

Introduction to Systems and Control Theory

Eva Zerz
Lehrstuhl D für Mathematik
RWTH Aachen
52062 Aachen
Germany

Abstract: These lecture notes give a completely self-contained introduction to the control theory of linear time-invariant systems. No prior knowledge is required apart from linear algebra and some basic familiarity with ordinary differential equations. Thus, the course is suited for students of mathematics in their second or third year, and for theoretically inclined engineering students. Because of its appealing simplicity and elegance, the behavioral approach has been adopted to a large extent. A short list of recommended text books on the subject has been added, as a suggestion for further reading.

Contents

1	Introduction	11
1.1	What is a system?	11
1.2	System properties: LTID systems	17
1.3	LTI differential equations	20
1.3.1	The homogeneous equation	20
1.3.2	The inhomogeneous equation	22
1.4	LTI difference equations	23
1.4.1	The homogeneous equation	23
1.4.2	The inhomogeneous equation	26
2	LTID systems: Basic facts	29
2.1	Representations	29
2.2	The fundamental principle	31
2.3	Elimination of latent variables	32
2.4	Inputs, outputs, and autonomous systems	34
2.5	Input-output representations	38

2.5.1	The continuous homogeneous equation	39
2.5.2	The continuous inhomogeneous equation	40
2.5.3	The discrete homogeneous equation ($T = \mathbb{N}$)	42
2.5.4	The discrete inhomogeneous equation	42
2.6	Reduction to first order	43
2.7	State	44
3	Stability	49
3.1	Stability of state space representations	53
3.2	Test for asymptotic stability	55
4	Reachability and controllability	59
4.1	Basic notions for state space systems	59
4.2	Controllable matrix pairs	72
4.3	Asymptotic controllability	78
4.4	Controllable behaviors	79
4.5	Non-linear systems and accessibility	82
5	Feedback control	93
5.1	Static state feedback	93
5.2	Feedback and controllability	94
5.3	Pole placement	96
5.4	Stabilization	98

<i>CONTENTS</i>	5
5.5 Feedback equivalence	100
5.6 Feedback for non-linear systems	106
5.7 Control as interconnection	109
6 Observability	113
6.1 Basic notions for state space systems	113
6.2 Observable matrix pairs	117
6.3 Asymptotic observability	120
6.4 Observable latent variable descriptions	121
6.5 Non-linear systems and zero-input observability	122
7 Observers	127
7.1 State observers	127
7.2 Pole placement	128
7.3 Detection	129
7.4 Compensators	130
8 Transfer matrices	131
8.1 Realization theory	132
8.2 Matrix fraction descriptions	137
8.3 Poles	145
8.4 Zeros	147

Appendices	150
A Distributions	151
B Jordan form	155
C Kronecker-Weierstraß form	157
D Smith form	159
E McMillan form	161
F An optimal control problem	163

Motivating example

A control problem you have, without doubt, already solved many times in your life is to balance a stick on your fingertip. The stick can be seen as a dynamical system. The upright position of the stick is an unstable equilibrium of the system. By moving your hand, you can (to some extent) control the system. Actually, you'll be trying to move your hand in a way that forces the stick to go back to the upright position when it is going to fall. The way you do this is by observing into which direction the stick is falling, and by reacting appropriately, that is, by moving the lower end of the stick into the same direction.

For simplicity, let us study the following closely related (but easier) problem: Suppose that the lower end of the stick is fixed, but one can apply an external torque to it. Thus we have an “inverted pendulum”.

Modelling: First, we need to set up a mathematical model describing the dynamical system under consideration. This step requires some background knowledge from physics. Also, one introduces simplifying assumptions at this stage: for instance, we'll assume that the mass of the stick is concentrated in a point that has distance l from its lower end. Thus, we are disregarding the actual shape and mass distribution of the stick. Secondly, we'll assume that there is no friction.

Let $\theta(t)$ be the angle between the pendulum and the vertical position at time t . Let $u(t)$ be the applied torque at time t . Let m and l be the mass and the length of the pendulum, respectively. Then we have

$$ml^2\ddot{\theta}(t) - mgl \sin(\theta(t)) = u(t)$$

for all $t \in \mathbb{R}$, where g is the gravitational acceleration.

Model analysis: Our model is a non-linear, second order, ordinary differential equation (ODE). Since m, l are, for physical reasons, positive real constants, the

equation can be made explicit, that is, it can be rewritten as

$$\ddot{\theta}(t) = \frac{g}{l} \sin(\theta(t)) + \frac{1}{ml^2} u(t).$$

Let us write $\omega := \sqrt{\frac{g}{l}}$ and $b := \frac{1}{ml^2}$, which are real and positive constants.

Transformation to standard form: To use basic facts from the theory of ODE, we bring the system into the standard form $\dot{x}(t) = f(t, x(t))$. This can be achieved by introducing

$$x(t) := \begin{bmatrix} \theta(t) \\ \dot{\theta}(t) \end{bmatrix}.$$

With this, we can rewrite our equation as

$$\begin{aligned} \dot{x}_1(t) &= x_2(t) \\ \dot{x}_2(t) &= \omega^2 \sin(x_1(t)) + bu(t). \end{aligned}$$

Provided that $u : \mathbb{R} \rightarrow \mathbb{R}$ is a continuous function, classical ODE theory implies that given $x(0) = x_0 \in \mathbb{R}^2$, the associated initial value problem has a unique solution $x : \mathbb{R} \rightarrow \mathbb{R}^2$, from which we get $\theta : \mathbb{R} \rightarrow \mathbb{R}$ via $\theta(t) = [1, 0]x(t)$.

Equilibria of the underlying uncontrolled system: Let us consider the system with $u \equiv 0$, that is, no external torque is acting. Then we have

$$\begin{aligned} \dot{x}_1(t) &= x_2(t) \\ \dot{x}_2(t) &= \omega^2 \sin(x_1(t)), \end{aligned}$$

which has the standard form $\dot{x}(t) = f(x(t))$, where

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}^2, \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \mapsto \begin{bmatrix} x_2 \\ \omega^2 \sin(x_1) \end{bmatrix}.$$

Such a system is called autonomous. The equilibria of $\dot{x}(t) = f(x(t))$ are the points $\bar{x} \in \mathbb{R}^2$ with $f(\bar{x}) = 0$. In the present case, these are the points of the form $\begin{bmatrix} k\pi \\ 0 \end{bmatrix}$, where $k \in \mathbb{Z}$.

Next, we study the stability of the equilibria. A sufficient condition for asymptotic stability of \bar{x} is that all the eigenvalues of

$$A = \left. \frac{\partial f}{\partial x} \right|_{x=\bar{x}}$$

have a negative real part. On the other hand, if A has an eigenvalue with a positive real part, then \bar{x} is an unstable equilibrium. We have

$$A = \left[\begin{array}{cc} 0 & 1 \\ \omega^2 \cos(x_1) & 0 \end{array} \right] \Big|_{x_1=k\pi} = \left[\begin{array}{cc} 0 & 1 \\ \omega^2(-1)^k & 0 \end{array} \right].$$

The characteristic polynomial of A is $\chi_A(s) = s^2 - \omega^2(-1)^k$. If k is even, we have

$$\lambda_{1,2} = \pm\omega$$

and we may conclude that that \bar{x} is unstable. If k is odd, we obtain

$$\lambda_{1,2} = \pm i\omega,$$

and hence the stability of \bar{x} cannot be decided based on the Jacobian matrix A . Nevertheless, using other methods, one can show that \bar{x} is stable for odd k (but not asymptotically stable). Physically, the equilibria with even k correspond to the upright position of the pendulum, and the equilibria with odd k correspond to the downright (hanging) position.

Linearization: We are interested in the behavior of the system near the equilibria. For this, we set

$$\begin{aligned} \tilde{x}_1 &:= x_1 - \bar{x}_1 = x_1 - k\pi \\ \tilde{x}_2 &:= x_2 - \bar{x}_2 = x_2. \end{aligned}$$

Thus, the new variables \tilde{x}_i describe the deviation of x_i from the equilibrium value \bar{x}_i . In terms of \tilde{x}_i , the system reads (omitting the argument t for simplicity)

$$\begin{aligned} \dot{\tilde{x}}_1 &= \dot{x}_1 = x_2 = \tilde{x}_2 \\ \dot{\tilde{x}}_2 &= \dot{x}_2 = \omega^2 \sin(x_1) + bu = \omega^2 \sin(\tilde{x}_1 + k\pi) + bu. \end{aligned}$$

Consider the Taylor series of $\sin(\tilde{x}_1 + k\pi)$ near zero, that is,

$$\sin(\tilde{x}_1 + k\pi) = \underbrace{\sin(k\pi)}_{=0} + \underbrace{\cos(k\pi)}_{=(-1)^k} \tilde{x}_1 + \text{higher order terms.}$$

Linearization means to disregard all higher powers of \tilde{x}_1 in the Taylor series expansion. This is justified when we assume that the deviation from the equilibrium is sufficiently small. We obtain

$$\begin{aligned} \dot{\tilde{x}}_1 &= \tilde{x}_2 \\ \dot{\tilde{x}}_2 &= \omega^2(-1)^k \tilde{x}_1 + bu, \end{aligned}$$

that is,

$$\dot{\tilde{x}} = \underbrace{\begin{bmatrix} 0 & 1 \\ \omega^2(-1)^k & 0 \end{bmatrix}}_{=A} \tilde{x} + \underbrace{\begin{bmatrix} 0 \\ b \end{bmatrix}}_{=B} u.$$

It should come as no surprise that A is nothing but the matrix $\frac{\partial f}{\partial x}|_{x=\bar{x}}$ computed earlier. We call $\dot{\tilde{x}} = A\tilde{x} + Bu$ the linearization of the system at \bar{x} .

Typical control questions: Consider $\bar{x} = 0 \in \mathbb{R}^2$ and

$$\begin{aligned} \dot{x}_1(t) &= x_2(t) \\ \dot{x}_2(t) &= \omega^2 \sin(x_1(t)) + bu(t). \end{aligned}$$

Suppose that $x_0 \in \mathbb{R}^2$ and $\tau > 0$ are given. Find (if possible) $u(\cdot)$ such that the solution of the corresponding initial value problem will satisfy $x(\tau) = \bar{x}$.

In terms of our pendulum, this means: given an initial angle $\theta(0)$ and an initial angular velocity $\dot{\theta}(0)$ (corresponding to a perturbation of the upright position), and given a time $\tau > 0$, find (if possible) a torque function $u(\cdot)$ such that the pendulum returns to its equilibrium position \bar{x} after time τ . If such a $u(\cdot)$ exists, one says that the system can be controlled from x_0 to \bar{x} in time τ .

Related questions: Is this possible for all x_0 that are “close enough” to \bar{x} ? Is this possible for all $\tau > 0$?

Another class of problems stems from the following consideration: Instead of requiring $x(\tau) = 0$ (i.e., the pendulum returns to the upright position in finite time τ), we may be satisfied with $\lim_{t \rightarrow \infty} x(t) = 0$ (i.e., the pendulum approaches the equilibrium asymptotically as time tends to infinity), but we would want this to work for all x_0 in a neighborhood of \bar{x} , with one and the same $u(\cdot)$. Does there exist $u(\cdot)$ such that for all x_0 that are close enough to \bar{x} , the solution of the resulting initial value problem satisfies $\lim_{t \rightarrow \infty} x(t) = 0$? If yes, how can we find such a function $u(\cdot)$?

The same questions are also relevant for the linearized system (where they are easier to test, as one might expect). Moreover, if the answer is positive for the linearization, can we conclude (under certain conditions) that the same is true for the original system?

The focus of these lecture notes is on the linear case. At some points, we will come back to the non-linear situation and we’ll discuss some facts about non-linear systems that can be derived from studying their linearizations.

Chapter 1

Introduction

1.1 What is a system?

The word “system” comes from the Greek word *συστημα*, which originally meant something like “to stand/put/place together.”

Here are some “definitions” I found surfing the web:

A system is

- ... a thing that has components which may act independently, but are connected somehow.
- ... a complex unity formed of many often diverse parts subject to a common plan or serving a common purpose.
- ... a collection of parts that interact with each other to function as a whole.
- ... a group of interrelated elements involved in a collective entity.
- ... a set of interrelated components related by flows of energy, material, or information.
- I will know one when I see it.

If that sounds vague, it's meant to.

We are interested in **dynamical** systems (Greek: *δυναμικός*; original meaning: powerful/strong; here: pertaining to power in motion, involving or causing action or change, opposed to static).

The components/parts/elements of a dynamical system evolve in time. Mathematically speaking, they are functions of time, and they will be called “**signals**.”

We start with some very general, abstract, and comprehensive definitions of a dynamical system. In the following section, we will introduce some important structural properties. Thus we will specialize step by step, until we arrive at a more concrete class of dynamical systems, which will then be investigated in detail.

Definition 1.1 A **dynamical system** Σ is determined by the following data:

- a set T , called the **time set**;
- a set W , called the **signal value set**;
- a set $\mathcal{B} \subseteq W^T$, called the **behavior**.

The set T is our mathematical model of **time**. We will deal exclusively with the following cases:

- $T = \mathbb{Z}$ or a subinterval, especially $\mathbb{N} := \{0, 1, 2, \dots\}$ (**discrete** time);
- $T = \mathbb{R}$ or a subinterval, especially $\mathbb{R}_+ := [0, \infty)$ (**continuous** time).

A **signal** w is a function of time, taking its values in the signal value set W . We write

$$w : T \rightarrow W, t \mapsto w(t).$$

The set W^T is the set of all functions from T to W , therefore it is the **set of all signals**.

Typically, not all signals in W^T can occur in our system (or at least, a system in which *anything* can happen would not be very interesting from the mathematical point of view). Usually, there will be a system **law** which is satisfied only by some signals.

The subset \mathcal{B} of W^T formalizes this law which governs the system. The signals $w \in \mathcal{B}$ are precisely those which are compatible with the system law, that is, they may occur in our system. We also call them **admissible**, and we write

$$\mathcal{B} = \{w \in W^T \mid w \text{ satisfies the system law}\}.$$

Typically, a signal has several, say q , components coming from the same set K (usually, $K = \mathbb{R}$). Then $W = K^q$, and a signal has the form

$$w : T \rightarrow K^q, t \mapsto w(t) = \begin{bmatrix} w_1(t) \\ \vdots \\ w_q(t) \end{bmatrix}.$$

In that case, we call w a **signal vector**, and W^T is the **set of all signal vectors**. Each component w_i of w is again called a **signal**.

This leads to a slight modification of our definition.

Definition 1.2 A dynamical system Σ is determined by the following data:

- a set T , called the **time set**;
- a set K , called the **signal value set**;
- a positive integer q , called the **number of signals**;
- a set $\mathcal{B} \subseteq (K^q)^T = (K^T)^q$, called the **behavior**.

Define $\mathcal{A} := K^T$. This will be called a **signal set**. Finally, we arrive at the following definition.

Definition 1.3 A dynamical system Σ is determined by the following data:

- a set \mathcal{A} , called the **signal set**;
- a positive integer q , called the **number of signals**;
- a set $\mathcal{B} \subseteq \mathcal{A}^q$, called the **behavior**.

Here, \mathcal{A} and q define the setting/mathematical framework for our description of the system: \mathcal{A} is the set of all signals, and \mathcal{A}^q is the set of all signal vectors with q components. The signal vectors $w \in \mathcal{B}$ are precisely those which are compatible with the system law, that is,

$$\mathcal{B} = \{w \in \mathcal{A}^q \mid w \text{ satisfies the system law}\}.$$

Remark 1.4 Prototypes of signal sets:

- the set of all functions from \mathbb{N} to \mathbb{R} ,

$$\mathcal{A} = \mathbb{R}^{\mathbb{N}}$$

(such functions are usually called sequences);

- the set of k times continuously differentiable functions from \mathbb{R} to \mathbb{R} (where $0 \leq k \leq \infty$)

$$\mathcal{A} = \mathcal{C}^k(\mathbb{R});$$

- the set of generalized functions or distributions

$$\mathcal{A} = \mathcal{D}'(\mathbb{R}).$$

Remark 1.5 The intermediary Definition 1.2 is actually superfluous. Any system according to 1.2 can put into the setting of Definition 1.1 by putting $W = K^q$ and into the setting of Definition 1.3 by putting $\mathcal{A} = K^T$.

A system according to Definition 1.1 can be transformed to Definition 1.3 only if $W = K^q$. However, this is not a serious restriction, because it is true for most systems of interest.

In that case, Definition 1.3 is indeed the most general one, because it encompasses distributions.

Example 1.6 The motion of two planets around the sun according to Kepler's laws. The time set is certainly continuous. Whether you choose \mathbb{R} or \mathbb{R}_+ as your time model, depends on your religious and/or scientific beliefs. Let's not touch these delicate issues and choose $T = \mathbb{R}$ for simplicity. The position of a planet in space is determined by its three real coordinates, therefore $K = \mathbb{R}$ and $q = 6$, or $W = \mathbb{R}^6$. If we let \mathcal{A} be the set of all functions from \mathbb{R} to \mathbb{R} , then $\mathcal{A}^q = W^T$. Being inhabitants of the Earth, we may suspect that $\mathcal{A} = \mathcal{C}^k(\mathbb{R})$ for some k which is large enough for comfort. Anyhow, we put

$$\mathcal{B} = \{w \in \mathcal{A}^6 \mid w \text{ satisfies Kepler's laws}\}.$$

Example 1.7 Suppose we have two signals x and b which are linked via

$$\dot{x}(t) = A(t)x(t) + b(t)$$

where $A : \mathbb{R} \rightarrow \mathbb{R}^{n \times n}, t \mapsto A(t)$ is a smooth map. From the theory of ordinary differential equations, we know that if $b : \mathbb{R} \rightarrow \mathbb{R}^n, t \mapsto b(t)$ is continuous, then the initial value problem (where $t_0 \in \mathbb{R}$ and $x_0 \in \mathbb{R}^n$ are arbitrary)

$$\begin{aligned}\dot{x}(t) &= A(t)x(t) + b(t) \\ x(t_0) &= x_0\end{aligned}$$

has a unique solution $x : \mathbb{R} \rightarrow \mathbb{R}^n, t \mapsto x(t)$ which is in $\mathcal{C}^1(\mathbb{R})^n$. Therefore, we may put $T = \mathbb{R}, W = \mathbb{R}^{2n}$ and

$$\mathcal{B} = \{(x, b) \in W^T \mid (x, b) \in \mathcal{C}^1(\mathbb{R})^n \times \mathcal{C}^0(\mathbb{R})^n \text{ and } \dot{x}(t) = A(t)x(t) + b(t)\}.$$

This behavior has the remarkable property

$$\forall b \in \mathcal{C}^0(\mathbb{R})^n \exists x \in \mathcal{C}^1(\mathbb{R})^n : (x, b) \in \mathcal{B}. \quad (1.1)$$

Later on, we will call signals b with this property **inputs**.

In order to avoid having to specify how many times a function is differentiable, it is convenient to work with distributions. We put $\mathcal{A} = \mathcal{D}'(\mathbb{R}), q = 2n$, and

$$\mathcal{B} = \{(x, b) \in \mathcal{A}^q \mid \dot{x} = A(t)x(t) + b(t)\}.$$

This is more general because we may now have discontinuous b , such as, e.g., the Heaviside function. Again, we have (compare with (1.1))

$$\forall b \in \mathcal{A}^n \exists x \in \mathcal{A}^n : (x, b) \in \mathcal{B}.$$

Take the scalar example

$$\dot{x}(t) = x(t) + b(t).$$

The classical solutions are

$$x(t) = e^t x_0 + \int_0^t e^{t-\tau} b(\tau) d\tau$$

where $x_0 \in \mathbb{R}$ is arbitrary. If b is the Heaviside function

$$b(t) = \begin{cases} 0 & \text{if } t < 0 \\ 1 & \text{if } t \geq 0 \end{cases}$$

then we obtain

$$x(t) = e^t x_0 + \begin{cases} 0 & \text{if } t < 0 \\ e^t - 1 & \text{if } t \geq 0 \end{cases}$$

which is \mathcal{C}^0 , but not \mathcal{C}^1 , and hence not a classical solution.

