mathbf{B}_P \\ \mathbf{RB}_P \end{bmatrix}, \mathbf{E}_R = (\mathbf{E}\mathbf{C}_Q^H + \mathbf{E}_B\mathbf{D}_Q^H)\mathbf{B}_P^H \quad (4.29)$$

Proof The first order variations of the extended loops can be calculated from (4.16) and can then be given by

$$\begin{aligned} \delta\mathbf{A}_i &= \mathbf{E}\delta\mathbf{A}_Q\mathbf{E}^H + \mathbf{E}\delta\mathbf{B}_Q\mathbf{E}_{B_i}^H + \mathbf{E}_{C_i}\delta\mathbf{C}_Q\mathbf{E}^H + \mathbf{E}_{C_i}\delta\mathbf{D}_Q\mathbf{E}_{B_i}^H + \mathbf{E}\delta\mathbf{R}\mathbf{E}_{R_i}^H \\ \delta\mathbf{B}_i &= (\mathbf{E}\delta\mathbf{B}_Q + \mathbf{E}_{C_i}\delta\mathbf{D}_Q + \mathbf{E}\delta\mathbf{R}\mathbf{B}_P\mathbf{D}_Q)\mathbf{D}_{yw_i} \\ \delta\mathbf{C}_i &= \mathbf{D}_{zu_i}(\delta\mathbf{C}_Q\mathbf{E}^H + \delta\mathbf{D}_Q\mathbf{E}_{B_i}^H) \\ \delta\mathbf{D}_i &= \mathbf{D}_{zu_i}\delta\mathbf{D}_Q\mathbf{D}_{yw_i} \end{aligned}$$

using (A.81) with $\mathbf{E}, \mathbf{E}_{B_i}, \mathbf{E}_{C_i}, \mathbf{E}_{R_i}$ as defined in (4.28). In the special case (4.3) the first order variations of the extended loops can be calculated from (4.17) in a similar way and are given by the same expressions but with $\mathbf{E}, \mathbf{E}_{B_i}, \mathbf{E}_{C_i}, \mathbf{E}_{R_i}$ as defined in (4.29). Then using (A.70) equation (4.21) can be written as

$$\begin{aligned} df \leq \operatorname{Re} \sum_i & ((\mathbf{E}^H \boldsymbol{\Omega}_{A_i} \mathbf{E}) \bullet \delta\mathbf{A}_Q + (\mathbf{E}^H (\boldsymbol{\Omega}_{A_i} \mathbf{E}_{B_i} + \boldsymbol{\Omega}_{B_i} \mathbf{D}_{yw_i}^H)) \bullet \delta\mathbf{B}_Q + \\ & + ((\mathbf{E}_{C_i}^H \boldsymbol{\Omega}_{A_i} + \mathbf{D}_{zu_i}^H \boldsymbol{\Omega}_{C_i}) \mathbf{E}) \bullet \delta\mathbf{C}_Q + \\ & + (\mathbf{E}_{C_i}^H (\boldsymbol{\Omega}_{A_i} \mathbf{E}_{B_i} + \boldsymbol{\Omega}_{B_i} \mathbf{D}_{yw_i}^H) + \mathbf{D}_{zu_i}^H (\boldsymbol{\Omega}_{C_i} \mathbf{E}_{B_i} + \boldsymbol{\Omega}_{D_i} \mathbf{D}_{yw_i}^H)) \bullet \delta\mathbf{D}_Q + \\ & + (\mathbf{E}^H (\boldsymbol{\Omega}_{A_i} \mathbf{E}_{R_i} + \boldsymbol{\Omega}_{B_i} \mathbf{D}_{yw_i}^H \mathbf{D}_Q^H \mathbf{B}_P^H)) \bullet \delta\mathbf{R}) + \boldsymbol{\Omega}_E \bullet \delta \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} \end{aligned} \quad (4.30)$$

The variable \mathbf{R} is constrained by the Sylvester equation (4.14) whose first order variation is given by

$$\delta\mathbf{R}\mathbf{A}_P = \delta\mathbf{A}_Q\mathbf{R} + \mathbf{A}_Q\delta\mathbf{R} + \delta\mathbf{B}_Q\mathbf{C}_P$$

so a directional variation for an arbitrary \mathbf{F}_i is given by

$$\mathbf{F}_i \bullet (\delta\mathbf{R}\mathbf{A}_P - \mathbf{A}_Q\delta\mathbf{R}) = \mathbf{F}_i \bullet (\delta\mathbf{A}_Q\mathbf{R} + \delta\mathbf{B}_Q\mathbf{C}_P)$$

which using (A.70) can be written in the form

$$(\mathbf{F}_i\mathbf{A}_P^H - \mathbf{A}_Q^H\mathbf{F}_i) \bullet \delta\mathbf{R} = (\mathbf{F}_i\mathbf{R}^H) \bullet \delta\mathbf{A}_Q + (\mathbf{F}_i\mathbf{C}_P^H) \bullet \delta\mathbf{B}_Q$$

Then if the \mathbf{F}_i solve the Sylvester equations (4.25) equation (4.30) can be written in the form (4.26) with $\boldsymbol{\Omega}_{A_Q}, \boldsymbol{\Omega}_{B_Q}, \boldsymbol{\Omega}_{C_Q}, \boldsymbol{\Omega}_{D_Q}$ given by (4.27).

4.2.8 Variation of the Parameter

In order to obtain the subgradient in (4.26) the directional variations $\Omega_{A_Q} \bullet \delta A_Q$, $\Omega_{B_Q} \bullet \delta B_Q$, $\Omega_{C_Q} \bullet \delta C_Q$, $\Omega_{D_Q} \bullet \delta D_Q$ of the parameter have to be calculated. This is done by applying the parametrization of stable observable pairs in (3.2.4) based on lossless systems. Thereby an infinitesimal variation $(\delta A_Q, \delta C_Q)$ can be calculated as a function of the variation δW of the stable observable pair parameter and thus an infinitesimal variation $\delta Q = (\delta A_Q, \delta B_Q, \delta C_Q, \delta D_Q)$ of the parameter Q and hence the directional variations required in (4.26) can be calculated as a function of δW and the variations $\delta B_Q, \delta D_Q$.

Proposition If the pair (A_{Q_0}, C_{Q_0}) defining the center of the current chart is stable and observable then it has a lossless completion (B_0, D_0) and the Stein and Sylvester equations

$$\begin{aligned}
A_{Q_0} L_A A_{Q_0}^H - L_A &= T A_{Q_0}^H \Omega_{A_Q} - A_{Q_0} T A_{Q_0}^H \Omega_{A_Q} G T A_{Q_0}^H - \\
&\quad - G^{-1} \Omega_{A_Q} W^H S^H B_0^H G_0 G^{-1} + \\
&\quad + G^{-1} G_0 B_0 M_A B_0^H G_0 G^{-1} + G^{-1} W^H N_A W G^{-1} \\
A_{Q_0} L_C A_{Q_0}^H - L_C &= T C_{Q_0}^H \Omega_{C_Q} - A_{Q_0} T C_{Q_0}^H \Omega_{C_Q} G T A_{Q_0}^H + \\
&\quad + G^{-1} G_0 B_0 M_C B_0^H G_0 G^{-1} + G^{-1} W^H N_C W G^{-1} \\
U^{\frac{1}{2}} M_A + M_A U^{\frac{1}{2}} &= V^{\frac{1}{2}} S^H B_0^H G_0 G^{-1} \Omega_{A_Q} W^H S^H \\
U^{\frac{1}{2}} M_C + M_C U^{\frac{1}{2}} &= V^{\frac{1}{2}} S^H D_0^H \Omega_{C_Q} W^H S^H \\
V^{\frac{1}{2}} N_A + N_A V^{\frac{1}{2}} &= S^H B_0^H G_0 G^{-1} \Omega_{A_Q} W^H S^H U^{\frac{1}{2}} \\
V^{\frac{1}{2}} N_C + N_C V^{\frac{1}{2}} &= S^H D_0^H \Omega_{C_Q} W^H S^H U^{\frac{1}{2}}
\end{aligned} \tag{4.31}$$

have unique solutions $L_A, L_C, M_A, M_C, N_A, N_C$ respectively and equation (4.26) can be written in the form

$$df \leq \text{Re}(\Omega_W \bullet \delta W + \Omega_{B_Q} \bullet \delta B_Q + \Omega_{D_Q} \bullet \delta D_Q) + \Omega_E \bullet \delta \begin{bmatrix} \alpha \\ \beta \end{bmatrix} \tag{4.32}$$

with

$$\begin{aligned}
\Omega_W &= S^H (B_0^H G_0 G^{-1} \Omega_{A_Q} + D_0^H \Omega_{C_Q}) + W (L_A + L_C + L_A^H + L_C^H) - \\
&\quad - (N_A + N_C + N_A^H + N_C^H) W G^{-1}
\end{aligned} \tag{4.33}$$

and

$$\begin{aligned}
S &= U^{-\frac{1}{2}} V^{-\frac{1}{2}} \\
T &= (G + W^H W)^{-1} \\
U &= B_0^H G_0 G^{-1} G_0 B_0 + D_0^H D_0 \\
V &= \mathbf{1} + W G^{-1} W^H
\end{aligned} \tag{4.34}$$

The Gram observability matrices G_0 of the chart center (A_{Q_0}, C_{Q_0}) and G of the current pair (A_Q, C_Q) are uniquely determined by the Stein equations

$$\begin{aligned}
A_{Q_0}^H G_0 A_{Q_0} - G_0 + C_{Q_0}^H C_{Q_0} &= \mathbf{0} \\
A_{Q_0}^H G A_{Q_0} - G + C_{Q_0}^H C_{Q_0} &= W^H W
\end{aligned} \tag{4.35}$$

respectively and $\mathbf{G} \leq \mathbf{G}_0$ must obey the condition

$$\mathbf{G} > 0 \quad (4.36)$$

otherwise the current chart is no longer valid and a new one has to be selected by centering on the balanced current pair ($\mathbf{A}_{Q_0} = \mathbf{J}\mathbf{A}_Q\mathbf{J}^{-1}$, $\mathbf{C}_{Q_0} = \mathbf{C}_Q\mathbf{J}^{-1}$) with \mathbf{J} defined by $\mathbf{G} = \mathbf{J}^H\mathbf{J}$ and resetting $\mathbf{W} = \mathbf{0}$ so that due to (4.35) there is $\mathbf{G} = \mathbf{G}_0 = \mathbf{1}$ and condition (4.36) again becomes valid.

Proof Equations (4.34), (4.35) and condition (4.36) are taken from (3.61), (3.62) and (3.63) respectively. Note that $\mathbf{G} > 0$ is required so that $\mathbf{T} > 0$, $\mathbf{V} > 0$ and since $(\mathbf{B}_0, \mathbf{D}_0)$ is a lossless completion also $\mathbf{U} > 0$. The first order variations of the stable observable pair $(\mathbf{A}_Q, \mathbf{C}_Q)$ can be calculated from (3.60) and can then be given by

$$\begin{aligned} \delta\mathbf{A}_Q &= \mathbf{A}_{Q_0}\delta\mathbf{T}\mathbf{G} + \mathbf{A}_{Q_0}\mathbf{T}\delta\mathbf{G} - \mathbf{G}^{-1}\delta\mathbf{G}\mathbf{G}^{-1}\mathbf{G}_0\mathbf{B}_0\mathbf{S}\mathbf{W} + \\ &\quad + \mathbf{G}^{-1}\mathbf{G}_0\mathbf{B}_0\delta\mathbf{S}\mathbf{W} + \mathbf{G}^{-1}\mathbf{G}_0\mathbf{B}_0\mathbf{S}\delta\mathbf{W} \\ \delta\mathbf{C}_Q &= \mathbf{C}_{Q_0}\delta\mathbf{T}\mathbf{G} + \mathbf{C}_{Q_0}\mathbf{T}\delta\mathbf{G} + \mathbf{D}_0\delta\mathbf{S}\mathbf{W} + \mathbf{D}_0\mathbf{S}\mathbf{W} \end{aligned}$$

using (A.81), (A.82). The first order variations of $\mathbf{G}, \mathbf{S}, \mathbf{T}, \mathbf{U}, \mathbf{V}$ can be calculated from (4.34) and (4.35). First

$$\mathbf{A}_{Q_0}^H\delta\mathbf{G}\mathbf{A}_{Q_0} = \delta\mathbf{G} + \delta\mathbf{W}^H\mathbf{W} + \mathbf{W}^H\delta\mathbf{W}$$

then since $\mathbf{S}^{-1} = \mathbf{V}^{\frac{1}{2}}\mathbf{U}^{\frac{1}{2}}$ we get

$$\delta\mathbf{S} + \mathbf{S}\delta\mathbf{V}^{\frac{1}{2}}\mathbf{U}^{\frac{1}{2}}\mathbf{S} + \mathbf{S}\mathbf{V}^{\frac{1}{2}}\delta\mathbf{U}^{\frac{1}{2}}\mathbf{S} = \mathbf{0}$$

and since $\mathbf{T}^{-1} = \mathbf{G} + \mathbf{W}^H\mathbf{W} = \mathbf{A}_{Q_0}^H\mathbf{G}\mathbf{A}_{Q_0} + \mathbf{C}_{Q_0}^H\mathbf{C}_{Q_0}$ we get

$$\delta\mathbf{T} + \mathbf{T}\mathbf{A}_{Q_0}^H\delta\mathbf{G}\mathbf{A}_{Q_0}\mathbf{T} = \mathbf{0}$$

and finally

$$\begin{aligned} \mathbf{U}^{\frac{1}{2}}\delta\mathbf{U}^{\frac{1}{2}} + \delta\mathbf{U}^{\frac{1}{2}}\mathbf{U}^{\frac{1}{2}} + \mathbf{B}_0^H\mathbf{G}_0\mathbf{G}^{-1}\delta\mathbf{G}\mathbf{G}^{-1}\mathbf{G}_0\mathbf{B}_0 &= \mathbf{0} \\ \mathbf{V}^{\frac{1}{2}}\delta\mathbf{V}^{\frac{1}{2}} + \delta\mathbf{V}^{\frac{1}{2}}\mathbf{V}^{\frac{1}{2}} + \mathbf{W}\mathbf{G}^{-1}\delta\mathbf{G}\mathbf{G}^{-1}\mathbf{W}^H &= \delta\mathbf{W}\mathbf{G}^{-1}\mathbf{W}^H + \mathbf{W}\mathbf{G}^{-1}\delta\mathbf{W}^H \end{aligned}$$

using (A.81), (A.82), (A.83). According to (A.77) the matrices $\mathbf{U}^{\frac{1}{2}} > 0$, $\mathbf{V}^{\frac{1}{2}} > 0$ are the unique positive matrix square roots of \mathbf{U}, \mathbf{V} . The first order directional

variations of these variables can then be given by

$$\begin{aligned} \operatorname{Re}(\Omega_{A_Q} \bullet \delta A_Q) &= \operatorname{Re}((T A_{Q_0}^H \Omega_{A_Q} - G^{-1} \Omega_{A_Q} W^H S^H B_0^H G_0 G^{-1}) \bullet \delta G + \\ &+ (B_0^H G_0 G^{-1} \Omega_{A_Q} W^H) \bullet \delta S + (A_{Q_0}^H \Omega_{A_Q} G) \bullet \delta T + \\ &+ (S^H B_0^H G_0 G^{-1} \Omega_{A_Q}) \bullet \delta W) \end{aligned} \quad (4.37)$$

$$\begin{aligned} \operatorname{Re}(\Omega_{C_Q} \bullet \delta C_Q) &= \operatorname{Re}((T C_{Q_0}^H \Omega_{C_Q}) \bullet \delta G + (D_0^H \Omega_{C_Q} W^H) \bullet \delta S + \\ &+ (C_{Q_0}^H \Omega_{C_Q} G) \bullet \delta T + (S^H D_0^H \Omega_{C_Q}) \bullet \delta W) \end{aligned} \quad (4.38)$$

$$\operatorname{Re}((A_{Q_0} L A_{Q_0}^H - L) \bullet \delta G) = \operatorname{Re}((W(L + L^H)) \bullet \delta W) \quad (4.39)$$

$$\operatorname{Re}(\Omega \bullet \delta S) = -\operatorname{Re}((V^{\frac{1}{2}} S^H \Omega S^H) \bullet \delta U^{\frac{1}{2}} + (S^H \Omega S^H U^{\frac{1}{2}}) \bullet \delta V^{\frac{1}{2}}) \quad (4.40)$$

$$\operatorname{Re}(\Omega \bullet \delta T) = -\operatorname{Re}((A_{Q_0} T \Omega T A_{Q_0}^H) \bullet \delta G) \quad (4.41)$$

$$\operatorname{Re}((U^{\frac{1}{2}} M + M U^{\frac{1}{2}}) \bullet \delta U^{\frac{1}{2}}) = -\operatorname{Re}((G^{-1} G_0 B_0 M B_0^H G_0 G^{-1}) \bullet \delta G) \quad (4.42)$$

$$\begin{aligned} \operatorname{Re}((V^{\frac{1}{2}} N + N V^{\frac{1}{2}}) \bullet \delta V^{\frac{1}{2}}) &= \operatorname{Re}(((N + N^H) W G^{-1}) \bullet \delta W - \\ &- (G^{-1} W^H N W G^{-1}) \bullet \delta G) \end{aligned} \quad (4.43)$$

using (A.69), (A.70) and particularly $\operatorname{Re}(A \bullet B) = \operatorname{Re}(\overline{B \bullet A}) = \operatorname{Re}(B \bullet A) = \operatorname{Re}(A^H \bullet B^H)$. Using (4.40), (4.41) on (4.37) leads to

$$\begin{aligned} \operatorname{Re}(\Omega_{A_Q} \bullet \delta A_Q) &= \operatorname{Re}((T A_{Q_0}^H \Omega_{A_Q} - G^{-1} \Omega_{A_Q} W^H S^H B_0^H G_0 G^{-1}) \bullet \delta G - \\ &- (A_{Q_0} T A_{Q_0}^H \Omega_{A_Q} G T A_{Q_0}^H) \bullet \delta G + (S^H B_0^H G_0 G^{-1} \Omega_{A_Q}) \bullet \delta W - \\ &- (V^{\frac{1}{2}} S^H B_0^H G_0 G^{-1} \Omega_{A_Q} W^H S^H) \bullet \delta U^{\frac{1}{2}} - \\ &- (S^H B_0^H G_0 G^{-1} \Omega_{A_Q} W^H S^H U^{\frac{1}{2}}) \bullet \delta V^{\frac{1}{2}}) \end{aligned}$$

According to (A.86) the Sylvester equations

$$\begin{aligned} U^{\frac{1}{2}} M_A + M_A U^{\frac{1}{2}} &= V^{\frac{1}{2}} S^H B_0^H G_0 G^{-1} \Omega_{A_Q} W^H S^H \\ V^{\frac{1}{2}} N_A + N_A V^{\frac{1}{2}} &= S^H B_0^H G_0 G^{-1} \Omega_{A_Q} W^H S^H U^{\frac{1}{2}} \end{aligned}$$

have unique solutions M_A, N_A since $U^{\frac{1}{2}} > 0, V^{\frac{1}{2}} > 0$. With (4.42) and (4.43) this leads to

$$\begin{aligned} \operatorname{Re}(\Omega_{A_Q} \bullet \delta A_Q) &= \operatorname{Re}((T A_{Q_0}^H \Omega_{A_Q} - G^{-1} \Omega_{A_Q} W^H S^H B_0^H G_0 G^{-1}) \bullet \delta G - \\ &- (A_{Q_0} T A_{Q_0}^H \Omega_{A_Q} G T A_{Q_0}^H) \bullet \delta G + (S^H B_0^H G_0 G^{-1} \Omega_{A_Q}) \bullet \delta W + \\ &+ (G^{-1} G_0 B_0 M_A B_0^H G_0 G^{-1} + G^{-1} W^H N_A W G^{-1}) \bullet \delta G - \\ &- ((N_A + N_A^H) W G^{-1}) \bullet \delta W) \end{aligned}$$

According to (A.88) the Stein equation

$$\begin{aligned} A_{Q_0} L_A A_{Q_0}^H - L_A &= T A_{Q_0}^H \Omega_{A_Q} - A_{Q_0} T A_{Q_0}^H \Omega_{A_Q} G T A_{Q_0}^H - \\ &- G^{-1} \Omega_{A_Q} W^H S^H B_0^H G_0 G^{-1} + \\ &+ G^{-1} G_0 B_0 M_A B_0^H G_0 G^{-1} + G^{-1} W^H N_A W G^{-1} \end{aligned}$$

has a unique solution L_A since A_{Q_0} is stable. With (4.39) this leads to

$$\begin{aligned} \operatorname{Re}(\Omega_{A_Q} \bullet \delta A_Q) &= \operatorname{Re}((S^H B_0^H G_0 G^{-1} \Omega_{A_Q} + W(L_A + L_A^H) - \\ &- (N_A + N_A^H) W G^{-1}) \bullet \delta W) \end{aligned} \quad (4.44)$$

Using (4.40), (4.41) on (4.38) leads to

$$\begin{aligned} \operatorname{Re}(\Omega_{C_Q} \bullet \delta C_Q) &= \operatorname{Re}((TC_{Q_0}^H \Omega_{C_Q}) \bullet \delta G - \\ &\quad - (A_{Q_0} TC_{Q_0}^H \Omega_{C_Q} G T A_{Q_0}^H) \bullet \delta G + (S^H D_0^H \Omega_{C_Q}) \bullet \delta W - \\ &\quad - (V^{\frac{1}{2}} S^H D_0^H \Omega_{C_Q} W^H S^H) \bullet \delta U^{\frac{1}{2}} - \\ &\quad - (S^H D_0^H \Omega_{C_Q} W^H S^H U^{\frac{1}{2}}) \bullet \delta V^{\frac{1}{2}}) \end{aligned}$$

According to (A.86) the Sylvester equations

$$\begin{aligned} U^{\frac{1}{2}} M_C + M_C U^{\frac{1}{2}} &= V^{\frac{1}{2}} S^H D_0^H \Omega_{C_Q} W^H S^H \\ V^{\frac{1}{2}} N_C + N_C V^{\frac{1}{2}} &= S^H D_0^H \Omega_{C_Q} W^H S^H U^{\frac{1}{2}} \end{aligned}$$

have unique solutions M_C, N_C since $U^{\frac{1}{2}} > 0, V^{\frac{1}{2}} > 0$. With (4.42) and (4.43) this leads to

$$\begin{aligned} \operatorname{Re}(\Omega_{C_Q} \bullet \delta C_Q) &= \operatorname{Re}((TC_{Q_0}^H \Omega_{C_Q}) \bullet \delta G - \\ &\quad - (A_{Q_0} TC_{Q_0}^H \Omega_{C_Q} G T A_{Q_0}^H) \bullet \delta G + (S^H D_0^H \Omega_{C_Q}) \bullet \delta W + \\ &\quad + (G^{-1} G_0 B_0 M_C B_0^H G_0 G^{-1} + G^{-1} W^H N_C W G^{-1}) \bullet \delta G - \\ &\quad - ((N_C + N_C^H) W G^{-1}) \bullet \delta W) \end{aligned}$$

According to (A.88) the Stein equation

$$\begin{aligned} A_{Q_0} L_C A_{Q_0}^H - L_C &= TC_{Q_0}^H \Omega_{C_Q} - A_{Q_0} TC_{Q_0}^H \Omega_{C_Q} G T A_{Q_0}^H + \\ &\quad + G^{-1} G_0 B_0 M_C B_0^H G_0 G^{-1} + G^{-1} W^H N_C W G^{-1} \end{aligned}$$

has a unique solution L_C since A_{Q_0} is stable. With (4.39) this leads to

$$\begin{aligned} \operatorname{Re}(\Omega_{C_Q} \bullet \delta C_Q) &= \operatorname{Re}((S^H D_0^H \Omega_{C_Q} + W(L_C + L_C^H) - \\ &\quad - (N_C + N_C^H) W G^{-1}) \bullet \delta W) \quad (4.45) \end{aligned}$$

The above Stein and Sylvester equations are collected in (4.31). Replacing the directional variations $\Omega_{A_Q} \bullet \delta A_Q$ and $\Omega_{C_Q} \bullet \delta C_Q$ in (4.26) by the terms in (4.44) and (4.45) respectively leads to the equation for the subgradient (4.32) with Ω_W given by (4.33).

4.2.9 Realness Constraint

To be practically usable, the parametrized controller K must remain real in the sense of (C.27) so that additional constraints have to be considered, according to the following idea, see also figure 4.2. If the plant P and the initial controller K_0 is real, then it is shown that a real initial parameter Q_0 can be found among the possible initial parameters Q_0 . It is then possible to keep the parameter Q real by satisfying some constraints. If the dynamic matrix A_P of the given plant P has no eigenvalues in common with the dynamic matrix A_Q of the current real parameter Q , then the reconstructed controller K is real.

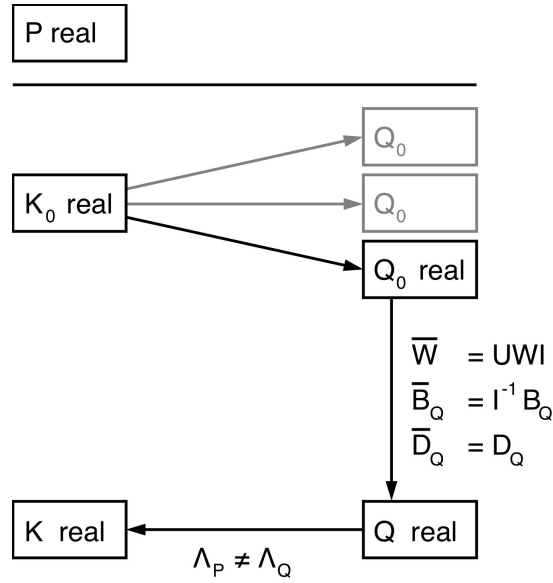


Figure 4.2: Principal idea behind the process to assure realness of the controller \mathbf{K} for a real plant \mathbf{P} . The spectra of the dynamic matrices $\mathbf{A}_P, \mathbf{A}_Q$ are denoted by Λ_P, Λ_Q respectively.

Proposition If the Riccati equation (4.12) has solutions and if the plant \mathbf{P} is real, then

$$\mathbf{K} \text{ real} \Rightarrow \text{exists a } \mathbf{Q} \text{ real} \quad (4.46)$$

and \mathbf{K}, \mathbf{Q} share the same realness transformation matrix \mathbf{I}

$$\begin{aligned} \begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix}} \\ \begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix}} \end{aligned} \quad (4.47)$$

Proof Plant \mathbf{P} and controller \mathbf{K} are real if and only if there exist \mathbf{I}_P, \mathbf{I} so that

$$\begin{aligned} \begin{bmatrix} \mathbf{A}_P & \mathbf{B}_P \\ \mathbf{C}_P & \mathbf{D}_P \end{bmatrix} &= \begin{bmatrix} \mathbf{I}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_P & \mathbf{B}_P \\ \mathbf{C}_P & \mathbf{D}_P \end{bmatrix}} \begin{bmatrix} \mathbf{I}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \\ \begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix}} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \end{aligned} \quad (4.48)$$

according to (C.28). Hence

$$\begin{bmatrix} \mathbf{A}_H & \mathbf{B}_H \\ \mathbf{C}_H & \mathbf{D}_H \end{bmatrix} = \begin{bmatrix} \mathbf{I}_H & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_H & \mathbf{B}_H \\ \mathbf{C}_H & \mathbf{D}_H \end{bmatrix}} \begin{bmatrix} \mathbf{I}_H & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \quad (4.49)$$

with $(\mathbf{A}_H, \mathbf{B}_H, \mathbf{C}_H, \mathbf{D}_H)$ being the realization of the loop \mathbf{H} as given in (2.19) and with

$$\mathbf{I}_H = \begin{bmatrix} \mathbf{I}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \quad (4.50)$$

which shows that the loop \mathbf{H} is real. With (4.48) and (4.50) and by using the substitutions in (4.13) equation (4.49) provides

$$\begin{aligned} \mathbf{A}_Q + \mathbf{R}\mathbf{B}_P\mathbf{C}_Q &= \mathbf{I}\overline{\mathbf{A}}_Q\mathbf{I}^{-1} + \mathbf{R}'\mathbf{B}_P\overline{\mathbf{C}}_Q\mathbf{I}^{-1} \\ \mathbf{B}_Q + \mathbf{R}\mathbf{B}_P\mathbf{D}_Q &= \mathbf{I}\overline{\mathbf{B}}_Q + \mathbf{R}'\mathbf{B}_P\overline{\mathbf{D}}_Q \\ \mathbf{C}_Q &= \overline{\mathbf{C}}_Q\mathbf{I}^{-1} \\ \mathbf{D}_Q &= \overline{\mathbf{D}}_Q \end{aligned}$$

with $\mathbf{R}' = \mathbf{I}\overline{\mathbf{R}}\mathbf{I}_P^{-1}$. This can be written in block matrix form

$$\begin{bmatrix} \mathbf{1} & (\mathbf{R} - \mathbf{R}')\mathbf{B}_P \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix}} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \quad (4.51)$$

Matrix \mathbf{R} is a solution of the Riccati equation (4.12) which can be written as the Sylvester equation (4.14) using the substitutions in (4.13). With (4.48) the conjugate of (4.14) provides

$$\mathbf{R}'\mathbf{A}_P = \mathbf{I}\overline{\mathbf{A}}_Q\mathbf{I}^{-1}\mathbf{R}' + \mathbf{I}\overline{\mathbf{B}}_Q\mathbf{C}_P$$

Subtracting this from (4.14) gives

$$(\mathbf{R} - \mathbf{R}')\mathbf{A}_P = \mathbf{A}_Q\mathbf{R} - \mathbf{I}\overline{\mathbf{A}}_Q\mathbf{I}^{-1}\mathbf{R}' + \mathbf{B}_Q\mathbf{C}_P - \mathbf{I}\overline{\mathbf{B}}_Q\mathbf{C}_P$$

Using (4.51) this simplifies to

$$\mathbf{A}_Q(\mathbf{R} - \mathbf{R}') = (\mathbf{R} - \mathbf{R}')(\mathbf{A}_P + \mathbf{B}_P(\mathbf{C}_Q\mathbf{R}' + \mathbf{D}_Q\mathbf{C}_P))$$

One solution of this Riccati equation is $\mathbf{R} = \mathbf{R}' = \mathbf{I}\overline{\mathbf{R}}\mathbf{I}_P^{-1}$. This solution leads to \mathbf{Q} real as can be seen in (4.51) with transformation matrix \mathbf{I} .