Example 1.8 Genetics/Gender-linked genes: An allele (a certain form of a gene) is located on the X-chromosome. Females have two X-chromosomes, males have one X- and one Y-chromosome. Let $p^f(i)$ be the frequency of the allele in the female gene pool of the i -th generation and let $p^m(i)$ be the same for the male gene pool. Since a son inherits his X-chromosome from the mother,

$$p^m(i+1) = p^f(i)$$

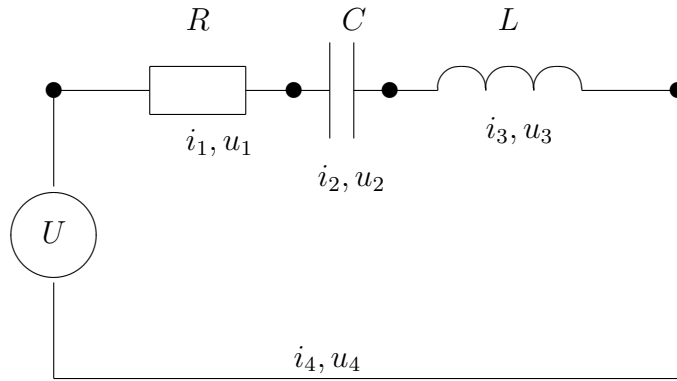
and since a daughter receives one X-chromosome from the father and one from the mother,

$$p^f(i+1) = \frac{1}{2}(p^f(i) + p^m(i)).$$

The time set is discrete. Both $T = \mathbb{Z}$ and $T = \mathbb{N}$ are suitable choices. We have two signals p^m and p^f , taking their values in \mathbb{R} , hence $W = \mathbb{R}^2$ and

$$\mathcal{B} = \left\{ \begin{bmatrix} p^m \\ p^f \end{bmatrix} \in W^T \mid \begin{bmatrix} p^m(i+1) \\ p^f(i+1) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} p^m(i) \\ p^f(i) \end{bmatrix} \text{ for all } i \in T \right\}.$$

Example 1.9 An electrical circuit.



This system involves 9 signals: U, i_1, \dots, i_4 , and u_1, \dots, u_4 . Kirchhoff's current law says that

$$i_1 = i_2 = i_3 = i_4 =: I$$

and Kirchhoff's voltage law says that

$$u_1 + u_2 + u_3 + u_4 = 0.$$

We have a free voltage source, $u_4 = U$. Moreover, we have the following constituent equations for the resistor (R), capacitor (C), and inductor (L), respectively:

$$u_1 = Ri_1 = RI, \quad C \frac{du_2}{dt} = i_2 = I, \quad L \frac{di_3}{dt} = L \frac{dI}{dt} = u_3.$$

Let us describe the system in terms of the signals I, U, u_1, u_2, u_3 (the other 4 signals are easily obtained from them.) We put $\mathcal{A} = \mathcal{D}'(\mathbb{R})$ or $\mathcal{D}'(\mathbb{R}_+)$, $q = 5$ and

$$\mathcal{B} = \{(I, U, u_1, u_2, u_3)^T \in \mathcal{A}^5 \mid u_1 + u_2 + u_3 + U = 0, u_1 = RI, C\dot{u}_2 = I, LI = u_3\}.$$

1.2 System properties: LTID systems

For the definition of linearity, we need to fix an underlying number field K . We will focus on $K = \mathbb{R}$, and therefore we give the definition for that case only.

Definition 1.10 A dynamical system $\Sigma = (\mathcal{A}, q, \mathcal{B})$ is called **linear** if \mathcal{A} is a real vector space, and \mathcal{B} is a subspace of \mathcal{A}^q .

The first requirement means that linear combinations of signals are again signals,

$$a_1, a_2 \in \mathcal{A}, \lambda_1, \lambda_2 \in \mathbb{R} \quad \Rightarrow \quad \lambda_1 a_1 + \lambda_2 a_2 \in \mathcal{A}$$

and the second requirement means that linear combinations of admissible signal vectors are again admissible signal vectors,

$$w_1, w_2 \in \mathcal{B}, \lambda_1, \lambda_2 \in \mathbb{R} \quad \Rightarrow \quad \lambda_1 w_1 + \lambda_2 w_2 \in \mathcal{B}.$$

We call this condition the **superposition principle**.

Remark 1.11 All signal spaces in 1.4 are real vector spaces.

Any \mathcal{A} of the form $\mathcal{A} = \mathbb{R}^T$ is a real vector space. For $a_1, a_2 \in \mathcal{A}$, $\lambda_1, \lambda_2 \in \mathbb{R}$, we have

$$\lambda_1 a_1 + \lambda_2 a_2 : T \rightarrow \mathbb{R}, t \mapsto \lambda_1 a_1(t) + \lambda_2 a_2(t).$$

More generally: If W is a real vector space, then so is W^T .

Definition 1.12 A dynamical system $\Sigma = (T, W, \mathcal{B})$ is called **linear** if W is a real vector space and \mathcal{B} is a subspace of W^T .

Definition 1.13 Let T be such that

$$t_1, t_2 \in T \quad \Rightarrow \quad t_1 + t_2 \in T. \quad (1.2)$$

For $\tau \in T$, we define the shift operator σ_τ by

$$\sigma_\tau : W^T \rightarrow W^T, w \mapsto \sigma_\tau w$$

where

$$(\sigma_\tau w)(t) = w(t + \tau).$$

A dynamical system $\Sigma = (T, W, \mathcal{B})$ is called **shift-invariant** (or: time-invariant) if for all $\tau \in T$

$$w \in \mathcal{B} \quad \Rightarrow \quad \sigma_\tau w \in \mathcal{B}.$$

Remark 1.14 Note that the time sets $T = \mathbb{R}$, $T = \mathbb{R}_+$, $T = \mathbb{Z}$, and $T = \mathbb{N}$ all satisfy (1.2). Recall that we always assume that T is one of these four sets. For $\mathcal{A} = \mathcal{D}'(\mathbb{R})$, the shift operator $\sigma_\tau : \mathcal{A} \rightarrow \mathcal{A}$ is defined by $(\sigma_\tau D)(\varphi) = D(\sigma_{-\tau}\varphi)$, which is motivated by the requirement $\sigma_\tau D_f = D_{\sigma_\tau f}$. Then shift-invariance can be defined as above, that is, by requiring that $\sigma_\tau \mathcal{B} \subseteq \mathcal{B}$ for all $\tau \in \mathbb{R}$.

Definition 1.15 A dynamical system $\Sigma = (T, W, \mathcal{B})$ is called a **differential (difference) system** if its time set is continuous (discrete) and its system law is given by differential (difference) equations.

This is the class of systems we will mainly study: linear, time-invariant (LTI) differential (difference) systems. The system laws will have the following form:

Differential systems: Systems of linear differential equations with constant coefficients. These can be put in the form

$$(R_d \frac{d^d}{dt^d} + \dots + R_1 \frac{d}{dt} + R_0)w = 0 \quad (1.3)$$

where $R_i \in \mathbb{R}^{p \times q}$ are real matrices. We define

$$R := R_d s^d + \dots + R_1 s + R_0.$$

Then R is a polynomial $p \times q$ matrix, and we may rewrite (1.3) in the concise form

$$R\left(\frac{d}{dt}\right)w = 0.$$

Difference systems: Systems of linear difference equations with constant coefficients. These can be put in the form

$$R_d w(t+d) + \dots + R_1 w(t+1) + R_0 w(t) = 0 \text{ for all } t \in T$$

where $R_i \in \mathbb{R}^{p \times q}$ are real matrices. Using the shift operator σ_τ , we may write

$$(R_d \sigma_d + \dots + R_1 \sigma_1 + R_0)w = 0. \quad (1.4)$$

We define $\sigma := \sigma_1$, then $\sigma_k = \sigma^k$ (k -fold application of σ). If we put again

$$R := R_d s^d + \dots + R_1 s + R_0$$

then we can write (1.4) in the concise form

$$R(\sigma)w = 0.$$

Remark 1.16 The examples discussed so far can be classified as follows.

Example 1.6: Kepler (1571-1630) formulated his laws in a non-differential way. The differential calculus was developed by Newton (1642-1727) and Leibniz (1646-1716). It was in fact Newton who came up with a differential equation for the motion of planets around the sun (the law of gravity). However, it is quite clear that Kepler's laws describe a non-linear and time-invariant system.

Example 1.7: The behavior given by $\dot{x}(t) = A(t)x(t) + b(t)$ is differential and linear. In general, it is not time-invariant. In fact, it is time-invariant if and only if A is a constant function.

Example 1.8: The system given by

$$w(t+1) = Aw(t) \text{ for all } t \in T$$

where

$$A = \begin{bmatrix} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

is a linear time-invariant difference system. Using the shift operator, it reads $\sigma w = Aw$ or

$$(\sigma I - A)w = 0.$$

The polynomial matrix R takes the form

$$R = sI - A = \begin{bmatrix} s & -1 \\ -\frac{1}{2} & s - \frac{1}{2} \end{bmatrix}.$$

Example 1.9: The system of the electrical circuit is a linear, time-invariant and differential. The system law reads

$$\begin{bmatrix} 0 & 1 & 1 & 1 & 1 \\ R & 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & C \frac{d}{dt} & 0 \\ L \frac{d}{dt} & 0 & 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} I \\ U \\ u_1 \\ u_2 \\ u_3 \end{bmatrix} = 0.$$

We may rewrite this in the form $\hat{R}(\frac{d}{dt})w = 0$ by putting

$$w = [I \quad U \quad u_1 \quad u_2 \quad u_3]^T$$

and

$$\hat{R} = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 \\ R & 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & Cs & 0 \\ Ls & 0 & 0 & 0 & -1 \end{bmatrix}.$$

1.3 LTI differential equations and their distributional solutions

Consider $P\left(\frac{d}{dt}\right)y = v$, where $P \in \mathbb{R}[s]^{p \times p}$ is non-singular, and $v \in \mathcal{D}'(T)^p$ for $T = \mathbb{R}$ or $T = \mathbb{R}_+$. Let y_p be one particular solution. Then any solution y has the form $y = y_p + y_h$, where y_h is a solution of the associated homogeneous equation $P\left(\frac{d}{dt}\right)y_h = 0$.

1.3.1 The homogeneous equation

Consider the homogeneous equation $P\left(\frac{d}{dt}\right)y = 0$. Assume that

$$P = P_d s^d + \dots + P_1 s + P_0$$

with coefficient matrices $P_i \in \mathbb{R}^{p \times p}$. After reduction to first order, we have

$$K\dot{\xi} = L\xi$$

where $K, L \in \mathbb{R}^{n \times n}$. Without loss of generality, we may put $n = dp$,

$$\xi = \begin{bmatrix} y \\ \frac{d}{dt}y \\ \vdots \\ \frac{d^{d-1}}{dt^{d-1}}y \end{bmatrix}, \quad K = \begin{bmatrix} I & & & \\ & \ddots & & \\ & & I & \\ & & & P_d \end{bmatrix} \quad \text{and} \quad L = \begin{bmatrix} 0 & I & & \\ \vdots & & \ddots & \\ 0 & & & I \\ -P_0 & -P_1 & \cdots & -P_{d-1} \end{bmatrix}.$$

Note that then $\det(P) = \det(sK - L)$, $\det(K) = \det(P_d)$, and $\det(L) = \pm \det(P_0)$.

Case 1: K is invertible, that is, after putting $A := K^{-1}L$, the system has the explicit form $\dot{\xi} = A\xi$. We say that ξ is a classical solution if it is \mathcal{C}^1 and $\dot{\xi} = A\xi$. It is well-known that

$$\Phi(t) := e^{At}$$

is a fundamental matrix for this system, that is, its n columns are a basis of the classical solution space. In other words, any classical solution ξ has the form $\xi(t) = e^{At}\xi_0$ for some $\xi_0 \in \mathbb{R}^n$. Note that the entries of Φ are \mathcal{C}^∞ . Hence every classical solution is also \mathcal{C}^∞ . It is known that the Wronski determinant $\det(\Phi(t)) = e^{\text{trace}(A)t} \neq 0$ for all $t \in T$, and thus also the entries of Φ^{-1} are in \mathcal{C}^∞ , in fact, $\Phi^{-1}(t) = e^{-At}$.

Theorem 1.17 The equation $\dot{\xi} = A\xi$ has no distributional solutions apart from the classical solutions, that is, $\xi(t) = e^{At}\xi_0$, where $\xi_0 \in \mathbb{R}^n$.

Proof: Let $\xi \in \mathcal{D}'(T)^n$ be a distributional solution. Set $\eta := \Phi^{-1}\xi$, then $\Phi\eta = \xi$. Differentiation yields

$$\dot{\xi} = \dot{\Phi}\eta + \Phi\dot{\eta} = A\Phi\eta + \Phi\dot{\eta}.$$

On the other hand, we have $\dot{\xi} = A\xi = A\Phi\eta$ by assumption. Therefore $\Phi\dot{\eta} = 0$ and thus $\dot{\eta} = 0$. According to Lemma A.1, this implies $\eta = \eta_0$, where η_0 is a constant vector. Thus

$$\xi = \Phi\eta = \Phi\eta_0$$

showing that the distributional solution ξ is indeed a classical solution. \square

Case 2: K is not invertible. Recall that $\det(sK - L) = \det(P)$ which is non-zero by assumption.

Theorem 1.18 (Kronecker-Weierstraß form) Let $K, L \in \mathbb{R}^{n \times n}$ be such that $\det(sK - L) \neq 0$. Then there exist non-singular matrices $U, V \in \mathbb{R}^{n \times n}$ such that

$$UKV = \begin{bmatrix} I_\nu & 0 \\ 0 & N \end{bmatrix} \quad \text{and} \quad ULV = \begin{bmatrix} A_1 & 0 \\ 0 & I_{n-\nu} \end{bmatrix}$$

where N is a nilpotent matrix, that is, $N^k = 0$ for some $k \in \mathbb{N}$. The number ν , and the nilpotency index of N , that is, the smallest integer κ such that $N^\kappa = 0$, are uniquely determined by K, L .

Since $K\dot{\xi} = L\xi$ is equivalent to $UKVV^{-1}\dot{\xi} = ULVV^{-1}\xi$, we set $x := V^{-1}\xi$ and obtain

$$\begin{aligned} \dot{x}_1 &= A_1 x_1 \\ Nx_2 &= x_2. \end{aligned}$$

The second equation implies (by repeated differentiation and multiplication by N) that $x_2 = 0$. Therefore, Case 2 can be reduced to Case 1.

Corollary 1.19 The equation $K\dot{\xi} = L\xi$ has no distributional solutions apart from the classical solutions, that is, $\xi(t) = Vx(t)$, where $x_1(t) = e^{A_1 t}x_{10}$ for some $x_{10} \in \mathbb{R}^\nu$ and $x_2(t) = 0$ for all t .

Summing up the two cases, we have the following result. We use that the entries of $\Phi(t) = e^{At}$ have the form

$$\Phi_{ij}(t) = \sum_{\lambda} a_{\lambda}(t)e^{\lambda t}$$

where $\lambda \in \mathbb{C}$ are the eigenvalues of A and $a_\lambda \in \mathbb{C}[t]$ are suitable polynomials. Note that

$$\det(P) = \det(sK - L) = \frac{1}{uv} \det(sI - A_1) \det(sN - I),$$

where $u = \det(U)$ and $v = \det(V)$ are non-zero constants. Since $\det(\lambda N - I) \neq 0$ for all $\lambda \in \mathbb{C}$, the zeros of $\det(P)$ are precisely the eigenvalues of A_1 .

Theorem 1.20 The system $P(\frac{d}{dt})y = 0$ has no distributional solutions apart from the classical solutions, that is, functions of the form

$$y(t) = \sum_{\lambda} a_{\lambda}(t)e^{\lambda t}$$

where $\lambda \in \mathbb{C}$ are the zeros of $\det(P)$, and $a_{\lambda} \in \mathbb{C}[t]^p$ are suitable polynomial vectors.

1.3.2 The inhomogeneous equation

Let $\mathcal{A} = \mathcal{D}'(T)$. Consider $P(\frac{d}{dt})y = v$, where $v \in \mathcal{A}^p$ is given. After reduction to first order, we obtain

$$K\dot{\xi} = L\xi + Mv$$

where $K, L \in \mathbb{R}^{n \times n}$ and $M \in \mathbb{R}^{n \times p}$.

Case 1: K is invertible, that is, after putting $A := K^{-1}L$, $B := K^{-1}M$, the system has the explicit form

$$\dot{\xi} = A\xi + Bv.$$

One uses the well-known “variation of constants” trick and sets

$$\xi = \Phi\psi$$

where $\Phi(t) = e^{At}$ is the fundamental matrix for the homogeneous equation, to obtain

$$\dot{\psi} = \Phi^{-1}Bv.$$

According to Theorem A.2, there exists a distribution ψ which satisfies this equation. Then $\xi = \Phi\psi$ is a particular solution of the inhomogeneous equation. Note that if v is \mathcal{C}^0 , that is, continuous, then ξ is \mathcal{C}^1 and hence it is a classical solution.

Case 2: K is not invertible. Without loss of generality, rewrite $K\dot{\xi} = L\xi + Mv$ as $UKVV^{-1}\dot{\xi} = ULVV^{-1}\xi + UMv$, where U and V are as in Theorem 1.18. Then, putting $\xi = Vx$ and $B = UM$, the equations read

$$\begin{aligned}\dot{x}_1 &= A_1x_1 + B_1v \\ Nx_2 &= x_2 + B_2v.\end{aligned}$$

Since $N^\kappa = 0$, this implies that

$$x_2 = -(B_2v + NB_2\dot{v} + \dots + N^{\kappa-1}B_2v^{(\kappa-1)}).$$

The equation for x_1 is again as in Case 1.

Theorem 1.21 Let $P \in \mathbb{R}[s]^{p \times p}$ be non-singular. For every $v \in \mathcal{A}^p$, there exists $y \in \mathcal{A}^p$ such that $P(\frac{d}{dt})y = v$.

This theorem is a special case of the so-called **fundamental principle**, which we will discuss later.

1.4 LTI difference equations

Consider $P(\sigma)y = v$, where $P \in \mathbb{R}[s]^{p \times p}$ is non-singular, and $v \in (\mathbb{R}^T)^p$ for $T = \mathbb{N}$ or $T = \mathbb{Z}$. Let y_p be one particular solution. Then any solution y has the form $y = y_p + y_h$, where y_h is a solution of the associated homogeneous equation $P(\sigma)y_h = 0$.

1.4.1 The homogeneous equation

Consider the homogeneous equation $P(\sigma)y = 0$. Assume that

$$P = P_d s^d + \dots + P_1 s + P_0$$

with coefficient matrices $P_i \in \mathbb{R}^{p \times p}$. After reduction to first order, we have

$$K\sigma\xi = L\xi,$$

that is

$$K\xi(t+1) = L\xi(t) \quad \text{for all } t \in T$$

where $K, L \in \mathbb{R}^{n \times n}$. We may choose K, L as in Section 1.3 if we set

$$\xi = \begin{bmatrix} y \\ \sigma y \\ \vdots \\ \sigma^{d-1}y \end{bmatrix}.$$

Time set $T = \mathbb{N}$

Case 1: K is invertible, that is, after putting $A := K^{-1}L$, the system has the explicit form $\sigma\xi = A\xi$. It is well-known that

$$\Phi(t) := A^t$$

is a fundamental matrix for this system, that is, its n columns are a basis of the solution space. In other words, any solution ξ has the form $\xi(t) = A^t\xi_0$ for some $\xi_0 \in \mathbb{R}^n$.