Proposition Realness of the parameter \mathbf{Q} is preserved

$$\mathbf{Q}_0 \text{ real} \Rightarrow \mathbf{Q} \text{ real} \quad (4.52)$$

if and only if the parameter \mathbf{Q} that is itself parametrized by $(\mathbf{W}, \mathbf{B}_Q, \mathbf{D}_Q)$ satisfies

$$\begin{aligned} \overline{\mathbf{W}} &= \mathbf{U}\mathbf{W}\mathbf{I} \\ \overline{\mathbf{B}}_Q &= \mathbf{I}^{-1}\mathbf{B}_Q \\ \overline{\mathbf{D}}_Q &= \mathbf{D}_Q \end{aligned} \quad (4.53)$$

with invertible \mathbf{I} uniquely defined by

$$\begin{bmatrix} \mathbf{A}_{Q_0} \\ \mathbf{C}_{Q_0} \end{bmatrix} \mathbf{I} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_{Q_0} \\ \mathbf{C}_{Q_0} \end{bmatrix}} \quad (4.54)$$

and unitary \mathbf{U} defined by

$$\mathbf{U} = \overline{\mathbf{B}}_0^H \mathbf{I}^H \mathbf{G}_0 \mathbf{B}_0 + \overline{\mathbf{D}}_0^H \mathbf{D}_0 \quad (4.55)$$

while the pair $(\mathbf{B}_0, \mathbf{D}_0)$ is a lossless completion (3.6) of the stable observable pair $(\mathbf{A}_{Q_0}, \mathbf{C}_{Q_0})$. Then \mathbf{Q}_0, \mathbf{Q} share the same realness transformation matrix \mathbf{I}

$$\begin{aligned} \begin{bmatrix} \mathbf{A}_{Q_0} & \mathbf{B}_{Q_0} \\ \mathbf{C}_{Q_0} & \mathbf{D}_{Q_0} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_{Q_0} & \mathbf{B}_{Q_0} \\ \mathbf{C}_{Q_0} & \mathbf{D}_{Q_0} \end{bmatrix}} \\ \begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix}} \end{aligned} \quad (4.56)$$

Proof Parameter \mathbf{Q}_0 realized by $(\mathbf{A}_{Q_0}, \mathbf{B}_{Q_0}, \mathbf{C}_{Q_0}, \mathbf{D}_{Q_0})$ is real if and only if there exists a matrix \mathbf{I} so that

$$\begin{bmatrix} \mathbf{A}_{Q_0} & \mathbf{B}_{Q_0} \\ \mathbf{C}_{Q_0} & \mathbf{D}_{Q_0} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_{Q_0} & \mathbf{B}_{Q_0} \\ \mathbf{C}_{Q_0} & \mathbf{D}_{Q_0} \end{bmatrix}} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1}$$

according to (C.28). The first block column in particular gives (4.54) hence condition (3.43) is satisfied and (3.44) can be applied so that

$$\begin{bmatrix} \mathbf{A}_Q \\ \mathbf{C}_Q \end{bmatrix} \mathbf{I} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_Q \\ \mathbf{C}_Q \end{bmatrix}} \Leftrightarrow \overline{\mathbf{W}} = \mathbf{U} \mathbf{W} \mathbf{I} \quad (4.57)$$

with unitary \mathbf{U} defined by (3.45). Now \mathbf{Q} is real if and only if there exists a matrix \mathbf{T} so that

$$\begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix} = \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix}} \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \quad (4.58)$$

according to (C.28). Since the pair $(\mathbf{A}_Q, \mathbf{C}_Q)$ is observable \mathbf{I} is uniquely defined by (4.57) according to (3.10). Comparing this to (4.58) shows that $\mathbf{T} = \mathbf{I}$ since \mathbf{I} is unique, so that realness of \mathbf{Q} follows from realness of \mathbf{Q}_0 if and only if the conditions (4.53) are satisfied.

Proposition If the dynamic matrix \mathbf{A}_P of the given plant \mathbf{P} has no eigenvalues in common with the dynamic matrix \mathbf{A}_Q of the current parameter \mathbf{Q} and if the plant \mathbf{P} is real, then

$$\mathbf{Q} \text{ real} \Rightarrow \mathbf{K} \text{ real} \quad (4.59)$$

and \mathbf{Q}, \mathbf{K} share the same realness transformation matrix \mathbf{I}

$$\begin{aligned} \begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix}} \\ \begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix}} \end{aligned} \quad (4.60)$$

Proof Plant \mathbf{P} and parameter \mathbf{Q} are real if and only if there exist \mathbf{I}_P, \mathbf{I} so that

$$\begin{aligned} \begin{bmatrix} \mathbf{A}_P & \mathbf{B}_P \\ \mathbf{C}_P & \mathbf{D}_P \end{bmatrix} &= \begin{bmatrix} \mathbf{I}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_P & \mathbf{B}_P \\ \mathbf{C}_P & \mathbf{D}_P \end{bmatrix}} \begin{bmatrix} \mathbf{I}_P & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \\ \begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix} &= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \overline{\begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix}} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \end{aligned} \quad (4.61)$$

according to (C.28). With (4.61) the conjugate of (4.14) provides

$$\overline{\mathbf{I}} \mathbf{I}_P^{-1} \mathbf{A}_P = \mathbf{A}_Q \overline{\mathbf{I}} \mathbf{I}_P^{-1} + \mathbf{B}_Q \mathbf{C}_P$$

Subtracting this from (4.14) gives

$$\mathbf{A}_Q (\mathbf{R} - \overline{\mathbf{I}} \mathbf{I}_P^{-1}) - (\mathbf{R} - \overline{\mathbf{I}} \mathbf{I}_P^{-1}) \mathbf{A}_P = \mathbf{0}$$

This is a Sylvester equation and according to (A.86) has a unique solution

$$\bar{\mathbf{R}} = \mathbf{I}^{-1} \mathbf{R} \mathbf{I}_P \quad (4.62)$$

if and only if $\lambda_{A_P} \neq \lambda_{A_Q}$ for all eigenvalues $\lambda_{A_P}, \lambda_{A_Q}$ of $\mathbf{A}_P, \mathbf{A}_Q$ respectively. With (4.61) and (4.62) it follows from (4.15) that

$$\bar{\mathbf{L}}_K = \mathbf{I}^{-1} \mathbf{L}_K \quad (4.63)$$

and with (4.63) further

$$\begin{bmatrix} \mathbf{A}_K & \mathbf{B}_K \\ \mathbf{C}_K & \mathbf{D}_K \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1} & \mathbf{L}_K \\ \mathbf{0} & \mathbf{E}_{KP} \end{bmatrix} \begin{bmatrix} \mathbf{A}_Q & \mathbf{B}_Q \\ \mathbf{C}_Q & \mathbf{D}_Q \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}$$

so that \mathbf{K} is real with transformation matrix \mathbf{I} .

4.3 Sufficient Conditions

The proposed multiobjective optimization method presented in the preceding sections roughly consists of first constructing a controller parameter \mathbf{Q} for a given initial controller \mathbf{K}_0 and constructing initial L1 norm parameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$, then calculating the local subgradient and thereby constantly improving the parameters $(\mathbf{Q}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ and finally reconstructing a final controller \mathbf{K}^* from the improved parameter \mathbf{Q} . While a detailed algorithm is presented in section 4.4, this section provides some sufficient conditions for these constructions, calculations and reconstructions to work.

4.3.1 Sufficient Condition for Parameter Construction

A necessary and sufficient condition for the ability to construct a parameter \mathbf{Q} from an initial controller \mathbf{K}_0 is the existence of a solution \mathbf{R} to the Riccati equation (4.12). A sufficient condition for the existence of \mathbf{R} can be given.

Proposition If the dynamic matrix \mathbf{A}_H of the closed loop \mathbf{H} formed by the given plant \mathbf{P} and a given controller \mathbf{K} is diagonalizable, then there exists a not necessarily unique \mathbf{R} so that a parameter \mathbf{Q} can be constructed from the given controller \mathbf{K} . If in addition to that \mathbf{H} is stable and \mathbf{K} is observable, then \mathbf{Q} is stable and observable.

Proof If \mathbf{A}_H is diagonalizable then according to (2.30) a not necessarily unique solution $\mathbf{R} = \mathbf{R}_{KP}$ to the Riccati equation (2.21), which is equal to the Riccati equation (4.12), exists. Then a parameter \mathbf{Q} can be constructed from controller \mathbf{K} by application of (4.13). If the closed loop \mathbf{H} is stable then \mathbf{K} is a stabilizing controller and thus \mathbf{Q} is stable due to (2.46). If \mathbf{K} is observable then $\mathbf{Q} = \mathbf{Q}_L$ is observable due to (2.45).

4.3.2 Sufficient Condition for Subgradient Calculation

A necessary and sufficient condition for the ability to write the subgradient (4.20) in the explicit form (4.32) is the existence of solutions \mathbf{F}_i to the Sylvester equations (4.25) and the existence of solutions $\mathbf{L}_A, \mathbf{L}_C, \mathbf{M}_A, \mathbf{M}_C, \mathbf{N}_A, \mathbf{N}_C$ to the Stein and Sylvester equations (4.31). A sufficient condition for the existence of $\mathbf{F}_i, \mathbf{L}_A, \mathbf{L}_C, \mathbf{M}_A, \mathbf{M}_C, \mathbf{N}_A, \mathbf{N}_C$ can be given.

Proposition If the dynamic matrix \mathbf{A}_P of the given plant P has no eigenvalues in common with the dynamic matrix \mathbf{A}_Q of the parameter Q and if Q is stable and observable, then there exist unique $\mathbf{F}_i, \mathbf{L}_A, \mathbf{L}_C, \mathbf{M}_A, \mathbf{M}_C, \mathbf{N}_A, \mathbf{N}_C$ and the subgradient of the cost function f can be given in the explicit form of (4.32) so that it can be calculated.

Proof The subgradient in (4.20) can be written in the form (4.21). Due to (A.86) the Sylvester equations (4.25) have unique solutions \mathbf{F}_i if and only if $\lambda_{A_P} \neq \lambda_{A_Q}$ for all eigenvalues $\lambda_{A_P}, \lambda_{A_Q}$ of $\mathbf{A}_P, \mathbf{A}_Q$ respectively. Then the subgradient (4.21) can be written in the form (4.26). If the pair $(\mathbf{A}_Q, \mathbf{C}_Q)$ of Q is stable and observable it can be set as the center of a chart $(\mathbf{A}_{Q_0} = \mathbf{A}_Q, \mathbf{C}_{Q_0} = \mathbf{C}_Q)$. Hence the Stein and Sylvester equations (4.31) have unique solutions and thus the subgradient (4.26) can be written in the form (4.32).

4.3.3 Sufficient Condition for Controller Reconstruction

A necessary and sufficient condition for the ability to reconstruct a final controller \mathbf{K}^* from a parameter Q is the existence of a solution \mathbf{R} to the Sylvester equation (4.14). A sufficient condition for the existence of \mathbf{R} can be given.

Proposition If the dynamic matrix \mathbf{A}_P of the given plant P has no eigenvalues in common with the dynamic matrix \mathbf{A}_Q of the current parameter Q , then there exists a unique \mathbf{R} so that a controller \mathbf{K} can be uniquely reconstructed from the current parameter Q .

Proof Due to (A.86) the Sylvester equation (4.14) has a unique solution \mathbf{R} if and only if $\lambda_{A_P} \neq \lambda_{A_Q}$ for all eigenvalues $\lambda_{A_P}, \lambda_{A_Q}$ of $\mathbf{A}_P, \mathbf{A}_Q$ respectively. Then a controller \mathbf{K} can be uniquely reconstructed from parameter Q by application of (4.15).

4.3.4 Sufficient Condition for L1 Parameter Construction

In order to construct initial L1 parameters α, β for extended loops \mathbf{H}_i controlled by an initial controller \mathbf{K}_0 , upper bounds of the L1 norms of the extended loops \mathbf{H}_i have to be found by solving the bilinear matrix inequalities in (4.6) which is difficult. A method for iteratively solving the bilinear matrix inequalities is to alternately fix α_i, β_i and \mathbf{Z}_i and to solve the resulting linear matrix inequalities at each step in order to decrease the upper bound γ_i . That way the linear matrix inequalities for fixed α_i, β_i are solved for \mathbf{Z}_i and provide a potentially better \mathbf{Z}_i for the following linear matrix inequality for fixed \mathbf{Z}_i and so on. A sufficient condition for the existence of initial $\alpha_i, \beta_i, \mathbf{Z}_i$ can be given so that the iterative process can be started, eventually resulting in the initial L1 parameters α, β .

Proposition If the extended loop \mathbf{H}_i is stable, then valid initial $\alpha_i, \beta_i, \mathbf{Z}_i$ can be derived from the spectral radius ρ_i of the extended loop dynamic matrix \mathbf{A}_i by selecting an $\alpha_i \in]0, 1 - \rho_i^2[$, solving the first linear matrix inequality in (4.6) for \mathbf{Z}_i and then solving the second linear matrix inequality in (4.6) for β_i . The spectral radii ρ_i can be collected in a radius vector ρ .

Proof The proof follows the ideas presented in [1] applied to the discrete case. If \mathbf{A}_i is stable then it has a spectral radius bounded by $\rho_i = \max_k |\lambda_k| < 1$ with λ_k denoting the eigenvalues of \mathbf{A}_i . Now select $\alpha_i \in]0, 1 - \rho_i^2[$ then $\frac{1}{\sqrt{1-\alpha_i}}\mathbf{A}_i$ is still stable with spectral radius $\frac{1}{\sqrt{1-\alpha_i}}\rho_i < 1$. Since $\frac{1}{\sqrt{1-\alpha_i}}\mathbf{A}_i$ is stable the Stein equation in \mathbf{Z}_i^{-1}

$$\frac{1}{1-\alpha_i}\mathbf{A}_i\mathbf{Z}_i^{-1}\mathbf{A}_i^H - \mathbf{Z}_i^{-1} + \frac{1}{\alpha_i}\mathbf{B}_i\mathbf{B}_i^H + \mathbf{E}_i = \mathbf{0}$$

with $\mathbf{E}_i > 0$ so that $\frac{1}{\alpha_i}\mathbf{B}_i\mathbf{B}_i^H + \mathbf{E}_i > 0$ has a unique solution $\mathbf{Z}_i^{-1} > 0$ according to (A.89). Thus there always exists a $\mathbf{Z}_i^{-1} > 0$ for any $\alpha_i \in]0, 1 - \rho_i^2[$ so that

$$\frac{1}{1-\alpha_i}\mathbf{A}_i\mathbf{Z}_i^{-1}\mathbf{A}_i^H - \mathbf{Z}_i^{-1} + \frac{1}{\alpha_i}\mathbf{B}_i\mathbf{B}_i^H < 0$$

which is equivalent to

$$\mathbf{Z}_i^{-1} - [\mathbf{A}_i \quad \mathbf{B}_i] \begin{bmatrix} \frac{1}{1-\alpha_i}\mathbf{Z}_i^{-1} & \mathbf{0} \\ \mathbf{0} & \frac{1}{\alpha_i}\mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A}_i^H \\ \mathbf{B}_i^H \end{bmatrix} > 0$$

Using (A.65) this is equivalent to

$$\begin{bmatrix} \mathbf{Z}_i^{-1} & \mathbf{A}_i & \mathbf{B}_i \\ \mathbf{A}_i^H & (1-\alpha_i)\mathbf{Z}_i & \mathbf{0} \\ \mathbf{B}_i^H & \mathbf{0} & \alpha_i\mathbf{1} \end{bmatrix} > 0$$

Using (A.64) this is equivalent to

$$\begin{bmatrix} \mathbf{Z}_i & \mathbf{Z}_i\mathbf{A}_i & \mathbf{Z}_i\mathbf{B}_i \\ \mathbf{A}_i^H\mathbf{Z}_i & (1-\alpha_i)\mathbf{Z}_i & \mathbf{0} \\ \mathbf{B}_i^H\mathbf{Z}_i & \mathbf{0} & \alpha_i\mathbf{1} \end{bmatrix} > 0$$

Using (A.65) this is equivalent to

$$\begin{bmatrix} (1-\alpha_i)\mathbf{Z}_i & \mathbf{0} \\ \mathbf{0} & \alpha_i\mathbf{1} \end{bmatrix} - \begin{bmatrix} \mathbf{A}_i^H\mathbf{Z}_i \\ \mathbf{B}_i^H\mathbf{Z}_i \end{bmatrix} \mathbf{Z}_i^{-1} [\mathbf{Z}_i\mathbf{A}_i \quad \mathbf{Z}_i\mathbf{B}_i] > 0$$

This finally gives

$$\begin{bmatrix} \mathbf{A}_i^H\mathbf{Z}_i\mathbf{A}_i + (\alpha_i - 1)\mathbf{Z}_i & \mathbf{A}_i^H\mathbf{Z}_i\mathbf{B}_i \\ \mathbf{B}_i^H\mathbf{Z}_i\mathbf{A}_i & \mathbf{B}_i^H\mathbf{Z}_i\mathbf{B}_i - \alpha_i\mathbf{1} \end{bmatrix} < 0$$

which is a slightly stricter version of the first matrix inequality in (4.6). This means that for any $\alpha_i \in]0, 1 - \rho_i^2[$ this linear matrix inequality in \mathbf{Z}_i has at least one solution $\mathbf{Z}_i > 0$ which can then be put as a constant matrix into the second matrix inequality in (4.6) that then becomes a linear matrix inequality in β_i .

4.4 Complete Algorithm

This section summarizes the steps to carry out a multiobjective optimization based on the parametrization of fixed order controllers and the parametrization of observable pairs with objectives formulated in terms of the L1, L2/H2, H_∞ norms. The proposed algorithm consists of the three stages initialization, optimization, finalization. An overview of this algorithm is presented in figure 4.3 while a more detailed description is provided below.

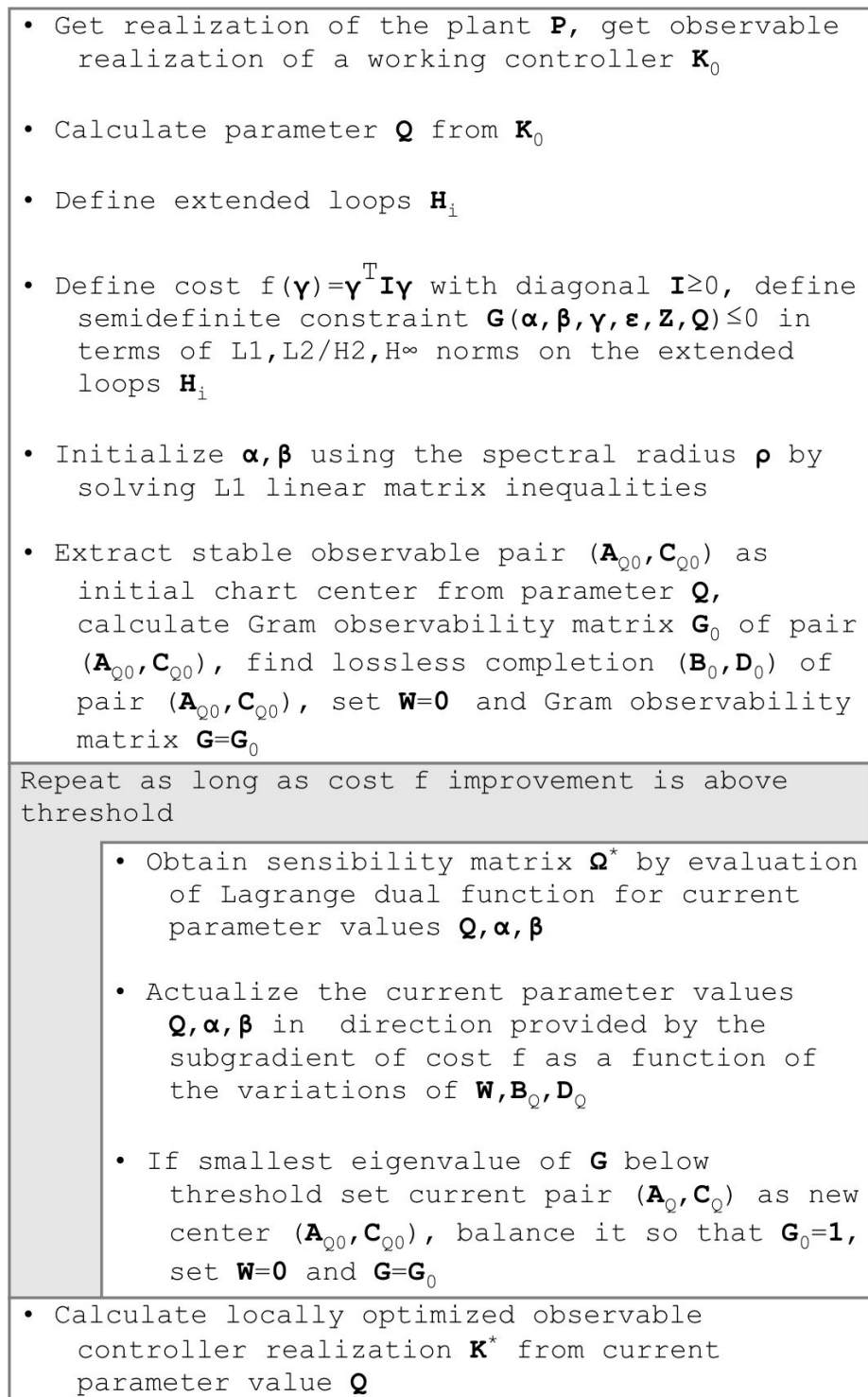


Figure 4.3: Overview of the optimization algorithm.

4.4.1 Initialization

In the initialization stage, the optimization problem is set up for a given plant \mathbf{P} that is to be optimally controlled by a controller \mathbf{K} in a feedback loop \mathbf{H} . The optimization objectives are defined in form of a scalar cost function and a negative semidefinite matrix constraint function. An initial parameter \mathbf{Q} is constructed from an admissible controller \mathbf{K}_0 .

- Get a realization of the plant \mathbf{P} and an observable realization of an initial working controller \mathbf{K}_0 .
- Define filters $\mathbf{P}_{zw_i}, \mathbf{P}_{yw_i}, \mathbf{P}_{zu_i}$ and thereby formally construct the extended loops \mathbf{H}_i as defined in section 4.1.1 for each optimization objective which allows to refine the objectives by selecting certain performance channels or by emphasizing certain frequencies.
- Define the optimization objectives in form of a scalar cost function $f(\boldsymbol{\gamma})$ of the form $f(\boldsymbol{\gamma}) = \boldsymbol{\gamma}^T \mathbf{I} \boldsymbol{\gamma}$ with diagonal $\mathbf{I} > 0$ and a constraint function $\mathbf{G}'(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\epsilon}, \mathbf{Z}, \mathbf{K}) \leq 0$ formulated in terms of the L1, L2/H2, H_∞ norms given in section 4.1.2 while assuring the initial working controller \mathbf{K}_0 to be admissible. It is admissible if there exist $\boldsymbol{\alpha} \geq 0, \boldsymbol{\beta} \geq 0, \boldsymbol{\gamma} > 0, \boldsymbol{\epsilon} > 0, \mathbf{Z} > 0$ so that $\mathbf{G}'(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\epsilon}, \mathbf{Z}, \mathbf{K}_0) \leq 0$. The optimization problem is now set up and can be written in the form (4.11).
- Try to construct initial L1 parameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$. It has been shown in section 4.3.4 that if the extended loop \mathbf{H}_i is stable then this is possible by first selecting an initial $\alpha_i \in]0, 1 - \rho_i^2[$ then calculating the initial \mathbf{Z}_i and finally the initial β_i , with ρ_i being the spectral radius of the extended loop dynamic matrix \mathbf{A}_i . Alternately fixing α_i, β_i and \mathbf{Z}_i and solving the resulting linear matrix inequalities at each step in order to decrease the upper L1 norm bound γ_i should then lead to usable initial L1 parameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$.
- Try to build \mathbf{Q} from \mathbf{K}_0 using the relations provided in section 4.2.1. If equation (4.12) has a solution \mathbf{R} then \mathbf{Q} can be calculated by (4.13). It has been shown in section 4.3.1 that if the closed loop dynamic matrix \mathbf{A}_H is diagonalizable then this is possible. Since \mathbf{K}_0 is required to be observable, also $\mathbf{Q} = \mathbf{Q}_L$ is observable due to (2.45). Further the extended loops \mathbf{H}_i need to be stable for the L1, L2/H2, H_∞ norms to remain bounded so that \mathbf{K}_0 is a stabilizing controller and thus \mathbf{Q} is stable due to (2.46). Hence \mathbf{Q} is stable and observable. The optimization problem is now ready for optimization and can be written in the form (4.18).
- Extract the stable observable pair $(\mathbf{A}_Q, \mathbf{C}_Q)$ from parameter \mathbf{Q} and set it as the center $(\mathbf{A}_{Q_0} = \mathbf{A}_Q, \mathbf{C}_{Q_0} = \mathbf{C}_Q)$ of the chart for the parametrization of stable observable pairs. Reset $\mathbf{W} = \mathbf{0}$ and let \mathbf{G}_0 be the Gram observability matrix of the center stable observable pair $(\mathbf{A}_{Q_0}, \mathbf{C}_{Q_0})$. Find a lossless completion pair $(\mathbf{B}_0, \mathbf{D}_0)$ of the center $(\mathbf{A}_{Q_0}, \mathbf{C}_{Q_0})$ as defined in (3.6).

4.4.2 Optimization

In the optimization stage, the parameter \mathbf{Q} along with α, β are improved, which means that they are moved from their initial positions so that the cost function f decreases and the constraint function \mathbf{G} remains negative semidefinite. The whole optimization stage is repeated until the cost function reaches a local minimum.

- Evaluate the Lagrange dual function of the cost function f for the current parameter value \mathbf{Q} and current values α, β and obtain the local sensitivity matrix $\mathbf{\Omega}^*$ as described in section 4.2.4.
- Actualize the parameter \mathbf{Q} and α, β so that f decreases, the decreasing direction being indicated by the subgradient of the cost function f as described in section 4.2.5. Try to calculate the subgradient as a function of the variation of the parameter state space matrices $\mathbf{B}_Q, \mathbf{D}_Q$ and the stable observable pair parameter \mathbf{W} as shown in the sections 4.2.6, 4.2.7, 4.2.8. It has been shown in section 4.3.2 that if the dynamic matrix \mathbf{A}_Q of the stable and observable parameter \mathbf{Q} has no eigenvalues in common with the dynamic matrix \mathbf{A}_P of the plant \mathbf{P} then this is possible.
- If the condition $\mathbf{G} > 0$ is close to be violated, center on the current balanced stable observable pair $(\mathbf{A}_{Q_0} = \mathbf{J}\mathbf{A}_Q\mathbf{J}^{-1}, \mathbf{C}_{Q_0} = \mathbf{C}_Q\mathbf{J}^{-1})$ with \mathbf{J} defined by $\mathbf{G} = \mathbf{J}^H\mathbf{J}$ and reset $\mathbf{W} = \mathbf{0}$ and $\mathbf{G}_0 = \mathbf{1}$. Find a lossless completion pair $(\mathbf{B}_0, \mathbf{D}_0)$ of the center $(\mathbf{A}_{Q_0}, \mathbf{C}_{Q_0})$ as defined in (3.6).

4.4.3 Finalization

In the finalization stage, an optimal controller \mathbf{K}^* is reconstructed from the improved parameter \mathbf{Q} .

- Try to build \mathbf{K}^* from \mathbf{Q} using the relations provided in section 4.2.1. If equation (4.14) has a solution \mathbf{R} then \mathbf{K}^* can be calculated by (4.15). It has been shown in section 4.3.3 that if the dynamic matrix \mathbf{A}_Q of the parameter \mathbf{Q} has no eigenvalues in common with the dynamic matrix \mathbf{A}_P of the plant \mathbf{P} then this is possible.

Chapter 5

Conclusion

In this report, a way to improve an existing controller with respect to a multi-objective specification, written in terms of the L1, L2/H2, H_∞ system norms, has been presented. Three key elements distinguish the proposed approach. First, a parametrization of fixed order controllers by a stable and observable controller parameter, based on the work in [13]. Second, a parametrization of all stable observable pairs that allows an infinitesimal adjustment of the controller parameter, based on the work in [2]. And third, the use of a sensitivity information to calculate the local subgradient of the optimization problem so that the controller parameter can be adjusted the right way, based on the work in [4].

Although many performance measures to quantify the specifications of a multiobjective optimization problem exist, the L1, L2/H2, H_∞ system norms are possibly the most important ones. They have been chosen in this report due to their practical relevance and also to illustrate the scope of the proposed approach. For a given system, the H2 and H_∞ norms stand exemplarily for linear matrix inequality constraints and the L1 norm stands exemplarily for a bilinear matrix inequality constraint. Despite these inherent differences, these constraints are treated in a unified way, as part of a more general semidefinite constraint. The necessity for the proposed approach to gain access to the sensitivity information currently limits the possible constraints to these three norms.

The proposed solution is based on the left controller parameter Q_L on which a parametrization of stable observable pairs is applied. However, the dual approach is also possible, based on the right controller parameter Q_R on which a parametrization of stable controllable pairs is applied. In practical applications one could imagine to use Q_L and the stable observable pair parametrization in case of K having more inputs than outputs since the output matrices C_K, C_{Q_L} are then of smaller dimension than the input matrices B_K, B_{Q_R} and thus the stable observable pair parametrization of (A_{Q_L}, C_{Q_L}) is then simpler than the stable controllable pair parametrization of (A_{Q_R}, B_{Q_R}) and to use Q_R and the stable controllable pair parametrization in the other case.

The Riccati equations (2.21), (2.22) may have multiple solutions R_{KP}, R_{PK} for a given plant P and controller K . The implications on the proposed algorithm are yet unknown. It may be possible that the descending paths of the algorithm are different for different R_{KP}, R_{PK} so that in the case of the de-

scending algorithm starting from a saddle point or a maximum of the cost function, the outcome of the local optimization may be totally different.

Appendix A

Algebra

A.1 Ring Structure

This section serves as the mathematical foundation on which the main text is built upon and provides in particular the coprime factorization and coprime completion that are the starting point in the search for parametrizations.

A.1.1 Rings and Units

This section introduces the algebraic structure of a ring and provides the laws that follow from the definitions. Some special rings are given and the notion of unit is defined.

Ring

A ring R is a set with two binary relations “+” called addition and “.” called multiplication with multiplication having priority over addition. Every element of a ring has an inverse with respect to addition, but not all elements have inverses with respect to multiplication. A ring is defined by the properties

$$R + R = R \quad (\text{completeness}) \quad (\text{A.1})$$

$$\forall a, b, c \in R : a + (b + c) = (a + b) + c \quad (\text{associativity}) \quad (\text{A.2})$$

$$\forall a, b \in R : a + b = b + a \quad (\text{commutativity}) \quad (\text{A.3})$$

$$\forall a \in R : a + 0 = 0 + a = a \quad (\text{identity}) \quad (\text{A.4})$$

$$\forall a \in R : a + (-a) = (-a) + a = a - a = 0 \quad (\text{inverse}) \quad (\text{A.5})$$

$$RR = R \quad (\text{completeness}) \quad (\text{A.6})$$

$$\forall a, b, c \in R : a(bc) = (ab)c \quad (\text{associativity}) \quad (\text{A.7})$$

$$\forall a \in R : a1 = 1a = a \quad (\text{identity}) \quad (\text{A.8})$$

$$\forall a, b, c \in R : a(b + c) = ab + ac, (a + b)c = ac + bc \quad (\text{distributivity}) \quad (\text{A.9})$$

Remark 1 Multiplication with the additive identity

$$\forall a \in R : a0 = 0a = 0 \quad (\text{A.10})$$

since $a = a1 = a(1 + 0) = a1 + a0 = a + a0 \Leftrightarrow a - a = 0 = a0$ and by analogy $a = 1a$ leads to $0 = 0a$.

Remark 2 Multiplication with an additive inverse

$$\forall a, b \in R : a(-b) = (-a)b = -ab \quad (\text{A.11})$$

since $ab + a(-b) = a[b + (-b)] = a0 = 0 \Leftrightarrow a(-b) = -ab$ and by analogy $ab + (-a)b$ leads to $(-a)b = -ab$.

Trivial Ring

The ring $R = \{0 = 1\}$ is called the *trivial ring*. In a nontrivial ring R that has other elements $a \in R : a \neq 0$ we always have

$$0 \neq 1 \Leftrightarrow \text{nontrivial ring} \quad (\text{A.12})$$

Proof Let $0 = 1$ and take an arbitrary $a \in R$ then $a = a1 = a0 = 0$ and thus $a = 0$ for all $a \in R$ so that R must be the trivial ring.

Commutative Ring

If in addition to laws (A.1) to (A.9) the elements of ring R further obey

$$\forall a, b \in R : ab = ba \quad (\text{A.13})$$

then R is said to be a *commutative ring*.

Integral Ring

If in addition to laws (A.1) to (A.9) there is no solution to

$$\forall a, b \in R \setminus \{0\} : ab = 0 \quad (\text{A.14})$$

then R has no zero divisors and is said to be an *integral ring*.

Remark On an integral ring, the equation $ax = a$ with $x \in R, a \in R \setminus \{0\}$ has a unique solution $x = 1$ since $ax - a = a(x - 1) = 0$ and $a \neq 0$.

Unit

If an element u in a ring R is invertible with respect to multiplication, then it is called a unit and obeys the rule

$$\forall u \in R : uu^{-1} = u^{-1}u = 1 \quad (\text{A.15})$$

This leads to a special property for units $u \in R$

$$Ru = uR = R \quad (\text{A.16})$$

Proof since on the one hand $Ru \subset R$ and $uR \subset R$ and on the other hand $R = R1 = Ru^{-1}u \subset Ru$ and $R = 1R = uu^{-1}R \subset uR$ and therefore being subsets of each other the two sets are equal and $Ru = uR = R$. This is in contrast to $Rr \subset R$, $rR \subset R$ for an arbitrary $r \in R$ due to multiplicative completeness (A.6).

Field

If every nonzero element $r \neq 0$ in a ring R is a unit then the ring R is called a *field*.

A.1.2 Ideals and Equivalence

This section introduces important subsets of a ring, the ideals. A measure for the size of ideals is given and the smallest ideals containing a number of elements are calculated. Finally, the notion of equivalence on a ring is defined and its relation to ideals is shown.

Ideal

An ideal I is an additive subgroup of a ring R that is invariant under multiplication by R . Since multiplication on a ring is not necessarily commutative, there exist left and right ideals. On a commutative ring they are both equal. A left ideal I is defined by the properties

$$I \subset R \quad (\text{subset}) \quad (\text{A.17})$$

$$I - I = I \quad (\text{additive subgroup}) \quad (\text{A.18})$$

$$RI = I \quad (\text{multiplicative invariance}) \quad (\text{A.19})$$

For a right ideal, equation (A.19) changes to $IR = I$. Equation (A.17) assures that I inherits the laws of the ring except completeness (A.1) and (A.6). Furthermore, since I is a subset, it not necessarily contains the identities 0 and 1 and the additive inverse $-a$ of an element $a \in I$. Equation (A.18) assures additive completeness of I and provides the ideal with the additive identity 0 and the additive inverse $-a$ of an element $a \in I$. With these two equations I becomes an additive subgroup of R . Equation (A.19) provides the multiplicative invariance property and since $I \subset R \Rightarrow II \subset RI = I$ assures multiplicative completeness.

Remark That means an ideal I is almost a ring. If I also contained the multiplicative identity 1, then I would indeed be a ring. Equation (A.20) will show that this ring is R .

Ideal containing a Unit

If any ideal I contains a unit $u \in R$ then the ideal is the ring itself

$$I = R \Leftrightarrow u \in I \quad (\text{A.20})$$

Proof Since $R = Ru = uR \subset I$ using (A.16) and (A.19) and $I \subset R$ due to (A.17) and therefore being subsets of each other the two sets are equal and $I = R$. The ideals 0 and R are called *trivial ideals*.

Intersection of Ideals

The intersection of two same sided ideals is again an ideal of the same side

$$I = I_1 \cap I_2 \tag{A.21}$$

Proof This is because first $I_1 \subset R$ and $I_2 \subset R$ so $I = I_1 \cap I_2 \subset R$ which verifies (A.17), second $I_1 - I_1 = I_1$ and $I_2 - I_2 = I_2$ so $I - I \subset I_1 - I_1 = I_1$ and $I - I \subset I_2 - I_2 = I_2$ so that $I - I \subset I_1 \cap I_2 = I$ and thus $I - I = I$ which verifies (A.18) and third $RI_1 = I_1$ and $RI_2 = I_2$ so $RI \subset RI_1 = I_1$ and $RI \subset RI_2 = I_2$ so that $RI \subset I_1 \cap I_2 = I$ and with $1 \in R$ this gives $RI = I$ which verifies (A.19). The same is true for right ideals.

Sum of Ideals

The sum of two same sided ideals is again an ideal of the same side

$$I = I_1 + I_2 \tag{A.22}$$

Proof This is because first $I_1 \subset R$ and $I_2 \subset R$ so $I = I_1 + I_2 \subset R + R = R$ which verifies (A.17), second $I_1 - I_1 = I_1$ and $I_2 - I_2 = I_2$ so $I - I = (I_1 + I_2) - (I_1 + I_2) = (I_1 - I_1) + (I_2 - I_2) = I_1 + I_2 = I$ which verifies (A.18) and third $RI_1 = I_1$ and $RI_2 = I_2$ so $RI \subset RI_1 + RI_2 = I_1 + I_2 = I$ and with $1 \in R$ this gives $RI = I$ which verifies (A.19). The same is true for right ideals.

Principal Ideal

A principal ideal $I(x)$ is an ideal with one additional property

$$Rx = I(x) \tag{generation} \tag{A.23}$$

This defines a left principal ideal. For a right principal ideal, this equation changes to $xR = I(x)$. This equation means that there is one element $x \in I(x)$ that generates the ideal. A principal ideal is therefore the set of all multiples of x by ring elements.

Size of an Ideal

Two same sided ideals can be compared in size with respect to mutual inclusion

$$I_1 \subset I_2 \Leftrightarrow I_1 \text{ smaller than } I_2 \tag{A.24}$$

and an ideal $I \neq R$ is said to be *maximal* if it is not contained in any other ideal of the same side.

Smallest Ideal containing two Ideals

The smallest ideal that contains two ideals $I_1, I_2 \subset R$ is the sum of the same sided ideals $I_1 + I_2$ generated by these ideals

$$I_1 + I_2 \text{ is the smallest ideal that contains } I_1, I_2 \tag{A.25}$$

Proof Let I be the smallest ideal that contains I_1, I_2 , thus $I \subset I_1 + I_2$. But I as an ideal verifies (A.18) and since $I_1 \subset I$ and $I_2 \subset I$ this leads to $I_1 + I_2 \subset I$ and therefore being subsets of each other the two sets are equal and $I = I_1 + I_2$.

Smallest Ideal containing one Element

The smallest ideal that contains an element $a \in R$ is the principal ideal $I(a)$ generated by this element

$$I(a) \text{ is the smallest ideal that contains } a \quad (\text{A.26})$$

Proof Let I be the smallest ideal that contains a , thus $I \subset I(a)$. But for left ideals $I(a) = Ra \subset RI = I$ and therefore being subsets of each other the two sets are equal and $I = I(a)$. The same is true for right ideals.

Smallest Ideal containing two Elements

The smallest ideal that contains two elements $a, b \in R$ is the sum of the same sided principal ideals $I(a) + I(b)$ generated by these elements

$$I(a) + I(b) \text{ is the smallest ideal that contains } a, b \quad (\text{A.27})$$

Proof Let I be the smallest ideal that contains a, b , thus $I \subset I(a) + I(b)$. But for left ideals $I(a) = Ra \subset RI = I$ and $I(b) = Rb \subset RI = I$ with (A.18) leads to $I(a) + I(b) \subset I$ and therefore being subsets of each other the two sets are equal and $I = I(a) + I(b)$. The same is true for right ideals.

Equivalence

Two elements $x, y \in R$ are said to be left equivalent if there exists a unit $u \in R$ for which $x = uy$. By analogy, they are right equivalent if there exists a unit $v \in R$ for which $x = yv$. Two equivalent elements generate the same principal ideal

$$I(x) = I(y) \Leftrightarrow x \sim y \quad (\text{equivalence}) \quad (\text{A.28})$$

Proof Left equivalence $x \sim y \Leftrightarrow x = uy$ leads to $I(x) = Rx = Ruy = Ry = I(y)$. To verify the reverse direction, let $x \approx y$. If $\exists a \in I(y) : a \notin I(x)$ then clearly $I(x) \neq I(y)$. If this is not the case then $\forall a \in I(y) : a \in I(x)$ and therefore $I(y) \subset I(x) \Rightarrow y \in Rx$. So $\exists m \in R : y = mx$, but m is not a unit since $x \approx y$ so that $I(m) \neq R$. Hence $\exists p \in R : p \notin I(m)$ and thus $\nexists q \in R : p = qm$. This allows the construction of $a = px$ so that $a \in I(x)$ but since $\nexists q \in R : p = qm$ this leads to $\nexists q \in R : a = px = qmx = qy$. This means that $a \notin I(y)$ so indeed $\exists a \in I(x) : a \notin I(y)$ and therefore $x \approx y \Rightarrow I(x) \neq I(y)$ which is equivalent to $x \sim y \Leftrightarrow I(x) = I(y)$. The same is true for right equivalence.

A.1.3 Division

This section introduces the notion of divisibility on a general ring, using the concept of ideals. It is shown how the least common multiple and the greatest common divisor of two ring elements can be found.

Divisibility

Since a principal ideal is the set of all multiples of its generator, any element of the ideal can be written as a product of the generator and a ring element. This introduces the notion of divisibility

$$a \in I(b) \Leftrightarrow a \text{ can be divided by } b \quad (\text{divisibility}) \quad (\text{A.29})$$

Proof Let $I(b)$ be a left ideal then $a \in I(b) = Rb$ and hence $\exists m \in R : a = mb$ and thus b is a right divisor of a . If it is a right ideal then, by analogy, its generator is a left divisor.

Remark A unit $u \in R$ is a divisor to every element of the ring since following equation (A.20) any ideal containing a unit is the ring itself. Especially $I(u) = R$ and thus every ring element is contained in this ideal and therefore divisible by its generator u .

Least Common Multiple

If the intersection of two principal ideals is itself a principal ideal, then the least common multiples are defined by

$$\text{The generators of } I(a) \cap I(b) \text{ are the least common multiples of } a, b \quad (\text{A.30})$$

Proof If the intersection of two principal ideals $I(a) \cap I(b)$ is itself a principal ideal, then we have $I(a) \cap I(b) = I(m)$. Since $I(m) \subset I(a), I(b)$ and especially $m \in I(a), I(b)$ divisibility (A.29) shows that m can be divided by a, b or in other words m is a common multiple of a, b . If m' is another common multiple of a, b of the same side as m , then $m' \in I(a), I(b)$ so that $m' \in I(a) \cap I(b) = I(m)$. This means m' is itself a multiple of m and therefore m is the least common multiple if $I(a) \cap I(b)$ is principal.

Remark Due to (A.28) all generators of $I(a) \cap I(b)$ and thus all least common multiples are equivalent.

Greatest Common Divisor

If the sum of two principal ideals is itself a principal ideal, then the greatest common divisors are defined by

$$\text{The generators of } I(a) + I(b) \text{ are the greatest common divisors of } a, b \quad (\text{A.31})$$

Proof If the sum of two principal ideals $I(a) + I(b)$ is itself a principal ideal, then we have $I(a) + I(b) = I(d)$. Since $I(a), I(b) \subset I(d)$ and especially $a, b \in I(d)$ divisibility (A.29) shows that d is a common divisor of a and b . If d' is another common divisor of a, b of the same side as d , then $a, b \in I(d')$. But following (A.27) the ideal $I(d) = I(a) + I(b)$ is the smallest ideal containing a, b , we have $I(d) \subset I(d')$ and thus d can be divided by d' and therefore d is the greatest common divisor if $I(a) + I(b)$ is principal.

Remark Due to (A.28) all generators of $I(a) + I(b)$ and thus all greatest common divisors are equivalent.

Coprime Elements

Two elements $a, b \in R$ are called *coprime* if their greatest common divisor is a unit $u \in R$

$$I(a) + I(b) = R \Leftrightarrow a, b \text{ are coprime} \quad (\text{A.32})$$

since $R = I(u)$. For a greatest common right divisor, they are right coprime with $I(a)$ and $I(b)$ being left ideals. They are left coprime in the other case.

Coprime Quotients

Since the greatest common divisor d of a, b contains all other common divisors, a, b divided by d do not share any common divisors and thus the quotients are coprime

$$a = md \text{ and } b = nd \Rightarrow m, n \text{ are right coprime} \quad (\text{A.33})$$

if d is the greatest common right divisor of a, b

If it is the greatest common left divisor, then we have $a = dm$ and $b = dn$ with m, n left coprime.

Proof Let m, n have a greatest common right divisor z so there are $v, w \in R$ with $m = vz$ and $n = wz$ so that $a = md = v(zd)$ and $b = nd = w(zd)$. So zd is a common right divisor of a, b so $a, b \in I(zd)$ and because (A.27) and (A.31) show that $I(d)$ is the smallest ideal to contain a, b thus $I(d) \subset I(zd)$. But $I(zd) = Rzd \subset RRd = Rd = I(d)$ and therefore being subsets of each other the two sets are equal and $I(d) = I(zd) \Leftrightarrow d \sim zd$ using (A.28). Therefore z must be a unit and thus m, n are right coprime. The same is true for left division and left coprimeness.

Bezout Identity

The greatest common divisor d of $a, b \in R$ is a generator of $I(a) + I(b)$ and thus $d \in I(a) + I(b)$. If the ideals are left sided then $d \in Ra + Rb$ and therefore

$$\exists x, y \in R : xa + yb = d, \text{ the greatest common right divisor of } a, b \quad (\text{A.34})$$

which is called the *Bezout identity*. For right ideals and left divisors, this equation changes to $ax + by = d$. On integral rings, the Bezout identity (A.34) can be reduced to

$$xv + yw = 1 \text{ with } x, y \text{ left coprime and } v, w \text{ right coprime} \quad (\text{A.35})$$

Proof According to (A.33) there are right coprime $v, w \in R$ with $a = vd$ and $b = wd$ so that $xvd + ywd = d$ leads to $xv + yw = 1$ on an integral ring. Let x, y have a greatest common left divisor z so there are $m, n \in R$ with $x = zm$ and $y = zn$ so that $zmv + znw = 1$. By $z(mv + nw)z = z$ we have $z(mv + nw) = (mv + nw)z = 1$ since the ring is integral. So the greatest common left divisor z is a unit and therefore x, y are left coprime.

A.1.4 Factorization

This section develops the factorization of elements on a special case of a ring, the principal commutative ring. On such a ring, it is shown that the factorization is essentially unique.

Principal Ring

A principal ring is a ring where every ideal is principal. So if in addition to laws (A.1) to (A.9) the elements and ideals of R further obey

$$\forall I \subset R, \exists x \in R: I = I(x) \quad (\text{A.36})$$

then R is said to be a *principal ring*.

Irreducible Elements

Every element $p \in R$ is at least contained in two principal ideals, namely the ideal $I(p)$ generated by p itself and the trivial ideal $I(1) = I(u) = R$. If p is not contained in any other principal ideal, then p is said to be *irreducible*. Hence, irreducible elements are only divisible by itself and the units $u \in R$. On a principal ring R the following equivalence holds

$$I(p) \text{ is maximal} \Leftrightarrow p \text{ is irreducible} \quad (\text{A.37})$$

Proof If on the one hand $p \in R$ is not irreducible, there is a $q \in R$ that divides p so that $p \in I(q)$ and using (A.26) this means $I(p) \subset I(q)$ and thus $I(p)$ is not maximal. If on the other hand $I(p)$ is not maximal then there is an ideal I so that $I(p) \subset I$ and since the ring is principal, I has a generator $q \in R$ so that $I(p) \subset I(q)$ and p can be divided by q and thus p is not irreducible.

Irreducibility and Coprimeness

For any element $x \in R$, the following equivalence holds if $p \in R$ is irreducible:

$$I(x) + I(p) = R \Leftrightarrow x \notin I(p) \quad (\text{A.38})$$

with the two ideals being same sided. In other words, if an irreducible $p \in R$ does not divide $x \in R$, then with (A.32) x, p are coprime and their greatest common divisor is a unit $u \in R$.

Proof If on the one hand $x \in I(p)$ then using (A.26) this leads to $I(x) \subset I(p)$ and thus $I(x) + I(p) = I(p) \neq R$ so that $I(x) + I(p) = R \Rightarrow x \notin I(p)$. If on the other hand $I(x) + I(p) \neq R$ then clearly $I(p) \subset I(x) + I(p)$ and because p is irreducible and with (A.37) $I(p)$ is maximal, this gives $I(x) + I(p) = I(p)$ and thus $I(x) \subset I(p)$ and especially $x \in I(p)$ so that $I(x) + I(p) = R \Leftarrow x \notin I(p)$.

Irreducibility and Divisibility

If a product ab of elements of a commutative ring $a, b \in R$ is not divisible by an irreducible $p \in R$, then neither of the factors a, b is divisible by p

$$ab = ba \notin I(p) \Leftrightarrow a, b \notin I(p) \quad (\text{A.39})$$

Proof If $a \in I(p)$ then with (A.26) this gives $I(a) \subset I(p)$ and because $I(ab) = Rab \subset Ra = I(a)$ we have $I(ab) \subset I(p)$ hence $ab \in I(p)$ and thus $ab \notin I(p) \Rightarrow a, b \notin I(p)$. Now let $ab \in I(p)$, then at least one of the factors a, b must be divisible by p to verify the reverse direction. If $a \notin I(p)$ then with (A.38) this gives $I(a) + I(p) = R$ and can be written as $\exists x, y \in R : xa + yb = 1$. Multiplying this by b leads to $xab + ypb = b$ and since $ab \in I(p)$ can be written as $\exists z \in R : ab = zp$ this gives $xzp + ypb = b$ so that $(xz + yb)p = b$ and thus $b \in I(p)$ which verifies $ab \notin I(p) \Leftarrow a, b \notin I(p)$.

Factorization on Principal Commutative Rings

Any element of a principal commutative ring $x \in R$ has an essentially unique factorization in irreducible elements

$$x = u \prod_{k=1}^n p_k \text{ where } u \in R \text{ is a unit} \quad (\text{A.40})$$

Proof This factorization is due to an algorithm that reduces x to a unit u . If x is not a unit, there must be a maximal ideal $I(p_1)$, which is principal since R is principal, with $x \in I(p_1)$ so that $\exists x_1 \in R : x = x_1 p_1$. If x_1 is not a unit, there must be a maximal ideal $I(p_2)$ with $x_1 \in I(p_2)$ so that $\exists x_2 \in R : x_1 = x_2 p_2$ and so on, until $x_n = u$ is a unit. Since the ideals $I(p_k)$ are maximal, with (A.37) the p_k are irreducible. It can be proven by Zorn's lemma that this algorithm and hence n in (A.40) is finite.

Now let $x = u \prod_{k=1}^m p_k = v \prod_{k=1}^n q_k$ be two factorizations of x . Since a q_j divides x we have $x = u \prod_{k=1}^m p_k \in I(q_j)$ and using (A.39) there is at least one p_i with $p_i \in I(q_j)$. Using (A.26) this gives $I(p_i) \subset I(q_j)$ but $I(p_i)$ is maximal so that $I(p_i) = I(q_j)$ and with (A.28) the two generators are equivalent $p_i \sim q_j$. So for every p_i there is an equivalent q_j and thus the number of irreducible factors is the same $m = n$ and the factorization (A.40) is essentially unique, with the word essentially referring to the possibility of exchanging a factor by an equivalent factor.

A.1.5 Extension to Matrices

This section treats matrices that are built from elements of a principal commutative integral ring.

Unit Matrices

A square matrix U built from elements of a principal commutative ring R is a unit if and only if its determinant is a unit in R

$$U \text{ is a unit in } R^{n,n} \Leftrightarrow \det U \text{ is a unit in } R \quad (\text{A.41})$$

Proof since on the one hand if U is a unit in $R^{n,n}$ then by definition (A.15) there exists a U^{-1} with $UU^{-1} = U^{-1}U = \mathbf{1}$ where $\mathbf{1}$ denotes the $R^{n,n}$ identity matrix. Therefore $\det(UU^{-1}) = \det U \det U^{-1} = \det(U^{-1}U) = \det U^{-1} \det U = \det \mathbf{1} = 1$ shows that $\det U$ is a unit in R . If on the other hand $\det U$ is a unit in R then $(\det U)^{-1} \in R$. Cramer's theorem states that

for every matrix $M \in R^{n,n}$ there is a unique adjoint matrix $\text{adj } M \in R^{n,n}$ that satisfies $M \text{adj } M = \text{adj } M M = \mathbf{1} \det M$. Since $(\det U)^{-1} \in R$, the matrix $\text{adj } U (\det U)^{-1} \in R^{n,n}$ is a valid inverse of U which therefore is a unit in $R^{n,n}$.

Bezout Identity for two Elements

The relations among two elements of a principal commutative integral ring $a, b \in R$ and their greatest common divisor $d \in R$ can be brought into the matrix representation of the Bezout identity

$$\exists \text{ unit } U \in R^{2,2} : U \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} d \\ 0 \end{bmatrix} \quad (\text{A.42})$$

with d being the greatest common divisor of $\{a, b\}$

Proof Since R is integral, the Bezout identity (A.34) can be reduced to (A.35), so for any two elements $a, b \in R$ with greatest common divisor $d \in R$ there exist $v, w, x, y \in R$ so that $xa + yb = d, xv + yw = 1, wa = vb$ with the third equation stemming from the fact that d is a common divisor of a, b so that $a = vd, b = wd$ hence $wa = wvd = vwd = vb$ since R is commutative. These three equations can then be brought into the matrix form (A.42)

$$\begin{bmatrix} x & y \\ \mp w & \pm v \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} d \\ 0 \end{bmatrix} \text{ with } \det \begin{bmatrix} x & y \\ \mp w & \pm v \end{bmatrix} = \pm xv \pm yw = \pm 1$$

Because of (A.41) the matrix U is a unit in $R^{2,2}$ since $\det U = \pm 1$ is a unit in R .

Bezout Identity for multiple Elements

This form of the Bezout identity can be generalized to a relation between a set of elements $\{a_1, a_2, \dots, a_n\} \subset R$ and their greatest common divisor $d \in R$

$$\exists \text{ unit } U \in R^{n,n} : U \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} d \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (\text{A.43})$$

with d being the greatest common divisor of $\{a_1, a_2, \dots, a_n\}$

Proof This is proven by induction. Following (A.42) this equation is correct for $n = 2$. Suppose (A.43) correct for $n - 1$ so that there exists a unit $U \in R^{n-1, n-1}$ with

$$U \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n-1} \end{bmatrix} = \begin{bmatrix} d \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & U \end{bmatrix} \begin{bmatrix} a_n \\ a_1 \\ a_2 \\ \vdots \\ a_{n-1} \end{bmatrix} = \begin{bmatrix} a_n \\ d \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

where d is the greatest common divisor of a_1, a_2, \dots, a_{n-1} . In order to get from $n - 1$ to n , equation (A.42) may be used for the greatest common divisor d' of

elements a_n, d and thus the greatest common divisor of all a_1, a_2, \dots, a_n so that there exists a unit $\mathbf{V} \in R^{2,2}$ with:

$$\mathbf{V} \begin{bmatrix} a_n \\ d \end{bmatrix} = \begin{bmatrix} d' \\ 0 \end{bmatrix} \Rightarrow \begin{bmatrix} \mathbf{V} & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n-2} \end{bmatrix} \begin{bmatrix} a_n \\ d \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} d' \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Hence the equation can be formulated for n :

$$\begin{bmatrix} \mathbf{V} & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n-2} \end{bmatrix} \begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{U} \end{bmatrix} \begin{bmatrix} a_n \\ a_1 \\ a_2 \\ \vdots \\ a_{n-1} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{V} & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n-2} \end{bmatrix} \begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{U} \end{bmatrix} \begin{bmatrix} \mathbf{0} & 1 \\ \mathbf{1}_{n-1} & \mathbf{0} \end{bmatrix}}_{\mathbf{U}'} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} d' \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

It is now left to verify that \mathbf{U}' is still a unit in $R^{n,n}$ what according to (A.41) is equal to $\det \mathbf{U}'$ being a unit in R :

$$\det \mathbf{U}' = \det \begin{bmatrix} \mathbf{V} & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n-2} \end{bmatrix} \det \begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{U} \end{bmatrix} \det \begin{bmatrix} \mathbf{0} & 1 \\ \mathbf{1}_{n-1} & \mathbf{0} \end{bmatrix} = (-1)^{n+1} \det \mathbf{V} \det \mathbf{U}$$

Since $(-1)^{n+1} = \pm 1$ is a unit in R and $\det \mathbf{U}, \det \mathbf{V}$ are units in R , $\det \mathbf{U}'$ is also a unit in R and thus the matrix \mathbf{U}' is a unit in $R^{n,n}$.

Equivalence to Triangular Matrix

Any matrix $\mathbf{M} \in R^{m,n}$ is left equivalent to an upper triangular matrix $\mathbf{T} \in R^{m,n}$ that has zeros below the diagonal

$$\exists \text{ unit } \mathbf{U} : \mathbf{U}\mathbf{M} = \mathbf{T} \text{ upper triangular} \quad (\text{A.44})$$

For right equivalence this equation changes to $\mathbf{M}^T \mathbf{U}^T = \mathbf{T}^T$ when transposed, so that any matrix \mathbf{M}^T is right equivalent to \mathbf{T}^T which is a lower triangular matrix with zeros above the diagonal.

Proof This is proven by induction. For $m = 1$ matrix \mathbf{M} is a single line and thus already an upper triangular matrix. For $n = 1$ matrix \mathbf{M} is a single column that according to (A.43) is left equivalent to an upper triangular matrix. Suppose (A.44) correct for $m-1, n-1$ so that there exists a unit $\mathbf{U} \in R^{m-1, m-1}$ with

$$\mathbf{U}\mathbf{M} = \mathbf{T} \text{ upper triangular}$$

For m, n the first column of a matrix $\mathbf{M}' \in R^{m,n}$ can be brought in upper triangular form using (A.43):

$$\mathbf{V}\mathbf{M}' = \begin{bmatrix} d & \mathbf{m} \\ \mathbf{0} & \mathbf{M} \end{bmatrix} \Rightarrow \underbrace{\begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{U} \end{bmatrix}}_{\mathbf{U}'} \mathbf{V}\mathbf{M}' = \begin{bmatrix} d & \mathbf{m} \\ \mathbf{0} & \mathbf{U}\mathbf{M} \end{bmatrix} = \underbrace{\begin{bmatrix} d & \mathbf{m} \\ \mathbf{0} & \mathbf{T} \end{bmatrix}}_{\mathbf{T}'}$$

with \mathbf{V} being a unit and $\det \mathbf{U}' = \det \mathbf{U} \det \mathbf{V}$ shows that \mathbf{U}' is a unit in $R^{m,m}$ so that \mathbf{M}' is equivalent to the resulting upper triangular matrix \mathbf{T}' and thus (A.44) is valid for m, n . It is interesting to note that due to (A.43) the entry d in \mathbf{T}' is the greatest common divisor of the first column of \mathbf{M}' .

Equivalence to Diagonal Matrix

Any matrix $M \in R^{m,n}$ is equivalent to a diagonal matrix $D \in R^{m,n}$ that has zeros below and above the diagonal

$$\exists \text{ units } U, V : UMV = D \text{ diagonal} \quad (\text{A.45})$$

Proof This is proven by induction. For $m = 1$ matrix M is a single line and thus by (A.43) already equivalent to the diagonal matrix $[d \ \mathbf{0}]$. For $n = 1$ matrix M is a single column and thus the transpose of a single line and thus also equivalent to a diagonal matrix. Suppose (A.45) correct for $m - 1, n - 1$ so that there exist units $U \in R^{m-1, m-1}$ and $V \in R^{n-1, n-1}$ with

$$UMV = D \text{ diagonal}$$

For m, n the first column of a matrix $M' \in R^{m,n}$ can be brought in upper triangular form using (A.43):

$$U_1 M' = \begin{bmatrix} d_1 & \mathbf{m}_1 \\ \mathbf{0} & M_1 \end{bmatrix}$$

with U_1 being a unit. If d_1 divides all elements of vector \mathbf{m}_1 so that $\mathbf{m}_1 = \mathbf{n}_1 d_1$ then

$$U_1 M' \underbrace{\begin{bmatrix} 1 & -\mathbf{n}_1 \\ \mathbf{0} & \mathbf{1} \end{bmatrix}}_{V_1} = \begin{bmatrix} d_1 & \mathbf{m}_1 \\ \mathbf{0} & M_1 \end{bmatrix} \begin{bmatrix} 1 & -\mathbf{n}_1 \\ \mathbf{0} & \mathbf{1} \end{bmatrix} = \begin{bmatrix} d_1 & \mathbf{0} \\ \mathbf{0} & M_1 \end{bmatrix}$$

with V_1 being a unit since $\det V_1 = 1$. If d_1 does not divide all elements of vector \mathbf{m}_1 so that $\mathbf{m}_1 \neq \mathbf{n}_1 d_1$ then again, using (A.44), there exists a unit $V_2 \in R^{n,n}$ with

$$V_2 (U_1 M')^T = \begin{bmatrix} d_2 & \mathbf{m}_2 \\ \mathbf{0} & M_2 \end{bmatrix}$$

where d_2 is the greatest common divisor of the first column of $(U_1 M')^T$ thus there exists a $r_2 \in R$ so that $d_1 = r_2 d_2$. Here the algorithm repeats until eventually a d_k divides all elements of \mathbf{m}_k so that $\mathbf{m}_k = \mathbf{n}_k d_k$. It can be proven by Zorn's lemma that this algorithm and hence k is finite so that there exist units $U' \in R^{m,m}$ and $V' \in R^{n,n}$ with

$$U' M' V' = \begin{bmatrix} d & \mathbf{0} \\ \mathbf{0} & M \end{bmatrix}$$

Since (A.45) is supposed to be true for $m - 1, n - 1$, this leads to

$$\underbrace{\begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & U \end{bmatrix}}_{U''} U' M' V' \underbrace{\begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & V \end{bmatrix}}_{V''} = \begin{bmatrix} d & \mathbf{0} \\ \mathbf{0} & UMV \end{bmatrix} = \underbrace{\begin{bmatrix} d & \mathbf{0} \\ \mathbf{0} & D \end{bmatrix}}_{D'}$$

with $\det U'' = \det U \det U'$ and $\det V'' = \det V \det V'$ shows that U'' and V'' are units in $R^{m,m}$ and $R^{n,n}$ respectively so that M' is equivalent to the resulting diagonal matrix D' and thus (A.45) is valid for m, n .