Case 2: K is not invertible. Using Theorem 1.18, the system can be rewritten as

$$\begin{aligned} \sigma x_1 &= A_1 x_1 \\ N\sigma x_2 &= x_2. \end{aligned}$$

The second equation implies that $x_2 = 0$, because $N^k = 0$ and thus $x_2(t) = Nx_2(t+1) = \dots = N^k x_2(t+k) = 0$ for all t . Therefore, Case 2 can be reduced to Case 1.

Lemma 1.22 The equation $K\sigma\xi = L\xi$ has the solutions $\xi(t) = Vx(t)$, where $x_1(t) = A_1^t x_{10}$ for some $x_{10} \in \mathbb{R}^r$ and $x_2(t) = 0$ for all t .

Summing up the two cases, we have the following result. We use that the entries of $\Phi(t) = A^t$, A invertible, have the form

$$\Phi_{ij}(t) = \sum_{\lambda} a_{\lambda}(t)\lambda^t$$

where $\lambda \in \mathbb{C}$ are the eigenvalues of A and $a_{\lambda} \in \mathbb{C}[t]$ are suitable polynomials. (If A is singular, the formula is still valid for all $t \geq n$.) Note that

$$\det(P) = \det(sK - L) = \frac{1}{uv} \det(sI - A_1) \det(sN - I),$$

where $u = \det(U)$ and $v = \det(V)$ are non-zero constants. Since $\det(\lambda N - I) \neq 0$ for all $\lambda \in \mathbb{C}$, the zeros of $\det(P)$ are precisely the eigenvalues of A_1 .

Theorem 1.23 The solutions of the system $P(\sigma)y = 0$ have the form

$$y(t) = \sum_{\lambda} a_{\lambda}(t)\lambda^t$$

for large enough t , where $\lambda \in \mathbb{C}$ are the zeros of $\det(P)$, and $a_{\lambda} \in \mathbb{C}[t]^p$ are suitable polynomial vectors.

Time set $T = \mathbb{Z}$

Case 1: K, L are both invertible. Then $A := K^{-1}L$ is also invertible, and thus $\Phi(t) = A^t$ is again a fundamental matrix (note that we need invertibility of A for $\Phi(t)$ to be defined also for negative values of t).

Case 2: L is invertible, but not K . As usual, we may transform the system into Kronecker-Weierstraß form

$$U_1KV_1 = \begin{bmatrix} I & 0 \\ 0 & N_1 \end{bmatrix} \quad \text{and} \quad U_1LV_1 = \begin{bmatrix} A_1 & 0 \\ 0 & I \end{bmatrix}$$

to obtain, for $x = V_1^{-1}\xi$,

$$\begin{aligned} \sigma x_1 &= A_1 x_1 \\ N_1 \sigma x_2 &= x_2 \end{aligned}$$

where A_1 is invertible. Again, the second equation implies that $x_2 = 0$. Thus we have reduced Case 2 to Case 1.

Case 3: K is invertible, but not L . We rewrite

$$K\xi(t+1) = L\xi(t) \quad \text{for all } t \in \mathbb{Z}$$

as

$$L\xi(\tau-1) = K\xi(\tau) \quad \text{for all } \tau \in \mathbb{Z}$$

and compute the Kronecker-Weierstraß form

$$U_2LV_2 = \begin{bmatrix} I & 0 \\ 0 & N_2 \end{bmatrix} \quad \text{and} \quad U_2KV_2 = \begin{bmatrix} A_2 & 0 \\ 0 & I \end{bmatrix}$$

where A_2 is invertible. Then we obtain the new system, for $y = V_2^{-1}\xi$,

$$\begin{aligned} y_1(\tau-1) &= A_2 y_1(\tau) \\ N_2 y_2(\tau-1) &= y_2(\tau). \end{aligned}$$

Then $y_2 = 0$ and $y_1(\tau) = A_2^{-\tau} y_{10}$ for some y_{10} . Note that

$$\det(P) = \det(sK - L) = \frac{1}{u_2 v_2} \det(sA_2 - I) \det(sI - N_2)$$

and thus, since zero is the only eigenvalue of the nilpotent matrix N_2 ,

$$\det(P)(\lambda) = 0 \quad \Leftrightarrow \quad \lambda = 0 \text{ or } \det(\lambda A_2 - I) = 0.$$

Thus the non-zero zeros of $\det(P)$ coincide with the eigenvalues of A_2^{-1} .

Case 4: K, L both singular. Then one uses the Kronecker-Weierstraß form as in Case 2 to obtain the values $\xi(t)$ for $t > 0$ and the Kronecker-Weierstraß form from Case 3 to determine the values $\xi(t)$ for $t < 0$. We omit the details.

Theorem 1.24 The solutions of $P(\sigma)y = 0$ have the form

$$y(t) = \sum_{\lambda} a_{\lambda}(t) \lambda^t$$

where $\lambda \in \mathbb{C}$ are the non-zero zeros of $\det(P)$ and $a_{\lambda} \in \mathbb{C}[t]^p$ are suitable polynomial vectors.

1.4.2 The inhomogeneous equation

For simplicity, we treat only the case $T = \mathbb{N}$, that is, $\mathcal{A} = \mathbb{R}^{\mathbb{N}}$. Consider $P(\sigma)y = v$, where $v \in \mathcal{A}^p$ is given. After reduction to first order, we obtain

$$K\sigma\xi = L\xi + Mv$$

where $K, L \in \mathbb{R}^{n \times n}$ and $M \in \mathbb{R}^{n \times p}$.

Case 1: K is invertible, that is, after putting $A := K^{-1}L$, $B := K^{-1}M$, the system has the form

$$\sigma\xi = A\xi + Bv.$$

Then

$$\xi(t) = \sum_{i=0}^{t-1} A^{t-i-1} Bv(i)$$

is a particular solution of the inhomogeneous equation.

Case 2: K is not invertible. Without loss of generality, rewrite $K\sigma\xi = L\xi + Mv$ as $UKVV^{-1}\sigma\xi = ULVV^{-1}\xi + UMv$, where U and V are as in Theorem 1.18. Then, putting $\xi = Vx$ and $B = UM$, the equations read

$$\begin{aligned}\sigma x_1 &= A_1 x_1 + B_1 v \\ N\sigma x_2 &= x_2 + B_2 v.\end{aligned}$$

Since $N^\kappa = 0$, this implies that

$$x_2 = -(B_2 v + NB_2 \sigma v + \dots + N^{\kappa-1} B_2 \sigma^{\kappa-1} v).$$

The equation for x_1 is again as in Case 1.

Theorem 1.25 Let $P \in \mathbb{R}[s]^{p \times p}$ be non-singular. For every $v \in \mathcal{A}^p$, there exists $y \in \mathcal{A}^p$ such that $P(\sigma)y = v$.

This is a discrete version of the fundamental principle.

Chapter 2

LTID systems: Basic facts

In the following, let R be a $p \times q$ polynomial matrix in the variable s , with real coefficients. We write

$$R \in \mathbb{R}[s]^{p \times q}.$$

From now on, we will restrict our discussion to the following **standard models**:

- The **continuous** standard model is

$$\mathcal{B} = \{w \in \mathcal{A}^q \mid R(\frac{d}{dt})w = 0\}$$

where $\mathcal{A} = \mathcal{D}'(T)$ with $T = \mathbb{R}$ or $T = \mathbb{R}_+$.

- The **discrete** standard model is

$$\mathcal{B} = \{w \in \mathcal{A}^q \mid R(\sigma)w = 0\}$$

where $\mathcal{A} = \mathbb{R}^T$ with $T = \mathbb{N}$ or $T = \mathbb{Z}$.

2.1 Representations

The polynomial matrix R is called a **representation** of \mathcal{B} . Note that once \mathcal{A} and q are fixed, the behavior \mathcal{B} is uniquely determined by R . Conversely, there are many polynomial matrices which represent the same behavior. To see this, we introduce the concept of a unimodular matrix.

Definition 2.1 A square polynomial matrix U is called **unimodular** if its determinant is a non-zero constant, that is, $\det(U) \in \mathbb{R} \setminus \{0\}$. This is equivalent to the existence of a polynomial matrix V such that

$$UV = VU = I.$$

Clearly, V is the inverse of U , which exists because U is non-singular, i.e., $\det(U) \neq 0$. Unimodularity is much stronger than non-singularity, the crucial point is that U possesses a *polynomial* (rather than a rational) inverse.

We observe that pre-multiplication by a unimodular matrix U does not change the behavior represented by R . More precisely, R and $\hat{R} = UR$ represent the same behavior, because

$$R\left(\frac{d}{dt}\right)w = 0 \quad \Rightarrow \quad U\left(\frac{d}{dt}\right)R\left(\frac{d}{dt}\right)w = \hat{R}\left(\frac{d}{dt}\right)w = 0$$

and

$$\hat{R}\left(\frac{d}{dt}\right)w = 0 \quad \Rightarrow \quad V\left(\frac{d}{dt}\right)\hat{R}\left(\frac{d}{dt}\right)w = R\left(\frac{d}{dt}\right)w = 0$$

where V is the polynomial inverse of U . The same holds if we replace $\frac{d}{dt}$ by σ .

Definition 2.2 A polynomial matrix R is called a **minimal** representation of \mathcal{B} if there exists no polynomial matrix which represents the same behavior and has a smaller number of rows.

The most important fact about polynomial matrices is stated next.

Theorem 2.3 (Smith form) For every polynomial matrix $R \in \mathbb{R}[s]^{p \times q}$ there exist unimodular matrices $U \in \mathbb{R}[s]^{p \times p}$ and $V \in \mathbb{R}[s]^{q \times q}$ such that

$$URV = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} \quad (2.1)$$

where $D \in \mathbb{R}[s]^{r \times r}$ is a non-singular diagonal matrix

$$D = \begin{bmatrix} d_1 & & \\ & \ddots & \\ & & d_r \end{bmatrix}$$

with $d_1 | d_2 | \dots | d_r$. This notation means that for $i = 1, \dots, r-1$, the polynomial d_i divides d_{i+1} , that is, $d_{i+1} = d_i e_i$ for some polynomial e_i . Clearly, the integer r is precisely the rank of the matrix R (over the quotient field $\mathbb{R}(s)$ of $\mathbb{R}[s]$). The matrix on the right hand side of (2.1) is called **Smith form** of R .

Corollary 2.4 For every polynomial matrix $R \in \mathbb{R}[s]^{p \times q}$ there exists a unimodular matrix $U \in \mathbb{R}[s]^{p \times p}$ such that

$$UR = \begin{bmatrix} R_1 \\ 0 \end{bmatrix}$$

where $R_1 \in \mathbb{R}[s]^{r \times q}$ has full row rank, that is, $\text{rank}(R_1) = r$.

Proof: Let U and V be as in Theorem 2.3, and set $W := V^{-1}$. Then we have

$$UR = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} W_1 \\ W_2 \end{bmatrix}$$

and we set $R_1 := DW_1$, which has full row rank, by construction. \square

Since R and R_1 represent the same behavior, we have the following conclusion.

Lemma 2.5 Any \mathcal{B} possesses a representation matrix with full row rank. If R is a minimal representation of \mathcal{B} , then R has full row rank.

Proof: The first statement follows directly from our considerations above. If $R \in \mathbb{R}[s]^{p \times q}$ is a representation with $\text{rank}(R) < p$, then there exists a representation with less than p rows, according to Corollary 2.4 (noting that $\text{rank}(R) = \text{rank}(R_1)$), and thus R cannot be minimal. \square

We will see later on that in fact, a representation R of \mathcal{B} is minimal if and only if R has full row rank.

2.2 The fundamental principle

We restate the fundamental principle for convenience.

Theorem 2.6 (Fundamental principle, square matrix version) Let $P \in \mathbb{R}[s]^{p \times p}$ be a non-singular polynomial matrix. Then

$$P\left(\frac{d}{dt}\right)y = v \quad \text{or} \quad P(\sigma)y = v$$

possesses a solution $y \in \mathcal{A}^p$ for every choice of the right hand side signal $v \in \mathcal{A}^p$.

The following variant for rectangular matrices with full row rank is an easy consequence.

Corollary 2.7 (Fundamental principle, non-square matrix version) Let $R \in \mathbb{R}[s]^{p \times q}$ be a full row rank polynomial matrix, that is, $\text{rank}(R) = p$. Then

$$R\left(\frac{d}{dt}\right)w = v \quad \text{or} \quad R(\sigma)w = v$$

possesses a solution $w \in \mathcal{A}^q$ for every choice of the right hand side signal $v \in \mathcal{A}^p$.

Proof: It suffices to do the continuous case. Since R has full row rank, its Smith form takes the form $URV = [D, 0]$. Thus we have $RV = [P, 0]$, where V is unimodular, and P is square and non-singular. Given v , we rewrite $R\left(\frac{d}{dt}\right)w = v$ as $(RVV^{-1})\left(\frac{d}{dt}\right)w = v$ or $P\left(\frac{d}{dt}\right)\tilde{w}_1 = v$, where $\tilde{w} := V^{-1}\left(\frac{d}{dt}\right)w$ is partitioned accordingly. From the square case, we know that there exists such a \tilde{w}_1 . The vector \tilde{w}_2 can be chosen arbitrarily, and we get w via $w = V\left(\frac{d}{dt}\right)\tilde{w}$. \square

The most important consequence of the fundamental principle is stated next.

2.3 Elimination of latent variables

Often, the components of a signal vector can be divided into two classes: a set of components we are truly interested in (called **manifest** variables), and a set of components that were introduced as auxiliary variables during the modelling process (called **latent** variables).

Remark 2.8 For instance, in the electrical circuit of Example 1.9, one is usually interested in the signals U and I because they represent the overall voltage and current that are relevant for this network. On the other hand, the voltages u_1, u_2, u_3 were introduced for determining the network equations, and those quantities are not really interesting by themselves. In fact, one would like to get rid of them in order to describe the relation between I and U alone. This process is called **elimination of latent variables**.

Consider

$$\mathcal{B} = \left\{ \left[\begin{array}{c} w \\ l \end{array} \right] \in \mathcal{A}^{q+r} \mid R\left(\frac{d}{dt}\right)w = M\left(\frac{d}{dt}\right)l \right\},$$

where $R \in \mathbb{R}[s]^{p \times q}$ and $M \in \mathbb{R}[s]^{p \times r}$ (we restrict to continuous time; the discrete case is analogous). Let's say that the components of w are the manifest variables, and l represents the latent variables. Then we are actually interested in the projection of \mathcal{B} onto the first q variables, that is, in the **system with latent variables** given by

$$\mathcal{B}_l = \{w \in \mathcal{A}^q \mid \exists l \in \mathcal{A}^r : R(\frac{d}{dt})w = M(\frac{d}{dt})l\}.$$

In other words, we do not care about the precise form of the latent variables l , only about their existence. The question arises whether we can write \mathcal{B}_l as a standard model, that is, whether we can find a polynomial matrix \hat{R} such that

$$\mathcal{B}_l = \hat{\mathcal{B}} = \{w \in \mathcal{A}^q \mid \hat{R}(\frac{d}{dt})w = 0\}.$$

The answer is yes, and moreover, there is an easy way to obtain the desired \hat{R} from the given R and M .

One solves the linear system of equations $\xi M = 0$, where $\xi \in \mathbb{R}[s]^{1 \times p}$. In other words, we compute the left kernel of M , over the polynomial ring. Using the Smith form, one can show that there exist $p - \text{rank}(M)$ linearly independent solutions that span this kernel. Thus, let $\xi_1, \dots, \xi_{p-\text{rank}(M)}$ be a generating system for the left kernel of M . Collecting these row vectors in a matrix X , we have constructed a matrix X which satisfies the following three conditions:

1. $XM = 0$;
2. any polynomial row vector ξ with $\xi M = 0$ can be written as a polynomial linear combination of the rows of X , that is, $\xi = \eta X$ for some polynomial row vector η ;
3. X has full row rank.

Lemma 2.9 Let X_1, X_2 be two matrices with the three properties from above. Then we must have $X_1 = UX_2$ for some unimodular matrix U .

Proof: Condition 1 implies $X_1M = 0$ and $X_2M = 0$. By Condition 2, each row of X_1 can be written as a polynomial linear combination of the rows of X_2 (and vice versa). This means that there exist polynomial matrices U and V such that

$$X_1 = UX_2 \quad \text{and} \quad X_2 = VX_1.$$

Then

$$(I - UV)X_1 = 0 \quad \text{and} \quad (I - VU)X_2 = 0.$$

Finally, Condition 3 implies that $UV = VU = I$, that is, U has a polynomial inverse, namely V , and thus U is unimodular. \square

Theorem 2.10 Let R, M be given polynomial matrices, with the same number of rows. Let X be as described above, and define $\hat{R} := XR$. Then

$$\exists l \in \mathcal{A}^r : R\left(\frac{d}{dt}\right)w = M\left(\frac{d}{dt}\right)l \Leftrightarrow \hat{R}\left(\frac{d}{dt}\right)w = 0$$

(analogously if $\frac{d}{dt}$ is replaced by σ), that is, using the notation from above, $\hat{\mathcal{B}} = \mathcal{B}_l$.

Proof: By Corollary 2.4, there exists a unimodular matrix V such that

$$VM = \begin{bmatrix} M_1 \\ 0 \end{bmatrix}$$

where M_1 has full row rank. Then

$$\exists l : R\left(\frac{d}{dt}\right)w = M\left(\frac{d}{dt}\right)l \Leftrightarrow$$

$$\exists l : V\left(\frac{d}{dt}\right)R\left(\frac{d}{dt}\right)w = V\left(\frac{d}{dt}\right)M\left(\frac{d}{dt}\right)l \Leftrightarrow$$

$$\exists l : \begin{bmatrix} R_1 \\ R_2 \end{bmatrix} \left(\frac{d}{dt}\right)w = \begin{bmatrix} M_1 \\ 0 \end{bmatrix} \left(\frac{d}{dt}\right)l.$$

Using Corollary 2.7, we obtain

$$\exists l : R\left(\frac{d}{dt}\right)w = M\left(\frac{d}{dt}\right)l \Leftrightarrow R_2\left(\frac{d}{dt}\right)w = 0.$$

It remains to establish a relation between R_2 and \hat{R} . Note that $\hat{X} := [0, I]V$ has the three properties given above. According to Lemma 2.9, there exists a unimodular matrix U such that

$$X = U\hat{X}.$$

Thus $\hat{R} = XR = U\hat{X}R = U[0, I]VR = UR_2$, which shows that

$$\hat{R}\left(\frac{d}{dt}\right)w = 0 \Leftrightarrow R_2\left(\frac{d}{dt}\right)w = 0$$

which completes the proof. \square

2.4 Inputs, outputs, and autonomous systems

Definition 2.11 Let the signal vector w be partitioned as

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix},$$

where $w_i \in \mathcal{A}^{q_i}$, where q_1 and q_2 are positive integers with $q_1 + q_2 = q$. The subvector w_1 is called a vector of **free variables** (or: inputs) of \mathcal{B} if it is unconstrained by the system law, that is,

$$\forall w_1 \in \mathcal{A}^{q_1} \exists w_2 \in \mathcal{A}^{q_2} : \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \in \mathcal{B}.$$

Free variables can be found as follows: We may assume that $R \in \mathbb{R}[s]^{p \times q}$ has full row rank, that is, $\text{rank}(R) = p \leq q$. Then there exists a $p \times p$ non-singular submatrix of R . Assume that $p < q$. We can always permute the columns of R (this just corresponds to renumbering our signal components) such that

$$R = \begin{bmatrix} -Q & P \end{bmatrix}$$

where P is such a non-singular $p \times p$ matrix. Writing

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

correspondingly, our system law $R(\frac{d}{dt})w = 0$ takes the form

$$P(\frac{d}{dt})y = Q(\frac{d}{dt})u. \quad (2.2)$$

From Theorem 2.6, we obtain that this equation has a solution $y \in \mathcal{A}^p$ for every choice of $u \in \mathcal{A}^m$, where $m := q - p$. Therefore, u is a vector of free variables of \mathcal{B} .