Smith Normal Form

Any matrix of elements of a principal commutative integral ring $M \in R^{m,n}$ is equivalent to a diagonal matrix $D \in R^{m,n}$ that has zeros below and above the diagonal and whose i th diagonal element $[D]_{ii} = d_i$ is a divisor of all lower diagonal elements $d_{i+\mathbb{N}}$

$$\begin{aligned} \exists \text{ units } U, V : UMV = D \text{ diagonal} & \quad (\text{A.46}) \\ \text{with } [D]_{ii} = d_i \text{ being a divisor of } d_{i+\mathbb{N}} & \end{aligned}$$

Proof Due to (A.45) matrix M is equivalent to a diagonal matrix D . It is left to verify the sequential divisibility of its elements. This is proven by induction. For $m = 1$ or $n = 1$, matrix M is equivalent to a diagonal matrix with only one diagonal element hence (A.46) is true. Suppose (A.46) correct for $m - 1, n - 1$ so that there exist units $U \in R^{m-1, m-1}$ and $V \in R^{n-1, n-1}$ with

$$UMV = D \text{ with } [D]_{ii} = d_i \text{ being a divisor of } d_{i+\mathbb{N}}$$

For m, n , using (A.45), there exist units $U' \in R^{m,m}$ and $V' \in R^{n,n}$ so that matrix $M \in R^{m,n}$ is equivalent to diagonal matrix $D' \in R^{m,n}$:

$$U'M'V' = \begin{bmatrix} d & \mathbf{0} \\ \mathbf{0} & D \end{bmatrix}$$

By using another unit $V'' \in R^{n,n}$, matrix D' can be brought into the form

$$D' \underbrace{\begin{bmatrix} 1 & \mathbf{0} \\ v & 1 \end{bmatrix}}_{V''} = \begin{bmatrix} d & \mathbf{0} \\ Dv & D \end{bmatrix}$$

with a vector $v \in R^{n-1}$ having identity elements $[v]_k = 1$. Since D is diagonal, Dv gives

$$Dv = \begin{bmatrix} d_1 \\ \vdots \\ d_l \end{bmatrix} \text{ with } l = \min(m-1, n-1)$$

Following (A.43) there exists a unit $U'' \in R^{m,m}$ that reduces the first column of $D'V''$ to upper triangular form:

$$\underbrace{\begin{bmatrix} u_1 & u_2 \\ u_3 & U_4 \end{bmatrix}}_{U''} \begin{bmatrix} d & \mathbf{0} \\ Dv & D \end{bmatrix} = \begin{bmatrix} \delta & n \\ \mathbf{0} & N \end{bmatrix} \Rightarrow u_1 d + u_2 Dv = \delta$$

with δ being the greatest common divisor of d and Dv that is the greatest common divisor of $\{d, d_1, \dots, d_l\}$. Since U'' is a unit, it has an inverse U''^{-1} so that this transformation is reversible:

$$\underbrace{\begin{bmatrix} u'_1 & u'_2 \\ u'_3 & U'_4 \end{bmatrix}}_{U''^{-1}} \begin{bmatrix} \delta & n \\ \mathbf{0} & N \end{bmatrix} = \begin{bmatrix} d & \mathbf{0} \\ Dv & D \end{bmatrix} \Rightarrow u'_1 \delta = d$$

With $\delta = u_1 d + u_2 Dv$ this gives $u'_1(u_1 d + u_2 Dv) = d$ hence $u'_1 u_2 = \mathbf{0}$ and $u'_1 u_1 = u_1 u'_1 = 1$ so that u_1, u'_1 are units since this must hold for all d and D .

But $u'_1\delta = d$ and u'_1 being a unit means that d, δ are equivalent greatest common divisors of $\{d, d_1, \dots, d_l\}$ so d divides all elements d_i of diagonal matrix \mathbf{D} so that we have

$$\mathbf{U}'\mathbf{M}'\mathbf{V}' = \mathbf{D}' \text{ with } [D']_{ii} = d_i \text{ being a divisor of } d_{i+\mathbb{N}}$$

and thus (A.46) is valid for m, n .

Bezout Identity for Matrices

Two matrices $\mathbf{A} \in R^{m,k}$ and $\mathbf{B} \in R^{n,k}$ with the same number of columns have a common square right divisor matrix $\mathbf{D} \in R^{k,k}$

$$\exists \mathbf{X}, \mathbf{Y} : \mathbf{X}\mathbf{A} + \mathbf{Y}\mathbf{B} = \mathbf{D} \tag{A.47}$$

with \mathbf{D} being the greatest common right divisor of \mathbf{A}, \mathbf{B}

which is the Bezout Identity for matrices. For left divisors this equation changes to $\mathbf{A}\mathbf{X} + \mathbf{B}\mathbf{Y} = \mathbf{D}$ for $\mathbf{A} \in R^{k,m}$ and $\mathbf{B} \in R^{k,n}$ having the same number of lines, so this is the transpose of (A.47).

Proof To verify the equation for right divisors, the matrices \mathbf{A}, \mathbf{B} can be collected in a vertical block matrix since they have the same number of columns. Then by applying (A.45) this leads to

$$\exists \text{ units } \mathbf{U}, \mathbf{V} : \begin{bmatrix} \mathbf{A} \\ \mathbf{B} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{U}_1 & \mathbf{U}_2 \\ \mathbf{U}_3 & \mathbf{U}_4 \end{bmatrix}}_{\mathbf{U}} \begin{bmatrix} \mathbf{D}' \\ \mathbf{0} \end{bmatrix} \mathbf{V}$$

with \mathbf{D}' diagonal. This gives

$$\mathbf{A} = \mathbf{U}_1\mathbf{D}'\mathbf{V}, \mathbf{B} = \mathbf{U}_3\mathbf{D}'\mathbf{V}$$

so that \mathbf{A}, \mathbf{B} have a common right divisor $\mathbf{D} = \mathbf{D}'\mathbf{V}$. Since \mathbf{U} is a unit, it is invertible over R :

$$\mathbf{U}^{-1} = \begin{bmatrix} \mathbf{X} & \mathbf{Y} \\ \mathbf{Z}_1 & \mathbf{Z}_2 \end{bmatrix}$$

This then provides the following representation:

$$\mathbf{U}^{-1} \begin{bmatrix} \mathbf{A} \\ \mathbf{B} \end{bmatrix} = \begin{bmatrix} \mathbf{X} & \mathbf{Y} \\ \mathbf{Z}_1 & \mathbf{Z}_2 \end{bmatrix} \begin{bmatrix} \mathbf{A} \\ \mathbf{B} \end{bmatrix} = \begin{bmatrix} \mathbf{D}' \\ \mathbf{0} \end{bmatrix} \mathbf{V} = \begin{bmatrix} \mathbf{D} \\ \mathbf{0} \end{bmatrix}$$

This equation includes the Bezout identity in the first line:

$$\mathbf{X}\mathbf{A} + \mathbf{Y}\mathbf{B} = \mathbf{D}$$

It is now left to verify that the common right divisor \mathbf{D} is also the greatest common right divisor of \mathbf{A}, \mathbf{B} . Let \mathbf{D}'' be another common right divisor of \mathbf{A}, \mathbf{B} , then there exist matrices \mathbf{M} and \mathbf{N} so that $\mathbf{A} = \mathbf{M}\mathbf{D}''$ and $\mathbf{B} = \mathbf{N}\mathbf{D}''$. Since $\mathbf{X}\mathbf{A} + \mathbf{Y}\mathbf{B} = \mathbf{D}$ this leads to $(\mathbf{X}\mathbf{M} + \mathbf{Y}\mathbf{N})\mathbf{D}'' = \mathbf{D}$. This shows that \mathbf{D}'' is itself a right divisor of \mathbf{D} and therefore \mathbf{D} is the greatest common right divisor of \mathbf{A}, \mathbf{B} .

Coprimeness

The pairs of matrices $\mathbf{A}_R \in R^{m,k}$, $\mathbf{B}_R \in R^{n,k}$ and $\mathbf{A}_L \in R^{k,m}$, $\mathbf{B}_L \in R^{k,n}$ are respectively right and left coprime if their respective greatest common right and left divisor is a unit matrix $\mathbf{U} \in R^{k,k}$:

$$\mathbf{A}_L \mathbf{A}_R + \mathbf{B}_L \mathbf{B}_R = \mathbf{U} \quad (\text{A.48})$$

which can be compared to (A.35) for ring elements.

Proof To verify the matrix version, comparison of (A.48) to (A.47) leads to \mathbf{U} being the greatest common right divisor of $\mathbf{A}_R, \mathbf{B}_R$ and by symmetry also to \mathbf{U} being the greatest common left divisor of $\mathbf{A}_L, \mathbf{B}_L$. Since their greatest common divisor is a unit matrix (A.41), the pairs $\mathbf{A}_R, \mathbf{B}_R$ and $\mathbf{A}_L, \mathbf{B}_L$ are right and left coprime respectively.

A.1.6 Extension to Fractions

This section extends the preceding equations for elements of a principal commutative integral ring to equations for fractions of two elements of such a ring.

Fractions on a Ring

A fraction $\frac{a}{b}$ of two elements of a principal commutative integral ring $a \in R, b \in R \setminus \{0\}$ is an element of the set $F = \frac{R}{R \setminus \{0\}}$. It can be proven using laws (A.1) to (A.9) and (A.13) and (A.14) that F is itself a commutative integral ring where every nonzero element is a unit and thus F is a field, in this case the field of fractions of R .

McMillan Normal Form

Any matrix $\mathbf{M} \in F^{m,n}$ is equivalent to a diagonal matrix $\mathbf{D} \in F^{m,n}$ that has zeros below and above the diagonal and whose i th diagonal element is $[\mathbf{D}]_{ii} = \frac{a_i}{b_i}$ with a_i, b_i coprime and where a_i divides all lower $a_{i+\mathbb{N}}$ and all lower $b_{i+\mathbb{N}}$ divide b_i

$$\exists \text{ units } \mathbf{U}, \mathbf{V} : \mathbf{U} \mathbf{M} \mathbf{V} = \mathbf{D} \text{ diagonal} \quad (\text{A.49})$$

$$\text{with } [\mathbf{D}]_{ii} = \frac{a_i}{b_i} \text{ where } a_i \text{ divides } a_{i+\mathbb{N}} \text{ and } b_{i+\mathbb{N}} \text{ divides } b_i$$

$$\text{and } a_i, b_i \in R \text{ are coprime}$$

with $\mathbf{U} \in R^{m,m}$ and $\mathbf{V} \in R^{n,n}$.

Proof By multiplying $\mathbf{M} \in F^{m,n}$ with the least common multiple d of all its denominators the resulting product $d\mathbf{M} \in R^{m,n}$ is a matrix with elements in R so that (A.46) applies. Hence there exist units $\mathbf{U} \in R^{m,m}$ and $\mathbf{V} \in R^{n,n}$ so that $\mathbf{U} d\mathbf{M} \mathbf{V} = \mathbf{D}$ is diagonal and $[\mathbf{D}]_{ii} = d_i \in R$ being a divisor of $d_{i+\mathbb{N}}$ so there exist $r_{i+\mathbb{N}} \in R$ with $d_{i+\mathbb{N}} = d_i r_{i+\mathbb{N}}$. Thus $\mathbf{U} \mathbf{M} \mathbf{V} = \frac{1}{d} \mathbf{D}$ is diagonal with $\frac{d_{i+\mathbb{N}}}{d} = \frac{d_i}{d} r_{i+\mathbb{N}}$. In general, a fraction can be reduced if its quotients are both divided by their greatest common divisor so that following (A.33) the reduced quotients are coprime. In this case the reduction of the quotients leads

to $\frac{d_i}{d} = \frac{a_i}{b_i}$ with $a_i, b_i \in R$ coprime and $\frac{d_{i+\mathbb{N}}}{d} = \frac{a_{i+\mathbb{N}}}{b_{i+\mathbb{N}}}$ with $a_{i+\mathbb{N}}, b_{i+\mathbb{N}}$ coprime. This gives $\frac{a_{i+\mathbb{N}}}{b_{i+\mathbb{N}}} = \frac{d_{i+\mathbb{N}}}{d} = \frac{d_i}{d} r_{i+\mathbb{N}} = \frac{a_i}{b_i} r_{i+\mathbb{N}}$ so that $a_{i+\mathbb{N}} b_i = a_i b_{i+\mathbb{N}} r_{i+\mathbb{N}}$. Since a_i, b_i are coprime a_i must be a divisor of $a_{i+\mathbb{N}}$ and since $a_{i+\mathbb{N}}, b_{i+\mathbb{N}}$ are coprime $b_{i+\mathbb{N}}$ must be a divisor of b_i .

Coprime Factorization

Any matrix $\mathbf{P} \in F^{m,n}$ has a non unique right and left coprime factorization with matrices $\mathbf{M}_R, \mathbf{M}_L \in R^{m,n}$, $\mathbf{N}_R \in R^{n,n}$, $\mathbf{N}_L \in R^{m,m}$ in the form of

$$\mathbf{P} = \mathbf{M}_R \mathbf{N}_R^{-1} = \mathbf{N}_L^{-1} \mathbf{M}_L \quad (\text{A.50})$$

with $\mathbf{M}_R, \mathbf{N}_R$ right coprime and $\mathbf{M}_L, \mathbf{N}_L$ left coprime

Proof Let \mathbf{P} be factorized using (A.49) hence there exist unit matrices \mathbf{U}, \mathbf{V} with $\mathbf{P} = \mathbf{U} \mathbf{D} \mathbf{V}$ so that $\mathbf{D} \in F^{m,n}$ is diagonal with $[\mathbf{D}]_{ii} = \frac{a_i}{b_i}$ and $a_i, b_i \in R$ are coprime. Matrix \mathbf{D} can be factorized easily since it is diagonal so that $\mathbf{D} = \mathbf{A} \mathbf{B}_R^{-1} = \mathbf{B}_L^{-1} \mathbf{A}$ with diagonal matrices $\mathbf{A} \in R^{m,n}$, $\mathbf{B}_R \in R^{n,n}$, $\mathbf{B}_L \in R^{m,m}$ of the shape $[\mathbf{A}]_{ii} = a_i$ and $[\mathbf{B}_{R,L}]_{ii} = \{b_i \text{ for } i \leq \min(m, n), 1 \text{ for } i > \min(m, n)\}$. Since a_i, b_i are coprime, using (A.34) there exist $x_i, y_i \in R$ so that $x_i a_i + y_i b_i = 1$. In matrix form this means there exist diagonal matrices $\mathbf{X} \in R^{n,m}$, $\mathbf{Y}_L \in R^{n,n}$, $\mathbf{Y}_R \in R^{m,m}$ of the shape $[\mathbf{X}]_{ii} = x_i$ and $[\mathbf{Y}_{L,R}]_{ii} = \{b_i \text{ for } i \leq \min(m, n), 1 \text{ for } i > \min(m, n)\}$ with $\mathbf{X} \mathbf{A} + \mathbf{Y}_L \mathbf{B}_R = \mathbf{1}$ and $\mathbf{A} \mathbf{X} + \mathbf{B}_L \mathbf{Y}_R = \mathbf{1}$. Following (A.48) this means that \mathbf{A}, \mathbf{B}_R are right coprime and \mathbf{A}, \mathbf{B}_L are left coprime.

With $\mathbf{D} = \mathbf{A} \mathbf{B}_R^{-1} = \mathbf{B}_L^{-1} \mathbf{A}$ the factorization $\mathbf{P} = \mathbf{U} \mathbf{D} \mathbf{V}$ can be written in the form

$$\mathbf{P} = (\mathbf{U} \mathbf{A})(\mathbf{V}^{-1} \mathbf{B}_R)^{-1} = (\mathbf{B}_L \mathbf{U}^{-1})^{-1} (\mathbf{A} \mathbf{V}) = \mathbf{M}_R \mathbf{N}_R^{-1} = \mathbf{N}_L^{-1} \mathbf{M}_L$$

with $\mathbf{M}_R = \mathbf{U} \mathbf{A}, \mathbf{N}_R = \mathbf{V}^{-1} \mathbf{B}_R$ right coprime since

$$\mathbf{1} = \mathbf{X} \mathbf{A} + \mathbf{Y}_L \mathbf{B}_R = \mathbf{X} \mathbf{U}^{-1} \mathbf{M}_R + \mathbf{Y}_L \mathbf{V} \mathbf{N}_R$$

and $\mathbf{M}_L = \mathbf{A} \mathbf{V}, \mathbf{N}_L = \mathbf{B}_L \mathbf{U}^{-1}$ left coprime since

$$\mathbf{1} = \mathbf{A} \mathbf{X} + \mathbf{B}_L \mathbf{Y}_R = \mathbf{M}_L \mathbf{V}^{-1} \mathbf{X} + \mathbf{N}_L \mathbf{U} \mathbf{Y}_R$$

using (A.48).

Coprime Completion

A given coprime factorization of $\mathbf{P} \in F^{m,n}$ as in (A.50) can be completed by a non unique set of matrices $\mathbf{X}_R, \mathbf{X}_L \in R^{n,m}$, $\mathbf{Y}_R \in R^{m,m}$, $\mathbf{Y}_L \in R^{n,n}$ in the form of

$$\mathbf{1} = \begin{bmatrix} \mathbf{X}_L & \mathbf{Y}_L \\ -\mathbf{N}_L & \mathbf{M}_L \end{bmatrix} \begin{bmatrix} \mathbf{M}_R & -\mathbf{Y}_R \\ \mathbf{N}_R & \mathbf{X}_R \end{bmatrix} \quad (\text{A.51})$$

with $\mathbf{X}_R, \mathbf{Y}_R$ right coprime and $\mathbf{X}_L, \mathbf{Y}_L$ left coprime

Furthermore all valid matrices $\mathbf{X}_i, \mathbf{Y}_i$ with $i \in \{R, L\}$ are a function of a free parameter matrix $\mathbf{Q} \in R^{n,m}$

$$\begin{aligned} \begin{bmatrix} \mathbf{X}_R \\ \mathbf{Y}_R \end{bmatrix} &= \begin{bmatrix} \mathbf{X}_{R0} & \mathbf{N}_R \\ \mathbf{Y}_{R0} & -\mathbf{M}_R \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \mathbf{Q} \end{bmatrix} && \text{with } \mathbf{Y}_R \text{ invertible} \\ \begin{bmatrix} \mathbf{X}_L & \mathbf{Y}_L \end{bmatrix} &= \begin{bmatrix} \mathbf{1} & \mathbf{Q} \end{bmatrix} \begin{bmatrix} \mathbf{X}_{L0} & \mathbf{Y}_{L0} \\ \mathbf{N}_L & -\mathbf{M}_L \end{bmatrix} && \text{with } \mathbf{Y}_L \text{ invertible} \end{aligned} \quad (\text{A.52})$$

with $\mathbf{X}_{i0}, \mathbf{Y}_{i0}$ being a set of valid initial matrices that satisfy (A.51) and allow the construction of all other solutions $\mathbf{X}_i, \mathbf{Y}_i$ using (A.52).

Proof According to (A.50) there exists a coprime factorization of \mathbf{P} in the form $\mathbf{P} = \mathbf{M}_R \mathbf{N}_R^{-1} = \mathbf{N}_L^{-1} \mathbf{M}_L$ so that $\mathbf{M}_L \mathbf{N}_R - \mathbf{N}_L \mathbf{M}_R = \mathbf{0}$. The factors $\mathbf{M}_R, \mathbf{N}_R$ are right coprime so with (A.48) there exist $\mathbf{X}'_L, \mathbf{Y}'_L$ and unit \mathbf{U}_L so that $\mathbf{X}'_L \mathbf{M}_R + \mathbf{Y}'_L \mathbf{N}_R = \mathbf{U}_L$. The factors $\mathbf{M}_L, \mathbf{N}_L$ are left coprime so with (A.48) there exist $\mathbf{X}'_R, \mathbf{Y}'_R$ and unit \mathbf{U}_R so that $\mathbf{M}_L \mathbf{X}_R + \mathbf{N}_L \mathbf{Y}_R = \mathbf{U}_R$. These equations can be written in block matrix form

$$\begin{bmatrix} \mathbf{X}'_L & \mathbf{Y}'_L \\ -\mathbf{N}_L & \mathbf{M}_L \end{bmatrix} \begin{bmatrix} \mathbf{M}_R & -\mathbf{Y}'_R \\ \mathbf{N}_R & \mathbf{X}'_R \end{bmatrix} = \begin{bmatrix} \mathbf{U}_L & \mathbf{Q}' \\ \mathbf{0} & \mathbf{U}_R \end{bmatrix}$$

with $\mathbf{Q}' = \mathbf{Y}'_L \mathbf{X}'_R - \mathbf{X}'_L \mathbf{Y}'_R$. The matrix on the right side is a unit matrix which can be inverted on R . This can be done without affecting the given factors $\mathbf{M}_i, \mathbf{N}_i$ with $i \in \{R, L\}$ by factorizing the unit matrix in the way

$$\begin{bmatrix} \mathbf{U}_L & \mathbf{Q}' \\ \mathbf{0} & \mathbf{U}_R \end{bmatrix} = \begin{bmatrix} \mathbf{U}_L & \mathbf{Q}_L \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{Q}_R \\ \mathbf{0} & \mathbf{U}_R \end{bmatrix}$$

so that $\mathbf{Q}' = \mathbf{U}_L \mathbf{Q}_R + \mathbf{Q}_L \mathbf{U}_R$. Hence

$$\begin{bmatrix} \mathbf{U}_L & \mathbf{Q}_L \\ \mathbf{0} & \mathbf{1} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{X}'_L & \mathbf{Y}'_L \\ -\mathbf{N}_L & \mathbf{M}_L \end{bmatrix} \begin{bmatrix} \mathbf{M}_R & -\mathbf{Y}'_R \\ \mathbf{N}_R & \mathbf{X}'_R \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{Q}_R \\ \mathbf{0} & \mathbf{U}_R \end{bmatrix}^{-1} = \mathbf{1}$$

which using (A.59) can be written in new variables in the form of (A.51). The new variables are then defined as

$$\begin{aligned} \mathbf{X}_R &= (\mathbf{X}'_R - \mathbf{N}_R \mathbf{Q}_R) \mathbf{U}_R^{-1} \\ \mathbf{Y}_R &= (\mathbf{Y}'_R + \mathbf{M}_R \mathbf{Q}_R) \mathbf{U}_R^{-1} \\ \mathbf{X}_L &= \mathbf{U}_L^{-1} (\mathbf{X}'_L + \mathbf{Q}_L \mathbf{N}_L) \\ \mathbf{Y}_L &= \mathbf{U}_L^{-1} (\mathbf{Y}'_L - \mathbf{Q}_L \mathbf{M}_L) \end{aligned}$$

It can be seen that there is a degree of freedom in $\mathbf{Q}_L, \mathbf{Q}_R$ since these two matrices of $R^{n,m}$ only obey one equation $\mathbf{U}_L \mathbf{Q}_R + \mathbf{Q}_L \mathbf{U}_R = \mathbf{Q}'$ of $R^{n,m}$. This degree of freedom can be extracted in the form of a free matrix \mathbf{Q} so that $\mathbf{U}_L \mathbf{Q}_R = \frac{1}{2} \mathbf{Q}' - \mathbf{U}_L \mathbf{Q} \mathbf{U}_R$ and $\mathbf{Q}_L \mathbf{U}_R = \frac{1}{2} \mathbf{Q}' + \mathbf{U}_L \mathbf{Q} \mathbf{U}_R$. This allows the new variables to be written in the form of (A.52).

It is now left to show that \mathbf{Q} is a ring matrix. The right coprime part of the solution (A.52) can be written as

$$\begin{bmatrix} -\mathbf{Y}_R \\ \mathbf{X}_R \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{M}_R & -\mathbf{Y}_{R0} \\ \mathbf{N}_R & \mathbf{X}_{R0} \end{bmatrix}}_{\mathbf{U}} \begin{bmatrix} \mathbf{Q} \\ \mathbf{1} \end{bmatrix}$$

But following (A.51) the matrix \mathbf{U} is a unit and hence $\mathbf{Q} \in R^{n,m}$.

A.2 Matrix Calculus

This section provides some important definitions and formulae for matrices that are used throughout this text.

A.2.1 Blocks

This section provides basic formulae for the treatment of block matrices.

Factorization of Block Matrices

A block matrix of the following shape can be factorized into

$$\begin{bmatrix} \mathbf{AEB} & \mathbf{AED} \\ \mathbf{CEB} & \mathbf{CED} \end{bmatrix} = \begin{bmatrix} \mathbf{A} \\ \mathbf{C} \end{bmatrix} \mathbf{E} \begin{bmatrix} \mathbf{B} & \mathbf{D} \end{bmatrix} \quad (\text{A.53})$$

Determinant of Sums

The determinant of a sum of matrices is given by

$$\det(\mathbf{A} + \mathbf{BCD}) = \det \mathbf{A} \det \mathbf{C} \det(\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B}) \quad (\text{A.54})$$

Inverse of Sums

The inverse of a sum of matrices is given by

$$(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B})^{-1}\mathbf{DA}^{-1} \quad (\text{A.55})$$

Useful Relations

Two useful relations between inverses of sums are

$$\begin{aligned} \mathbf{A}(\mathbf{1} + \mathbf{BA})^{-1} &= (\mathbf{1} + \mathbf{AB})^{-1}\mathbf{A} \\ \mathbf{AB}(\mathbf{1} + \mathbf{AB})^{-1} &= \mathbf{1} - (\mathbf{1} + \mathbf{AB})^{-1} \end{aligned} \quad (\text{A.56})$$

Schur Decomposition

A block matrix can be decomposed in its block Schur form

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{CA}^{-1} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{CA}^{-1}\mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{A}^{-1}\mathbf{B} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \quad (\text{A.57})$$

if \mathbf{A} is invertible.

Determinant of Block Matrices

The determinant of a block matrix is given by

$$\det \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{cases} \det \mathbf{A} \det(\mathbf{D} - \mathbf{CA}^{-1}\mathbf{B}); & \exists \mathbf{A}^{-1} \\ \det \mathbf{D} \det(\mathbf{A} - \mathbf{BD}^{-1}\mathbf{C}); & \exists \mathbf{D}^{-1} \end{cases} \quad (\text{A.58})$$

which follows from (A.57).

Inverse of Block Matrices

The inverse of a block matrix is given by

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{cases} \begin{bmatrix} \mathbf{X} & -\mathbf{A}^{-1}\mathbf{B}\mathbf{Y} \\ -\mathbf{Y}\mathbf{C}\mathbf{A}^{-1} & \mathbf{Y} \end{bmatrix}; & \exists \mathbf{A}^{-1}, \exists \mathbf{D}^{-1} \\ \begin{bmatrix} \mathbf{X} & -\mathbf{X}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{X} & \mathbf{Y} \end{bmatrix}; & \exists \mathbf{A}^{-1}, \exists \mathbf{D}^{-1} \\ \begin{bmatrix} \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{B}\mathbf{Y}\mathbf{C}\mathbf{A}^{-1} & -\mathbf{A}^{-1}\mathbf{B}\mathbf{Y} \\ -\mathbf{Y}\mathbf{C}\mathbf{A}^{-1} & \mathbf{Y} \end{bmatrix}; & \exists \mathbf{A}^{-1} \\ \begin{bmatrix} \mathbf{X} & -\mathbf{X}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{X} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{X}\mathbf{B}\mathbf{D}^{-1} \end{bmatrix}; & \exists \mathbf{D}^{-1} \end{cases} \quad (\text{A.59})$$

with $\mathbf{X} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}$ and $\mathbf{Y} = (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}$ which follows from (A.57).

A.2.2 Spectra

This section provides formulae for eigenvalues and definiteness of matrices.

Eigenvalues

If λ is an eigenvalue of $\mathbf{A} \in \mathbb{C}^{n,n}$ that is $\det(\lambda\mathbf{1} - \mathbf{A}) = 0$ then

Matrix	Eigenvalue
\mathbf{A}^T	λ
$\overline{\mathbf{A}}$	$\overline{\lambda}$
\mathbf{A}^{-1}	λ^{-1}
$\alpha\mathbf{A}$	$\alpha\lambda$
$\beta\mathbf{1} + \mathbf{A}$	$\beta + \lambda$

(A.60)

Proof This is due to $\det(\lambda\mathbf{1} - \mathbf{A}) = 0 \Leftrightarrow \det(\lambda\mathbf{1} - \mathbf{A})^T = \det(\lambda\mathbf{1} - \mathbf{A}^T) = 0 \Leftrightarrow \det(\lambda\mathbf{1} - \overline{\mathbf{A}}) = \det(\overline{\lambda}\mathbf{1} - \overline{\mathbf{A}}) = 0 \Leftrightarrow \det(\lambda\mathbf{A}(\mathbf{A}^{-1} - \lambda^{-1}\mathbf{1})) = \lambda^n \det \mathbf{A} \det(\mathbf{A}^{-1} - \lambda^{-1}\mathbf{1}) = 0 \Leftrightarrow \det(\lambda^{-1}\mathbf{1} - \mathbf{A}^{-1}) = 0 \Leftrightarrow \alpha^n \det(\lambda\mathbf{1} - \mathbf{A}) = \det(\alpha\lambda\mathbf{1} - \alpha\mathbf{A}) = 0 \Leftrightarrow \det((\beta + \lambda)\mathbf{1} - (\beta\mathbf{1} + \mathbf{A})) = 0$.

Spectral Decomposition of Normal Matrices

A matrix \mathbf{A} is normal if and only if $\mathbf{A}\mathbf{A}^H = \mathbf{A}^H\mathbf{A}$. Then there always exists a spectral decomposition of \mathbf{A} in the form of

$$\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^H = \sum_k \lambda_k \mathbf{a}_k \mathbf{a}_k^H \quad (\text{A.61})$$

with unitary \mathbf{U} so that $\mathbf{U}^{-1} = \mathbf{U}^H$. The λ_k are the eigenvalues and the \mathbf{a}_k the eigenvectors of \mathbf{A} with $\mathbf{a}_i^H \mathbf{a}_i = 1$ and $\mathbf{a}_i^H \mathbf{a}_{j \neq i} = 0$.

Remark Hermitian matrices $\mathbf{A} = \mathbf{A}^H$ are normal since $\mathbf{A}\mathbf{A}^H = \mathbf{A}^2 = \mathbf{A}^H\mathbf{A}$.

Definiteness of Hermitian Matrices

A hermitian matrix $\mathbf{A} = \mathbf{A}^H$ is called positive semidefinite if

$$\mathbf{x}^H \mathbf{A} \mathbf{x} \geq 0, \forall \mathbf{x} \neq \mathbf{0} \Leftrightarrow \lambda_k \geq 0 \quad (\text{A.62})$$

for all its eigenvalues λ_k . A matrix is called positive definite if “ \geq ” can be replaced by “ $>$ ”. A matrix \mathbf{B} is called negative (semi)definite if $-\mathbf{B}$ is positive (semi)definite.

Proof Since \mathbf{A} is hermitian it always has a spectral decomposition $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{U}^H$ due to (A.61). Then we have $\mathbf{x}^H \mathbf{A} \mathbf{x} = (\mathbf{U}^H \mathbf{x}) \mathbf{A} (\mathbf{U}^H \mathbf{x}) = \mathbf{y}^H \mathbf{D} \mathbf{y} = \sum_k \lambda_k |\mathbf{y}_k|^2 \geq 0, \forall \mathbf{y} \neq \mathbf{0} \Leftrightarrow \lambda_k \geq 0$ for all eigenvalues λ_k . Finally since $\mathbf{y} = \mathbf{U}^H \mathbf{x}$ with unitary \mathbf{U} we have $\mathbf{y} \neq \mathbf{0} \Leftrightarrow \mathbf{x} \neq \mathbf{0}$ so that (A.62) follows.

Remark For a hermitian matrix $\mathbf{A} = \mathbf{A}^H$ the notation $\mathbf{A} = \mathbf{0}$ is equivalent to $\mathbf{A} = \mathbf{0}$. On the one hand $\mathbf{A} = \mathbf{0} \Rightarrow \mathbf{x}^H \mathbf{A} \mathbf{x} = 0, \forall \mathbf{x}$ and on the other hand $\mathbf{x}^H \mathbf{A} \mathbf{x} = 0, \forall \mathbf{x} \Rightarrow \mathbf{x} \mathbf{U} \mathbf{D} \mathbf{U}^H \mathbf{x} = (\mathbf{U}^H \mathbf{x})^H \mathbf{D} (\mathbf{U}^H \mathbf{x}) = \mathbf{y}^H \mathbf{D} \mathbf{y} = \sum_k \lambda_k |\mathbf{y}_k|^2 = 0, \forall \mathbf{y} \Leftrightarrow \lambda_k = 0$ for all eigenvalues λ_k so that $\mathbf{D} = \mathbf{0}$ and thus $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{U}^H = \mathbf{0}$.