Next, we are interested in the maximal number of free variables of a system \mathcal{B} . For this, we define the **input-dimension** of \mathcal{B} by

$$\text{idim}(\mathcal{B}) := \max\{k \in \mathbb{N} \mid \exists \Pi \text{ such that } \mathcal{B} \rightarrow \mathcal{A}^k, w \mapsto [I_k, 0]\Pi w \text{ is surjective}\},$$

where Π is a permutation matrix. Note that the mapping $\mathcal{B} \rightarrow \mathcal{A}^k$ is simply the projection of w onto a subvector consisting of k components. Thus, the input-dimension is the largest integer k such that there exist k components of $w \in \mathcal{B}$ that are free (i.e., unconstrained by the system law). We first observe that surjectivity of $\mathcal{B} \rightarrow \mathcal{A}^k, w \mapsto Sw$, where $S = [I_k, 0]\Pi$, amounts to saying

$$\forall \xi \in \mathcal{A}^k \exists w \in \mathcal{A}^q : \begin{bmatrix} I \\ 0 \end{bmatrix} \xi = \begin{bmatrix} S \\ R \end{bmatrix} (\frac{d}{dt})w. \quad (2.3)$$

As a by-product of introducing the concept of input-dimension, we can now prove the characterization of minimality of representations that was announced earlier.

Theorem 2.12 Any two representations of \mathcal{B} have the same rank. A representation of \mathcal{B} is minimal if and only if it has full row rank.

Proof: Let $R \in \mathbb{R}[s]^{p \times q}$ be a representation of \mathcal{B} . By Lemma 2.5, we may assume that R has full row rank p . If $p < q$, then without loss of generality, let $R = [-Q, P]$, where $P \in \mathbb{R}[s]^{p \times p}$ is non-singular. We have seen above that $u \in \mathcal{A}^m$ is a vector of free variables of \mathcal{B} . Thus $\text{idim}(\mathcal{B}) \geq m = q - \text{rank}(R)$. Clearly, this inequality holds also $p = q$, that is, for $m = 0$.

Conversely, using the elimination of latent variables, we find that

$$\exists w \in \mathcal{A}^q : \begin{bmatrix} I \\ 0 \end{bmatrix} \xi = \begin{bmatrix} S \\ R \end{bmatrix} \left(\frac{d}{dt}\right)w \Leftrightarrow X_1 \left(\frac{d}{dt}\right)\xi = 0,$$

where $X = [X_1, X_2]$ is a polynomial matrix whose rows generate the left kernel of $\begin{bmatrix} S \\ R \end{bmatrix}$. Thus

$$\forall \xi \in \mathcal{A}^k \exists w \in \mathcal{A}^q : \begin{bmatrix} I \\ 0 \end{bmatrix} \xi = \begin{bmatrix} S \\ R \end{bmatrix} \left(\frac{d}{dt}\right)w \Leftrightarrow \forall \xi \in \mathcal{A}^k : X_1 \left(\frac{d}{dt}\right)\xi = 0.$$

However, this is true if and only if $X_1 = 0$. Thus, in view of (2.3), surjectivity of $\mathcal{B} \rightarrow \mathcal{A}^k$, $w \mapsto Sw$ is equivalent to saying that $[X_1, X_2] \begin{bmatrix} S \\ R \end{bmatrix} = 0$ implies that $X_1 = 0$. Since we may assume, without loss of generality, that R has full row rank, this is also equivalent to saying that $\begin{bmatrix} S \\ R \end{bmatrix} \in \mathbb{R}[s]^{(k+p) \times q}$ has full row rank. This implies that $k + p \leq q$. Therefore we must have $\text{idim}(\mathcal{B}) \leq q - \text{rank}(R)$.

Combining this with the inequality from above, we have $\text{idim}(\mathcal{B}) = q - \text{rank}(R)$, showing that the rank of a representation matrix is an invariant of \mathcal{B} , i.e., it does not depend on the specific choice of R . From this, we get that a full row rank representation must be minimal (the converse direction was already shown in Lemma 2.5). \square

Thus, in spite of the seemingly complicated definition, we find that the input-dimension is given by the simple formula $\text{idim}(\mathcal{B}) = q - \text{rank}(R)$, which is just the “number of free variables” of a linear system of equations that one would expect from linear algebra.

Definition 2.13 A system law in the form of equation (2.2), where P is square and non-singular, is called an **input-output representation** of \mathcal{B} . One calls $p = \text{rank}(R)$ the number of outputs (or output-dimension), and $m = q - p$ the number of inputs (or input-dimension). The signal subvector u is called input, and the signal subvector y is called output.

It is important to note that in general, one and the same behavior may have several different input-output representations. This is due to the fact that there may be more than one way to select a non-singular $p \times p$ submatrix of R , where $p = \text{rank}(R)$.

Example 2.14 The electrical circuit from Example 1.9 admits 5 different input-output representations. The matrix

$$\hat{R} = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 \\ R & 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & Cs & 0 \\ Ls & 0 & 0 & 0 & -1 \end{bmatrix}.$$

has rank 4 and every 4×4 submatrix is non-singular (the parameters R, C, L are supposed to be positive). Each of the 5 signals I, U, u_1, u_2, u_3 can take the role of the input, and then the remaining 4 signals will be outputs. However, from the physical point of view, one would usually consider U as the input and I, u_1, u_2, u_3 as outputs.

Example 2.15 The system from Example 1.8 does not have any free variables. If

$$R(\sigma)w = \begin{bmatrix} \sigma & -1 \\ -\frac{1}{2} & \sigma - \frac{1}{2} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = 0$$

then

$$\begin{bmatrix} \sigma - \frac{1}{2} & 1 \\ \frac{1}{2} & \sigma \end{bmatrix} \begin{bmatrix} \sigma & -1 \\ -\frac{1}{2} & \sigma - \frac{1}{2} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} \sigma^2 - \frac{1}{2}\sigma - \frac{1}{2} & 0 \\ 0 & \sigma^2 - \frac{1}{2}\sigma - \frac{1}{2} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = 0$$

and thus for $i = 1, 2$:

$$(\sigma^2 - \frac{1}{2}\sigma - \frac{1}{2})w_i = 0.$$

This implies that if $w_i(0)$ and $w_i(1)$ are given, all $w_i(t)$ for $t \in T$ are uniquely determined. In particular, w_i is not free. In fact, considering the characteristic equation

$$\lambda^2 - \frac{1}{2}\lambda - \frac{1}{2} = 0$$

with its solutions $\lambda_1 = 1$ and $\lambda_2 = -\frac{1}{2}$, we must have

$$w_i(t) = a_1\lambda_1^t + a_2\lambda_2^t = a_1 + a_2(-\frac{1}{2})^t$$

where the constants a_1, a_2 can be computed from $w_i(0), w_i(1)$. In fact, the original equations show that it suffices to know $w_i(0)$ for $i = 1, 2$ to obtain a unique solution.

Systems without inputs are called autonomous (Greek: *αυτονομος*; under its own law/self-governed/independent).

Definition 2.16 A system which has no free variables is called **autonomous**. In other words, \mathcal{B} is autonomous if and only if $\text{idim}(\mathcal{B}) = 0$.

Example 2.17 A system given by the scalar equation $p(\frac{d}{dt})w = 0$, where $0 \neq p \in \mathbb{R}[s]$, is autonomous. Any w satisfying $p(\frac{d}{dt})w = 0$ must be of the form

$$w(t) = \sum_{\lambda} a_{\lambda}(t)e^{\lambda t}$$

where $\lambda \in \mathbb{C}$ are the zeros of p , and a_{λ} are polynomials. Thus w is certainly constrained by the system law. Similarly, $p(\sigma)w = 0$ implies

$$w(t) = \sum_{\lambda} a_{\lambda}(t)\lambda^t$$

in the discrete case.

Lemma 2.18 \mathcal{B} is autonomous if and only if it has a square non-singular representation matrix.

Proof: If \mathcal{B} is autonomous, then $\text{idim}(\mathcal{B}) = 0$, that is, any representation matrix has full column rank. On the other hand, we can always find a representation with full row rank. Combining this, we get that there exists a square representation matrix with full rank, that is, a non-singular matrix. Conversely, let \mathcal{B} be represented by a square non-singular matrix, say, $P \in \mathbb{R}[s]^{q \times q}$ with $\text{rank}(P) = q$. Then $\text{idim}(\mathcal{B}) = q - q = 0$, that is, \mathcal{B} is autonomous. \square

2.5 Input-output representations

Let $P(\frac{d}{dt})y = Q(\frac{d}{dt})u$ be a system law in input-output form and assume that $u \in \mathcal{A}^m$ is given. From the fundamental principle, we know that there exist outputs belonging to this input. What can be said about the set of these outputs? Let y_p be one particular output that belongs to u , that is, $P(\frac{d}{dt})y_p = Q(\frac{d}{dt})u$. Let y be another output that belongs to u , then $P(\frac{d}{dt})(y - y_p) = 0$, that is, $y_h := y - y_p$

is an element of the autonomous system defined by $P(\frac{d}{dt})y_h = 0$. Therefore, any output y that belongs to u can be written in the form

$$y = y_p + y_h$$

where y_p is one particular solution of the inhomogeneous equation $P(\frac{d}{dt})y = Q(\frac{d}{dt})u$ and y_h is an arbitrary solution of the corresponding homogeneous equation.

2.5.1 The continuous homogeneous equation

Consider $P(\frac{d}{dt})y = 0$, where $P \in \mathbb{R}[s]^{p \times p}$ is non-singular. Let $0 \neq \det(P) \in \mathbb{R}[s]$ be its determinant and let

$$\Lambda = \{\lambda \in \mathbb{C} \mid \det(P)(\lambda) = 0\}$$

denote the set of its zeros. Any solution y of $P(\frac{d}{dt})y = 0$ has the form

$$y(t) = \sum_{\lambda \in \Lambda} a_\lambda(t) e^{\lambda t}$$

where $a_\lambda \in \mathbb{C}[t]^p$.

In the **scalar** case ($p = 1$), the degree of each polynomial a_λ is at most $\mu(\lambda) - 1$, where $\mu(\lambda)$ is the multiplicity of λ as a zero of P , that is, according to the fundamental theorem of algebra,

$$P = c \prod_{\lambda \in \Lambda} (s - \lambda)^{\mu(\lambda)}$$

where c is the leading coefficient of P . Therefore, each y satisfying $P(\frac{d}{dt})y = 0$ is uniquely determined by the

$$\sum_{\lambda \in \Lambda} \mu(\lambda) = \deg(P)$$

coefficients of these polynomials. In particular, the dimension of the solution space of the scalar equation $P(\frac{d}{dt})y = 0$ is precisely the degree of P , that is, the order of the differential equation $P(\frac{d}{dt})y = 0$.

In the **general** case ($p \geq 1$), we use the Smith form $UPV = D = \text{diag}(d_1, \dots, d_p)$. Since $P(\frac{d}{dt})y = 0$ is equivalent to $D(\frac{d}{dt})\tilde{y} = 0$, where $y = V(\frac{d}{dt})\tilde{y}$, we consider

$$\mathcal{B} = \{y \in \mathcal{A}^p \mid P(\frac{d}{dt})y = 0\} \quad \text{and} \quad \tilde{\mathcal{B}} = \{\tilde{y} \in \mathcal{A}^p \mid D(\frac{d}{dt})\tilde{y} = 0\}$$

and we conclude from the scalar case that

$$\dim \tilde{\mathcal{B}} = \sum_{i=1}^p \deg(d_i) = \deg \prod_{i=1}^p d_i = \deg \det(D) = \deg \det(P).$$

The map $\tilde{y} \mapsto y = V(\frac{d}{dt})\tilde{y}$ provides a vector space isomorphism between $\tilde{\mathcal{B}}$ and \mathcal{B} and thus, their dimensions coincide.

Theorem 2.19 \mathcal{B} is autonomous if and only if it is finite-dimensional as a real vector space. More precisely, if \mathcal{B} is represented by the square non-singular matrix P , then $\dim \mathcal{B} = \deg \det(P)$.

Remark 2.20 Note that our signals are supposed to be real-valued (recall that although we work with distributional solutions in the continuous case, the solutions of an autonomous system are classical functions). Then \mathcal{B} is always a real vector space. However, we still have to admit complex coefficients in the polynomials a_λ , because $\lambda \in \Lambda$ and hence $e^{\lambda t}$ will be complex, in general. However, since the coefficients of P are supposed to be real, we have

$$\lambda \in \Lambda \quad \Rightarrow \quad \bar{\lambda} \in \Lambda,$$

that is, the zeros of $\det(P)$ come in pairs of complex conjugate numbers. For our signals to have real values, we must have

$$a_{\bar{\lambda}} = \overline{a_\lambda}$$

and thus $\deg \det(P)$ is really the vector space dimension of \mathcal{B} over \mathbb{R} and not \mathbb{C} .

2.5.2 The continuous inhomogeneous equation

Suppose that $P(\frac{d}{dt})h^{(i)} = Q(\frac{d}{dt})\delta^{(i)}$, where $\delta^{(i)} \in \mathcal{A}^m$ has δ in the i -th position and zeros everywhere else. Here, $\delta \in \mathcal{A} = \mathcal{D}'(T)$ denotes the Dirac delta distribution. Define

$$h = [h^{(1)} \quad \dots \quad h^{(m)}] \in \mathcal{A}^{p \times m}.$$

Then h is called **impulse response** (or: fundamental solution), because its columns can be seen as the system's response (output) to an input which is an impulse (delta distribution).

Then a particular output y belonging to the input u is given by the distributional **convolution** (if it exists)

$$y = h * u$$

that is, in the classical case (u, h locally integrable functions, i.e., regular distributions)

$$y(t) = \int_{-\infty}^{\infty} h(t - \tau)u(\tau)d\tau.$$

The existence of the convolution is guaranteed by any of the following conditions:

- Both h and u have their support in $[0, \infty)$. This means that they assign zero to every test function whose support is in $(-\infty, 0)$. If h, u are continuous functions, this means that $h(t) = 0$ and $u(t) = 0$ for all $t < 0$. Then we have

$$y(t) = \int_0^t h(t - \tau)u(\tau)d\tau,$$

which is an integral over a compact interval, which exists due to the assumption of continuity. Similarly, the convolution is always well-defined if $T = \mathbb{R}_+$.

- If h or u has compact support, then $h * u$ exists.
- Let u be a bounded function, and let h be a (globally) integrable function. Then

$$\|y\|_{\infty} \leq \|h\|_1 \cdot \|u\|_{\infty}$$

and thus, a bounded input $u \in L^{\infty}$ leads to a bounded output $y \in L^{\infty}$ provided that $h \in L^1$. This is known as “bounded input, bounded output (BIBO) stability”.

Example: Consider $\dot{x}(t) = Ax(t) + b(t)$. We take b as the input, and x as the output. Then $P = sI - A$ and $Q = I$. An impulse response is given by

$$h(t) = \begin{cases} e^{At} & \text{if } t \geq 0 \\ 0 & \text{if } t < 0 \end{cases}$$

for both $T = \mathbb{R}$ and $T = \mathbb{R}_+$. If $T = \mathbb{R}$, assume that u is a continuous function with $u(t) = 0$ for $t < 0$. Then we have

$$y(t) = (h * u)(t) = \int_0^t e^{A(t-\tau)}u(\tau)d\tau = \int_0^t e^{A\tau}u(t - \tau)d\tau$$

as a particular output $y = x$ belonging to the input $u = b$.

2.5.3 The discrete homogeneous equation ($T = \mathbb{N}$)

Consider $P(\sigma)y = 0$, where $P \in \mathbb{R}[s]^{p \times p}$ is non-singular. Let $0 \neq \det(P) \in \mathbb{R}[s]$ be its determinant and let

$$\Lambda = \{\lambda \in \mathbb{C} \mid \det(P)(\lambda) = 0\}$$

denote the set of its zeros. Any solution y of $P(\sigma)y = 0$ has the form

$$y(t) = \sum_{\lambda \in \Lambda} a_\lambda(t) \lambda^t$$

where $a_\lambda \in \mathbb{C}[t]^p$, for all t that are large enough (recall that this restriction is due to the possible presence of 0 in Λ). Therefore, it does not suffice, in the discrete case, to count the possible choices of the coefficients of the a_λ . Nevertheless, we can argue as follows.

In the **scalar** case ($p = 1$), we have $P(\sigma)y = 0$, where $P = p_d s^d + \dots + p_1 s + p_0$ for some $d \geq 0$, $p_i \in \mathbb{R}$, $p_d \neq 0$. Thus the difference equation reads

$$p_d y(t+d) + \dots + p_1 y(t+1) + p_0 y(t) = 0 \quad \text{for all } t \in \mathbb{N}.$$

The first d values $y(0), \dots, y(d-1)$ of y are unconstrained by the system law, whereas $y(t)$ for $t \geq d$ is uniquely determined by the previous values. Therefore, each y satisfying $P(\sigma)y = 0$ is uniquely determined by choosing $d = \deg(P)$ real parameters (namely, the values $y(0), \dots, y(d-1) \in \mathbb{R}$).

In the **general** case ($p \geq 1$), we use again the Smith form $UPV = D = \text{diag}(d_1, \dots, d_p)$, to obtain the following result which is completely analogous to its continuous counterpart.

Theorem 2.21 \mathcal{B} is autonomous if and only if it is a finite-dimensional vector space. More precisely, if \mathcal{B} is represented by the square non-singular matrix P , then $\dim \mathcal{B} = \deg \det(P)$.

2.5.4 The discrete inhomogeneous equation

Suppose that $P(\sigma)h^{(i)} = Q(\sigma)\delta^{(i)}$, where $\delta^{(i)} \in \mathcal{A}^m$ has δ in the i -th position and zeros everywhere else. Here, $\delta \in \mathcal{A} = \mathbb{R}^T$ denotes the sequence that takes the value 1 at time zero, and the value zero everywhere else. Define

$$h = [h^{(1)} \quad \dots \quad h^{(m)}] \in \mathcal{A}^{p \times m}.$$

Then h is called **impulse response** (or: fundamental solution), because its columns can be seen as the system's response (output) to an input which is an impulse (discrete version of delta function).

Then a particular output y belonging to the input u is given by the discrete **convolution** (if it exists)

$$y = h * u$$

that is,

$$y(t) = \sum_{i=-\infty}^{\infty} h(t-i)u(i).$$

We identify $\mathbb{R}^{\mathbb{N}}$ with the set of sequences $a \in \mathbb{R}^{\mathbb{Z}}$ with $a(t) = 0$ for all $t < 0$. The existence of the convolution is guaranteed by any of the following conditions:

- Let $h(t) = 0$ and $u(t) = 0$ for all $t < 0$. Then

$$y(t) = \sum_{i=0}^t h(t-i)u(i),$$

a finite sum. Thus the convolution always exists for $T = \mathbb{N}$.

- At least one of h , u has compact support.
- u is bounded, and h is summable (“BIBO stability”).

Example: Consider $x(t+1) = Ax(t) + b(t)$. We take b as the input, and x as the output. Then $P = sI - A$ and $Q = I$. An impulse response is given by $h(t) = A^{t-1}$ for $t \geq 1$ and $h(t) = 0$ for all $t \leq 0$. If $T = \mathbb{N}$, or $T = \mathbb{Z}$ and $u(t) = 0$ for $t < 0$, we have

$$y(t) = (h * u)(t) = \sum_{i=0}^t h(t-i)u(i) = \sum_{i=0}^{t-1} h(t-i)u(i) = \sum_{i=0}^{t-1} A^{t-i-1}u(i)$$

as a particular output $y = x$ belonging to the input $u = b$.

2.6 Reduction to first order

A polynomial matrix $R \in \mathbb{R}[s]^{p \times q}$ can be written in the form

$$R = R_d s^d + \dots + R_1 s + R_0$$

where $R_i \in \mathbb{R}^{p \times q}$. We may assume that R_d is not the zero matrix. Our system law takes the form

$$(R_d \frac{d^d}{dt^d} + \dots + R_1 \frac{d}{dt} + R_0)w = 0$$

or

$$(R_d \sigma^d + \dots + R_1 \sigma + R_0)w = 0$$

see (1.3) and (1.4). If we put

$$\xi = \begin{bmatrix} w \\ \frac{d}{dt}w \\ \vdots \\ \frac{d^{d-1}}{dt^{d-1}}w \end{bmatrix} \quad \text{or} \quad \xi = \begin{bmatrix} w \\ \sigma w \\ \vdots \\ \sigma^{d-1}w \end{bmatrix}$$

in the continuous or discrete case, respectively, then we can rewrite the system law as

$$K \frac{d}{dt}\xi = L\xi \quad \text{or} \quad K\sigma\xi = L\xi$$

where

$$K = \begin{bmatrix} I_q & & & \\ & \ddots & & \\ & & I_q & \\ & & & R_d \end{bmatrix} \quad \text{and} \quad L = \begin{bmatrix} 0 & I_q & & \\ \vdots & & \ddots & \\ 0 & & & I_q \\ -R_0 & -R_1 & \cdots & -R_{d-1} \end{bmatrix}.$$

Putting $n = dq$ and $k = (d-1)q + p$, we have $K, L \in \mathbb{R}^{k \times n}$ and

$$R(\frac{d}{dt})w = 0 \quad \Leftrightarrow \quad \exists \xi \in \mathcal{A}^n : \begin{cases} K \frac{d}{dt}\xi = L\xi \\ w = [I_q, 0, \dots, 0]\xi \end{cases} \quad (2.4)$$

and similarly in the discrete case. This shows that reduction to first order is nothing but a special way of introducing latent variables. Another way of doing this is discussed in the next section.

2.7 State

Theorem 2.22 Let $R \in \mathbb{R}[s]^{p \times q}$. There exists an integer $n \in \mathbb{N}$ and real matrices $K \in \mathbb{R}^{n \times n}$, $L \in \mathbb{R}^{n \times q}$, $M \in \mathbb{R}^{p \times n}$, $N \in \mathbb{R}^{p \times q}$ such that the system law $R(\frac{d}{dt})w = 0$ has a first order latent variable representation of the form

$$R(\frac{d}{dt})w = 0 \quad \Leftrightarrow \quad \exists x \in \mathcal{A}^n : \begin{cases} \frac{d}{dt}x = Kx + Lw \\ 0 = Mx + Nw \end{cases} \quad (2.5)$$

and similarly for σ instead of $\frac{d}{dt}$.

Proof: Let $R = R_d s^d + \dots + R_1 s + R_0$ with $R_i \in \mathbb{R}^{p \times q}$. We put $n = dp$ and

$$K = \begin{bmatrix} 0 & \cdots & \cdots & 0 \\ I_p & \ddots & & \vdots \\ & \ddots & \ddots & \vdots \\ & & I_p & 0 \end{bmatrix} \in \mathbb{R}^{n \times n} \quad L = \begin{bmatrix} R_0 \\ \vdots \\ \vdots \\ R_{d-1} \end{bmatrix} \in \mathbb{R}^{n \times q}$$

$$M = \begin{bmatrix} 0 & \cdots & 0 & I_p \end{bmatrix} \in \mathbb{R}^{p \times n} \quad N = R_d \in \mathbb{R}^{p \times q}.$$

Then

$$\begin{bmatrix} \frac{d}{dt} I_n - K \\ -M \end{bmatrix} x = \begin{bmatrix} L \\ N \end{bmatrix} w$$

can be pre-multiplied by $U(\frac{d}{dt})$, where U is the unimodular matrix of size $n + p = (d + 1)p$

$$U = \begin{bmatrix} I_p & sI_p & \cdots & s^d I_p \\ & \ddots & \ddots & \vdots \\ & & \ddots & sI_p \\ & & & I_p \end{bmatrix}$$

to obtain the equivalent equation

$$\begin{bmatrix} 0 \\ -I_n \end{bmatrix} x = \begin{bmatrix} R \\ * \end{bmatrix} \left(\frac{d}{dt}\right) w$$

where the $*$ denotes a polynomial matrix whose precise form is not important here, because we only need that such an x exists if and only if $R(\frac{d}{dt})w = 0$. \square

Remark 2.23 Comparing (2.4) and (2.5), we note that in (2.4), we have to deal with an implicit equation $K\dot{\xi} = L\xi$, and moreover, K, L are not square, in general. Thus we cannot use the Kronecker-Weierstraß form from Theorem 1.18, which works only for square matrices. The non-square Kronecker form is much more complicated, and not treated here. On the other hand, we have an explicit equation $\dot{x} = Kx + Lw$ in (2.5). Together with $0 = Mx + Nw$, this yields a so-called semi-explicit system. Another difference is the size n of the latent variable ξ or x , respectively: ξ has $n_1 = dq$ components, whereas x has $n_2 = dp$ components. Recall that we may assume, without loss of generality, that R has full row rank $p \leq q$, and thus $n_2 \leq n_1$. Finally, the construction of (2.5) can easily be modified by first applying it to each row of R separately, and then combining the results. This may lead to even “smaller” (with respect to the number n of latent variables x_i) representations. In fact, if $d_i \in \mathbb{N}$ is the highest power of s appearing in the i -th row of R (assuming that R contains no zero row), we will get a representation of size $\sum_{i=1}^p d_i$ instead of $pd = p \max_i \{d_i\}$. Summing up, we see that (2.5) is usually preferable to the “naive” reduction to first order from (2.4).

For every input-output representation $P(\frac{d}{dt})y = Q(\frac{d}{dt})u$, we can find, according to Theorem 2.22, a first order representation of the form

$$P(\frac{d}{dt})y = Q(\frac{d}{dt})u \quad \Leftrightarrow \quad \exists x \in \mathcal{A}^n : \begin{cases} \frac{d}{dt}x &= Kx + L_1u + L_2y \\ 0 &= Mx + N_1u + N_2y \end{cases} \quad (2.6)$$

and similarly for σ instead of $\frac{d}{dt}$.

A particularly important case arises when N_2 is non-singular. We will see later on (in Chapter 8) that it is always possible to choose a representation and an input-output decomposition such that this is true. Then we can solve the second equation for y to obtain

$$y = -N_2^{-1}(Mx + N_1u)$$

and plug that into the first equation. We get

$$\frac{d}{dt}x = (K - L_2N_2^{-1}M)x + (L_1 - L_2N_2^{-1}N_1)u.$$

Setting

$$\begin{aligned} A &= K - L_2N_2^{-1}M \\ B &= L_1 - L_2N_2^{-1}N_1 \\ C &= -N_2^{-1}M \\ D &= -N_2^{-1}N_1 \end{aligned}$$

we have

$$P(\frac{d}{dt})y = Q(\frac{d}{dt})u \quad \Leftrightarrow \quad \exists x \in \mathcal{A}^n : \begin{cases} \frac{d}{dt}x &= Ax + Bu \\ y &= Cx + Du. \end{cases}$$

The explicit equations

$$\begin{aligned} \dot{x} &= Ax + Bu & \text{or} & & \sigma x &= Ax + Bu \\ y &= Cx + Du & & & y &= Cx + Du \end{aligned}$$

are called **state space representations**. We call u the input, y the output, and x the **state** of this system law. This notion comes from the following observation in the classical case, where u is a continuous function (for the continuous time sets $T = \mathbb{R}$ or $T = \mathbb{R}_+$): If $x(t_0)$ is known for some $t_0 \in T$, and if u is known on some interval $T_{01} := [t_0, t_1] \cap T$, where $t_1 > t_0$, then x , and thus y , are uniquely determined everywhere in T_{01} . This is due to the solution formulas

$$x(t) = e^{A(t-t_0)}x(t_0) + \int_{t_0}^t e^{A(t-\tau)}Bu(\tau)d\tau$$

or

$$x(t) = A^{t-t_0}x(t_0) + \sum_{i=t_0}^{t-1} A^{t-1-i}Bu(i)$$

respectively. Thus x represents the system's memory in the sense that $x(t_0)$ contains all the information about the "past" needed for determining the "future" (provided that the future input is given). Roughly speaking, the "history" of the system up to time t_0 is stored in $x(t_0)$, which is therefore called the system's "state" at time t_0 .

Let us write $\varphi(t_1, t_0, x_0, u)$ for the state of the system at time t_1 provided that the state at time t_0 was x_0 and that the input function is u . Then the **state transition map** φ has the following properties for all $t_2 > t_1 > t_0 \in T$, $x_0, x_1, x_2 \in \mathbb{R}^n$, $u, u_1, u_2 \in \mathcal{A}^m$:

Consistency: $\varphi(t_0, t_0, x_0, u) = x_0$.

Causality: If $u_1(t) = u_2(t)$ for all $t \in T_{01}$, then $\varphi(t_1, t_0, x_0, u_1) = \varphi(t_1, t_0, x_0, u_2)$.

Semigroup property: If

$$\varphi(t_1, t_0, x_0, u) = x_1 \quad \text{and} \quad \varphi(t_2, t_1, x_1, u) = x_2$$

then

$$\varphi(t_2, t_0, x_0, u) = x_2.$$

Chapter 3

Stability

Stability is concerned with the behavior of signals on the non-negative time set

$$T_+ := [0, \infty) \cap T$$

in particular, in the limit as time tends to infinity. In the continuous case, $T_+ = \mathbb{R}_+$, and in the discrete case, $T_+ = \mathbb{N}$.

Let $P(\frac{d}{dt})y = Q(\frac{d}{dt})u$ or $P(\sigma)y = Q(\sigma)u$ be input-output representations of the system law. Let y_1, y_2 be two outputs belonging to the same input. What can be said about the size of their difference $y_1 - y_2$? Recall that $y_1 - y_2$ satisfies the homogeneous equation $P(\frac{d}{dt})y_h = 0$ or $P(\sigma)y_h = 0$. Therefore, it is a smooth function of time in the continuous case, and $y_h(t) \in \mathbb{R}^p$ is well-defined (in general, this makes no sense for distributions). For the following definition, let $\|\cdot\|$ denote a norm on \mathbb{R}^p , e.g., the Euclidean norm.

Definition 3.1 Let \mathcal{B} be represented by $P(\frac{d}{dt})y = Q(\frac{d}{dt})u$ or $P(\sigma)y = Q(\sigma)u$, respectively. The input-output representation is called **stable** if any two outputs y_1, y_2 belonging to the same input u satisfy

$$\|y_1(t) - y_2(t)\| \leq M \quad \text{for all } t \in T_+$$

for some constant M which is independent of t (but may depend on the specific choice of y_1, y_2). It is called **asymptotically stable** if we have additionally

$$\lim_{t \rightarrow \infty} \|y_1(t) - y_2(t)\| = 0.$$

Stability means: $P(\frac{d}{dt})y = 0$ implies that y is bounded on T_+ . Asymptotic stability means: Moreover, $P(\frac{d}{dt})y = 0$ implies that $\lim_{t \rightarrow \infty} \|y(t)\| = 0$. Since

these notions depend only on P and not on Q , we have the following modified definition.

Definition 3.2 An autonomous system is stable if all its signals are bounded on T_+ ; and asymptotically stable if additionally, all its signals tend to zero as time tends to infinity.

Let $P \in \mathbb{R}[s]^{p \times p}$ be non-singular. In the scalar case ($p = 1$), the solutions to $P(\frac{d}{dt})y = 0$ or $P(\sigma)y = 0$ have the form (for large enough t , in the discrete case)

$$y(t) = \sum_{\lambda} a_{\lambda}(t)e^{\lambda t} \quad \text{or} \quad y(t) = \sum_{\lambda} a_{\lambda}(t)\lambda^t,$$

where $\lambda \in \mathbb{C}$ are the zeros of P , and $a_{\lambda} \in \mathbb{C}[t]$. Moreover, $\deg(a_{\lambda}) \leq \mu(\lambda) - 1$. Thus we can see that stability depends on the location of the zeros λ in the complex plane, and on their multiplicities $\mu(\lambda)$.

Remark 3.3 Let $\lambda \in \mathbb{C}$. Consider the function $f : \mathbb{R}_+ \rightarrow \mathbb{R}$, $t \mapsto |a(t)e^{\lambda t}| = |a(t)|e^{\operatorname{Re}(\lambda)t}$, where a is a non-zero polynomial. It grows without bound if and only if we have either $\operatorname{Re}(\lambda) > 0$, or $\operatorname{Re}(\lambda) = 0$ and $\deg(a) \geq 1$. It is constant if and only if $\operatorname{Re}(\lambda) = 0$ and $\deg(a) = 0$. It tends to zero if and only if $\operatorname{Re}(\lambda) < 0$.

Consider the function $f : \mathbb{N} \rightarrow \mathbb{R}$, $t \mapsto |a(t)\lambda^t| = |a(t)||\lambda|^t$, where a is a non-zero polynomial. It grows without bound if and only if either $|\lambda| > 1$, or $|\lambda| = 1$ and $\deg(a) \geq 1$. It is constant if and only if $|\lambda| = 1$ and $\deg(a) = 0$. It tends to zero if and only if $|\lambda| < 1$.

Thus we can characterize (asymptotic) stability in the scalar case as follows: The system given by the scalar equation $P(\frac{d}{dt})y = 0$ or $P(\sigma)y = 0$, where $0 \neq P \in \mathbb{R}[s]$ and $\Lambda = \{\lambda \in \mathbb{C} \mid P(\lambda) = 0\}$, is

continuous-time asymptotically stable if and only if $\operatorname{Re}(\lambda) < 0$ for all $\lambda \in \Lambda$;

continuous-time stable if and only if $\operatorname{Re}(\lambda) \leq 0$ for all $\lambda \in \Lambda$, and if $\operatorname{Re}(\lambda) = 0$, then λ is simple, that is, $\mu(\lambda) = 1$;

discrete-time asymptotically stable if and only if $|\lambda| < 1$ for all $\lambda \in \Lambda$;

discrete-time stable if and only if $|\lambda| \leq 1$ for all $\lambda \in \Lambda$, and if $|\lambda| = 1$, then λ is simple, that is, $\mu(\lambda) = 1$.

Now we generalize this to the matrix case ($p \geq 1$). For this, we need the notion of semi-simple zeros.

Definition 3.4 Let P be a square and non-singular polynomial matrix. A zero λ of $\det(P)$ is called a **semi-simple** zero of P if the multiplicity of λ as a zero of $\det(P)$ (called the algebraic multiplicity $\mu(\lambda)$ of λ) equals the dimension of the kernel of $P(\lambda)$ (called the geometric multiplicity $\nu(\lambda)$ of λ).

Note that in the scalar case ($p = 1$), the geometric multiplicity is always one, and thus, a zero is semi-simple if and only if it is simple. In the general case, however, we only have that simple zeros are semi-simple (this follows from $1 \leq \nu(\lambda) \leq \mu(\lambda)$ which holds for all zeros due to the Smith form). The converse is not necessarily true, as can be seen from the following example.

Example 3.5 Consider

$$P = \begin{bmatrix} s & 0 \\ 0 & s \end{bmatrix} \quad \text{and} \quad \hat{P} = \begin{bmatrix} s & 1 \\ 0 & s \end{bmatrix}.$$

Then $\det(P) = \det(\hat{P}) = s^2$. The only zero of s^2 is $\lambda = 0$, and its algebraic multiplicity equals 2 in both cases. However,

$$P(0) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \text{and} \quad \hat{P}(0) = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix},$$

and hence, the geometric multiplicity of $\lambda = 0$ as a zero of P is 2, and its geometric multiplicity as a zero of \hat{P} equals 1. Thus $\lambda = 0$ is a semi-simple (but not simple) zero of P . On the other hand, $\lambda = 0$ is not a semi-simple zero of \hat{P} .

Theorem 3.6 Let P be a square and non-singular polynomial matrix. The autonomous system represented by P is

continuous-time asymptotically stable if and only if the zeros of $\det(P)$ have a negative real part;

continuous-time stable if and only if the zeros of $\det(P)$ have a non-positive real part and moreover, each zero λ with $\operatorname{Re}(\lambda) = 0$ is semi-simple;

discrete-time asymptotically stable if and only if the zeros of $\det(P)$ have modulus less than one;

discrete-time stable if and only if the zeros of $\det(P)$ have modulus less than or equal to one, and moreover, each zero λ with $|\lambda| = 1$ is semi-simple.

Proof: We do only the continuous case; the discrete case is completely analogous. Let $UPV = D = \text{diag}(d_1, \dots, d_p)$ be the Smith form of P . Consider

$$\tilde{\mathcal{B}} = \{\tilde{y} \in \mathcal{A}^p \mid D(\frac{d}{dt})\tilde{y} = 0\}.$$

Then P and D have (up to a non-zero constant) the same determinant and hence the same determinantal zeros, with the same algebraic multiplicities. Moreover, since $U(\lambda)P(\lambda)V(\lambda) = D(\lambda)$, where $U(\lambda)$ and $V(\lambda)$ are non-singular complex matrices, also the geometric multiplicities coincide. Therefore it suffices to prove the statement for $\tilde{\mathcal{B}}$, which is (asymptotically) stable if and only if \mathcal{B} is.

Asymptotic stability of $\tilde{\mathcal{B}}$ clearly amounts to the requirement that all zeros of the polynomials d_i have a negative real part.

For stability, let λ be a zero of one of the d_i with $\text{Re}(\lambda) = 0$. We may assume, without loss of generality, that $d_1 \mid \dots \mid d_p$. Then we have

$$D = \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix}$$

with $D_1(\lambda)$ non-singular, and $D_2(\lambda) = 0_k$ (the $k \times k$ zero matrix). Moreover,

$$D_2 = \begin{bmatrix} (s - \lambda)^{l_1} p_1 & & \\ & \ddots & \\ & & (s - \lambda)^{l_k} p_k \end{bmatrix}$$

where $1 \leq l_1 \leq \dots \leq l_k$, and p_i are polynomials with $p_i(\lambda) \neq 0$. Then the algebraic multiplicity of λ equals $l_1 + \dots + l_k$, and the geometric multiplicity of λ equals k . The two multiplicities coincide if and only if $l_i = 1$ for all i , that is, if $s - \lambda$ enters linearly in each diagonal entry of D_2 . This means that λ is a semi-simple zero of D if and only if it is a simple zero of all diagonal entries of D_2 . Equivalently, the polynomial coefficient vector a_λ of $e^{\lambda t}$ in the representation $\tilde{y}(t) = \sum a_{\tilde{\lambda}} e^{\tilde{\lambda} t}$ can only be a constant vector. \square

Example 3.7 The Smith forms of the matrices in Example 3.5 are

$$D = \begin{bmatrix} s & 0 \\ 0 & s \end{bmatrix} \quad \text{and} \quad \hat{D} = \begin{bmatrix} 1 & 0 \\ 0 & s^2 \end{bmatrix}$$

respectively. The system given by $P(\frac{d}{dt})y = 0$ is stable (its solutions are precisely the constant vectors $y(t) = y(0)$), whereas the system given by $\hat{P}(\frac{d}{dt})y = 0$ admits the solution

$$y(t) = \begin{bmatrix} y_1(0) - ty_2(0) \\ y_2(0) \end{bmatrix}$$

which is unbounded whenever $y_2(0) \neq 0$.