Definiteness of Hermitian Block Matrices

A necessary condition for a hermitian block matrix can be given by

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^H & \mathbf{C} \end{bmatrix} \geq 0 \Rightarrow \mathbf{A} \geq 0 \text{ and } \mathbf{C} \geq 0 \quad (\text{A.63})$$

The condition also holds if “ \geq ” is replaced by “ $>$ ”.

Proof It follows from (A.62) that

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^H & \mathbf{C} \end{bmatrix} \geq 0 \Leftrightarrow \mathbf{x}^H \mathbf{A} \mathbf{x} + 2 \operatorname{Re}(\mathbf{x}^H \mathbf{B} \mathbf{y}) + \mathbf{y}^H \mathbf{C} \mathbf{y} \geq 0$$

for all $[\mathbf{x} \ \mathbf{y}]^T \neq \mathbf{0}$. Let $\mathbf{x} \neq \mathbf{0}, \mathbf{y} = \mathbf{0}$ then $\mathbf{x}^H \mathbf{A} \mathbf{x} \geq 0$ for all $\mathbf{x} \neq \mathbf{0}$ so that $\mathbf{A} \geq 0$ and let $\mathbf{x} = \mathbf{0}, \mathbf{y} \neq \mathbf{0}$ then $\mathbf{y}^H \mathbf{C} \mathbf{y} \geq 0$ for all $\mathbf{y} \neq \mathbf{0}$ so that $\mathbf{C} \geq 0$. The proof for positive definiteness is similar.

Congruence

Congruence preserves definiteness

$$\mathbf{A} \geq 0 \Leftrightarrow \mathbf{T}^H \mathbf{A} \mathbf{T} \geq 0 \quad (\text{A.64})$$

for invertible \mathbf{T} .

Proof It follows from (A.62) that $\mathbf{A} \geq 0 \Leftrightarrow \mathbf{x}^H \mathbf{A} \mathbf{x} \geq 0, \forall \mathbf{x}$ and with $\mathbf{x} = \mathbf{T} \mathbf{y}$ we have $\mathbf{x}^H \mathbf{A} \mathbf{x} = \mathbf{y}^H \mathbf{T}^H \mathbf{A} \mathbf{T} \mathbf{y}, \forall \mathbf{y} \Leftrightarrow \mathbf{T}^H \mathbf{A} \mathbf{T} \geq 0$. The same is also true for positive and negative (semi)definiteness.

Schur Congruence

If $\mathbf{C} > 0$ then the following is true

$$\begin{aligned} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^H & \mathbf{C} \end{bmatrix} > 0 &\Leftrightarrow \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^H > 0 \\ \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^H & \mathbf{C} \end{bmatrix} \geq 0 &\Leftrightarrow \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^H \geq 0 \end{aligned} \quad (\text{A.65})$$

which is an application of the congruence.

Proof Rearranging the matrix gives

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^H & \mathbf{C} \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{1} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{C} & \mathbf{B}^H \\ \mathbf{B} & \mathbf{A} \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{1} \\ \mathbf{1} & \mathbf{0} \end{bmatrix}$$

Then according to (A.57), the block matrix can be factorized

$$\begin{bmatrix} \mathbf{C} & \mathbf{B}^H \\ \mathbf{B} & \mathbf{A} \end{bmatrix} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{B}\mathbf{C}^{-1} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^H \end{bmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{C}^{-1}\mathbf{B}^H \\ \mathbf{0} & \mathbf{1} \end{bmatrix}$$

Due to (A.64) this leads to

$$\begin{aligned} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^H & \mathbf{C} \end{bmatrix} > 0 &\Leftrightarrow \begin{bmatrix} \mathbf{C} & \mathbf{B}^H \\ \mathbf{B} & \mathbf{A} \end{bmatrix} > 0 \Leftrightarrow \begin{bmatrix} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^H \end{bmatrix} > 0 \\ \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^H & \mathbf{C} \end{bmatrix} \geq 0 &\Leftrightarrow \begin{bmatrix} \mathbf{C} & \mathbf{B}^H \\ \mathbf{B} & \mathbf{A} \end{bmatrix} \geq 0 \Leftrightarrow \begin{bmatrix} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^H \end{bmatrix} \geq 0 \end{aligned}$$

Hence (A.65) follows since $\mathbf{C} > 0$.

S-Lemma

The so called S-lemma states

$$\begin{aligned} \exists \sigma \geq 0 : \mathbf{Q} \geq \sigma \mathbf{P} &\text{ is equivalent to} \\ \forall \mathbf{x} : \mathbf{x}^H \mathbf{P} \mathbf{x} \geq 0 &\Rightarrow \mathbf{x}^H \mathbf{Q} \mathbf{x} \geq 0 \end{aligned} \quad (\text{A.66})$$

if there exists a \mathbf{y} so that $\mathbf{y}^H \mathbf{P} \mathbf{y} > 0$.

A.2.3 Operations

This section provides formulae for special operations on matrices.

Geometric Series

The inverse of $\mathbf{1} - \mathbf{A}$ can be expressed in form of the geometric series

$$(\mathbf{1} - \mathbf{A})^{-1} = \sum_{k=0}^{\infty} \mathbf{A}^k \text{ if } |\lambda_A| < 1 \quad (\text{A.67})$$

for all eigenvalues λ_A of \mathbf{A} .

Proof Let $\mathbf{S} = \sum_{k=0}^{\infty} \mathbf{A}^k$ which only converges if $|\lambda_A| < 1$. Then \mathbf{S} satisfies the equation $\mathbf{A}\mathbf{S} = \mathbf{S} - \mathbf{1}$ so that (A.67) follows.

Scalar Product of Matrices

The trace of the hermitian product defines a scalar product

$$\mathbf{A} \bullet \mathbf{B} = \text{tr}(\mathbf{A}^H \mathbf{B}) = \sum_i \sum_j \bar{a}_{ij} b_{ij} \quad (\text{A.68})$$

Proof This is a scalar product since $\mathbf{A} \bullet \mathbf{B} = \sum_i \sum_j \bar{a}_{ij} b_{ij} = \sum_i \sum_j \overline{a_{ij} b_{ij}} = \overline{\mathbf{B} \bullet \mathbf{A}}$ and $\mathbf{C} \bullet (\alpha \mathbf{A} + \beta \mathbf{B}) = \alpha \text{tr}(\mathbf{C}^H \mathbf{A}) + \beta \text{tr}(\mathbf{C}^H \mathbf{B}) = \alpha \mathbf{C} \bullet \mathbf{A} + \beta \mathbf{C} \bullet \mathbf{B}$ and $\mathbf{A} \bullet \mathbf{A} = \sum_i \sum_j \bar{a}_{ij} a_{ij} = \sum_i \sum_j |a_{ij}|^2 \geq 0$.

General Properties

Aside from that, the scalar product has the properties

$$\mathbf{A} \bullet \mathbf{B} = \mathbf{B}^H \bullet \mathbf{A}^H = \mathbf{A}^T \bullet \mathbf{B}^T \quad (\text{A.69})$$

$$\mathbf{A} \bullet (\mathbf{B}\mathbf{C}) = (\mathbf{B}^H \mathbf{A}) \bullet \mathbf{C} = (\mathbf{A}\mathbf{C}^H) \bullet \mathbf{B} \quad (\text{A.70})$$

$$\alpha \bullet \text{tr} \mathbf{B} = \alpha \mathbf{1} \bullet \mathbf{B} \quad (\text{A.71})$$

$$[\mathbf{A} \ \mathbf{B}] \bullet [\mathbf{C} \ \mathbf{D}] = \mathbf{A} \bullet \mathbf{C} + \mathbf{B} \bullet \mathbf{D} \quad (\text{A.72})$$

Proof The first line is due to $\mathbf{A} \bullet \mathbf{B} = \text{tr}(\mathbf{A}^H \mathbf{B}) = \text{tr}(\mathbf{B}\mathbf{A}^H) = \mathbf{B}^H \bullet \mathbf{A}^H = \text{tr}(\mathbf{B}\mathbf{A}^H)^T = \mathbf{A}^T \bullet \mathbf{B}^T$. The second line follows from $\mathbf{A} \bullet (\mathbf{B}\mathbf{C}) = \text{tr}(\mathbf{A}^H \mathbf{B}\mathbf{C}) = \text{tr}((\mathbf{B}^H \mathbf{A})^H \mathbf{C}) = (\mathbf{B}^H \mathbf{A}) \bullet \mathbf{C} = \text{tr}((\mathbf{A}\mathbf{C}^H)^H \mathbf{B}) = (\mathbf{A}\mathbf{C}^H) \bullet \mathbf{B}$. The third line is due to $\alpha \bullet \text{tr} \mathbf{B} = \text{tr}(\bar{\alpha} \text{tr} \mathbf{B}) = \bar{\alpha} \text{tr} \mathbf{B} = \text{tr}(\bar{\alpha} \mathbf{B}) = \alpha \mathbf{1} \bullet \mathbf{B}$.

Properties for Hermitian Matrices

For hermitian matrices, the scalar product has the additional properties

$$\mathbf{A} = \mathbf{A}^H, \mathbf{B} = \mathbf{B}^H \Rightarrow \mathbf{A} \bullet \mathbf{B} \in \mathbb{R} \quad (\text{A.73})$$

$$\mathbf{A} = \mathbf{A}^H \Rightarrow \begin{cases} \mathbf{A} \bullet (\mathbf{B} + \mathbf{B}^H) = 2 \text{Re}(\mathbf{A} \bullet \mathbf{B}) \\ \mathbf{A} \bullet (\mathbf{B} - \mathbf{B}^H) = 2 \text{Im}(\mathbf{A} \bullet \mathbf{B}) \end{cases} \quad (\text{A.74})$$

$$\mathbf{A} = \mathbf{A}^H \geq 0, \mathbf{B} = \mathbf{B}^H \geq 0 \Rightarrow \mathbf{A} \bullet \mathbf{B} \geq 0 \quad (\text{A.75})$$

Proof For hermitian matrices \mathbf{A}, \mathbf{B} the scalar product is commutative since $\mathbf{A} \bullet \mathbf{B} = \mathbf{B}^H \bullet \mathbf{A}^H = \mathbf{B} \bullet \mathbf{A} = \overline{\mathbf{A} \bullet \mathbf{B}}$ due to (A.69) and the definition of the scalar product. Thus being equal to its conjugate, the scalar product of hermitian matrices is real. The next equality follows from $\mathbf{A} \bullet (\mathbf{B} \pm \mathbf{B}^H) = \mathbf{A} \bullet \mathbf{B} \pm \mathbf{A} \bullet \mathbf{B}^H = \mathbf{A} \bullet \mathbf{B} \pm \mathbf{B} \bullet \mathbf{A}^H = \mathbf{A} \bullet \mathbf{B} \pm \mathbf{B} \bullet \mathbf{A} = \mathbf{A} \bullet \mathbf{B} \pm \overline{\mathbf{A} \bullet \mathbf{B}}$ which gives the double of the real and imaginary parts of $\mathbf{A} \bullet \mathbf{B}$ respectively. Since \mathbf{A} is hermitian it always has a spectral decomposition $\mathbf{A} = \sum_k \lambda_k \mathbf{a}_k \mathbf{a}_k^H$ due to (A.61). Applying this on (A.68) gives $\mathbf{A} \bullet \mathbf{B} = \text{tr}(\mathbf{A}^H \mathbf{B}) = \sum_k \lambda_k \text{tr}(\mathbf{a}_k \mathbf{a}_k^H \mathbf{B}) = \sum_k \lambda_k \text{tr}(\mathbf{a}_k^H \mathbf{B} \mathbf{a}_k) = \sum_k \lambda_k \mathbf{a}_k^H \mathbf{B} \mathbf{a}_k \geq 0$ since $\lambda_k \geq 0$ due to $\mathbf{A} \geq 0$ and $\mathbf{a}_k^H \mathbf{B} \mathbf{a}_k \geq 0$ due to (A.62) and $\mathbf{B} \geq 0$.

Matrix Square Root

A matrix \mathbf{B} is a square root of matrix \mathbf{A} if

$$\mathbf{A} = \mathbf{B}^2 \Leftrightarrow \mathbf{B} = \mathbf{A}^{\frac{1}{2}} \Leftrightarrow \mathbf{U}\mathbf{B}\mathbf{U}^H = (\mathbf{U}\mathbf{A}\mathbf{U}^H)^{\frac{1}{2}} \quad (\text{A.76})$$

$$\mathbf{A} = \mathbf{B}^2 \geq 0 \Rightarrow \mathbf{B} = \mathbf{A}^{\frac{1}{2}} \geq 0 \text{ is unique} \quad (\text{A.77})$$

$$\mathbf{A} = \mathbf{B}^H\mathbf{B} \geq 0 \Rightarrow \mathbf{B} = \mathbf{U}\mathbf{A}^{\frac{1}{2}} \quad (\text{A.78})$$

with \mathbf{U} unitary, that is $\mathbf{U}^H\mathbf{U} = \mathbf{1}$.

Proof The first line is the definition of the matrix square root followed by a property that can be verified by $\mathbf{U}\mathbf{A}\mathbf{U}^H = \mathbf{U}\mathbf{B}^2\mathbf{U}^H = \mathbf{U}\mathbf{B}\mathbf{U}^H\mathbf{U}\mathbf{B}\mathbf{U}^H$ for unitary \mathbf{U} . The second line can be verified since for hermitian matrices and thus positive semidefinite matrices there always exists a spectral decomposition which can be written in the form of $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^H$ with diagonal $\mathbf{D} \geq 0$ and unitary \mathbf{V} . Then due to (A.76) this leads to $\mathbf{B} = \mathbf{A}^{\frac{1}{2}} = (\mathbf{V}\mathbf{D}\mathbf{V}^H)^{\frac{1}{2}} = \mathbf{V}\mathbf{D}^{\frac{1}{2}}\mathbf{V}^H$ and since $\mathbf{D} \geq 0$ is diagonal we have $[\mathbf{D}]_{ii} = d_i \geq 0$ and thus its square root $[\mathbf{D}^{\frac{1}{2}}]_{ii} = d_i^{\frac{1}{2}}$ is unique and positive semidefinite. The third line again makes use of the spectral decomposition $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^H$ which can be written as $\mathbf{A} = \mathbf{V}\mathbf{D}^{\frac{1}{2}}\mathbf{D}^{\frac{1}{2}}\mathbf{V}^H = \mathbf{V}\mathbf{D}^{\frac{1}{2}}\mathbf{V}^H\mathbf{U}^H\mathbf{U}\mathbf{V}\mathbf{D}^{\frac{1}{2}}\mathbf{V}^H$ so that $\mathbf{B} = \mathbf{U}\mathbf{V}\mathbf{D}^{\frac{1}{2}}\mathbf{V}^H = \mathbf{U}\mathbf{A}^{\frac{1}{2}}$.

Pseudoinverse Matrix

A matrix $\mathbf{A} \in \mathbb{C}^{m,n}$ with $m \geq n$ and $\text{rk } \mathbf{A} = n$ has a unique pseudoinverse matrix $\mathbf{A}^+ \in \mathbb{C}^{n,m}$ with $\text{rk } \mathbf{A}^+ = n$ so that

$$\begin{aligned} \mathbf{A}^+\mathbf{A} &= \mathbf{1} \\ (\mathbf{A}\mathbf{A}^+)^H &= \mathbf{A}\mathbf{A}^+ \\ \text{with } \mathbf{A}^+ &= (\mathbf{A}^H\mathbf{A})^{-1}\mathbf{A}^H \end{aligned} \quad (\text{A.79})$$

This is the full column rank case of the Moore Penrose pseudoinverse.

Orthogonal Completion

A matrix $\mathbf{A} \in \mathbb{C}^{m,n}$ with $m > n$ and $\text{rk } \mathbf{A} = n$ has an orthogonal completion $\mathbf{A}_\perp \in \mathbb{C}^{m,m-n}$ with $\text{rk } \mathbf{A}_\perp = m - n$ so that

$$\begin{aligned} \mathbf{A}^H\mathbf{A}_\perp &= \mathbf{0} \\ [\mathbf{A} \quad \mathbf{A}_\perp] \begin{bmatrix} \mathbf{A}^+ \\ \mathbf{A}_\perp^+ \end{bmatrix} &= \mathbf{1} \\ \text{parametrized by } \mathbf{A}_\perp &= \mathbf{A}_\perp\mathbf{0}\mathbf{T} \end{aligned} \quad (\text{A.80})$$

with \mathbf{T} invertible and the pseudoinverse given by (A.79).

Proof Let $\mathbf{A} = [\mathbf{a}_1 \ \cdots \ \mathbf{a}_n]$. Since \mathbf{A} has full rank, all $\mathbf{a}_k \in \mathbb{C}^m$ are linearly independent. Then $m - n$ linearly independent vectors $\mathbf{b}_l \in \mathbb{C}^m$ can be found so that $\mathbf{a}_k^H\mathbf{b}_l = 0; \forall k \in \mathbb{N}_n, \forall l \in \mathbb{N}_{m-n}$, which means that they are all orthogonal on every vector \mathbf{a}_k . The vectors \mathbf{b}_l can then be stacked to form a

matrix $\mathbf{A}_\perp = [\mathbf{b}_1 \ \cdots \ \mathbf{b}_{m-n}] \in \mathbb{C}^{m, m-n}$ with $\text{rk } \mathbf{A}_\perp = m - n$ so that due to their orthogonality on the vectors \mathbf{a}_k we get

$$\mathbf{A}^H \mathbf{A}_\perp = \begin{bmatrix} \mathbf{a}_1^H \mathbf{b}_1 & \cdots & \mathbf{a}_1^H \mathbf{b}_{m-n} \\ \vdots & & \vdots \\ \mathbf{a}_n^H \mathbf{b}_1 & \cdots & \mathbf{a}_n^H \mathbf{b}_{m-n} \end{bmatrix} = \mathbf{0}$$

This property does not depend on the ordering of the column vectors \mathbf{b}_k in \mathbf{A}_\perp so that an invertible matrix \mathbf{T} parametrizes the orthogonal completions in the form of $\mathbf{A}_\perp = \mathbf{A}_\perp \mathbf{T}$. Since matrices $\mathbf{A}, \mathbf{A}_\perp$ have full rank, according to (A.79) there exist pseudoinverse matrices $\mathbf{A}^+ = (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H, \mathbf{A}_\perp^+ = (\mathbf{A}_\perp^H \mathbf{A}_\perp)^{-1} \mathbf{A}_\perp^H$ so that $\mathbf{A}^+ \mathbf{A} = \mathbf{1}, \mathbf{A}_\perp^+ \mathbf{A}_\perp = \mathbf{1}$. Along with the orthogonality $\mathbf{A}^H \mathbf{A}_\perp = \mathbf{0}$, this can be written in the form

$$\begin{bmatrix} \mathbf{A}^+ \\ \mathbf{A}_\perp^+ \end{bmatrix} [\mathbf{A} \ \mathbf{A}_\perp] = \mathbf{1}$$

so that (A.80) follows.

Matrix Variations

The product rule for matrices is given by

$$\delta(\mathbf{A}\mathbf{B}) = \delta\mathbf{A}\mathbf{B} + \mathbf{A}\delta\mathbf{B} \quad (\text{A.81})$$

The variation of an inverse matrix is given by

$$\delta(\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \delta\mathbf{A} \mathbf{A}^{-1} \quad (\text{A.82})$$

The variation of a matrix square root is determined by

$$\delta\mathbf{B} = \delta(\mathbf{B}^{\frac{1}{2}}) \mathbf{B}^{\frac{1}{2}} + \mathbf{B}^{\frac{1}{2}} \delta(\mathbf{B}^{\frac{1}{2}}) \quad (\text{A.83})$$

Proof The product rule for vectors can be deduced from the scalar product rule

$$\delta(\mathbf{a}^T \mathbf{b}) = \delta \left(\sum_{i=1}^k a_i b_i \right) = \sum_{i=1}^k \delta a_i b_i + \sum_{i=1}^k a_i \delta b_i = \delta \mathbf{a}^T \mathbf{b} + \mathbf{a}^T \delta \mathbf{b}$$

so that the product rule for matrices follows

$$\begin{aligned} \delta(\mathbf{A}\mathbf{B}) &= \delta \left(\begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix} [\mathbf{b}_1 \ \cdots \ \mathbf{b}_n] \right) = \delta \begin{bmatrix} \mathbf{a}_1^T \mathbf{b}_1 & \cdots & \mathbf{a}_1^T \mathbf{b}_n \\ \vdots & & \vdots \\ \mathbf{a}_m^T \mathbf{b}_1 & \cdots & \mathbf{a}_m^T \mathbf{b}_n \end{bmatrix} = \\ &= \begin{bmatrix} \delta \mathbf{a}_1^T \mathbf{b}_1 + \mathbf{a}_1^T \delta \mathbf{b}_1 & \cdots & \delta \mathbf{a}_1^T \mathbf{b}_n + \mathbf{a}_1^T \delta \mathbf{b}_n \\ \vdots & & \vdots \\ \delta \mathbf{a}_m^T \mathbf{b}_1 + \mathbf{a}_m^T \delta \mathbf{b}_1 & \cdots & \delta \mathbf{a}_m^T \mathbf{b}_n + \mathbf{a}_m^T \delta \mathbf{b}_n \end{bmatrix} = \delta\mathbf{A}\mathbf{B} + \mathbf{A}\delta\mathbf{B} \end{aligned}$$

For the identity matrix, the variation is zero

$$\delta(\mathbf{A}\mathbf{A}^{-1}) = \delta\mathbf{1} = \mathbf{0}$$

so that the variation of an inverse matrix follows

$$\delta \mathbf{A} \mathbf{A}^{-1} + \mathbf{A} \delta(\mathbf{A}^{-1}) = \delta(\mathbf{A} \mathbf{A}^{-1}) = \mathbf{0} \Leftrightarrow \delta(\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \delta \mathbf{A} \mathbf{A}^{-1}$$

The variation of a squared matrix is given by the product rule

$$\delta(\mathbf{A}^2) = \delta \mathbf{A} \mathbf{A} + \mathbf{A} \delta \mathbf{A}$$

so that the variation of a matrix square root follows by setting $\mathbf{A} = \mathbf{B}^{\frac{1}{2}}$.

A.3 Matrix Equations

This section discusses some important matrix equations and provides their solutions or solution existence conditions.

A.3.1 Underdetermined Linear Equations

The linear equation

$$\mathbf{X} \mathbf{A} = \mathbf{B} \tag{A.84}$$

for $\mathbf{A} \in \mathbb{C}^{m,n}$ with $m \geq n$ and $\text{rk } \mathbf{A} = n$ has solutions

$$\mathbf{X} = \mathbf{B} \mathbf{A}^+ + \mathbf{Y}(\mathbf{1} - \mathbf{A} \mathbf{A}^+) \tag{A.85}$$

with free parameter \mathbf{Y} .

Proof Define $\mathbf{X} \mathbf{A}_\perp = \mathbf{C}$ using the orthogonal completion \mathbf{A}_\perp of \mathbf{A} . This together with (A.84) gives

$$\mathbf{X} \begin{bmatrix} \mathbf{A} & \mathbf{A}_\perp \end{bmatrix} = \begin{bmatrix} \mathbf{B} & \mathbf{C} \end{bmatrix} \Leftrightarrow \mathbf{X} = \begin{bmatrix} \mathbf{B} & \mathbf{C} \end{bmatrix} \begin{bmatrix} \mathbf{A}^+ \\ \mathbf{A}_\perp^+ \end{bmatrix}$$

according to (A.80), so that $\mathbf{X} = \mathbf{B} \mathbf{A}^+ + \mathbf{C} \mathbf{A}_\perp^+$. This is a solution of (A.84) since $\mathbf{X} \mathbf{A} = \mathbf{B} \mathbf{A}^+ \mathbf{A} + \mathbf{C} \mathbf{A}_\perp^+ \mathbf{A} = \mathbf{B}$ which does not depend on the choice of \mathbf{C} so that instead of $\mathbf{C} = \mathbf{X} \mathbf{A}_\perp$ also any other $\mathbf{C} = \mathbf{Y} \mathbf{A}_\perp$ with an arbitrary compatible matrix \mathbf{Y} leads to a valid solution \mathbf{X} . Hence $\mathbf{X} = \mathbf{B} \mathbf{A}^+ + \mathbf{Y} \mathbf{A}_\perp \mathbf{A}_\perp^+ = \mathbf{B} \mathbf{A}^+ + \mathbf{Y}(\mathbf{1} - \mathbf{A} \mathbf{A}^+)$ using (A.80).

A.3.2 Sylvester Equations

The Sylvester matrix equation

$$\mathbf{A} \mathbf{X} + \mathbf{X} \mathbf{B} + \mathbf{C} = \mathbf{0} \tag{A.86}$$

with square matrices \mathbf{A}, \mathbf{B} has a unique solution \mathbf{X} if and only if $\lambda_A + \lambda_B \neq 0$ for all eigenvalues λ_A, λ_B of \mathbf{A}, \mathbf{B} respectively. For the important special case $\mathbf{B} = \mathbf{A}^H$ there is an additional property

$$\mathbf{B} = \mathbf{A}^H \text{ and } \mathbf{C} \geq 0 \Rightarrow \mathbf{X} \geq 0 \tag{A.87}$$

The same is true if $\mathbf{C} > 0$ then $\mathbf{X} > 0$.

A.3.3 Stein Equations

The Sylvester matrix equation

$$\mathbf{A}\mathbf{X}\mathbf{B} - \mathbf{X} + \mathbf{C} = \mathbf{0} \quad (\text{A.88})$$

with square matrices \mathbf{A}, \mathbf{B} has a unique solution \mathbf{X} if and only if $\lambda_A \lambda_B \neq 1$ for all eigenvalues λ_A, λ_B of \mathbf{A}, \mathbf{B} respectively. For the important special case $\mathbf{B} = \mathbf{A}^H$ there is an additional property

$$\mathbf{B} = \mathbf{A}^H \text{ and } \mathbf{C} \geq 0 \Rightarrow \mathbf{X} \geq 0 \quad (\text{A.89})$$

The same is true if $\mathbf{C} > 0$ then $\mathbf{X} > 0$.

Appendix B

Analysis

B.1 Complex Analysis

This section provides some important theorems of complex analysis.

B.1.1 Analytic Functions

A function f is analytic on a domain $D \subset \mathbb{C}$ if it is equal to its unique Taylor series representation

$$f(z) = \sum_{k=0}^{\infty} c_k(z - z_0)^k \Leftrightarrow f \text{ analytic on } D \quad (\text{B.1})$$

for any $z_0 \in D$.

B.1.2 Isolated Singularity

A function f has an isolated singularity at $z_0 \in \mathbb{C}$ if it is analytic on a domain $D \setminus z_0 \subset \mathbb{C}$. Then it is equal to its unique Laurent series representation

$$f(z) = \sum_{k=-\infty}^{\infty} c_k(z - z_0)^k \Leftrightarrow f \text{ analytic on } D \setminus z_0 \quad (\text{B.2})$$

B.1.3 Residue

The residue of a function f with an isolated singularity at $z_0 \in \mathbb{C}$ is defined as

$$\text{Res } f = \frac{1}{2\pi j} \oint_C f(z) dz = c_{-1} \quad (\text{B.3})$$

where C is a circle around z_0 and c_{-1} is the coefficient of $(z - z_0)^{-1}$ in the Laurent series representation of f as in (B.2).

Proof Parametrizing the circle C by $C = \{z : z = z_0 + re^{j2\pi t}; \forall t \in [0, 1], r > 0\}$ so that $dz = 2\pi j(z - z_0)dt$ and then integrating the Laurent series representation of f leads to

$$\frac{1}{2\pi j} \oint_C f(z)dz = \frac{1}{2\pi j} \sum_{k=-\infty}^{\infty} c_k \oint_C (z - z_0)^k dz = \sum_{k=-\infty}^{\infty} c_k \oint_C (re^{j2\pi t})^{k+1} dt = c_{-1}$$

B.1.4 Identity Theorem

If a function f that is analytic on an open and connected domain $D \subset \mathbb{C}$ vanishes on a small region $R \subset D$, then it also vanishes on the entire domain D :

$$f(R) = 0 \Leftrightarrow f(D) = 0, \text{ if } f \text{ analytic on } D \text{ and } R \subset D \quad (\text{B.4})$$

and $f(D)$ is then the analytic continuation of $f(R)$.

Proof From $R \subset D$ the direction $f(D) = 0 \Rightarrow f(R) = 0$ follows immediately. To verify the reverse direction, assume $f(R) = 0$. Then due to (B.1) the Taylor series on $z \in R$ around an arbitrary $z_0 \in R$ is given by $f(z) = \sum_{k=0}^{\infty} \frac{1}{k!} \frac{\partial^k f(z_0)}{\partial^k z} (z - z_0)^k = 0$ for all $z \in R$. This can be derived so that $\frac{\partial^l f(z)}{\partial^l z} = \sum_{k=0}^{\infty} \frac{1}{k!} \frac{\partial^{k+l} f(z_0)}{\partial^{k+l} z} (z - z_0)^k = 0$ for all $z \in R$ and for all $l \in \mathbb{Z}_+$ so that we get $\frac{\partial^k f(z_0)}{\partial^k z} = 0$ for all $k \in \mathbb{Z}_+$. Since f is analytic on D the Taylor series on $z \in D$ around z_0 is the same, so that $f(z) = \sum_{k=0}^{\infty} \frac{1}{k!} \frac{\partial^k f(z_0)}{\partial^k z} (z - z_0)^k = 0$ since the gradients are all zero in z_0 and thus (B.4) is verified.

B.1.5 Open Mapping Theorem

If a function f is analytic on an open and connected domain $D \subset \mathbb{C}$ then either $f(D)$ is also open or f is constant:

$$D \text{ open} \Rightarrow \text{either } f(D) \text{ open or } f \text{ constant, if } f \text{ analytic on } D \quad (\text{B.5})$$

Proof Define an infinitesimally small circular region $R(z_0) = \{z : z - z_0 = \epsilon e^{j\phi}\}$ around $z_0 \in D$. Since $R(z_0)$ is infinitesimally small and f is analytic, the Taylor series of f around z_0 can be reduced to $f(R(z_0)) = f(z_0) + \frac{\partial f(z_0)}{\partial z} (z - z_0) = f(z_0) + \frac{\partial f(z_0)}{\partial z} \epsilon e^{j\phi}$ and with $\frac{\partial f(z_0)}{\partial z} = g e^{j\phi_0}$ it can be seen that $f(R(z_0)) - f(z_0) = g \epsilon e^{j(\phi + \phi_0)}$. On the one hand if $g \neq 0$, $f(R(z_0))$ also defines an infinitesimally small circular region around $f(z_0)$ and therefore $f(z_0)$ cannot be a boundary point. This reasoning is valid for every $z_0 \in D$ so that $f(D)$ must be open. On the other hand if $g = 0$ then $f(R(z_0)) - f(z_0) = 0$ and using (B.4) we get $f(D) - f(z_0) = 0$ and thus f is constant.

B.1.6 Maximum Modulus Theorem

If a function f is analytic on an open and connected domain $D \subset \mathbb{C}$ then its modulus $|f|$ is maximal at the closure ∂D of the domain D :

$$\max_{z \in D} |f(z)| = \max_{z \in \partial D} |f(z)|, \text{ if } f \text{ analytic on } D \quad (\text{B.6})$$

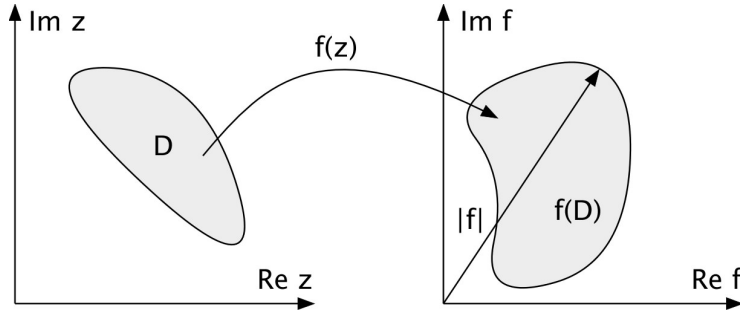


Figure B.1: The modulus of f has its extrema on the closure of $f(D)$.

Proof As can be seen in figure B.1, the modulus $|f|$ has its maximum on the boundary ∂D and (B.5) states that either $f(D)$ is open or f is constant. So either a $z \in \partial D$ provides a maximal $|f|$ and (B.6) is verified or f is constant and (B.6) is also true.

B.2 Norms

This section introduces the norms as a universal measure of length for vectors, matrices, signals and systems.

B.2.1 Norms for Constants

The matrix and induced matrix norms measure the length of constant matrices and transformation matrices.

The p Norm

For a constant matrix $\mathbf{A} \in \mathbb{C}^{m,n}$ the p norm is defined by

$$\|\mathbf{A}\|_p = \sqrt[p]{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^p} \quad (\text{B.7})$$

Usually $p \in \{1, 2, \infty\}$ and in the case $p = \infty$ the p norm is given by

$$\|\mathbf{A}\|_\infty = \max_{i \in \mathbb{N}_m, j \in \mathbb{N}_n} |a_{ij}| \quad (\text{B.8})$$

The i_p Norm

For a linear transformation $\mathbf{y} = \mathbf{A}\mathbf{x}$ with $\mathbf{A} \in \mathbb{C}^{m,n}$ the i_p norm is defined by

$$\|\mathbf{A}\|_{i_p} = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{y}\|_p}{\|\mathbf{x}\|_p} = \max_{\|\mathbf{x}\|_p=1} \|\mathbf{A}\mathbf{x}\|_p \quad (\text{B.9})$$

and the ip norm is induced by the matrix p norm. Usually $p \in \{1, 2, \infty\}$ and in these cases, the ip norm is given by

$$\|\mathbf{A}\|_{i1} = \max_{j \in \mathbb{N}_n} \sum_{i=1}^m |a_{ij}| \quad (\text{B.10})$$

$$\|\mathbf{A}\|_{i2} = \sigma_{\max} \quad (\text{B.11})$$

$$\|\mathbf{A}\|_{i\infty} = \max_{i \in \mathbb{N}_m} \sum_{j=1}^n |a_{ij}| \quad (\text{B.12})$$

whereby σ_{\max} is the largest singular value of \mathbf{A} .

B.2.2 Norms for Functions with Real Variables

The Lebesgue norm measures the length of a function that depends on a real or integer variable.

The Lp Norm

For continuous function f_c and discrete function f_d defined on domains $D_c \subset \mathbb{R}, D_d \subset \mathbb{Z}$ respectively, the Lp norm is defined by

$$\begin{aligned} \|f_c(D_c)\|_{Lp} &= \sqrt[p]{\int_{D_c} \|f_c(t)\|^p dt} \\ \|f_d[D_d]\|_{Lp} &= \sqrt[p]{\sum_{D_d} \|f_d[k]\|^p} \end{aligned} \quad (\text{B.13})$$

Any matrix norm of $f_c(t), f_d[k]$ can be used. Usually $p \in \{1, 2, \infty\}$ and in the case $p = \infty$ the Lp norm is given by

$$\begin{aligned} \|f_c(D_c)\|_{L\infty} &= \sup_{t \in D_c} \|f_c(t)\| \\ \|f_d[D_d]\|_{L\infty} &= \sup_{k \in D_d} \|f_d[k]\| \end{aligned} \quad (\text{B.14})$$

B.2.3 Norms for Functions with Complex Variables

The Hardy norm measures the length of a function that depends on a complex variable.

The Hp Norm

For a function F analytic on a domain $D \subset \mathbb{C}$ the Hp norm is defined by

$$\|F(D)\|_{Hp} = \|F(\partial D)\|_{Lp} \quad (\text{B.15})$$

Proof This is due to the matrix version of (B.6) which states that $\|F\|$ is largest on the boundary ∂D and with the boundary parametrized by a $t \in \mathbb{R}$, the problem of finding the largest $\|F\|$ over a complex domain reduces to the Lp problem of finding the largest $\|F\|$ over a real domain.

Example In the case of F_c analytic on $D_c = \{s \in \mathbb{C} : \operatorname{Re} s \geq 0\}$ and F_d analytic on $D_d = \{z \in \mathbb{C} : |z| \geq 1\}$ the Hp norm is given by

$$\begin{aligned} \|F_c(D_c)\|_{Hp} &= \|F_c(\partial D_c)\|_{Lp} = \|F_c(j2\pi\mathbb{R})\|_{Lp} \\ \|F_d(D_d)\|_{Hp} &= \|F_d(\partial D_d)\|_{Lp} = \|F_d(e^{j2\pi\mathbb{R}})\|_{Lp} \end{aligned} \quad (\text{B.16})$$

The examples are important in the study of continuous (index c) and discrete (index d) stable systems, where the system transfer function is analytic in D_c and D_d respectively. The factor 2π appearing in these definitions stems from the normalization condition provided by the Parseval Identity (C.9).

B.3 Optimization

This section provides some helpful results in the domain of constrained optimization, with the constraint given in the form of a matrix inequality. Details can be found in [4].

B.3.1 Problem Formulation

An optimization problem of minimizing a scalar function f under matrix inequality constraint $\mathbf{G} = \mathbf{G}^H$ may be formulated in the following way

$$f(\mathbf{x}^*) = \inf_{\mathbf{x} \in F} f(\mathbf{x}) \text{ with } F = \{\mathbf{x} \in X : \mathbf{G}(\mathbf{x}) \leq 0\} \quad (\text{B.17})$$

with \mathbf{x} being the vector of variables over a space X . The set F is the feasible set, the set of all \mathbf{x} that satisfy the constraint \mathbf{G} . A maximization of f can be formulated as a minimization of $-f$.

B.3.2 Lagrange Function

The constrained problem (B.17) can be formulated as an equivalent unconstrained problem

$$f(\mathbf{x}^*) = \inf_{\mathbf{x} \in X} \sup_{\mathbf{Y} \geq 0} L(\mathbf{x}, \mathbf{Y}) \quad (\text{B.18})$$

with the Lagrange multiplier $\mathbf{Y} = \mathbf{Y}^H \geq 0$ and the Lagrange function defined as

$$L(\mathbf{x}, \mathbf{Y}) = f(\mathbf{x}) + \mathbf{Y} \bullet \mathbf{G}(\mathbf{x}) \quad (\text{B.19})$$

Proof This stems from the fact that

$$\sup_{\mathbf{Y} \geq 0} L(\mathbf{x}, \mathbf{Y}) = f(\mathbf{x}) + \begin{cases} 0; & \mathbf{G}(\mathbf{x}) \leq 0 \\ \infty; & \text{else} \end{cases}$$

since in the case $\mathbf{G}(\mathbf{x}) \leq 0$ the upper limit follows by applying (A.75) on the hermitian matrices $\mathbf{Y} \geq 0$, $-\mathbf{G} \geq 0$ and in the other case there is no upper limit to the scalar product $\mathbf{Y} \bullet \mathbf{G}$. This means that in order to get $L(\mathbf{x}, \mathbf{Y})$ minimal, \mathbf{x} must be in the feasible set F so that $\sup_{\mathbf{Y} \geq 0} L(\mathbf{x}, \mathbf{Y}) = f(\mathbf{x})$ and hence (B.18) is equivalent to (B.17) in the sense that they share the same minimizing \mathbf{x}^* and the same minimal value $f(\mathbf{x}^*)$.

B.3.3 Primal and Dual Problem

To the primal optimization problem (B.17) and equivalently (B.18), a dual problem in the form of

$$g(\mathbf{Y}^*) = \sup_{\mathbf{Y} \geq 0} \inf_{\mathbf{x} \in X} L(\mathbf{x}, \mathbf{Y}) \quad (\text{B.20})$$

can be assigned due to the symmetry to (B.18). In general

$$f(\mathbf{x}^*) \geq g(\mathbf{Y}^*) \Leftrightarrow \text{weak duality} \quad (\text{B.21})$$

which is called weak duality. In the case

$$f(\mathbf{x}^*) = g(\mathbf{Y}^*) \Leftrightarrow \text{strong duality} \quad (\text{B.22})$$

which is called strong duality, the Lagrange function L has a saddle point at $L(\mathbf{x}^*, \mathbf{Y}^*) = f(\mathbf{x}^*) = g(\mathbf{Y}^*)$. This case is of high practical importance which is due to the large number of problems that can be found to exhibit strong duality and the performance of the solvers that specialize on these problems. It also allows the parameter perturbation analysis.

B.3.4 First Order Perturbation Analysis

The problem (B.17) may depend on parameters \mathbf{z} that can be made explicit in the form of

$$f(\mathbf{x}^*, \mathbf{z}) = \inf_{\mathbf{x} \in F(\mathbf{z})} f(\mathbf{x}, \mathbf{z}) \text{ with } F(\mathbf{z}) = \{\mathbf{x} \in X : \mathbf{G}(\mathbf{x}, \mathbf{z}) \leq 0\} \quad (\text{B.23})$$

in order to allow the analysis of the effects of small parameter changes on the optimal value $f(\mathbf{x}^*, \mathbf{z})$. In the case of strong duality (B.22) an upper limit on the first order effects on the optimal value can be given by

$$f(\mathbf{x}^*, \mathbf{z}) - f(\mathbf{x}_0^*, \mathbf{z}_0) = df \leq \delta L = \frac{\partial L(\mathbf{x}_0^*, \mathbf{Y}_0^*, \mathbf{z}_0)}{\partial \mathbf{z}} \delta \mathbf{z} \quad (\text{B.24})$$

with \mathbf{x}_0^* and \mathbf{Y}_0^* being the solution of the primal and dual problems for a nominal parameter vector \mathbf{z}_0 and \mathbf{x}^* and \mathbf{Y}^* being the solution of the primal and dual problems for a perturbed parameter vector $\mathbf{z} = \mathbf{z}_0 + \delta \mathbf{z}$. The parameter dependent Lagrange function is defined as

$$L(\mathbf{x}, \mathbf{Y}, \mathbf{z}) = f(\mathbf{x}, \mathbf{z}) + \mathbf{Y} \bullet \mathbf{G}(\mathbf{x}, \mathbf{z}) \quad (\text{B.25})$$

in analogy to (B.19).

Proof This upper limit can be verified by using strong duality (B.22) and (B.20) to get

$$f(\mathbf{x}^*, \mathbf{z}) = L(\mathbf{x}^*, \mathbf{Y}^*, \mathbf{z}) = \inf_{\mathbf{x} \in X} L(\mathbf{x}, \mathbf{Y}^*, \mathbf{z}) \leq L(\mathbf{x}^*, \mathbf{Y}^*, \mathbf{z})$$

For $\mathbf{x}_0^* \in X$ this leads to

$$\begin{aligned} f(\mathbf{x}^*, \mathbf{z}) &\leq L(\mathbf{x}_0^*, \mathbf{Y}^*, \mathbf{z}) = \\ &= \underbrace{L(\mathbf{x}_0^*, \mathbf{Y}_0^*, \mathbf{z}_0)}_{f(\mathbf{x}_0^*, \mathbf{z}_0)} + \underbrace{\frac{\partial L(\mathbf{x}_0^*, \mathbf{Y}_0^*, \mathbf{z}_0)}{\partial \mathbf{Y}} \delta \mathbf{Y}}_{\delta \mathbf{Y} \bullet \mathbf{G}(\mathbf{x}_0^*, \mathbf{z}_0)} + \frac{\partial L(\mathbf{x}_0^*, \mathbf{Y}_0^*, \mathbf{z}_0)}{\partial \mathbf{z}} \delta \mathbf{z} \end{aligned}$$

with $\mathbf{Y}^* = \mathbf{Y}_0^* + \delta\mathbf{Y}$. An upper limit for the scalar product can be given by

$$\delta\mathbf{Y} \bullet \mathbf{G}(\mathbf{x}_0^*, z_0) = \underbrace{\overbrace{\mathbf{Y}^* \bullet \mathbf{G}(\mathbf{x}_0^*, z_0)}^{\geq 0}}_{\leq 0} - \underbrace{\mathbf{Y}_0^* \bullet \mathbf{G}(\mathbf{x}_0^*, z_0)}_{=0}$$

The first term is non positive due to (A.75) applied on the hermitian matrices $\mathbf{Y}^* \geq 0, -\mathbf{G}(\mathbf{x}_0^*, z_0) \geq 0$. The second term is zero because with strong duality we have

$$f(\mathbf{x}_0^*, z_0) = L(\mathbf{x}_0^*, \mathbf{Y}_0^*, z_0) = f(\mathbf{x}_0^*, z_0) + \mathbf{Y}_0^* \bullet \mathbf{G}(\mathbf{x}_0^*, z_0)$$

so that $\delta\mathbf{Y} \bullet \mathbf{G}(\mathbf{x}_0^*, z_0) \leq 0$ and (B.24) follows. A more rigorous proof and further results can be found in [4].

Appendix C

Systems

C.1 Linear Systems

This section provides some important characteristics of linear time invariant systems.

C.1.1 System Representations

A linear time invariant system with input $\mathbf{u} \in \mathbb{R}^m$, output $\mathbf{y} \in \mathbb{R}^n$ and state $\mathbf{x} \in \mathbb{C}^p$ may be time continuous or discrete. Continuous (discrete) systems can be represented by a set of linear first order differential (difference) equations in the following standard form

$$\begin{aligned}\dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) \\ \mathbf{y}(t) &= \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)\end{aligned}\tag{C.1}$$

for continuous systems and

$$\begin{aligned}\mathbf{x}[k+1] &= \mathbf{A}\mathbf{x}[k] + \mathbf{B}\mathbf{u}[k] \\ \mathbf{y}[k] &= \mathbf{C}\mathbf{x}[k] + \mathbf{D}\mathbf{u}[k]\end{aligned}\tag{C.2}$$

for discrete systems. These are the state space representations in time domain of these systems. The matrices \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{D} are often called *dynamic*, *input*, *output* and *throughput* matrices respectively. These systems being linear, there exist integral transforms

$$\mathbf{x}(s) = \int_{\mathbb{R}_+} \mathbf{x}(t)e^{-st} dt\tag{C.3}$$

for continuous systems and

$$\mathbf{x}(z) = \sum_{\mathbb{Z}_+} \mathbf{x}[k]z^{-k}\tag{C.4}$$

for discrete systems, that transform these systems to a set of linear algebraic equations

$$\mathbf{y}(s) = \mathbf{P}(s)\mathbf{u}(s)\tag{C.5}$$

for continuous systems and

$$\mathbf{y}(z) = \mathbf{P}(z)\mathbf{u}(z) \quad (\text{C.6})$$

for discrete systems. These are the transfer function representations in frequency domain of these systems. The matrix of transfer functions \mathbf{P} is called the transfer matrix. The link between state space and transfer function representation can be given by

$$\mathbf{P}(x) = \mathbf{C}(x\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D} =: \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \hline \mathbf{C} & \mathbf{D} \end{array} \right] \quad (\text{C.7})$$

with $x \in \{s, z\}$ for continuous and discrete systems respectively. There is also a link between a continuous and a discrete system given by

$$\mathbf{x}[k] = \mathbf{x}(kT) \text{ and } z = e^{sT} \quad (\text{C.8})$$

with T being the sample time of the discrete system.

Proof This can be seen by approximating the integral in (C.3) as a sum

$$\mathbf{x}(s) \approx \sum_{\mathbb{Z}_+} \mathbf{x}(kT)e^{-skT}$$

and comparing this to (C.4).

C.1.2 Parseval Theorem

The transformations (C.3) and (C.4) relate the L2 norm of continuous $\mathbf{x}_c(t)$ and discrete $\mathbf{x}_d[k]$ to the H2 norm of their respective spectra $\mathbf{x}_c(s)$ and $\mathbf{x}_d(z)$ in the form of

$$\begin{aligned} \|\mathbf{x}_c(t)\|_{L_2} &= \|\mathbf{x}_c(s)\|_{H_2} \\ \|\mathbf{x}_d[k]\|_{L_2} &= \|\mathbf{x}_d(z)\|_{H_2} \end{aligned} \quad (\text{C.9})$$

which is called the *Parseval theorem*.

Proof This can be verified for a discrete signal $\mathbf{x}[k]$. According to (B.16) the square of the H2 norm of $\mathbf{x}(z)$ is given by

$$\|\mathbf{x}(z)\|_{H_2}^2 = \|\mathbf{x}(e^{j2\pi t})\|_{L_2}^2 = \int_0^1 \|\mathbf{x}(e^{j2\pi t})\|^2 dt$$

The 2 norm (B.7) is used so that $\|\mathbf{x}(e^{j2\pi t})\|^2 = \mathbf{x}^H(e^{j2\pi t})\mathbf{x}(e^{j2\pi t})$. With the transformation (C.4) this leads to

$$\begin{aligned} \|\mathbf{x}(z)\|_{H_2}^2 &= \int_0^1 \left(\sum_{k=0}^{\infty} \mathbf{x}^H[k]e^{j2\pi kt} \right) \left(\sum_{l=0}^{\infty} \mathbf{x}[l]e^{-j2\pi lt} \right) dt = \\ &= \sum_{k=0}^{\infty} \sum_{l=0}^{\infty} \mathbf{x}^H[k]\mathbf{x}[l] \int_0^1 e^{j2\pi(k-l)t} dt \end{aligned}$$

Note that

$$\int_0^1 e^{j2\pi kt} dt = \begin{cases} 1; & k = 0 \\ 0; & k \in \mathbb{Z} \setminus \{0\} \end{cases}$$

therefore

$$\|\mathbf{x}(z)\|_{H_2}^2 = \sum_{k=0}^{\infty} \mathbf{x}^H[k] \mathbf{x}[k] = \|\mathbf{x}[k]\|_{L_2}^2$$

and (C.9) follows for $\mathbf{x}_d = \mathbf{x}$. The proof for continuous systems is similar.

C.1.3 L2/H2 System Norm

The L2/H2 norm of a discrete system \mathbf{P} as in (C.2) can be given explicitly by

$$\|\mathbf{P}[k]\|_{L_2} = \|\mathbf{P}(z)\|_{H_2} = \sqrt{\text{tr}(\mathbf{B}^H \mathbf{G} \mathbf{B} + \mathbf{D}^H \mathbf{D})} \quad (\text{C.10})$$

with \mathbf{G} being the Gram observability matrix that according to (C.34) solves

$$\mathbf{A}^H \mathbf{G} \mathbf{A} - \mathbf{G} + \mathbf{C}^H \mathbf{C} = \mathbf{0} \quad (\text{C.11})$$

using the 2 norm as in (B.7). This can also be stated in inequality form by

$$\|\mathbf{P}[k]\|_{L_2} = \|\mathbf{P}(z)\|_{H_2} < \gamma \quad (\text{C.12})$$

if and only if there exists a $\mathbf{Z} > 0$ so that

$$\begin{aligned} \mathbf{A}^H \mathbf{Z} \mathbf{A} - \mathbf{Z} + \mathbf{C}^H \mathbf{C} &< 0 \\ \text{tr}(\mathbf{B}^H \mathbf{G} \mathbf{B} + \mathbf{D}^H \mathbf{D}) &< \gamma^2 \end{aligned} \quad (\text{C.13})$$

Proof Due to (C.9) we have $\|\mathbf{P}[k]\|_{L_2} = \|\mathbf{P}(z)\|_{H_2}$. The impulse response $\mathbf{P}[k]$ of the discrete system (C.2) is given by

$$\mathbf{P}[k] = \begin{cases} \mathbf{D}; & k = 0 \\ \mathbf{C} \mathbf{A}^{k-1} \mathbf{B}; & k > 0 \end{cases}$$

so that the L2 norm can be calculated directly

$$\begin{aligned} \|\mathbf{P}[k]\|_{L_2}^2 &= \sum_{k=0}^{\infty} \|\mathbf{P}[k]\|^2 = \text{tr}(\mathbf{D}^H \mathbf{D}) + \sum_{k=0}^{\infty} \text{tr}(\mathbf{B}^H \mathbf{A}^{Hk} \mathbf{C}^H \mathbf{C} \mathbf{A}^k \mathbf{B}) = \\ &= \text{tr}(\mathbf{D}^H \mathbf{D} + \mathbf{B}^H \left(\sum_{k=0}^{\infty} \mathbf{A}^{Hk} \mathbf{C}^H \mathbf{C} \mathbf{A}^k \right) \mathbf{B}) \end{aligned}$$

according to (B.13) while using the 2 norm (B.7). It can be easily verified that the series $\mathbf{G} = \sum_{k=0}^{\infty} \mathbf{A}^{Hk} \mathbf{C}^H \mathbf{C} \mathbf{A}^k$ solves the Stein equation (C.11) so that the result (C.10) is obtained.

Now assume \mathbf{A} stable and (C.12). Let $\mathbf{Z}(\epsilon) = \sum_{k=0}^{\infty} \mathbf{A}^{Hk} (\mathbf{C}^H \mathbf{C} + \epsilon \mathbf{1}) \mathbf{A}^k$ so that $\mathbf{Z}(0) = \mathbf{G}$. Since $\mathbf{Z}(\epsilon)$ is a continuous function there exists a small enough $\epsilon_0 > 0$ so that $\text{tr}(\mathbf{B}^H \mathbf{Z}(\epsilon_0) \mathbf{B} + \mathbf{D}^H \mathbf{D}) < \gamma^2$. It can be easily verified that $\mathbf{Z}(\epsilon_0)$ solves the Stein equation $\mathbf{A}^H \mathbf{Z}(\epsilon_0) \mathbf{A} - \mathbf{Z}(\epsilon_0) + \mathbf{C}^H \mathbf{C} + \epsilon_0 \mathbf{1} = \mathbf{0}$. Due to \mathbf{A} stable and $\mathbf{C}^H \mathbf{C} + \epsilon_0 \mathbf{1} > 0$ we get $\mathbf{Z}(\epsilon_0) > 0$ according to (A.89).

Thus $\mathbf{A}^H \mathbf{Z}(\epsilon_0) \mathbf{A} - \mathbf{Z}(\epsilon_0) + \mathbf{C}^H \mathbf{C} < 0$ so that (C.13) follows. In order to verify the reverse direction assume (C.13). Then there exists an $\mathbf{E} > 0$ so that $\mathbf{A}^H \mathbf{Z} \mathbf{A} - \mathbf{Z} + \mathbf{C}^H \mathbf{C} + \mathbf{E} = \mathbf{0}$. Since $\mathbf{C}^H \mathbf{C} + \mathbf{E} > 0$ and $\mathbf{Z} > 0$ the dynamic matrix \mathbf{A} is stable and the L2/H2 norm is given by (C.10) with Gram observability matrix \mathbf{G} being the unique positive definite solution of (C.11). Thus $\mathbf{C}^H \mathbf{C} = \mathbf{G} - \mathbf{A}^H \mathbf{G} \mathbf{A}$ so that $\mathbf{A}^H (\mathbf{Z} - \mathbf{G}) \mathbf{A} - (\mathbf{Z} - \mathbf{G}) + \mathbf{E} = \mathbf{0}$. Since \mathbf{A} is stable and $\mathbf{E} > 0$ this Stein equation has a unique positive definite solution $\mathbf{Z} - \mathbf{G} > 0$ so that $\gamma^2 > \text{tr}(\mathbf{B}^H \mathbf{Z} \mathbf{B} + \mathbf{D}^H \mathbf{D}) \geq \text{tr}(\mathbf{B}^H \mathbf{G} \mathbf{B} + \mathbf{D}^H \mathbf{D})$ and (C.12) follows.

Remark A dual result to (C.10), (C.11) can be obtained using the Gram controllability matrix. The L2/H2 norm of \mathbf{P} as in (C.2) can be given by

$$\|\mathbf{P}[k]\|_{L_2} = \|\mathbf{P}(z)\|_{H_2} = \sqrt{\text{tr}(\mathbf{B}^H \mathbf{G} \mathbf{B} + \mathbf{D}^H \mathbf{D})}$$

with \mathbf{G} being the Gram controllability matrix that according to (C.32) solves

$$\mathbf{A} \mathbf{G} \mathbf{A}^H - \mathbf{G} + \mathbf{B} \mathbf{B}^H = \mathbf{0}$$

using the 2 norm as in (B.7).

C.1.4 Star Norm

The concept of the star norm used in section 4.1.2 as an upper bound on the L1 norm was developed in [10] and thoroughly treated in [1] for the continuous case. The discrete case was treated in [9], [12]. An alternative proof of the discrete case matrix inequalities based on the concept given in [9] is presented in the following.

Proposition Consider $\mathbf{z} = \mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{y}$ and let X, Y, Z be the spaces of $\mathbf{x}, \mathbf{y}, \mathbf{z}$ respectively and define

$$\begin{aligned} X_R &= \{\mathbf{x} : \mathbf{x}^H \mathbf{R} \mathbf{x} \leq 1\} \\ Y_S &= \{\mathbf{y} : \mathbf{y}^H \mathbf{S} \mathbf{y} \leq 1\} \\ Z_T &= \{\mathbf{z} : \mathbf{z}^H \mathbf{T} \mathbf{z} \leq 1\} \end{aligned} \quad (\text{C.14})$$

with $\mathbf{R} \geq 0, \mathbf{S} \geq 0, \mathbf{T} \geq 0$ but $\mathbf{R} \neq 0, \mathbf{S} \neq 0$. Then

$$\begin{aligned} \forall \mathbf{x} \in X_R, \forall \mathbf{y} \in Y_S : \mathbf{z} \in Z_T &\Leftrightarrow \\ \Leftrightarrow \forall \mathbf{x} \in X, \forall \mathbf{y} \in Y : \mathbf{x}^H \mathbf{R} \mathbf{x} \leq \mathbf{y}^H \mathbf{S} \mathbf{y} &\Rightarrow \mathbf{z}^H \mathbf{T} \mathbf{z} \leq \mathbf{y}^H \mathbf{S} \mathbf{y} \end{aligned} \quad (\text{C.15})$$

$$\Leftrightarrow \forall \mathbf{x} \in X, \forall \mathbf{y} \in Y : \mathbf{y}^H \mathbf{S} \mathbf{y} \leq \mathbf{x}^H \mathbf{R} \mathbf{x} \Rightarrow \mathbf{z}^H \mathbf{T} \mathbf{z} \leq \mathbf{x}^H \mathbf{R} \mathbf{x} \quad (\text{C.16})$$

$$\Leftrightarrow \forall \mathbf{x} \in X, \forall \mathbf{y} \in Y, \exists \alpha \in [0, 1] : \mathbf{z}^H \mathbf{T} \mathbf{z} \leq \alpha \mathbf{x}^H \mathbf{R} \mathbf{x} + (1 - \alpha) \mathbf{y}^H \mathbf{S} \mathbf{y} \quad (\text{C.17})$$

$$\Leftrightarrow \exists \alpha \in [0, 1] : \begin{bmatrix} \mathbf{A}^H \mathbf{T} \mathbf{A} - \alpha \mathbf{R} & \mathbf{A}^H \mathbf{T} \mathbf{B} \\ \mathbf{B}^H \mathbf{T} \mathbf{A} & \mathbf{B}^H \mathbf{T} \mathbf{B} - (1 - \alpha) \mathbf{S} \end{bmatrix} \leq 0 \quad (\text{C.18})$$

Proof On the one hand suppose $\exists \mathbf{x} \in X, \exists \mathbf{y} \in Y : \mathbf{x}^H \mathbf{R} \mathbf{x} \leq \mathbf{y}^H \mathbf{S} \mathbf{y} \Rightarrow \mathbf{z}^H \mathbf{T} \mathbf{z} > \mathbf{y}^H \mathbf{S} \mathbf{y}$. Hence there exists a $\beta > 0$ so that $\mathbf{z}^H \mathbf{T} \mathbf{z} > \beta^2 > \mathbf{y}^H \mathbf{S} \mathbf{y}$. Now define $\bar{\mathbf{x}} = \mathbf{x}/\beta, \bar{\mathbf{y}} = \mathbf{y}/\beta, \bar{\mathbf{z}} = \mathbf{z}/\beta$ then $\bar{\mathbf{z}} = \mathbf{A} \bar{\mathbf{x}} + \mathbf{B} \bar{\mathbf{y}}$ and $\bar{\mathbf{z}}^H \mathbf{T} \bar{\mathbf{z}} > 1$ as well

as $\bar{\mathbf{x}}^H \mathbf{R} \bar{\mathbf{x}} < 1$ and $\bar{\mathbf{y}}^H \mathbf{S} \bar{\mathbf{y}} < 1$. Thus $\exists \bar{\mathbf{x}} \in X_R, \exists \bar{\mathbf{y}} \in Y_S : \bar{\mathbf{z}} \notin Z_T$ which verifies sufficiency in (C.15).

On the other hand suppose $\forall \mathbf{x} \in X, \forall \mathbf{y} \in Y : \mathbf{x}^H \mathbf{R} \mathbf{x} \leq \mathbf{y}^H \mathbf{S} \mathbf{y} \Rightarrow \mathbf{z}^H \mathbf{T} \mathbf{z} \leq \mathbf{y}^H \mathbf{S} \mathbf{y}$. Using $\mathbf{z} = \mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{y}$ this translates into

$$\begin{aligned} & \forall \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \in X \times Y : \\ & \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^H \begin{bmatrix} -\mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \geq 0 \Rightarrow \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^H \begin{bmatrix} -\mathbf{A}^H \mathbf{T} \mathbf{A} & -\mathbf{A}^H \mathbf{T} \mathbf{B} \\ \mathbf{B}^H \mathbf{T} \mathbf{A} & \mathbf{S} - \mathbf{B}^H \mathbf{T} \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \geq 0 \end{aligned} \quad (\text{C.19})$$

Since $\mathbf{R} \geq 0, \mathbf{S} \geq 0$ but $\mathbf{S} \neq 0$ there

$$\exists \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \in X \times Y : \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^H \begin{bmatrix} -\mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} > 0 \quad (\text{C.20})$$

Due to (C.20) the S-lemma (A.66) can be applied on (C.19) which then becomes equivalent to

$$\exists \alpha \geq 0 : \begin{bmatrix} \mathbf{A}^H \mathbf{T} \mathbf{A} - \alpha \mathbf{R} & \mathbf{A}^H \mathbf{T} \mathbf{B} \\ \mathbf{B}^H \mathbf{T} \mathbf{A} & \mathbf{B}^H \mathbf{T} \mathbf{B} - (1 - \alpha) \mathbf{S} \end{bmatrix} \leq 0$$

With $\mathbf{B}^H \mathbf{T} \mathbf{B} \geq 0$ and $\mathbf{S} \geq 0$ and by applying (A.63) it can be deduced from the lower right block that $\alpha \leq 1$ so that the right hand side of (C.18) follows. This matrix inequality form is equivalent to the right hand side of (C.17) from where it can be seen that $\forall \mathbf{x} \in X_R, \forall \mathbf{y} \in Y_S : \mathbf{z} \in Z_T$ which verifies necessity in (C.15). Condition (C.16) then follows from the symmetry in (C.17) with respect to \mathbf{x}, \mathbf{R} and \mathbf{y}, \mathbf{S} .