3.1 Stability of state space representations

Consider the state space equations

$$\begin{aligned} \dot{x} &= Ax + Bu & \text{or} & & \sigma x &= Ax + Bu \\ y &= Cx + Du & & & y &= Cx + Du. \end{aligned} \quad (3.1)$$

We call them **stable** if two states x_1, x_2 belonging to the same input u satisfy

$$\|x_1(t) - x_2(t)\| \leq M \quad \text{for all } t \in T_+$$

for some constant M . Then

$$\|y_1(t) - y_2(t)\| = \|C(x_1(t) - x_2(t))\| \leq M_1 \quad \text{for all } t \in T_+.$$

We call the state space equations **asymptotically stable** if additionally

$$\lim_{t \rightarrow \infty} \|x_1(t) - x_2(t)\| = 0.$$

Then

$$\lim_{t \rightarrow \infty} \|y_1(t) - y_2(t)\| = 0.$$

Note that here, $\|\cdot\|$ is used to denote both a norm on \mathbb{R}^n and a norm on \mathbb{R}^p . One may think of the respective Euclidean norms, for instance.

Stability means that $\dot{x} = Ax$ or $\sigma x = Ax$ imply that x is bounded on T_+ . Asymptotic stability means that additionally, $\lim_{t \rightarrow \infty} \|x(t)\| = 0$.

Note that (asymptotic) stability of a state-space system implies (asymptotic) stability of the associated input-output system $\{[u^T, y^T]^T \mid \exists x \text{ with (3.1)}\}$. However, the converse is not true in general.

The equations $\dot{x} = Ax$ and $\sigma x = Ax$ are special cases of autonomous systems, namely with

$$P = sI - A.$$

Then the zeros of $\det(P)$ are precisely the eigenvalues of A . An eigenvalue λ of A is called semi-simple if it is a semi-simple zero of $sI - A$.

Corollary 3.8 Let A be a square real matrix. The autonomous system represented by $\dot{x} = Ax$ or $\sigma x = Ax$, respectively, is

continuous-time asymptotically stable if and only if the eigenvalues of A have a negative real part;

continuous-time stable if and only if the eigenvalues of A have a non-positive real part and moreover, each eigenvalue λ with $\operatorname{Re}(\lambda) = 0$ is semi-simple;

discrete-time asymptotically stable if and only if the eigenvalues of A have modulus less than one;

discrete-time stable if and only if the eigenvalues of A have modulus less than or equal to one, and moreover, each eigenvalue λ with $|\lambda| = 1$ is semi-simple.

Note that the solution of $\dot{x} = Ax$ or $\sigma x = Ax$ is

$$x(t) = e^{At}x_0 \quad \text{or} \quad x(t) = A^t x_0$$

respectively, where $x(0) = x_0$. Using these solutions formulas, we can sharpen the notion of stability as follows.

Corollary 3.9 Consider the autonomous state space system $\dot{x} = Ax$ or $\sigma x = Ax$. For $t \in T_+$ and $x_0 \in \mathbb{R}^n$, set $\Phi(t) = e^{At}$ in the continuous case, and $\Phi(t) = A^t$ in the discrete case. The following are equivalent:

1. The system is stable, that is,

$$\forall x_0 \exists M(x_0) \forall t : \quad \|\Phi(t)x_0\| \leq M(x_0).$$

- 2.

$$\exists N \forall x_0, t : \quad \|\Phi(t)x_0\| \leq N\|x_0\|.$$

3. The system is stable in the sense of Lyapunov, that is,

$$\forall \varepsilon > 0 \exists \delta > 0 \forall x_0, t : \quad \|x_0\| \leq \delta \Rightarrow \|\Phi(t)x_0\| \leq \varepsilon.$$

Proof: “1 \Rightarrow 2”: By assumption, there exist M_1, \dots, M_n such that $\|\Phi(t)e_i\| \leq M_i$, where e_i denotes the i -th standard basis vector. For an arbitrary $x_0 = \sum a_i e_i$, this implies

$$\|\Phi(t)x_0\| = \left\| \sum a_i \Phi(t)e_i \right\| \leq \sum |a_i| \|\Phi(t)e_i\| \leq \sum |a_i| M_i \leq \max\{M_i\} \|x_0\|_1,$$

from which the result follows, since all norms on \mathbb{R}^n are equivalent.

“2 \Rightarrow 1”: Set $M(x_0) := N\|x_0\|$.

“2 \Rightarrow 3”: Let $\varepsilon > 0$ be given. Choose $\delta := \frac{\varepsilon}{N}$. If $\|x_0\| \leq \delta$, then

$$\|\Phi(t)x_0\| \leq N\|x_0\| \leq N\delta = \varepsilon.$$

“3 \Rightarrow 2”: There exists $\delta_1 > 0$ such that

$$\|x_0\| \leq \delta_1 \quad \Rightarrow \quad \|\Phi(t)x_0\| \leq 1.$$

Then we have for all $x_0 \neq 0$

$$\|\Phi(t)x_0\| = \frac{\|x_0\|}{\delta_1} \|\Phi(t)\frac{\delta_1 x_0}{\|x_0\|}\| \leq \frac{1}{\delta_1} \|x_0\|,$$

and thus we set $N = \frac{1}{\delta_1}$. For $x_0 = 0$, the statement is anyhow obvious. \square

The crucial point is that the constant N in Condition 2 is not only independent of t , but also independent of the choice of the solution x of $\dot{x} = Ax$ or $\sigma x = Ax$, respectively. In other words, Condition 2 amounts to the uniform boundedness of the matrix norms $\|e^{At}\|$ or $\|A^t\|$ on T_+ . Note that Conditions 1 and 3 are both trivially implied by Condition 2, whereas for the respective converse directions, we have heavily relied on the crucial fact that the mapping from $x(0) = x_0$ to $x(t) = \Phi(t)x_0$ is linear.

Lyapunov stability can also be defined as follows: For all $\varepsilon > 0$, there exists $\delta > 0$ such that $\|x_1(0) - x_2(0)\| \leq \delta$ implies

$$\|x_1(t) - x_2(t)\| \leq \varepsilon \quad \text{for all } t \in T_+,$$

where $x_i(t)$ for $i = 1, 2$ is the solution at time t when starting in $x_i(0)$. (For the system class under consideration, this definition is equivalent to the one from above due to linearity.) Roughly speaking: If the initial values are close to each other, then the solutions will remain close to each other for all $t > 0$.

3.2 Test for asymptotic stability

Let $A \in \mathbb{R}^{n \times n}$ be given. We wish to test whether A is (discrete- or continuous-time) asymptotically stable. The naive way to do this is to compute all the eigenvalues λ of A and to check whether they satisfy $|\lambda| < 1$ or $\operatorname{Re}(\lambda) < 0$, respectively. However, when n is large, this is quite a difficult and computationally expensive task. Moreover, it seems to be a waste of time and effort to determine the precise location of the eigenvalues in the complex plane, when all we want to know is whether they are contained in a specific region (the open unit disc, or the open left half plane, respectively).

One way out is to compute the characteristic polynomial

$$\chi_A(s) := \det(sI - A) = s^n + a_{n-1}s^{n-1} + \dots + a_1s + a_0.$$

Then there are certain criteria, in terms of the coefficients a_i of χ_A , which make it possible to determine whether A is asymptotically stable or not. The best known of these criteria is the **Routh-Hurwitz test** for continuous time, and the **Schur-Cohn test** for discrete time.

For instance, if $n = 2$, we have $\chi_A(s) = s^2 + a_1s + a_0$, where $a_1 = -\text{trace}(A)$ and $a_0 = \det(A)$. The matrix $A \in \mathbb{R}^{2 \times 2}$ is continuous-time asymptotically stable if and only if $a_1 > 0$ and $a_0 > 0$. In the general case ($n \geq 3$), however, the Routh-Hurwitz criterion is much more complex.

Here, we will adopt another approach, based on **Lyapunov equations**. Recall that a symmetric matrix $P \in \mathbb{R}^{n \times n}$ is called **positive semi-definite** (written $P \geq 0$) if we have $x^T Px \geq 0$ for all $x \in \mathbb{R}^n$. Note that then also for $x \in \mathbb{C}^n$, the number $x^* Px$ is real and non-negative. Here, x^* denotes the Hermitian transpose of x . This is due to the fact that if $x = a + ib$ for some $a, b \in \mathbb{R}^n$, then

$$x^* Px = (a + ib)^* P(a + ib) = (a^T - ib^T)P(a + ib) = a^T Pa + b^T Pb,$$

because $a^T Pb = (a^T Pb)^T = b^T P^T a = b^T Pa$.

Theorem 3.10 The following are equivalent:

1. The matrix A is continuous-time asymptotically stable.
2. For every $Q \geq 0$, there exists $P \geq 0$ such that

$$A^T P + PA + Q = 0. \tag{3.2}$$

3. There exists $P \geq 0$ such that

$$A^T P + PA + I = 0.$$

Equation (3.2) is called **continuous-time Lyapunov equation**.

Proof: “1 \Rightarrow 2”: If A is asymptotically stable, it makes sense to define

$$P := \int_0^\infty e^{A^T t} Q e^{At} dt.$$

This is clearly positive semi-definite if Q is, and it satisfies the Lyapunov equation, because

$$A^T P + PA = \int_0^\infty \frac{d}{dt}(e^{A^T t} Q e^{At}) dt = e^{A^T t} Q e^{At} \Big|_0^\infty = -Q.$$

“2 \Rightarrow 3” is a special case ($Q = I$).

“3 \Rightarrow 1”: Let λ be an eigenvalue of A and let $0 \neq x \in \mathbb{C}^n$ be an associated eigenvector, that is, $Ax = \lambda x$. Pre-multiply $A^T P + PA + I = 0$ by x^* and post-multiply by x to get

$$x^* A^T P x + x^* P A x + x^* x = (\bar{\lambda} + \lambda) x^* P x + x^* x = 0.$$

Thus

$$2\operatorname{Re}(\lambda) x^* P x = -x^* x.$$

The right hand side of this equation is negative, and $x^* P x \geq 0$. This can only be true if $\operatorname{Re}(\lambda) < 0$. \square

For the sake of completeness, we give also the discrete version of this theorem.

Theorem 3.11 The following are equivalent:

1. The matrix A is discrete-time asymptotically stable.
2. For every $Q \geq 0$, there exists $P \geq 0$ such that

$$A^T P A - P + Q = 0. \tag{3.3}$$

3. There exists $P \geq 0$ such that

$$A^T P A - P + I = 0.$$

Equation (3.3) is called **discrete-time Lyapunov equation**.

Proof: “1 \Rightarrow 2”: If A is asymptotically stable, it makes sense to define

$$P := \sum_{i=0}^{\infty} (A^T)^i Q A^i.$$

This is clearly positive semi-definite if Q is, and it satisfies the Lyapunov equation, because

$$A^T P A - P = \sum_{i=0}^{\infty} (A^T)^{i+1} Q A^{i+1} - \sum_{i=0}^{\infty} (A^T)^i Q A^i = -Q.$$

“2 \Rightarrow 3” is a special case ($Q = I$).

“3 \Rightarrow 1”: Let λ be an eigenvalue of A and let $0 \neq x \in \mathbb{C}^n$ be an associated eigenvector, that is, $Ax = \lambda x$. Pre-multiply $A^T P A - P + I = 0$ by x^* and post-multiply by x to get

$$x^* A^T P A x - x^* P x + x^* x = (\bar{\lambda} \lambda - 1) x^* P x + x^* x = 0.$$

Thus

$$(|\lambda|^2 - 1) x^* P x = -x^* x.$$

The right hand side of this equation is negative, and $x^* P x \geq 0$. This can only be true if $|\lambda|^2 - 1 < 0$, that is, $|\lambda| < 1$. \square

Remark 3.12 If A is (continuous- or discrete-time) asymptotically stable, then the solutions of (3.2) or (3.3), respectively, are uniquely determined. In the continuous case, this can be seen as follows (the discrete case is analogous): Let P_1, P_2 be two solutions of (3.2), and let $P := P_1 - P_2$. Then $A^T P + P A = 0$. Now consider the matrix-valued function $f(t) := e^{A^T t} P e^{A t}$. Since

$$\frac{d}{dt} f(t) = e^{A^T t} (A^T P + P A) e^{A t} = 0$$

for all t , the function f is actually constant, that is, $f(t) = f(0) = P$ for all t . Now letting t tend to infinity, we see that P must be zero, because $\lim_{t \rightarrow \infty} f(t) = 0$.

The solution formulas $P = \int_0^{\infty} e^{A^T t} Q e^{A t} dt$ and $P = \sum_{i=0}^{\infty} (A^T)^i Q A^i$ are only of theoretical interest. In practice, one solves the linear matrix equations

$$A^T P + P A + Q = 0 \quad \text{or} \quad A^T P A - P + Q = 0$$

(these are n^2 linear equations for n^2 unknowns; exploiting symmetry, this reduces to $\frac{1}{2}n(n+1)$ unknowns), and then tests whether $P \geq 0$. The solution formulas show that we have a positive definite solution P provided that Q is positive definite (to test a symmetric matrix for positive definiteness is easier than to test it for positive semi-definiteness).

Chapter 4

Reachability and controllability

4.1 Basic notions for state space systems

Consider the state space equations

$$\begin{aligned} \dot{x} &= Ax + Bu & \text{or} & & \sigma x &= Ax + Bu \\ y &= Cx + Du & & & y &= Cx + Du, \end{aligned}$$

where $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{p \times n}$ and $D \in \mathbb{R}^{p \times m}$.

Let \mathcal{U} denote the space of **admissible input functions** (until now, we had $\mathcal{U} = \mathcal{D}'(T)^m$ in the continuous and $\mathcal{U} = (\mathbb{R}^T)^m$ in the discrete case, where $T = \mathbb{R}$, \mathbb{R}_+ and $T = \mathbb{N}, \mathbb{Z}$ are our usual time sets). In this section, we restrict to piecewise continuous input functions u in the continuous case, then the state function x is piecewise \mathcal{C}^1 (in particular, it is a classical function). One calls $X = \mathbb{R}^n$ the **state space** of these systems. Its elements are called **states**. The state transition map

$$\varphi : \{(t, t_0) \in T^2 \mid t \geq t_0\} \times X \times \mathcal{U} \rightarrow X, \quad (t, t_0, x_0, u) \mapsto \varphi(t, t_0, x_0, u)$$

yields the state at time t if the state at time t_0 was x_0 , and the input function was u . More concretely, we have

$$\varphi(t, t_0, x_0, u) = e^{A(t-t_0)}x_0 + \int_{t_0}^t e^{A(t-\tau)}Bu(\tau)d\tau \quad (4.1)$$

and

$$\varphi(t, t_0, x_0, u) = A^{t-t_0}x_0 + \sum_{i=t_0}^{t-1} A^{t-i-1}Bu(i). \quad (4.2)$$

The state transition maps are **consistent**, i.e.,

$$\varphi(t, t, x, u) = x \quad (4.3)$$

for all $t \in T$, $x \in X$, $u \in \mathcal{U}$; and they are **strictly causal**, that is, if $u_1(t) = u_2(t)$ for all $t_0 \leq t < t_1$ (note the strict inequality!), then

$$\varphi(t_1, t_0, x, u_1) = \varphi(t_1, t_0, x, u_2) \quad (4.4)$$

for all $x \in X$. Moreover, we have the **semigroup property**:

$$\varphi(t_2, t_1, \varphi(t_1, t_0, x, u), u) = \varphi(t_2, t_0, x, u) \quad (4.5)$$

for all $t_0 \leq t_1 \leq t_2 \in T$, $x \in X$, $u \in \mathcal{U}$. As a consequence, we get the following **concatenation property**: Let $t_1 \in T$, and let $u_1, u_2 \in \mathcal{U}$ be two admissible inputs. Define their concatenation at time t_1 by

$$u(t) = \begin{cases} u_1(t) & \text{if } t < t_1 \\ u_2(t) & \text{if } t_1 \leq t \end{cases}$$

which is again an admissible input. Then

$$\varphi(t_2, t_0, x, u) = \varphi(t_2, t_1, \varphi(t_1, t_0, x, u_1), u_2)$$

for all $x \in X$, and all $t_0 \leq t_1 \leq t_2 \in T$.

Due to **linearity**, we have

$$\varphi(t, t_0, \lambda_1 x_1 + \lambda_2 x_2, \lambda_1 u_1 + \lambda_2 u_2) = \lambda_1 \varphi(t, t_0, x_1, u_1) + \lambda_2 \varphi(t, t_0, x_2, u_2) \quad (4.6)$$

for all $t \geq t_0 \in T$, $\lambda_1, \lambda_2 \in \mathbb{R}$, $x_1, x_2 \in X$, $u_1, u_2 \in \mathcal{U}$. In particular,

$$\varphi(t, t_0, 0, 0) = 0$$

for all $t \geq t_0 \in T$. One says that the zero state is an **equilibrium** of the system when the zero input function is applied.

Due to **time-invariance**, we have

$$\varphi(t, t_0, x_0, u) = \varphi(t - \tau, t_0 - \tau, x_0, \sigma_\tau u) \quad (4.7)$$

for all $x_0 \in X$, $u \in \mathcal{U}$, and all $t \geq t_0 \in T$ and all $\tau \in T$ with the property that $t - \tau, t_0 - \tau \in T$. Here, σ_τ is the shift operator defined by $\sigma_\tau u(s) = u(s + \tau)$ for all s . In particular,

$$\varphi(t, t_0, x_0, u) = \varphi(t - t_0, 0, x_0, \sigma_{t_0} u).$$

Reachability and controllability are concerned with the following questions: Is it possible to steer the system, by a suitable choice of the input function, from one particular state x_0 (which is determined by the past of the system and may be thought of as “unwanted”) to another particular state x_1 (which is prescribed by us and thus “desired”)? How long does the transition from x_0 to x_1 take? Can we find a concrete formula for an input function that forces the system to go from x_0 to x_1 ? To formulate and answer these questions, some new concepts need to be introduced.

Definition 4.1 Let $t_0 \in T$ be fixed. One says that the state $x_1 \in X$

can be reached from $x_0 \in X$ in time $\tau \in T$ ($\tau \geq 0$) if there exists $u \in \mathcal{U}$ such that

$$\varphi(t_0 + \tau, t_0, x_0, u) = x_1.$$

Equivalently, we say that **x_0 can be controlled to x_1 in time τ** .

can be reached from $x_0 \in X$ if this holds for at least one $\tau \geq 0$. Equivalently, **x_0 can be controlled to x_1** .

We say that the system is

completely reachable from $x_0 \in X$ if any $x_1 \in X$ is reachable from x_0 .

completely controllable to $x_1 \in X$ if any $x_0 \in X$ can be controlled to x_1 .

completely reachable (controllable) if x_1 can be reached from x_0 (or: x_0 can be controlled to x_1) for all $x_0, x_1 \in X$.

The choice of the starting time t_0 plays no role since we are dealing with time-invariant systems: We have

$$\varphi(t_0 + \tau, t_0, x_0, u) = \varphi(\tau, 0, x_0, \sigma_{t_0} u).$$

Therefore, if there exists a $t_0 \in T$ such that x_1 is reachable from x_0 (when starting at time t_0) then this is true for any other starting time, e.g., for $t_0 = 0$. Thus we can often choose $t_0 = 0$ for simplicity.

Let $\tau \in T$, $\tau \geq 0$, and $x_0, x_1 \in X$. Define the following sets:

$$\begin{aligned} \mathcal{R}(\tau, x_0) &:= \{x \in X \mid x \text{ is reachable from } x_0 \text{ in time } \tau\} \\ \mathcal{C}(\tau, x_1) &:= \{x \in X \mid x \text{ is controllable to } x_1 \text{ in time } \tau\}. \end{aligned}$$

Moreover, we let $\mathcal{R}(\tau) := \mathcal{R}(\tau, 0)$ denote the set of states that are reachable from $x_0 = 0$ in time τ , and $\mathcal{C}(\tau) := \mathcal{C}(\tau, 0)$ is the set of states that are controllable to $x_1 = 0$ in time τ . Finally,

$$\mathcal{R} := \bigcup_{\tau \geq 0} \mathcal{R}(\tau) \quad \text{and} \quad \mathcal{C} := \bigcup_{\tau \geq 0} \mathcal{C}(\tau)$$

are the set of states that are reachable from zero, and the set of states that are controllable to zero, respectively. The system is completely reachable from zero if and only if $\mathcal{R} = X$, and it is completely controllable to zero if and only if $\mathcal{C} = X$.

Theorem 4.2 Let $s, t \in T$, $0 \leq s \leq t$. We have

1. $\mathcal{R}(s) \subseteq \mathcal{R}(t)$ and $\mathcal{C}(s) \subseteq \mathcal{C}(t)$;
2. $\mathcal{R}(t), \mathcal{C}(t), \mathcal{R}, \mathcal{C}$ are subspaces of $X = \mathbb{R}^n$;
3. There exists $\tau^* \in T$, $\tau^* \geq 0$ such that

$$\mathcal{R} = \mathcal{R}(\tau) \quad \text{and} \quad \mathcal{C} = \mathcal{C}(\tau) \quad \text{for all } \tau \geq \tau^*.$$

Proof:

1. Let $x \in \mathcal{R}(s)$. Then there exists a input function u such that

$$\varphi(t, t-s, 0, u) = x.$$

On the other hand,

$$\varphi(t-s, 0, 0, 0) = 0.$$

Let \tilde{u} be the input function defined by

$$\tilde{u}(\tau) = \begin{cases} 0 & \text{if } \tau < t-s \\ u(\tau) & \text{if } t-s \leq \tau. \end{cases}$$

We have

$$\varphi(t, 0, 0, \tilde{u}) = \varphi(t, t-s, \varphi(t-s, 0, 0, 0), u) = \varphi(t, t-s, 0, u) = x$$

which shows that $x \in \mathcal{R}(t)$. The statement for \mathcal{C} is analogous.

2. We show that $\mathcal{R}(t)$ is a vector space, the statement for \mathcal{R} , $\mathcal{C}(t)$, \mathcal{C} is analogous. We have $0 \in \mathcal{R}(t)$, because

$$\varphi(t, 0, 0, 0) = 0.$$

Let $x_1, x_2 \in \mathcal{R}(t)$, and $\lambda_1, \lambda_2 \in \mathbb{R}$. We need to show that $\lambda_1 x_1 + \lambda_2 x_2 \in \mathcal{R}(t)$. By assumption, there exist input functions u_i such that

$$\varphi(t, 0, 0, u_i) = x_i$$

for $i = 1, 2$. Then

$$\varphi(t, 0, 0, \lambda_1 u_1 + \lambda_2 u_2) = \lambda_1 x_1 + \lambda_2 x_2$$

and hence $\lambda_1 x_1 + \lambda_2 x_2 \in \mathcal{R}(t)$.

3. Consider a strictly increasing sequence

$$0 = \tau_0 < \tau_1 < \tau_2 < \dots$$

in T with $\lim_{i \rightarrow \infty} \tau_i = \infty$. By part 1,

$$\mathcal{R}(\tau_0) \subseteq \mathcal{R}(\tau_1) \subseteq \mathcal{R}(\tau_2) \subseteq \dots$$

By part 2, this is a sequence of subspaces of $X = \mathbb{R}^n$, with

$$\dim \mathcal{R}(\tau_0) \leq \dim \mathcal{R}(\tau_1) \leq \dim \mathcal{R}(\tau_2) \leq \dots \leq n.$$

This is an increasing sequence of integers less than or equal to n . Such a sequence must become stationary, that is, there exists i_0 such that

$$\dim \mathcal{R}(\tau_i) = \dim \mathcal{R}(\tau_{i_0}) \quad \text{for all } i \geq i_0.$$

We use the following fact from linear algebra: If a vector space is contained in another vector space of the same finite dimension, then the two vector spaces must be the same. Thus

$$\mathcal{R}(\tau_i) = \mathcal{R}(\tau_{i_0}) \quad \text{for all } i \geq i_0.$$

For any $\tau \in T$, $\tau \geq \tau_{i_0}$, there exists $j \geq i_0$ such that $\tau \leq \tau_j$. Then $\mathcal{R}(\tau_{i_0}) \subseteq \mathcal{R}(\tau) \subseteq \mathcal{R}(\tau_j) = \mathcal{R}(\tau_{i_0})$. We conclude that

$$\mathcal{R}(\tau) = \mathcal{R}(\tau_{i_0}) \quad \text{for all } \tau \geq \tau_{i_0}$$

and thus

$$\mathcal{R} = \bigcup_{\tau \geq 0} \mathcal{R}(\tau) = \bigcup_{\tau \geq \tau_{i_0}} \mathcal{R}(\tau) = \mathcal{R}(\tau_{i_0}).$$

Put $\tau^* := \tau_{i_0}$, then for $\tau \geq \tau^*$, we have

$$\mathcal{R} \supseteq \mathcal{R}(\tau) \supseteq \mathcal{R}(\tau^*) = \mathcal{R}$$

and thus $\mathcal{R}(\tau) = \mathcal{R}$ for all $\tau \geq \tau^*$. □

Corollary 4.3 In discrete time,

$$\mathcal{R}(n) = \mathcal{R} \quad \text{and} \quad \mathcal{C}(n) = \mathcal{C},$$

where n is the dimension of the state space. In continuous time,

$$\mathcal{R}(\varepsilon) = \mathcal{R} \quad \text{and} \quad \mathcal{C}(\varepsilon) = \mathcal{C}$$

for every $\varepsilon > 0$.

Remark 4.4 This is probably the first time we encounter a significant difference between the continuous and discrete cases. In a discrete system, if x can be reached from x_0 at all, then it can also be reached in time n , where n is the dimension of the state space. In a continuous system, if x can be reached from x_0 at all, then it can also be reached in an arbitrarily small time ε . This is counter-intuitive at first sight: In a “real world” system, it certainly takes “some time” to change from one state to another. The reason is that we admit arbitrarily large input values here, i.e., we make the optimistic assumption that we can put as much “energy” as we like into the system. In a real world system, there are constraints which limit the size of the admissible inputs, and this has the consequence that the transition from one state to another cannot be done arbitrarily fast in practice.

Proof: For discrete time, we use that we know from the previous proof that

$$\mathcal{R}(0) \subseteq \mathcal{R}(1) \subseteq \mathcal{R}(2) \subseteq \dots \tag{4.8}$$

becomes stationary, that is, there exists i_0 such that $\mathcal{R}(i) = \mathcal{R}(i_0)$ for all $i \geq i_0$ and then $\mathcal{R}(i_0) = \mathcal{R}$. We have to show that this happens for some $i_0 \leq n$. Then we are finished, because $\mathcal{R} = \mathcal{R}(i_0) \subseteq \mathcal{R}(n) \subseteq \mathcal{R}$ yields the desired result. Considering the dimensions $d_i := \dim \mathcal{R}(i) \leq n$, we have

$$0 = d_0 \leq d_1 \leq d_2 \leq \dots \leq d_n \leq d_{n+1}.$$

These inequalities cannot all be strict, i.e., we must have $d_i = d_{i+1}$ and hence

$$\mathcal{R}(i) = \mathcal{R}(i+1) \tag{4.9}$$

for some $i \leq n$. The claim is that then we may put $i_0 = i$, that is, the first equality in (4.8) will already yield stationarity. Thus we have to show that (4.9) implies

$$\mathcal{R}(i) = \mathcal{R}(i+k) \quad \text{for all } k \geq 0.$$

We do this by induction on k . The statement is trivial for $k = 0$. Let's assume that the statement is true for k . We need to show it for $k + 1$. The inclusion

$$\mathcal{R}(i) \subseteq \mathcal{R}(i + k + 1)$$

is clear. For the converse, let $x \in \mathcal{R}(i + k + 1)$. This means that there exists an input function u such that

$$\varphi(i + k + 1, 0, 0, u) = x.$$

Set

$$x_1 := \varphi(i + k, 0, 0, u).$$

Then

$$x = \varphi(i + k + 1, i + k, x_1, u) = \varphi(i + 1, i, x_1, \sigma^k u)$$

and $x_1 \in \mathcal{R}(i + k)$, which equals $\mathcal{R}(i)$ by the inductive assumption. Thus there exists an input function u_1 with

$$\varphi(i, 0, 0, u_1) = x_1.$$

Let u_2 be defined by

$$u_2(\tau) = \begin{cases} u_1(\tau) & \text{if } \tau < i \\ (\sigma^k u)(\tau) & \text{if } i \leq \tau. \end{cases}$$

Then

$$\varphi(i + 1, 0, 0, u_2) = \varphi(i + 1, i, \varphi(i, 0, 0, u_1), \sigma^k u) = \varphi(i + 1, i, x_1, \sigma^k u) = x$$

which shows that $x \in \mathcal{R}(i + 1)$. Finally, (4.9) implies that $x \in \mathcal{R}(i)$ as desired. For continuous time, let $\varepsilon > 0$ be given. Consider

$$\mathcal{R}(0) \subseteq \mathcal{R}\left(\frac{\varepsilon}{n}\right) \subseteq \mathcal{R}\left(\frac{2\varepsilon}{n}\right) \subseteq \dots$$

and apply the same argument as for discrete time. □

Corollary 4.5 The following are equivalent: The system is

1. completely reachable (controllable);
2. completely reachable from zero, that is, $\mathcal{R} = X$.

Proof: It is clear that statement 1 implies 2. For the converse direction, let $x_0, x_1 \in X$ be given. We wish to show that x_1 can be reached from x_0 .

In continuous time, pick any $\tau^* = \varepsilon > 0$. In discrete time, choose $\tau^* = n$, where n is the dimension of X . Define $x := x_1 - \varphi(\tau^*, 0, x_0, 0)$. By assumption, $x \in \mathcal{R} = \mathcal{R}(\tau^*)$, that is, there exists an input function $u \in \mathcal{U}$ with

$$x = \varphi(\tau^*, 0, 0, u).$$

This can be rewritten as

$$x_1 = \varphi(\tau^*, 0, x_0, 0) + \varphi(\tau^*, 0, 0, u) = \varphi(\tau^*, 0, x_0, u)$$

showing that x_1 can be reached from x_0 . \square

Until now, we have not used the specific form of the state transition maps from (4.1) and (4.2), but only their properties (4.3)–(4.7). This will change now.

Corollary 4.6 Consider $\dot{x} = Ax + Bu$ or $\sigma x = Ax + Bu$. In the discrete case, assume that A is invertible. Then the following are equivalent: The system is

1. completely reachable (controllable);
2. completely reachable from zero, that is, $\mathcal{R} = X$;
3. completely controllable to zero, that is, $\mathcal{C} = X$.

Proof: We only need to prove “3 \Rightarrow 1”. Let $x_0, x_1 \in X$ be given. We wish to show that x_1 can be reached from x_0 .

In continuous time, pick $\varepsilon > 0$ and define $x := x_0 - e^{-A\varepsilon}x_1$. By assumption, $x \in \mathcal{C} = \mathcal{C}(\varepsilon)$, that is, there exists an input function $u \in \mathcal{U}$ with

$$0 = \varphi(\varepsilon, 0, x, u) = e^{A\varepsilon}x + \int_0^\varepsilon e^{A(\varepsilon-\tau)}Bu(\tau)d\tau.$$

Plugging in for x , this can be rewritten as

$$x_1 = e^{A\varepsilon}x_0 + \int_0^\varepsilon e^{A(\varepsilon-\tau)}Bu(\tau)d\tau = \varphi(\varepsilon, 0, x_0, u)$$

showing that x_1 can be reached from x_0 .

If A is invertible, an analogous argument can be applied in the discrete case. \square

Remark 4.7 Here we have another difference between continuous and discrete systems. The reason is that e^{At} is always an invertible matrix, whereas its discrete counterpart A^t is invertible if and only if A is. Therefore, we have to make this additional assumption in the discrete case. Without it, complete controllability to zero may be strictly weaker than complete controllability. Take for instance $\sigma x = Ax + Bu$ with

$$A = \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

Let $x_0 = x(0) \in \mathbb{R}^2$ be any given initial state. Put $u(0) = -x_1(0)$, then

$$x_1(1) = x_1(0) + u(0) = 0 \quad \text{and} \quad x_2(1) = x_1(0) + u(0) = 0,$$

that is, $x(1) = 0$. This shows that any x_0 can be controlled to zero (in time 1). Hence the system is completely controllable to zero, or $\mathcal{C} = \mathcal{C}(1) = \mathbb{R}^2$. However, it is not completely controllable. If we start in $x_0 = 0$, then every state $x(t)$ will satisfy $x_1(t) = x_2(t)$. This shows that $\mathcal{R} \neq \mathbb{R}^2$, i.e., the system is not completely reachable from zero, and thus it is not completely reachable (controllable).

The next goal is to determine conditions for (complete) reachability/controlability in terms of the matrices A, B . For this, we define the **controllability Gramians** (named after the Danish mathematician J. P. Gram, 1850–1916)

$$W(t) = \int_0^t e^{A\tau} B B^T e^{A^T \tau} d\tau \quad \text{or} \quad W(t) = \sum_{i=0}^{t-1} A^i B B^T (A^T)^i \in \mathbb{R}^{n \times n} \quad (4.10)$$

and the **Kalman reachability/controlability matrix** (named after the Hungarian control scientist R. E. Kalman, 1930–)

$$K = [B \quad AB \quad A^2B \quad \dots \quad A^{n-1}B] \in \mathbb{R}^{n \times nm}. \quad (4.11)$$

Remark 4.8 The definitions (4.10) and (4.11) are motivated by the following observations: Consider $\sigma x = Ax + Bu$. We know that $\mathcal{R} = \mathcal{R}(n)$. Moreover, $x \in \mathcal{R}$ if and only if there exists $u \in \mathcal{U}$, that is, a sequence of input vectors $u(0), \dots, u(n-1) \in \mathbb{R}^m$ such that

$$\begin{aligned} x &= \sum_{i=0}^{n-1} A^{n-1-i} B u(i) = A^{n-1} B u(0) + \dots + A B u(n-2) + B u(n-1) \\ &= [B \quad AB \quad \dots \quad A^{n-1} B] \begin{bmatrix} u(n-1) \\ u(n-2) \\ \vdots \\ u(0) \end{bmatrix} =: K v. \end{aligned}$$

Therefore, x is reachable from zero if and only if the equation $x = Kv$ possesses a solution $v \in \mathbb{R}^{nm}$. This is the case if and only if $x \in \text{im}(K)$. So we have $\mathcal{R} = \text{im}(K)$. Moreover, note that $W(n) = KK^T$. Since $\text{im}(K) = \text{im}(KK^T)$ holds for any real matrix K , we also have $\mathcal{R} = \text{im}(W(n))$.

If $x \in \mathcal{R}$, we therefore have $x = W(n)z = KK^Tz$ for some $z \in \mathbb{R}^n$. Then $v^* := K^Tz$ is a special solution of $Kv = x$. This corresponds to

$$\begin{bmatrix} u(n-1) \\ u(n-2) \\ \vdots \\ u(0) \end{bmatrix} = \begin{bmatrix} B^T \\ B^T A^T \\ \vdots \\ B^T (A^T)^{n-1} \end{bmatrix} z.$$

In other words, $u(t) = B^T (A^T)^{n-t-1} z$ is a special input function that steers the system from state 0 to state x in time n .

Interestingly, the analogous statement is also valid for continuous systems, although the proof is a bit more involved. This is the content of the main theorem of this section, which is stated next.

Theorem 4.9 Consider $\dot{x} = Ax + Bu$ or $\sigma x = Ax + Bu$. In continuous time, let $\tau^* = \varepsilon > 0$ be arbitrary. In discrete time, let $\tau^* = n$, where n is the dimension of the state space. We have

$$\mathcal{R} = \mathcal{R}(\tau^*) = \text{im}(W(\tau^*)) = \text{im}(K).$$

Therefore, the following are equivalent:

1. $\dot{x} = Ax + Bu$ or $\sigma x = Ax + Bu$ is (completely) reachable/controllable;
2. $W(\tau^*)$ is non-singular;
3. K has full row rank.

Moreover in that case, an input function which steers the system from state 0 to state x in time τ^* is given by

$$u(t) = B^T e^{A^T(\tau^*-t)} W(\tau^*)^{-1} x \quad \text{or} \quad u(t) = B^T (A^T)^{\tau^*-t-1} W(\tau^*)^{-1} x.$$

Note that this special input function is smooth in the continuous case (although only piecewise continuity has been required at the beginning).

Since $W(\tau^*)$ is always positive semi-definite due to its form, condition 2 from above is also equivalent to: $W(\tau^*)$ is positive definite, that is, $x^T W(\tau^*) x > 0$ for all $0 \neq x \in \mathbb{R}^n$. Then we write $W(\tau^*) > 0$.

Remark 4.10 It is worth noting that we obtain the same, purely algebraic condition for reachability/controllability both in continuous and in discrete time, namely, $\text{rank}(K) = n$, where n is the dimension of the state space. In that case, it is not ambiguous to simply say that the matrix pair (A, B) is reachable/controllable (this notion is independent of the time set).

The well-known Hamilton-Cayley theorem, i.e.,

$$\chi_A(A) = A^n + a_{n-1}A^{n-1} + \dots + a_1A + a_0I = 0$$

implies the n -th power (and hence all higher powers) of an $n \times n$ matrix A is a linear combination of the first n powers of A , that is, $A^0 = I, A^1 = A, \dots, A^{n-1}$. Therefore we have

$$\begin{aligned} \text{im}(K) &= K\mathbb{R}^{nm} = \text{span}\{A^i b_j \mid i = 0, \dots, n-1, j = 1, \dots, m\} \\ &= \text{span}\{A^i b_j \mid i \in \mathbb{N}, j = 1, \dots, m\} \end{aligned}$$

where $b_j \in \mathbb{R}^n$ are the columns of B , that is, $B = [b_1, \dots, b_m]$. This will be used in the proof of Theorem 4.9.

The result

$$\mathcal{R} = \text{im}(K) = \text{im}(B) + A \text{im}(B) + \dots + A^{n-1} \text{im}(B)$$

can also be formulated as follows: \mathcal{R} is the smallest A -invariant (that is, $A\mathcal{R} \subseteq \mathcal{R}$) subspace of the state space X that contains $\text{im}(B)$.