Definition Consider the discrete system (C.2) written shortly as

$$\begin{aligned} \mathbf{x}_+ &= \mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{u} \\ \mathbf{y} &= \mathbf{C} \mathbf{x} + \mathbf{D} \mathbf{u} \end{aligned} \quad (\text{C.21})$$

and let U, X, Y be the input, state and output spaces. Define

$$\begin{aligned} U_1 &= \{\mathbf{u} : \mathbf{u}^H \mathbf{u} \leq 1\} \\ X_Z &= \{\mathbf{x} : \mathbf{x}^H \mathbf{Z} \mathbf{x} \leq 1\} \\ Y_Z &= \{\mathbf{y} : \mathbf{y}^H \frac{1}{\gamma^2} \mathbf{y} \leq 1\} \end{aligned} \quad (\text{C.22})$$

with $\mathbf{Z} \geq 0, \gamma > 0$ but $\mathbf{Z} \neq 0$. Then the ellipsoid

$$\begin{aligned} X_Z \text{ is inescapable} &\Leftrightarrow \forall \mathbf{u} \in U_1, \forall \mathbf{x} \in X_Z : \mathbf{x}_+ \in X_Z \\ Y_Z \text{ is inescapable} &\Leftrightarrow \forall \mathbf{u} \in U_1, \forall \mathbf{x} \in X_Z : \mathbf{y} \in Y_Z \end{aligned} \quad (\text{C.23})$$

Proposition Consider the discrete system (C.21) and the ellipsoids as defined in (C.22). Then the ellipsoids X_Z, Y_Z are inescapable if and only if

$$\begin{aligned} & \exists \alpha \in [0, 1] : \begin{bmatrix} \mathbf{A}^H \mathbf{Z} \mathbf{A} - (1 - \alpha) \mathbf{Z} & \mathbf{A}^H \mathbf{Z} \mathbf{B} \\ \mathbf{B}^H \mathbf{Z} \mathbf{A} & \mathbf{B}^H \mathbf{Z} \mathbf{B} - \alpha \mathbf{1} \end{bmatrix} \leq 0 \\ & \exists \beta \in [0, \gamma^2] : \begin{bmatrix} \mathbf{C}^H \mathbf{C} - \beta \mathbf{Z} & \mathbf{C}^H \mathbf{D} \\ \mathbf{D}^H \mathbf{C} & \mathbf{D}^H \mathbf{D} - (\gamma^2 - \beta) \mathbf{1} \end{bmatrix} \leq 0 \end{aligned} \quad (\text{C.24})$$

Proof From the definition (C.23) and the definitions (C.22) the ellipsoids X_Z, Y_Z are inescapable if and only if

$$\begin{aligned} \forall \mathbf{x} \in X, \forall \mathbf{u} \in U, \exists \alpha' \in [0, 1] : \mathbf{x}_+^H \mathbf{Z} \mathbf{x}_+ &\leq \alpha' \mathbf{x}^H \mathbf{Z} \mathbf{x} + (1 - \alpha') \mathbf{u}^H \mathbf{u} \\ \forall \mathbf{x} \in X, \forall \mathbf{u} \in U, \exists \beta' \in [0, 1] : \mathbf{y}^H \frac{1}{\gamma^2} \mathbf{y} &\leq \beta' \mathbf{x}^H \mathbf{Z} \mathbf{x} + (1 - \beta') \mathbf{u}^H \mathbf{u} \end{aligned}$$

according to (C.17). With (C.18) and the substitutions $\alpha = 1 - \alpha', \beta = \gamma^2 \beta'$ the matrix inequalities in (C.24) follow.

C.1.5 Realness

A physical signal, continuous $\mathbf{x}_c(t)$ or discrete $\mathbf{x}_d[k]$, in time domain is a measurable real quantity, hence $\mathbf{x}_c(t), \mathbf{x}_d[k] \in \mathbb{R}$. In frequency domain, this implies

$$\begin{aligned} \overline{\mathbf{x}_c(s)} &= \mathbf{x}_c(\bar{s}) \\ \overline{\mathbf{x}_d(z)} &= \mathbf{x}_d(\bar{z}) \end{aligned} \tag{C.25}$$

A linear time invariant continuous or discrete system $\mathbf{P}(x)$ with $x \in \{s, z\}$ and with realization $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ is real if it maps real valued input signals onto real valued output signals. This is true if and only if

$$\overline{\mathbf{P}(x)} = \mathbf{P}(\bar{x}) \tag{C.26}$$

or equivalently if and only if

$$\left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{array} \right] = \left[\begin{array}{c|c} \overline{\mathbf{A}} & \overline{\mathbf{B}} \\ \overline{\mathbf{C}} & \overline{\mathbf{D}} \end{array} \right] \tag{C.27}$$

which means that these realizations are equivalent, see (C.35) and hence

$$\exists \mathbf{T} : \left[\begin{array}{cc} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{array} \right] = \left[\begin{array}{cc} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{array} \right] \overline{\left[\begin{array}{cc} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{array} \right]} \left[\begin{array}{cc} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{array} \right]^{-1} \tag{C.28}$$

according to (C.36).

Proof The condition (C.25) for signals can be verified with the transformations (C.3) and (C.4) so that

$$\begin{aligned} \overline{\mathbf{x}_c(s)} &= \int_{\mathbb{R}_+} \overline{\mathbf{x}_c(t)} e^{-\bar{s}t} dt = \int_{\mathbb{R}_+} \mathbf{x}_c(t) e^{-\bar{s}t} dt = \mathbf{x}_c(\bar{s}) \\ \overline{\mathbf{x}_d(z)} &= \sum_{\mathbb{Z}_+} \overline{\mathbf{x}_d[k]} \bar{z}^{-k} = \sum_{\mathbb{Z}_+} \mathbf{x}_d[k] \bar{z}^{-k} = \mathbf{x}_d(\bar{z}) \end{aligned}$$

because $\mathbf{x}_c(t) = \overline{\mathbf{x}_c(t)}, \mathbf{x}_d[k] = \overline{\mathbf{x}_d[k]}$ since they are real. With the same reasoning a system $\mathbf{P}(x)$ continuous or discrete with $x \in \{s, z\}$, is real if and only if for all real inputs $\mathbf{u}(x) = \mathbf{u}(\bar{x})$ it produces real outputs $\mathbf{y}(x) = \overline{\mathbf{P}(x)\mathbf{u}(x)} = \overline{\mathbf{P}(x)\mathbf{u}(\bar{x})} = \mathbf{P}(\bar{x})\mathbf{u}(\bar{x}) = \mathbf{y}(\bar{x})$. Hence $\overline{\mathbf{P}(x)} = \mathbf{P}(\bar{x})$ with $x \in \{s, z\}$. Using the link (C.7) this then leads to (C.27).

C.1.6 Stability

A linear time invariant system with dynamic matrix \mathbf{A} is said to be stable if and only if

$$\begin{aligned} \operatorname{Re} s < 0 \quad \forall s : \det(s\mathbf{1} - \mathbf{A}) = 0 & \quad \text{for continuous systems} \\ |z| < 1 \quad \forall z : \det(z\mathbf{1} - \mathbf{A}) = 0 & \quad \text{for discrete systems} \end{aligned} \quad (\text{C.29})$$

Furthermore is the system stable if and only if for any $\mathbf{Q} > 0$ there exists a unique $\mathbf{P} > 0$ which is the solution to

$$\begin{aligned} \mathbf{A}^H \mathbf{P} + \mathbf{P} \mathbf{A} + \mathbf{Q} = \mathbf{0} & \quad \text{for continuous systems} \\ \mathbf{A}^H \mathbf{P} \mathbf{A} - \mathbf{P} + \mathbf{Q} = \mathbf{0} & \quad \text{for discrete systems} \end{aligned} \quad (\text{C.30})$$

which are called Lyapunov matrix equations.

Proof This is due to the solutions of (C.1), (C.2) being respectively given by

$$\begin{aligned} \mathbf{y}(t) &= \mathbf{C} e^{\mathbf{A}t} \left(\mathbf{x}(0) + \int_0^t e^{-\mathbf{A}\tau} \mathbf{B} \mathbf{u}(\tau) d\tau \right) + \mathbf{D} \mathbf{u}(t) \quad \text{for continuous systems} \\ \mathbf{y}[k] &= \mathbf{C} \mathbf{A}^k \left(\mathbf{x}[0] + \sum_{\kappa=0}^{k-1} \mathbf{A}^{-(\kappa+1)} \mathbf{B} \mathbf{u}[\kappa] \right) + \mathbf{D} \mathbf{u}[k] \quad \text{for discrete systems} \end{aligned}$$

so that the solution only converges and thus the output only remains bounded for bounded inputs if (C.29) is satisfied. Now on the one hand assume \mathbf{A} stable then according to (A.87), (A.89) the Sylvester and Stein equations in (C.30) have a unique solution $\mathbf{P} > 0$ for any $\mathbf{Q} > 0$. On the other hand assume (C.30) and define $V(\mathbf{x}) = \mathbf{x}^H \mathbf{P} \mathbf{x}$ which is a Lyapunov function candidate due to $\mathbf{P} > 0$ then

$$\begin{aligned} \dot{V}(\mathbf{x}) &= \mathbf{x}^H (\mathbf{A}^H \mathbf{P} + \mathbf{P} \mathbf{A}) \mathbf{x} = -\mathbf{x}^H \mathbf{Q} \mathbf{x} \quad \text{for continuous systems} \\ \Delta V(\mathbf{x}) &= \mathbf{x}^H (\mathbf{A}^H \mathbf{P} \mathbf{A} - \mathbf{P}) \mathbf{x} = -\mathbf{x}^H \mathbf{Q} \mathbf{x} \quad \text{for discrete systems} \end{aligned}$$

so that $V(\mathbf{x})$ is a Lyapunov function since $\mathbf{Q} > 0$ and \mathbf{A} is stable according to Lyapunov stability theory.

C.1.7 Controllability

A linear time invariant system with dynamic matrix $\mathbf{A} \in \mathbb{C}^{p,p}$ and input matrix $\mathbf{B} \in \mathbb{C}^{p,m}$ is said to be controllable if the future state \mathbf{x} can be controlled by system input \mathbf{u} which is true if and only if

$$\mathbf{A}, \mathbf{B} \text{ controllable} \Leftrightarrow \begin{cases} \operatorname{rnk} [\mathbf{B} \quad \mathbf{A}\mathbf{B} \quad \dots \quad \mathbf{A}^{p-1}\mathbf{B}] = p \\ \mathbf{B}^T \mathbf{x} \neq \mathbf{0} \quad \forall \mathbf{x} : \mathbf{A}^T \mathbf{x} = \lambda \mathbf{x} \end{cases} \quad (\text{C.31})$$

Furthermore is the system controllable if and only if the Gram controllability matrix \mathbf{G} which is the unique solution to

$$\begin{aligned} \mathbf{A}\mathbf{G} + \mathbf{G}\mathbf{A}^H + \mathbf{B}\mathbf{B}^H &= \mathbf{0} \quad \text{for continuous systems} \\ \mathbf{A}\mathbf{G}\mathbf{A}^H - \mathbf{G} + \mathbf{B}\mathbf{B}^H &= \mathbf{0} \quad \text{for discrete systems} \end{aligned} \quad (\text{C.32})$$

exists and $\mathbf{G} > 0$.

C.1.8 Observability

A linear time invariant system with dynamic matrix $\mathbf{A} \in \mathbb{C}^{p,p}$ and output matrix $\mathbf{C} \in \mathbb{C}^{n,p}$ is said to be observable if the past state \mathbf{x} can be observed by system output \mathbf{y} which is true if and only if

$$\mathbf{A}, \mathbf{C} \text{ observable} \Leftrightarrow \begin{cases} \text{rk} \begin{bmatrix} \mathbf{C}^T & \mathbf{A}^T \mathbf{C}^T & \dots & \mathbf{A}^{T^{p-1}} \mathbf{C}^T \end{bmatrix} = p \\ \mathbf{C}\mathbf{x} \neq \mathbf{0} \forall \mathbf{x} : \mathbf{A}\mathbf{x} = \lambda\mathbf{x} \end{cases} \quad (\text{C.33})$$

Furthermore is the system observable if and only if the Gram observability matrix \mathbf{G} which is the unique solution to

$$\begin{aligned} \mathbf{A}^H \mathbf{G} + \mathbf{G}\mathbf{A} + \mathbf{C}^H \mathbf{C} &= \mathbf{0} && \text{for continuous systems} \\ \mathbf{A}^H \mathbf{G}\mathbf{A} - \mathbf{G} + \mathbf{C}^H \mathbf{C} &= \mathbf{0} && \text{for discrete systems} \end{aligned} \quad (\text{C.34})$$

exists and $\mathbf{G} > 0$.

C.1.9 Equivalence

Two state space representations are equivalent if their transfer functions are identical:

$$(\mathbf{A}_1, \mathbf{B}_1, \mathbf{C}_1, \mathbf{D}_1) \sim (\mathbf{A}_2, \mathbf{B}_2, \mathbf{C}_2, \mathbf{D}_2) \Leftrightarrow \left[\begin{array}{c|c} \mathbf{A}_1 & \mathbf{B}_1 \\ \hline \mathbf{C}_1 & \mathbf{D}_1 \end{array} \right] = \left[\begin{array}{c|c} \mathbf{A}_2 & \mathbf{B}_2 \\ \hline \mathbf{C}_2 & \mathbf{D}_2 \end{array} \right] \quad (\text{C.35})$$

with $x \in \{s, z\}$ for continuous and discrete systems respectively. Equivalent state space representations therefore describe the same system.

C.1.10 Change of Basis

The basis in which the states of a system in state space representation are described can be changed. However, stability, controllability, observability and the transfer matrix \mathbf{P} are preserved under basis transformation and the transformed state space representation is equivalent to the original one.

$$\mathbf{P} = \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \hline \mathbf{C} & \mathbf{D} \end{array} \right] = \left[\begin{array}{c|c} \mathbf{T}\mathbf{A}\mathbf{T}^{-1} & \mathbf{T}\mathbf{B} \\ \hline \mathbf{C}\mathbf{T}^{-1} & \mathbf{D} \end{array} \right] \quad (\text{C.36})$$

This also means that two realizations are equivalent according to (C.35) if and only if there exists an invertible \mathbf{T} so that (C.36) is true.

Proof The equivalence follows from (C.1) or (C.2) if the original state \mathbf{x} is replaced by the transformed state $\boldsymbol{\xi} = \mathbf{T}\mathbf{x}$. The invariances follow by reexamination of (C.29), (C.31), (C.33) and (C.7) with the new state space representation.

C.1.11 Kalman Decomposition

A change of basis can be used to recover the controllable and observable eigenmodes of the system. This special basis transformation is called the Kalman

decomposition that transforms any system \mathbf{P} into the form

$$\mathbf{P} = \left[\begin{array}{cccc|c} \mathbf{A}_{c\bar{o}} & \mathbf{A}_{12} & \mathbf{A}_{13} & \mathbf{A}_{14} & \mathbf{B}_{c\bar{o}} \\ \mathbf{0} & \mathbf{A}_{co} & \mathbf{0} & \mathbf{A}_{24} & \mathbf{B}_{co} \\ \mathbf{0} & \mathbf{0} & \mathbf{A}_{\bar{c}\bar{o}} & \mathbf{A}_{34} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{A}_{\bar{c}o} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{C}_{co} & \mathbf{0} & \mathbf{C}_{\bar{c}o} & \mathbf{D} \end{array} \right] \text{ with} \quad (C.37)$$

controllable and observable $\left[\begin{array}{c|c} \mathbf{A}_{co} & \mathbf{B}_{co} \\ \hline \mathbf{C}_{co} & \mathbf{D} \end{array} \right]$

controllable $\left[\begin{array}{cc|c} \mathbf{A}_{c\bar{o}} & \mathbf{A}_{12} & \mathbf{B}_{c\bar{o}} \\ \mathbf{0} & \mathbf{A}_{co} & \mathbf{B}_{co} \\ \hline \mathbf{0} & \mathbf{C}_{co} & \mathbf{D} \end{array} \right]$

observable $\left[\begin{array}{cc|c} \mathbf{A}_{co} & \mathbf{A}_{24} & \mathbf{B}_{co} \\ \mathbf{0} & \mathbf{A}_{\bar{c}o} & \mathbf{0} \\ \hline \mathbf{C}_{co} & \mathbf{C}_{\bar{c}o} & \mathbf{D} \end{array} \right]$

This decomposes the system's total number of states p into

$$p = p_{c\bar{o}} + p_{co} + p_{\bar{c}\bar{o}} + p_{\bar{c}o} \quad (C.38)$$

where $p_{c\bar{o}}, p_{co}, p_{\bar{c}\bar{o}}, p_{\bar{c}o}$ are the respective dimensions of $\mathbf{A}_{c\bar{o}}, \mathbf{A}_{co}, \mathbf{A}_{\bar{c}\bar{o}}, \mathbf{A}_{\bar{c}o}$. Regarding input output behavior, the system \mathbf{P} can be reduced in state dimension to its controllable and observable subsystem by cutting off the uncontrollable and unobservable parts so that

$$\mathbf{P} = \left[\begin{array}{c|c} \mathbf{A}_{co} & \mathbf{B}_{co} \\ \hline \mathbf{C}_{co} & \mathbf{D} \end{array} \right] \quad (C.39)$$

since the input controls only controllable parts and the output observes only observable parts. The resulting state space representation can not be further reduced without changing the system \mathbf{P} . Therefore a controllable and observable state space representation is called *minimal*.

C.1.12 McMillan Degree

The state dimension p_{co} of this minimal state space representation of the system \mathbf{P} is the number of controllable and observable eigenmodes. Since controllability and observability is invariant under change of basis (C.36), this number is also a characteristic property of \mathbf{P} regardless of the state space representation chosen and is called the McMillan degree δ of the system

$$\delta(\mathbf{P}) = p_{co} \leq p \quad (C.40)$$

where p is the state dimension of any state space representation of \mathbf{P} . This inequality reflects the fact that a controllable and observable state space representation has state dimension $p = p_{co}$ and is therefore minimal.

The combined McMillan degree $\delta(\mathbf{P}_1, \mathbf{P}_2)$ of two systems $\mathbf{P}_1, \mathbf{P}_2$ can not exceed

$$\delta(\mathbf{P}_1, \mathbf{P}_2) \leq \delta(\mathbf{P}_1) + \delta(\mathbf{P}_2) \quad (C.41)$$

but may be lower, depending on the actual state space representations and the connection between them.

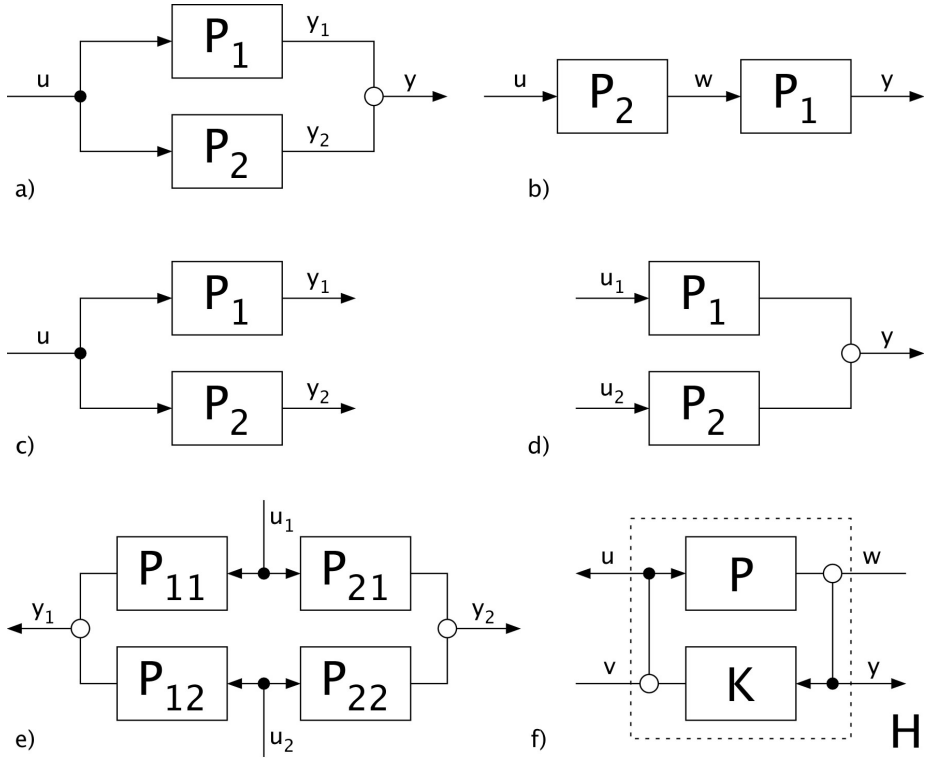


Figure C.1: The different system connections: a) parallel b) serial c) common input d) common output e) bridge f) feedback loop.

Proof This stems from the fact that using (C.37) we can reduce the state space representations of P_1, P_2 to its minimal form so that the combined state vector has the dimension $p = \delta(P_1) + \delta(P_2)$. From (C.40) we have $\delta(P_1, P_2) \leq p$ so that (C.41) is obtained.

C.2 Connections

This section provides some formulae for operations on systems and the various possible connections between systems to form super systems as illustrated in figure C.1.

C.2.1 Flow Inversion

The signal flow through a system

$$P = \left[\begin{array}{c|c} A & B \\ \hline C & D \end{array} \right]$$

may be inverted if $\exists D^{-1}$ and the inverted system is then given by

$$P^{-1} = \left[\begin{array}{c|c} A - BD^{-1}C & BD^{-1} \\ \hline -D^{-1}C & D^{-1} \end{array} \right] \quad (\text{C.42})$$

with McMillan degree

$$\delta(\mathbf{P}) = \delta(\mathbf{P}^{-1}) \quad (\text{C.43})$$

Proof This follows from the application of (A.55) on (C.7). To verify the preservation of the McMillan degree let \mathbf{P}^{-1} be unobservable so that $\forall \mathbf{x} : (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})\mathbf{x} = \lambda\mathbf{x}$ we have $\mathbf{D}^{-1}\mathbf{C}\mathbf{x} = \mathbf{0}$ using (C.33). This leads to $(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})\mathbf{x} = \mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ and $\mathbf{C}\mathbf{x} = \mathbf{0}$ so that if \mathbf{P}^{-1} is unobservable then \mathbf{P} is unobservable or alternatively if \mathbf{P} is observable then \mathbf{P}^{-1} is observable. The same is true for controllability, so that if \mathbf{P} is minimal then \mathbf{P}^{-1} is minimal. This means there may be cases where \mathbf{P} is not minimal with \mathbf{P}^{-1} minimal. Since \mathbf{P}^{-1} has the same state dimension than \mathbf{P} following (C.42) this leads to $\delta(\mathbf{P}) \leq \delta(\mathbf{P}^{-1})$ and therefore $\delta(\mathbf{P}) \leq \delta(\mathbf{P}^{-1}) \leq \delta((\mathbf{P}^{-1})^{-1}) = \delta(\mathbf{P})$ so that (C.43) follows.

C.2.2 Parallel Connection

Two systems

$$\mathbf{P}_i = \left[\begin{array}{c|c} \mathbf{A}_i & \mathbf{B}_i \\ \hline \mathbf{C}_i & \mathbf{D}_i \end{array} \right]; i \in \{1, 2\}$$

in parallel connection give

$$\mathbf{P}_1 + \mathbf{P}_2 = \left[\begin{array}{cc|c} \mathbf{A}_1 & \mathbf{0} & \mathbf{B}_1 \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{B}_2 \\ \hline \mathbf{C}_1 & \mathbf{C}_2 & \mathbf{D}_1 + \mathbf{D}_2 \end{array} \right] \quad (\text{C.44})$$

with McMillan degree

$$\delta(\mathbf{P}_1 + \mathbf{P}_2) \leq \delta(\mathbf{P}_1) + \delta(\mathbf{P}_2) \quad (\text{C.45})$$

Proof This is due to $\mathbf{y} = \mathbf{y}_1 + \mathbf{y}_2 = \mathbf{C}_1\mathbf{x}_1 + \mathbf{C}_2\mathbf{x}_2 + (\mathbf{D}_1 + \mathbf{D}_2)\mathbf{u}$. The McMillan degree can have maximal range $\delta(\mathbf{P}_1 + \mathbf{P}_2) \in [0, \delta(\mathbf{P}_1) + \delta(\mathbf{P}_2)]$. The lower bound is attained if for example $\mathbf{P}_2 = -\mathbf{P}_1$ so that $\delta(\mathbf{P}_1 + \mathbf{P}_2) = \delta(\mathbf{0}) = 0$.

C.2.3 Serial Connection

Two systems

$$\mathbf{P}_i = \left[\begin{array}{c|c} \mathbf{A}_i & \mathbf{B}_i \\ \hline \mathbf{C}_i & \mathbf{D}_i \end{array} \right]; i \in \{1, 2\}$$

in serial connection give

$$\mathbf{P}_1\mathbf{P}_2 = \left[\begin{array}{cc|c} \mathbf{A}_1 & \mathbf{B}_1\mathbf{C}_2 & \mathbf{B}_1\mathbf{D}_2 \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{B}_2 \\ \hline \mathbf{C}_1 & \mathbf{D}_1\mathbf{C}_2 & \mathbf{D}_1\mathbf{D}_2 \end{array} \right] \quad (\text{C.46})$$

with McMillan degree

$$\delta(\mathbf{P}_1\mathbf{P}_2) \leq \delta(\mathbf{P}_1) + \delta(\mathbf{P}_2) \quad (\text{C.47})$$

Proof This is due to $\mathbf{w} = \mathbf{C}_2\mathbf{x}_2 + \mathbf{D}_2\mathbf{u}$ so that $\dot{\mathbf{x}}_1 = \mathbf{A}_1\mathbf{x}_1 + \mathbf{B}_1(\mathbf{C}_2\mathbf{x}_2 + \mathbf{D}_2\mathbf{u})$ and $\mathbf{y} = \mathbf{C}_1\mathbf{x}_1 + \mathbf{D}_1(\mathbf{C}_2\mathbf{x}_2 + \mathbf{D}_2\mathbf{u})$. The McMillan degree can have maximal range $\delta(\mathbf{P}_1\mathbf{P}_2) \in [0, \delta(\mathbf{P}_1) + \delta(\mathbf{P}_2)]$. The lower bound is attained if for example $\mathbf{P}_2 = \mathbf{P}_1^{-1}$ so that $\delta(\mathbf{P}_1\mathbf{P}_2) = \delta(\mathbf{1}) = 0$.

C.2.4 Feedback Loop

A plant P and controller K

$$P = \left[\begin{array}{c|c} \mathbf{A}_P & \mathbf{B}_P \\ \hline \mathbf{C}_P & \mathbf{D}_P \end{array} \right]; \quad K = \left[\begin{array}{c|c} \mathbf{A}_K & \mathbf{B}_K \\ \hline \mathbf{C}_K & \mathbf{D}_K \end{array} \right]$$

may be put in feedback loop if $\exists \mathbf{E} = (\mathbf{1} - \mathbf{D}_P \mathbf{D}_K)^{-1}$, $\exists \mathbf{F} = (\mathbf{1} - \mathbf{D}_K \mathbf{D}_P)^{-1}$ then the closed loop system is given by

$$\begin{aligned} \begin{bmatrix} \mathbf{y} \\ \mathbf{u} \end{bmatrix} &= \mathbf{H} \begin{bmatrix} \mathbf{v} \\ \mathbf{w} \end{bmatrix} \quad \text{with } \mathbf{H} = \begin{bmatrix} -\mathbf{K} & \mathbf{1} \\ \mathbf{1} & -\mathbf{P} \end{bmatrix}^{-1} = \\ &= \left[\begin{array}{cc|cc} \mathbf{A}_P + \mathbf{B}_P \mathbf{F} \mathbf{D}_K \mathbf{C}_P & \mathbf{B}_P \mathbf{F} \mathbf{C}_K & \mathbf{B}_P \mathbf{F} & \mathbf{B}_P \mathbf{F} \mathbf{D}_K \\ \mathbf{B}_K \mathbf{E} \mathbf{C}_P & \mathbf{A}_K + \mathbf{B}_K \mathbf{E} \mathbf{D}_P \mathbf{C}_K & \mathbf{B}_K \mathbf{E} \mathbf{D}_P & \mathbf{B}_K \mathbf{E} \\ \hline \mathbf{E} \mathbf{C}_P & \mathbf{E} \mathbf{D}_P \mathbf{C}_K & \mathbf{E} \mathbf{D}_P & \mathbf{E} \\ \mathbf{F} \mathbf{D}_K \mathbf{C}_P & \mathbf{F} \mathbf{C}_K & \mathbf{F} & \mathbf{F} \mathbf{D}_K \end{array} \right] \end{aligned} \quad (\text{C.48})$$

with McMillan degree

$$\delta(\mathbf{H}) = \delta(\mathbf{P}) + \delta(\mathbf{K}) \quad (\text{C.49})$$

Proof The result follows from $\mathbf{y} = \mathbf{C}_P \mathbf{x}_P + \mathbf{D}_P \mathbf{u} + \mathbf{w}$ and $\mathbf{u} = \mathbf{C}_K \mathbf{x}_K + \mathbf{D}_K \mathbf{y} + \mathbf{v}$ so that $\mathbf{y} = \mathbf{E}(\mathbf{C}_P \mathbf{x}_P + \mathbf{D}_P \mathbf{C}_K \mathbf{x}_K + \mathbf{w} + \mathbf{D}_P \mathbf{v})$ and $\mathbf{u} = \mathbf{F}(\mathbf{C}_K \mathbf{x}_K + \mathbf{D}_K \mathbf{C}_P \mathbf{x}_P + \mathbf{v} + \mathbf{D}_K \mathbf{w})$ leads to the above result by making use of (A.56). To verify the McMillan degree let the state space representation of P, K be minimal and assume H unobservable. Using inobservability from (C.33) on (C.48) there exist $\mathbf{x}_P, \mathbf{x}_K$ so that

$$\left[\begin{array}{cc|c} \mathbf{A}_P + \mathbf{B}_P \mathbf{F} \mathbf{D}_K \mathbf{C}_P & \mathbf{B}_P \mathbf{F} \mathbf{C}_K & \mathbf{x}_P \\ \mathbf{B}_K \mathbf{E} \mathbf{C}_P & \mathbf{A}_K + \mathbf{B}_K \mathbf{E} \mathbf{D}_P \mathbf{C}_K & \mathbf{x}_K \\ \mathbf{E} \mathbf{C}_P & \mathbf{E} \mathbf{D}_P \mathbf{C}_K & \mathbf{0} \\ \mathbf{F} \mathbf{D}_K \mathbf{C}_P & \mathbf{F} \mathbf{C}_K & \mathbf{0} \end{array} \right] \begin{bmatrix} \mathbf{x}_P \\ \mathbf{x}_K \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{x}_P \\ \mathbf{x}_K \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}$$

This simplifies to

$$\left[\begin{array}{cc|c} \mathbf{A}_P & \mathbf{0} & \mathbf{x}_P \\ \mathbf{0} & \mathbf{A}_K & \mathbf{x}_K \\ \mathbf{C}_P & \mathbf{D}_P \mathbf{C}_K & \mathbf{0} \\ \mathbf{D}_K \mathbf{C}_P & \mathbf{C}_K & \mathbf{0} \end{array} \right] \begin{bmatrix} \mathbf{x}_P \\ \mathbf{x}_K \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{x}_P \\ \mathbf{x}_K \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}$$

But P, K are supposed to have minimal state space representations thus in addition to that we have P, K observable

$$\begin{aligned} \mathbf{C}_P \mathbf{x}_P &\neq \mathbf{0} \quad \forall \mathbf{x}_P : \mathbf{A}_P \mathbf{x}_P = \lambda \mathbf{x}_P \\ \mathbf{C}_K \mathbf{x}_K &\neq \mathbf{0} \quad \forall \mathbf{x}_K : \mathbf{A}_K \mathbf{x}_K = \lambda \mathbf{x}_K \end{aligned}$$

If λ is an eigenvalue of \mathbf{A}_P but not \mathbf{A}_K then $\mathbf{x}_K = \mathbf{0}$. This further simplifies the inobservability equation

$$\left[\begin{array}{c} \mathbf{A}_P \\ \mathbf{C}_P \\ \mathbf{D}_K \mathbf{C}_P \end{array} \right] \mathbf{x}_P = \lambda \begin{bmatrix} \mathbf{x}_P \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}$$

But $C_P \mathbf{x}_P = \mathbf{0}$ contradicts the observability of \mathbf{P} so \mathbf{H} is observable. The same is true if λ is an eigenvalue of \mathbf{A}_K but not \mathbf{A}_P . If λ is a common eigenvalue of $\mathbf{A}_P, \mathbf{A}_K$ then the inobservability equation gives

$$\begin{aligned} C_P \mathbf{x}_P + D_P C_K \mathbf{x}_K &= \mathbf{0} \\ D_K C_P \mathbf{x}_P + C_K \mathbf{x}_K &= \mathbf{0} \end{aligned}$$

This leads to

$$\begin{aligned} (\mathbf{1} - D_P D_K) C_P \mathbf{x}_P &= \mathbf{E} C_P \mathbf{x}_P = \mathbf{0} \Leftrightarrow C_P \mathbf{x}_P = \mathbf{0} \\ (\mathbf{1} - D_K D_P) C_K \mathbf{x}_K &= \mathbf{F} C_K \mathbf{x}_K = \mathbf{0} \Leftrightarrow C_K \mathbf{x}_K = \mathbf{0} \end{aligned}$$

since \mathbf{E}, \mathbf{F} are required to be invertible to form the feedback loop. This contradicts the observability of \mathbf{P} and \mathbf{K} so \mathbf{H} is observable. So \mathbf{P}, \mathbf{K} minimal leads to \mathbf{H} observable. The same is true for controllability so that if \mathbf{P}, \mathbf{K} are minimal then \mathbf{H} is minimal. This means there may be cases where one of \mathbf{P}, \mathbf{K} is not minimal with \mathbf{H} minimal. Since the state dimension of \mathbf{H} is the sum of the state dimensions of \mathbf{P}, \mathbf{K} following (C.48) this leads to $\delta(\mathbf{H}) \geq \delta(\mathbf{P}) + \delta(\mathbf{K})$. But from (C.41) we see that $\delta(\mathbf{H}) \leq \delta(\mathbf{P}) + \delta(\mathbf{K})$ so that (C.49) follows.

C.2.5 Common Input

Two systems

$$\mathbf{P}_i = \left[\begin{array}{c|c} \mathbf{A}_i & \mathbf{B}_i \\ \hline \mathbf{C}_i & \mathbf{D}_i \end{array} \right]; \quad i \in \{1, 2\}$$

with common input give

$$\left[\begin{array}{c} \mathbf{P}_1 \\ \mathbf{P}_2 \end{array} \right] = \left[\begin{array}{cc|c} \mathbf{A}_1 & \mathbf{0} & \mathbf{B}_1 \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{B}_2 \\ \hline \mathbf{C}_1 & \mathbf{0} & \mathbf{D}_1 \\ \mathbf{0} & \mathbf{C}_2 & \mathbf{D}_2 \end{array} \right] \quad (\text{C.50})$$

with McMillan degree

$$\max(\delta(\mathbf{P}_1), \delta(\mathbf{P}_2)) \leq \delta \left(\left[\begin{array}{c} \mathbf{P}_1 \\ \mathbf{P}_2 \end{array} \right] \right) \leq \delta(\mathbf{P}_1) + \delta(\mathbf{P}_2) \quad (\text{C.51})$$

Proof To verify the lower bound of the McMillan degree, let $\mathbf{P}_1, \mathbf{P}_2$ be minimal or reduce them to their minimal subsystems. Then the connected system is observable:

$$\begin{aligned} \left[\begin{array}{cc} \mathbf{C}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_2 \end{array} \right] \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{C}_1 \mathbf{x}_1 \\ \mathbf{C}_2 \mathbf{x}_2 \end{bmatrix} \neq \mathbf{0} \\ \forall \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} : \left[\begin{array}{cc} \mathbf{A}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{array} \right] \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_1 \mathbf{x}_1 \\ \mathbf{A}_2 \mathbf{x}_2 \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \end{aligned}$$

If λ is an eigenvalue of \mathbf{A}_1 but not \mathbf{A}_2 then $\mathbf{x}_2 = \mathbf{0}$ however $\mathbf{C}_1 \mathbf{x}_1 \neq \mathbf{0}$ from \mathbf{P}_1 being observable suffices to assure observability. The same is true if λ is an eigenvalue of \mathbf{A}_2 but not \mathbf{A}_1 . If λ is a common eigenvalue of $\mathbf{A}_1, \mathbf{A}_2$ then

$C_1 \mathbf{x}_1 \neq \mathbf{0}$ from \mathbf{P}_1 being observable and $C_2 \mathbf{x}_2 \neq \mathbf{0}$ from \mathbf{P}_2 being observable. But the connected system is not always controllable:

$$\begin{aligned} \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} &= \mathbf{B}_1^T \mathbf{x}_1 + \mathbf{B}_2^T \mathbf{x}_2 \stackrel{?}{\neq} \mathbf{0} \\ \forall \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} : \begin{bmatrix} \mathbf{A}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_1^T \mathbf{x}_1 \\ \mathbf{A}_2^T \mathbf{x}_2 \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \end{aligned}$$

If λ is an eigenvalue of \mathbf{A}_1 but not \mathbf{A}_2 then $\mathbf{x}_2 = \mathbf{0}$ and $\mathbf{B}_1^T \mathbf{x}_1 \neq \mathbf{0}$ from \mathbf{P}_1 being controllable suffices to assure controllability. The same is true if λ is an eigenvalue of \mathbf{A}_2 but not \mathbf{A}_1 . If λ is a common eigenvalue of $\mathbf{A}_1, \mathbf{A}_2$ then $\mathbf{B}_1^T \mathbf{x}_1 + \mathbf{B}_2^T \mathbf{x}_2 = \mathbf{0}$ is possible. In this case, the common λ is an uncontrollable eigenmode of the connected system and reduces the McMillan degree by one. If it is a multiple eigenvalue of both matrices then it reduces the McMillan degree by its multiplicity. Now let $\delta(\mathbf{P}_1) \leq \delta(\mathbf{P}_2)$ then the lower bound of the McMillan degree of the connected system can be attained if all eigenvalues of \mathbf{A}_1 are also eigenvalues of \mathbf{A}_2 . Then the McMillan degree of the unconnected system $\delta(\mathbf{P}_1) + \delta(\mathbf{P}_2)$ can be reduced by the degree of the smaller system $\delta(\mathbf{P}_1)$ through connection. The result is the McMillan degree of the larger system $\delta(\mathbf{P}_2)$ of the two or in general $\max(\delta(\mathbf{P}_1), \delta(\mathbf{P}_2))$.

Remark An important special case is $\mathbf{A}_i = \mathbf{A}, \mathbf{B}_i = \mathbf{B}; i \in \{1, 2\}$ then

$$\begin{bmatrix} \mathbf{P}_1 \\ \mathbf{P}_2 \end{bmatrix} = \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \hline \mathbf{C}_1 & \mathbf{D}_1 \\ \mathbf{C}_2 & \mathbf{D}_2 \end{array} \right] \quad (\text{C.52})$$

C.2.6 Common Output

Two systems

$$\mathbf{P}_i = \left[\begin{array}{c|c} \mathbf{A}_i & \mathbf{B}_i \\ \hline \mathbf{C}_i & \mathbf{D}_i \end{array} \right]; i \in \{1, 2\}$$

with common output give

$$[\mathbf{P}_1 \quad \mathbf{P}_2] = \left[\begin{array}{cc|cc} \mathbf{A}_1 & \mathbf{0} & \mathbf{B}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{0} & \mathbf{B}_2 \\ \hline \mathbf{C}_1 & \mathbf{C}_2 & \mathbf{D}_1 & \mathbf{D}_2 \end{array} \right] \quad (\text{C.53})$$

with McMillan degree

$$\max(\delta(\mathbf{P}_1), \delta(\mathbf{P}_2)) \leq \delta([\mathbf{P}_1 \quad \mathbf{P}_2]) \leq \delta(\mathbf{P}_1) + \delta(\mathbf{P}_2) \quad (\text{C.54})$$

Proof To verify the lower bound of the McMillan degree, let $\mathbf{P}_1, \mathbf{P}_2$ be minimal or reduce them to their minimal subsystems. Then the connected system is controllable:

$$\begin{aligned} \begin{bmatrix} \mathbf{B}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{B}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{B}_1^T \mathbf{x}_1 \\ \mathbf{B}_2^T \mathbf{x}_2 \end{bmatrix} \neq \mathbf{0} \\ \forall \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} : \begin{bmatrix} \mathbf{A}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_1^T \mathbf{x}_1 \\ \mathbf{A}_2^T \mathbf{x}_2 \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \end{aligned}$$

If λ is an eigenvalue of \mathbf{A}_1 but not \mathbf{A}_2 then $\mathbf{x}_2 = \mathbf{0}$ however $\mathbf{B}_1^T \mathbf{x}_1 \neq \mathbf{0}$ from \mathbf{P}_1 being controllable suffices to assure controllability. The same is true if λ is an eigenvalue of \mathbf{A}_2 but not \mathbf{A}_1 . If λ is a common eigenvalue of $\mathbf{A}_1, \mathbf{A}_2$ then $\mathbf{B}_1^T \mathbf{x}_1 \neq \mathbf{0}$ from \mathbf{P}_1 being controllable and $\mathbf{B}_2^T \mathbf{x}_2 \neq \mathbf{0}$ from \mathbf{P}_2 being controllable. But the connected system is not always observable:

$$\begin{aligned} [\mathbf{C}_1 \quad \mathbf{C}_2] \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} &= \mathbf{C}_1 \mathbf{x}_1 + \mathbf{C}_2 \mathbf{x}_2 \stackrel{?}{\neq} \mathbf{0} \\ \forall \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} : \begin{bmatrix} \mathbf{A}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{A}_1 \mathbf{x}_1 \\ \mathbf{A}_2 \mathbf{x}_2 \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \end{aligned}$$

If λ is an eigenvalue of \mathbf{A}_1 but not \mathbf{A}_2 then $\mathbf{x}_2 = \mathbf{0}$ and $\mathbf{C}_1 \mathbf{x}_1 \neq \mathbf{0}$ from \mathbf{P}_1 being observable suffices to assure observability. The same is true if λ is an eigenvalue of \mathbf{A}_2 but not \mathbf{A}_1 . If λ is a common eigenvalue of $\mathbf{A}_1, \mathbf{A}_2$ then $\mathbf{C}_1 \mathbf{x}_1 + \mathbf{C}_2 \mathbf{x}_2 = \mathbf{0}$ is possible. In this case, the common λ is an unobservable eigenmode of the connected system and reduces the McMillan degree by one. If it is a multiple eigenvalue of both matrices then it reduces the McMillan degree by its multiplicity. Now let $\delta(\mathbf{P}_1) \leq \delta(\mathbf{P}_2)$ then the lower bound of the McMillan degree of the connected system can be attained if all eigenvalues of \mathbf{A}_1 are also eigenvalues of \mathbf{A}_2 . Then the McMillan degree of the unconnected system $\delta(\mathbf{P}_1) + \delta(\mathbf{P}_2)$ can be reduced by the degree of the smaller system $\delta(\mathbf{P}_1)$ through connection. The result is the McMillan degree of the larger system $\delta(\mathbf{P}_2)$ of the two or in general $\max(\delta(\mathbf{P}_1), \delta(\mathbf{P}_2))$.

Remark An important special case is $\mathbf{A}_i = \mathbf{A}, \mathbf{C}_i = \mathbf{C}; i \in \{1, 2\}$ then

$$[\mathbf{P}_1 \quad \mathbf{P}_2] = \left[\begin{array}{c|cc} \mathbf{A} & \mathbf{B}_1 & \mathbf{B}_2 \\ \hline \mathbf{C} & \mathbf{D}_1 & \mathbf{D}_2 \end{array} \right] \quad (\text{C.55})$$

C.2.7 Bridge

Four systems

$$\mathbf{P}_{ij} = \left[\begin{array}{c|c} \mathbf{A}_{ij} & \mathbf{B}_{ij} \\ \hline \mathbf{C}_{ij} & \mathbf{D}_{ij} \end{array} \right]; i, j \in \{1, 2\}$$

in bridge connection give

$$[\mathbf{P}_{11} \quad \mathbf{P}_{12} \\ \mathbf{P}_{21} \quad \mathbf{P}_{22}] = \left[\begin{array}{cccc|cc} \mathbf{A}_{11} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{B}_{11} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{12} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{B}_{12} \\ \mathbf{0} & \mathbf{0} & \mathbf{A}_{21} & \mathbf{0} & \mathbf{B}_{21} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{A}_{22} & \mathbf{0} & \mathbf{B}_{22} \\ \hline \mathbf{C}_{11} & \mathbf{C}_{12} & \mathbf{0} & \mathbf{0} & \mathbf{D}_{11} & \mathbf{D}_{12} \\ \mathbf{0} & \mathbf{0} & \mathbf{C}_{21} & \mathbf{C}_{22} & \mathbf{D}_{21} & \mathbf{D}_{22} \end{array} \right] \quad (\text{C.56})$$

with McMillan degree

$$\begin{aligned} \max(\delta(\mathbf{P}_{11}), \delta(\mathbf{P}_{12}), \delta(\mathbf{P}_{21}), \delta(\mathbf{P}_{22})) &\leq \\ &\leq \delta \left(\begin{bmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} \\ \mathbf{P}_{21} & \mathbf{P}_{22} \end{bmatrix} \right) \leq \\ &\leq \delta(\mathbf{P}_{11}) + \delta(\mathbf{P}_{12}) + \delta(\mathbf{P}_{21}) + \delta(\mathbf{P}_{22}) \quad (\text{C.57}) \end{aligned}$$

Proof The bridge can be obtained by first connecting $\mathbf{P}_{11}, \mathbf{P}_{12}$ and $\mathbf{P}_{21}, \mathbf{P}_{22}$ each to common outputs using (C.53) and (C.54). And then connecting the two resulting $[\mathbf{P}_{11} \ \mathbf{P}_{12}]$ and $[\mathbf{P}_{21} \ \mathbf{P}_{22}]$ to a common input using (C.50) and (C.51).

Remark An important special case is $\mathbf{A}_{ij} = \mathbf{A}, \mathbf{B}_{ij} = \mathbf{B}_j, \mathbf{C}_{ij} = \mathbf{C}_i; i, j \in \{1, 2\}$ then

$$\begin{bmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} \\ \mathbf{P}_{21} & \mathbf{P}_{22} \end{bmatrix} = \left[\begin{array}{c|cc} \mathbf{A} & \mathbf{B}_1 & \mathbf{B}_2 \\ \mathbf{C}_1 & \mathbf{D}_{11} & \mathbf{D}_{12} \\ \mathbf{C}_2 & \mathbf{D}_{21} & \mathbf{D}_{22} \end{array} \right] \quad (\text{C.58})$$

C.3 Factorizations

This section provides formulae for the important right and left factorization of systems.

C.3.1 Right Factorization

A system

$$\mathbf{P} = \mathbf{X}\mathbf{Y}^{-1} = \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{array} \right]$$

has a right factorization in the form of

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \left[\begin{array}{c|c} \mathbf{A} + \mathbf{B}\mathbf{F} & \mathbf{B}\mathbf{U} \\ \mathbf{C} + \mathbf{D}\mathbf{F} & \mathbf{D}\mathbf{U} \\ \mathbf{F} & \mathbf{U} \end{array} \right] \quad (\text{C.59})$$

with free parameters \mathbf{F}, \mathbf{U} and \mathbf{U} invertible. Its McMillan degree is bounded by

$$\delta \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \right) \geq \delta(\mathbf{P}) \quad (\text{C.60})$$

see [7] for details.

Proof Start with two systems having common input

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \left[\begin{array}{c|c} \mathbf{A}' & \mathbf{B}' \\ \mathbf{C}_X & \mathbf{D}_X \\ \mathbf{C}_Y & \mathbf{D}_Y \end{array} \right]$$

with \mathbf{D}_Y invertible in order to allow the inversion of \mathbf{Y} . Using flow inversion (C.42) and serial connection (C.46) this leads to

$$\mathbf{X}\mathbf{Y}^{-1} = \left[\begin{array}{c|c} \mathbf{A}' & -\mathbf{B}'\mathbf{D}_Y^{-1}\mathbf{C}_Y \\ \mathbf{0} & \mathbf{A}' - \mathbf{B}'\mathbf{D}_Y^{-1}\mathbf{C}_Y \\ \mathbf{C}_X & -\mathbf{D}_X\mathbf{D}_Y^{-1}\mathbf{C}_Y \end{array} \middle| \begin{array}{c} \mathbf{B}'\mathbf{D}_Y^{-1} \\ \mathbf{B}'\mathbf{D}_Y^{-1} \\ \mathbf{D}_X\mathbf{D}_Y^{-1} \end{array} \right]$$

Applying the transformation \mathbf{T} as in (C.36) and cutting off the uncontrollable part leads to

$$\mathbf{X}\mathbf{Y}^{-1} = \left[\begin{array}{c|c} \mathbf{A}' - \mathbf{B}'\mathbf{D}_Y^{-1}\mathbf{C}_Y & \mathbf{B}'\mathbf{D}_Y^{-1} \\ \mathbf{C}_X - \mathbf{D}_X\mathbf{D}_Y^{-1}\mathbf{C}_Y & \mathbf{D}_X\mathbf{D}_Y^{-1} \end{array} \right] \text{ with } \mathbf{T} = \begin{bmatrix} \mathbf{1} & -\mathbf{1} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}$$

Comparing this to

$$P = XY^{-1} = \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \hline \mathbf{C} & \mathbf{D} \end{array} \right]$$

provides a solution:

$$\mathbf{B} = \mathbf{B}'\mathbf{D}_Y^{-1} \Rightarrow \mathbf{A} = \mathbf{A}' - \mathbf{B}\mathbf{C}_Y, \quad \mathbf{D} = \mathbf{D}_X\mathbf{D}_Y^{-1} \Rightarrow \mathbf{C} = \mathbf{C}_X - \mathbf{D}\mathbf{C}_Y$$

with free parameters $\mathbf{F} = \mathbf{C}_Y$, $\mathbf{U} = \mathbf{D}_Y^{-1}$ and \mathbf{U} invertible since \mathbf{D}_Y is supposed to be invertible, this immediately gives (C.59). To verify the McMillan degree let (C.59) be uncontrollable so that by using incontrollability from (C.31) on (C.59) there exists \mathbf{x} so that

$$[\mathbf{A} + \mathbf{B}\mathbf{F} \quad \mathbf{B}\mathbf{U}]^T \mathbf{x} = \lambda [\mathbf{x} \quad \mathbf{0}]$$

This simplifies to

$$[\mathbf{A} \quad \mathbf{B}]^T \mathbf{x} = \lambda [\mathbf{x} \quad \mathbf{0}]$$

since \mathbf{U} is invertible. This means that if (C.59) is uncontrollable then \mathbf{P} is uncontrollable or alternatively if \mathbf{P} is controllable then (C.59) is controllable. Now let (C.59) be unobservable so that by using inobservability from (C.33) on (C.59) there exists \mathbf{x} so that

$$\begin{bmatrix} \mathbf{A} + \mathbf{B}\mathbf{F} \\ \mathbf{C} + \mathbf{D}\mathbf{F} \\ \mathbf{F} \end{bmatrix} \mathbf{x} = \lambda \begin{bmatrix} \mathbf{x} \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}$$

This simplifies to

$$\begin{bmatrix} \mathbf{A} \\ \mathbf{C} \end{bmatrix} \mathbf{x} = \lambda \begin{bmatrix} \mathbf{x} \\ \mathbf{0} \end{bmatrix}$$

This means that if (C.59) is unobservable then \mathbf{P} is unobservable or alternatively if \mathbf{P} is observable then (C.59) is observable. Therefore if \mathbf{P} is minimal then (C.59) is minimal. This means there may be cases where \mathbf{P} is not minimal with (C.59) minimal. Since (C.59) has the same state dimensions than \mathbf{P} , (C.60) follows.

C.3.2 Left Factorization

A system

$$P = Y^{-1}X = \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \hline \mathbf{C} & \mathbf{D} \end{array} \right]$$

has a left factorization in the form of

$$[\mathbf{X} \quad \mathbf{Y}] = \left[\begin{array}{c|cc} \mathbf{A} + \mathbf{L}\mathbf{C} & \mathbf{B} + \mathbf{L}\mathbf{D} & \mathbf{L} \\ \hline \mathbf{V}\mathbf{C} & \mathbf{V}\mathbf{D} & \mathbf{V} \end{array} \right] \quad (\text{C.61})$$

with free parameters \mathbf{L} , \mathbf{V} and \mathbf{V} invertible. Its McMillan degree is bounded by

$$\delta([\mathbf{X} \quad \mathbf{Y}]) \geq \delta(\mathbf{P}) \quad (\text{C.62})$$

see [7] for details.

Proof Start with two systems having common output

$$[\mathbf{X} \quad \mathbf{Y}] = \left[\begin{array}{c|cc} \mathbf{A}' & \mathbf{B}'_X & \mathbf{B}_Y \\ \hline \mathbf{C}' & \mathbf{D}_X & \mathbf{D}_Y \end{array} \right]$$

with \mathbf{D}_Y invertible in order to allow the inversion of \mathbf{Y} . Using flow inversion (C.42) and serial connection (C.46) this leads to

$$\mathbf{Y}^{-1}\mathbf{X} = \left[\begin{array}{cc|c} \mathbf{A}' - \mathbf{B}_Y\mathbf{D}_Y^{-1}\mathbf{C}' & -\mathbf{B}_Y\mathbf{D}_Y^{-1}\mathbf{C}' & -\mathbf{B}_Y\mathbf{D}_Y^{-1}\mathbf{D}_X \\ \mathbf{0} & \mathbf{A}' & \mathbf{B}_X \\ \hline \mathbf{D}_Y^{-1}\mathbf{C}' & \mathbf{D}_Y^{-1}\mathbf{C}' & \mathbf{D}_Y^{-1}\mathbf{D}_X \end{array} \right]$$

Applying the transformation \mathbf{T} as in (C.36) and cutting off the unobservable part leads to

$$\mathbf{Y}^{-1}\mathbf{X} = \left[\begin{array}{c|c} \mathbf{A}' - \mathbf{B}_Y\mathbf{D}_Y^{-1}\mathbf{C}' & \mathbf{B}_X - \mathbf{B}_Y\mathbf{D}_Y^{-1}\mathbf{D}_X \\ \hline \mathbf{D}_Y^{-1}\mathbf{C}' & \mathbf{D}_Y^{-1}\mathbf{D}_X \end{array} \right] \text{ with } \mathbf{T} = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{1} & \mathbf{1} \end{bmatrix}$$

Comparing this to

$$\mathbf{P} = \mathbf{Y}^{-1}\mathbf{X} = \left[\begin{array}{c|c} \mathbf{A} & \mathbf{B} \\ \hline \mathbf{C} & \mathbf{D} \end{array} \right]$$

provides a solution:

$$\mathbf{C} = \mathbf{D}_Y^{-1}\mathbf{C}' \Rightarrow \mathbf{A} = \mathbf{A}' - \mathbf{B}_Y\mathbf{C}, \quad \mathbf{D} = \mathbf{D}_Y^{-1}\mathbf{D}_X \Rightarrow \mathbf{B} = \mathbf{B}_X - \mathbf{B}_Y\mathbf{D}$$

with free parameters $\mathbf{L} = \mathbf{B}_Y$, $\mathbf{V} = \mathbf{D}_Y^{-1}$ and \mathbf{V} invertible since \mathbf{D}_Y is supposed to be invertible, this immediately gives (C.61). To verify the McMillan degree let (C.61) be unobservable so that by using inobservability from (C.33) on (C.61) there exists \mathbf{x} so that

$$\begin{bmatrix} \mathbf{A} + \mathbf{LC} \\ \mathbf{VC} \end{bmatrix} \mathbf{x} = \lambda \begin{bmatrix} \mathbf{x} \\ \mathbf{0} \end{bmatrix}$$

This simplifies to

$$\begin{bmatrix} \mathbf{A} \\ \mathbf{C} \end{bmatrix} \mathbf{x} = \lambda \begin{bmatrix} \mathbf{x} \\ \mathbf{0} \end{bmatrix}$$

since \mathbf{V} is invertible. This means that if (C.61) is unobservable then \mathbf{P} is unobservable or alternatively if \mathbf{P} is observable then (C.61) is observable. Now let (C.61) be uncontrollable so that by using incontrollability from (C.31) on (C.61) there exists \mathbf{x} so that

$$[\mathbf{A} + \mathbf{LC} \quad \mathbf{B} + \mathbf{LD} \quad \mathbf{L}]^T \mathbf{x} = \lambda [\mathbf{x} \quad \mathbf{0} \quad \mathbf{0}]$$

This simplifies to

$$[\mathbf{A} \quad \mathbf{B}]^T \mathbf{x} = \lambda [\mathbf{x} \quad \mathbf{0}]$$

This means that if (C.61) is uncontrollable then \mathbf{P} is uncontrollable or alternatively if \mathbf{P} is controllable then (C.61) is controllable. Therefore if \mathbf{P} is minimal then (C.61) is minimal. This means there may be cases where \mathbf{P} is not minimal with (C.61) minimal. Since (C.61) has the same state dimensions than \mathbf{P} , (C.62) follows.

List of Figures

2.1	A linear time invariant system \mathbf{P} with input \mathbf{u} and output \mathbf{y} . . .	6
2.2	The plant \mathbf{P} and the controller \mathbf{K} form the control loop.	8
2.3	Decomposition of the feedback loop by extraction of the free parameter \mathbf{Q}	12
2.4	Realization of a linear time invariant system with input \mathbf{u} , output \mathbf{y} and system state \mathbf{x}	13
4.1	The plant \mathbf{P} is extended by filters $\mathbf{P}_{zw_i}, \mathbf{P}_{yw_i}, \mathbf{P}_{zu_i}$ to form the extended plant \mathbf{P}_i , which together with the controller \mathbf{K} forms the extended loop \mathbf{H}_i	50
4.2	Principal idea behind the process to assure realness of the controller \mathbf{K} for a real plant \mathbf{P} . The spectra of the dynamic matrices $\mathbf{A}_P, \mathbf{A}_Q$ are denoted by $\mathbf{\Lambda}_P, \mathbf{\Lambda}_Q$ respectively.	65
4.3	Overview of the optimization algorithm.	71
B.1	The modulus of f has its extrema on the closure of $f(D)$	105
C.1	The different system connections: a) parallel b) serial c) common input d) common output e) bridge f) feedback loop.	120

Bibliography

- [1] J. Abedor, K. Nagpal, K. Poolla: “A Linear Matrix Inequality Approach to Peak to Peak Gain Minimization”, *International Journal of Robust and Nonlinear Control*, 1996
- [2] D. Alpay, L. Baratchart, A. Gombani: “On the Differential Structure of Matrix-valued Rational Inner Functions”, *Operator Theory: Advances and Applications*, Vol. 73, Pp. 30-66, 1994
- [3] P. Apkarian, D. Noll: “Nonsmooth H_∞ Synthesis”, *IEEE Transactions on Automatic Control*, Vol. 51, No. 1, 2006
- [4] J. F. Bonnans, A. Shapiro: “Perturbation Analysis of Optimization Problems”, *Springer Series in Operations Research*, 2000
- [5] S. Boyd, L. Vandenberghe: “Convex Optimization”, *Cambridge University Press*, 2004
- [6] P. A. Fuhrmann: “A Polynomial Approach to Linear Algebra”, *Springer*, 1996
- [7] P. A. Fuhrmann, R. J. Ober: “State Space Formulas for coprime Factorizations”, *Operator theory: Advances and Applications: Contributions to Operator Theory and its Applications, The T. A. Ando Birthday Volume*, pages 39–75. *Birkh?auser Verlag*, 1993
- [8] S. Gumussoy, D. Henrion, M. Millstone, M. L. Overton: “Multiobjective Robust Control with HIFOO 2.0”, *Proceedings of the 6th IFAC Symposium on Robust Control Design*, 2009
- [9] X. F. Ji, H. Y. Su, J. Chu: “Peak to Peak Gain Minimization for Uncertain Linear Discrete Systems: A Matrix Inequality Approach”, *Acta Automatica Sinica*, Vol. 33, No. 7, 2007
- [10] K. Nagpal, J. Abedor, K. Poolla: “An LMI Approach to Peak-to-Peak Gain Minimization: Filtering and Control”, *Proceedings of the American Control Conference*, 1994
- [11] M. Olivi: “Parametrization of Rational Lossless Matrices with Applications to Linear Systems Theory”, *Habilitation at Université de Nice Sophia Antipolis*, 2010

- [12] B. T. Polyak, A. V. Nazin, M. V. Topunov, S. A. Nazin: “Rejection of Bounded Disturbances via Invariant Ellipsoids Technique”, Proceedings of the 45th IEEE Conference on Decision and Control, 2006
- [13] M. Vidyasagar: “Control System Synthesis: A Factorization Approach”, MIT Press, 1985
- [14] K. Zhou, J. C. Doyle, K. Glover: “Robust and Optimal Control”, Prentice Hall, 1996

Contrôle Multiobjectif d'ordre réduit

Résumé : Cette thèse présente une méthode permettant d'améliorer un régulateur existant relativement à des spécifications multiples formulées en termes de normes de système L_* , L_2/H_2 et H_∞ . Trois éléments clés distinguent l'approche proposée. Premièrement, un paramétrage des régulateurs à ordre fixé au moyen d'un facteur stable et observable de la boucle fermée. Deuxièmement, un paramétrage minimal de toutes les paires stables observables fondé sur la théorie des systèmes conservatifs et permettant un ajustement infinitésimal du régulateur à améliorer. Et troisièmement, l'utilisation d'une information de sensibilité pour calculer le sous gradient local du critère d'optimisation et construire une suite minimisante.

Mots clés : Contrôle multi-objectifs, Systèmes conservatifs, Paramétrage du régulateur, Factorisation

Reduced Order Multiobjective Control

Abstract: A way to improve an existing controller with respect to a multiobjective specification, written in terms of the L_* , L_2/H_2 and H_∞ system norms, is presented. Three key elements distinguish the proposed approach. First, a parametrization of fixed order controllers by a stable and observable controller parameter. Second, a parametrization of all stable observable pairs based on the theory of lossless systems that allows an infinitesimal adjustment of the controller parameter. And third, the use of a sensitivity information to calculate the local subgradient of the optimization problem so that the controller parameter can be adjusted the right way.

Keywords: Multiobjective Control, Lossless Systems, Controller Parametrization, Factorization