Proof: Due to Remark 4.8, we only need to do the continuous case. Then $\tau^* = \varepsilon > 0$ is arbitrary. We start with showing that $\mathcal{R}(\varepsilon) \subseteq \text{im}(K)$: If $x \in \mathcal{R}(\varepsilon)$, then there exists $u \in \mathcal{U}$ such that

$$x = \int_0^\varepsilon e^{A\tau} B u(\varepsilon - \tau) d\tau = \int_0^\varepsilon \sum_{i=0}^{\infty} \frac{\tau^i}{i!} A^i B u(\varepsilon - \tau) d\tau \in \text{im}(K).$$

Secondly, we wish to show that $\text{im}(K) = \text{im}(W(\varepsilon))$. From linear algebra, we know that it is equivalent to prove that $\text{im}(K)^\perp = \text{im}(W(\varepsilon))^\perp$. Let $x \in \text{im}(K)^\perp$, that is, $\langle x, y \rangle = 0$ for all $y \in \text{im}(K)$. Then $x^T K z = 0$ for all $z \in \mathbb{R}^{nm}$, which means that $x^T K = 0$ and hence $x^T A^i B = 0$ for all i . Then also $x^T e^{At} B = x^T \sum_{i=0}^{\infty} \frac{t^i}{i!} A^i B = 0$ and thus $x^T W(\varepsilon) = 0$. This shows that $x \in \text{im}(W(\varepsilon))^\perp$. Conversely, let $x \in \text{im}(W(\varepsilon))^\perp$, then

$$x^T W(\varepsilon) = \int_0^\varepsilon x^T e^{A\tau} B B^T e^{A^T \tau} d\tau = 0.$$

Post-multiplying this by x , we obtain

$$\int_0^\varepsilon \|B^T e^{A^T \tau} x\|^2 d\tau = 0.$$

We conclude that the smooth function $f(\tau) = B^T e^{A^T \tau} x$ is the constant zero function. Then also all its derivatives are zero. Evaluating them at $\tau = 0$, we obtain

$$B^T x = 0, \quad B^T A^T x = 0, \quad B^T (A^T)^2 x = 0, \quad \dots$$

that is, $x^T K = 0$ and hence $x \in \text{im}(K)^\perp$.

Finally, we need to show that $\text{im}(W(\varepsilon)) \subseteq \mathcal{R}(\varepsilon)$. Let $x \in \text{im}(W(\varepsilon))$, then there exists $z \in \mathbb{R}^n$ such that

$$x = W(\varepsilon)z.$$

Set $u(t) = B^T e^{A^T(\varepsilon-t)} z$. Then

$$\begin{aligned} \varphi(\varepsilon, 0, 0, u) &= \int_0^\varepsilon e^{A(\varepsilon-\tau)} B u(\tau) d\tau \\ &= \int_0^\varepsilon e^{A(\varepsilon-\tau)} B B^T e^{A^T(\varepsilon-\tau)} z d\tau \\ &= \int_0^\varepsilon e^{A\tau} B B^T e^{A^T \tau} z d\tau = W(\varepsilon)z = x \end{aligned}$$

which shows that $x \in \mathcal{R}(\varepsilon)$. □

Remark 4.11 One can show that the special input functions given in Theorem 4.9 are optimal in the sense that the **energy** associated with them, that is

$$E(u) = \int_0^{\tau^*} \|u(\tau)\|^2 d\tau \quad \text{or} \quad E(u) = \sum_{i=0}^{\tau^*-1} \|u(i)\|^2$$

is minimal among the energies of all u with $\varphi(\tau^*, 0, 0, u) = x$. This minimal energy (for controlling the system from 0 to x in time τ^*) is given by $E_{\min}(\tau^*, x) = x^T W(\tau^*)^{-1} x$. This shows that the smaller t is, the more energy is needed to do the transition from 0 to x in time t (compare this with Remark 4.4). More precisely, in continuous time, if $0 < s < t$, then $W(t) - W(s) > 0$, which implies $W(s)^{-1} - W(t)^{-1} > 0$, and hence

$$E_{\min}(s, x) = x^T W(s)^{-1} x > x^T W(t)^{-1} x = E_{\min}(t, x) \quad \text{for all } x \neq 0.$$

Again, the statement is easy to prove in discrete time (then we put $\tau^* = n$): Let v be defined as in Remark 4.8, then $E(u) = \|v\|^2$. We need to find the solution v

of $Kv = x$ which makes $\|v\|^2$ minimal. The special input from Theorem 4.9 corresponds to $v^* = K^T(KK^T)^{-1}x$. Let v be another solution of $Kv = x$, then $v = v^* + v_0$, where v_0 solves the homogeneous equation $Kv_0 = 0$. Then

$$\|v\|^2 = \|v^* + v_0\|^2 = \|v^*\|^2 + 2\langle v^*, v_0 \rangle + \|v_0\|^2.$$

However, $\langle v^*, v_0 \rangle = \langle K^T(KK^T)^{-1}x, v_0 \rangle = \langle (KK^T)^{-1}x, Kv_0 \rangle = 0$. Hence

$$\|v\|^2 = \|v^*\|^2 + \|v_0\|^2 \geq \|v^*\|^2$$

which shows that v^* is indeed the solution of $Kv = x$ that has minimum norm. Finally,

$$E_{\min}(n, x) = \|v^*\|^2 = \|K^T(KK^T)^{-1}x\|^2 = x^T(KK^T)^{-1}x = x^TW(n)^{-1}x.$$

For continuous time, see Appendix F.

Corollary 4.12 Consider $\dot{x} = Ax + Bu$ or $\sigma x = Ax + Bu$. In the discrete case, assume that A is invertible. Then $\mathcal{C} = \mathcal{R}$.

In general, we only have $\mathcal{R} \subseteq \mathcal{C}$, see for example Remark 4.7.

Proof: Let τ^* be as usual. We have $x \in \mathcal{C} = \mathcal{C}(\tau^*)$ if and only if there exists $u \in \mathcal{U}$ such that

$$\varphi(\tau^*, 0, x, u) = \varphi(\tau^*, 0, x, 0) + \varphi(\tau^*, 0, 0, u) = 0,$$

that is, $\varphi(\tau^*, 0, x, 0) = \varphi(\tau^*, 0, 0, -u)$. Thus $x \in \mathcal{C}$ if and only if $\varphi(\tau^*, 0, x, 0) \in \mathcal{R}$. In continuous time, this means

$$x \in \mathcal{C} \quad \Leftrightarrow \quad e^{A\tau^*}x \in \mathcal{R}. \quad (4.12)$$

Since $e^{A\tau^*}$ is invertible, this shows that $e^{A\tau^*}\mathcal{C} = \mathcal{R}$ and hence $\dim(\mathcal{C}) = \dim(\mathcal{R})$. In discrete time,

$$x \in \mathcal{C} \quad \Leftrightarrow \quad A^{\tau^*}x \in \mathcal{R}. \quad (4.13)$$

If A is invertible, we can argue as in the continuous case to see that \mathcal{R} and \mathcal{C} have the same dimension. Thus it suffices to show that $\mathcal{R} \subseteq \mathcal{C}$. If $x \in \mathcal{R} = \mathcal{R}(\tau^*)$, then $A^{\tau^*}x$ and $e^{A\tau^*}x$ are also in \mathcal{R} (this is due to the A -invariance of \mathcal{R}). According to (4.12) and (4.13), this implies $x \in \mathcal{C}$. \square

4.2 Controllable matrix pairs

Let $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times m}$. We say that the matrix pair (A, B) is **controllable** if the associated Kalman controllability matrix

$$K = [B \quad AB \quad \cdots \quad A^{n-1}B]$$

has full row rank, that is, $\text{rank}(K) = n$.

If a state space system $\dot{x} = Ax + Bu$ is subject to a coordinate transform $x = Tz$, where $T \in \mathbb{R}^{n \times n}$ is invertible, then we get

$$\dot{z} = T^{-1}ATz + T^{-1}Bu.$$

Discrete systems behave analogously. We say that the matrix pair $(T^{-1}AT, T^{-1}B)$ is **similar** to the matrix pair (A, B) . Of course, a coordinate transform should not change structural system properties such as stability and controllability. Indeed, similar matrices have the same eigenvalues, and the ranks of the Kalman controllability matrices of similar matrix pairs coincide.

The following result is limited to the **single-input** case, that is, $m = 1$. Then B is a single column vector. In that case, we simply write b instead of B . The associated Kalman controllability matrix

$$K = [b \quad Ab \quad \cdots \quad A^{n-1}b]$$

is then a square matrix.

Theorem 4.13 Let $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$, and let (A, b) be a controllable matrix pair. Then there exists an invertible matrix $T \in \mathbb{R}^{n \times n}$ such that

$$\tilde{A} := T^{-1}AT = \begin{bmatrix} 0 & \cdots & 0 & -a_0 \\ 1 & & & -a_1 \\ & \ddots & & \vdots \\ & & 1 & -a_{n-1} \end{bmatrix} \quad \text{and} \quad \tilde{b} := T^{-1}b = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

The numbers a_i are precisely the coefficients of the characteristic polynomial, that is,

$$\chi_A(s) = \chi_{T^{-1}AT}(s) = s^n + a_{n-1}s^{n-1} + \cdots + a_1s + a_0.$$

This is called the **controllability form** of (A, b) . Moreover, there exists an invertible matrix $T_1 \in \mathbb{R}^{n \times n}$ such that

$$T_1^{-1}AT_1 = \begin{bmatrix} 0 & 1 & & \\ \vdots & & \ddots & \\ 0 & & & 1 \\ -a_0 & -a_1 & \cdots & -a_{n-1} \end{bmatrix} \quad \text{and} \quad T_1^{-1}b = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}.$$

This is called the **controller form** of (A, b) . The coefficients a_i are the same as with the controllability form.

Proof: According to the Hamilton-Cayley theorem,

$$A^n = -a_{n-1}A^{n-1} - \dots - a_1A - a_0I.$$

Thus we have

$$\begin{aligned} AK &= [Ab \ A^2b \ \dots \ A^nb] \\ &= [b \ Ab \ \dots \ A^{n-1}b] \begin{bmatrix} 0 & \dots & 0 & -a_0 \\ 1 & & & -a_1 \\ & \ddots & & \vdots \\ & & 1 & -a_{n-1} \end{bmatrix} = K\tilde{A}. \end{aligned}$$

Moreover, we have

$$b = [b \ Ab \ \dots \ A^{n-1}b] \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = K\tilde{b}.$$

Since (A, b) is controllable, K is invertible. Thus we simply put $T = K$. This completes the proof for the controllability form. For the controller form, the construction is a bit more involved, and it is omitted here. \square

A matrix pair (A, b) is in controllability form if and only if its Kalman matrix is the identity matrix. If a scalar input-output representation

$$y^{(n)} + a_{n-1}y^{(n-1)} + \dots + a_1\dot{y} + a_0y = u$$

is reduced to first order in the usual way, i.e., via putting $x = [y, \dot{y}, \dots, y^{(n-1)}]^T$, then the resulting state space system is precisely in controller form.

We return to the general **multi-input** case, and we give another result about transforming a given matrix pair into some special form via a similarity transform (this corresponds to a coordinate transform in the state space).

Theorem 4.14 (Kalman controllability decomposition) Let $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times m}$. Let K be the associated Kalman controllability matrix and let $r := \text{rank}(K)$. Then there exists an invertible matrix $T \in \mathbb{R}^{n \times n}$ such that

$$\tilde{A} := T^{-1}AT = \begin{bmatrix} A_1 & A_2 \\ 0 & A_3 \end{bmatrix} \quad \text{and} \quad \tilde{B} := T^{-1}B = \begin{bmatrix} B_1 \\ 0 \end{bmatrix}$$

where $A_1 \in \mathbb{R}^{r \times r}$, $B_1 \in \mathbb{R}^{r \times m}$ is a controllable matrix pair.

Remark 4.15 The theorem says that by a suitable coordinate transform, namely $x = Tz$, the given system $\dot{x} = Ax + Bu$ can be put into the form

$$\dot{z}_1 = A_1 z_1 + A_2 z_2 + B_1 u \quad (4.14)$$

$$\dot{z}_2 = A_3 z_2 \quad (4.15)$$

where (A_1, B_1) is controllable. The second equation is certainly not controllable (it is completely decoupled from the input, hence we cannot influence z_2 by the choice of the input u). In fact, it is an autonomous equation. If we start in $z_2(0) = 0$, then we have $z_2(t) = 0$ for all t . Then the first equation becomes

$$\dot{z}_1 = A_1 z_1 + B_1 u$$

which is controllable. Thus the states that are reachable from zero in the system (4.14), (4.15) are precisely those of the form $\begin{pmatrix} \zeta \\ 0 \end{pmatrix}$, where $\zeta \in \mathbb{R}^r$ is arbitrary. In other words, the reachable space of (4.14), (4.15) takes the simple form

$$\{z \in \mathbb{R}^n \mid z \text{ can be reached from } 0\} = \mathbb{R}^r \times \{0\}.$$

This can be used to determine the reachable space of the original system, because

$$\{x \in \mathbb{R}^n \mid x \text{ can be reached from } 0\} = T(\mathbb{R}^r \times \{0\}).$$

Note that if the original (A, B) is controllable, then $r = n$, and the Kalman controllability decomposition becomes trivial. Thus the interesting case arises when (A, B) itself is not controllable.

Proof: Let $v_1, \dots, v_r \in \mathbb{R}^n$ be a basis of \mathcal{R} . Let $w_1, \dots, w_{n-r} \in \mathbb{R}^n$ be such that $v_1, \dots, v_r, w_1, \dots, w_{n-r}$ is a basis of \mathbb{R}^n . Then

$$T := \begin{bmatrix} V & W \end{bmatrix} := \begin{bmatrix} v_1 & \cdots & v_r & w_1 & \cdots & w_{n-r} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

is an invertible matrix. Due to the A -invariance of \mathcal{R} , the columns of AV are again in \mathcal{R} . Thus they can be written as linear combinations of the vectors v_i , that is,

$$AV = VA_1$$

for some matrix $A_1 \in \mathbb{R}^{r \times r}$. On the other hand, the columns of AW are in \mathbb{R}^n and thus they can be written as linear combinations of v_i, w_j , that is,

$$AW = VA_2 + WA_3$$

for some matrices A_2, A_3 of appropriate sizes. Summing up, we have

$$AT = \begin{bmatrix} AV & AW \end{bmatrix} = \begin{bmatrix} VA_1 & VA_2 + WA_3 \end{bmatrix} = \begin{bmatrix} V & W \end{bmatrix} \begin{bmatrix} A_1 & A_2 \\ 0 & A_3 \end{bmatrix} = T\tilde{A}.$$

Since $\text{im}(B) \subseteq \mathcal{R}$, the columns of B are linear combinations of the vectors v_i , that is,

$$B = VB_1 = \begin{bmatrix} V & W \end{bmatrix} \begin{bmatrix} B_1 \\ 0 \end{bmatrix} = T\tilde{B}$$

for some matrix $B_1 \in \mathbb{R}^{r \times m}$. It remains to show that (A_1, B_1) is a controllable matrix pair. The Kalman controllability matrix associated to (\tilde{A}, \tilde{B}) is

$$\tilde{K} = \begin{bmatrix} \tilde{B} & \tilde{A}\tilde{B} & \cdots & \tilde{A}^{n-1}\tilde{B} \end{bmatrix} = \begin{bmatrix} B_1 & A_1B_1 & \cdots & A_1^{n-1}B_1 \\ 0 & 0 & \cdots & 0 \end{bmatrix}$$

and this has rank r , just like the Kalman controllability matrix of the original matrix pair. Therefore,

$$r = \text{rank} \begin{bmatrix} B_1 & A_1B_1 & \cdots & A_1^{n-1}B_1 \end{bmatrix}.$$

Due to the Hamilton-Cayley theorem, this implies that

$$r = \text{rank} \begin{bmatrix} B_1 & A_1B_1 & \cdots & A_1^{r-1}B_1 \end{bmatrix}$$

which shows that (A_1, B_1) is controllable. \square

Remark 4.16 Let χ_A and $\text{spec}(A)$ denote the characteristic polynomial and the spectrum of A , respectively. In a Kalman controllability decomposition, we clearly have

$$\chi_A = \chi_{A_1} \cdot \chi_{A_3}$$

and thus

$$\text{spec}(A) = \text{spec}(A_1) \cup \text{spec}(A_3).$$

One calls χ_{A_3} the **uncontrollable part** of the characteristic polynomial of A with respect to B , and $\lambda \in \text{spec}(A_3)$ an **uncontrollable mode** of (A, B) . Of course, it has to be verified that these notions do not depend on the specific choice of the Kalman decomposition (which is non-unique, in general, since its construction involves several choices): In fact, the matrix A_1 is the matrix representation of $A|_{\mathcal{R}} : \mathcal{R} \rightarrow \mathcal{R}$ with respect to the chosen basis of \mathcal{R} . Thus, A_1 depends on the choice of the basis, but χ_{A_1} does not. Therefore, this holds also for $\chi_{A_3} = \chi_A / \chi_{A_1}$.

Example 4.17 Consider the matrices from Remark 4.7,

$$A = \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

The Kalman controllability matrix is

$$K = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

which has rank $r = 1$. We choose the vector $v_1 := B$ as a basis of $\mathcal{R} = \text{im}(K)$. If we choose

$$T := \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

then

$$T^{-1}AT = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \quad \text{and} \quad T^{-1}B = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

Thus we see that $\chi_A(s) = (s-1)s$ with uncontrollable part s , and hence 0 is an uncontrollable mode of (A, B) .

There is also a direct way to characterize the uncontrollable modes of a matrix pair.

Theorem 4.18 Let $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, and $\lambda \in \mathbb{C}$. The following are equivalent:

1. λ is an uncontrollable mode of (A, B) ;
2. $\text{rank} \begin{bmatrix} \lambda I - A & B \end{bmatrix} < n$.

Proof: Both conditions are invariant under similarity transforms. For the first condition, this follows from the argument given in Remark 4.16, and for the second condition, let $\tilde{A} = T^{-1}AT$ and $\tilde{B} = T^{-1}B$. Since

$$\begin{bmatrix} \lambda I - \tilde{A} & \tilde{B} \end{bmatrix} = T^{-1} \begin{bmatrix} \lambda I - A & B \end{bmatrix} \begin{bmatrix} T & 0 \\ 0 & I \end{bmatrix},$$

we have

$$\text{rank} \begin{bmatrix} \lambda I - \tilde{A} & \tilde{B} \end{bmatrix} = \text{rank} \begin{bmatrix} \lambda I - A & B \end{bmatrix}.$$

Thus we may assume, without loss of generality, that a Kalman controllability decomposition has already been performed, i.e.,

$$A = \begin{bmatrix} A_1 & A_2 \\ 0 & A_3 \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} B_1 \\ 0 \end{bmatrix}$$

where (A_1, B_1) is controllable. Then we have

$$H(\lambda) := \begin{bmatrix} \lambda I - A & B \end{bmatrix} = \begin{bmatrix} \lambda I - A_1 & -A_2 & B_1 \\ 0 & \lambda I - A_3 & 0 \end{bmatrix}.$$

If λ is an uncontrollable mode, then it is an eigenvalue of A_3 . Thus it makes $\lambda I - A_3$ singular, showing that the rank of $H(\lambda)$ cannot be full, that is, it must be less than n . Conversely, assume that the rank of $H(\lambda)$ is not full. Then there exists a vector $x \neq 0$ such that $xH(\lambda) = 0$, that is,

$$\begin{aligned} x_1(\lambda I - A_1) &= 0 \\ -x_1A_2 + x_2(\lambda I - A_3) &= 0 \\ x_1B_1 &= 0. \end{aligned}$$

The first and third equations imply that

$$x_1B_1 = 0, \quad x_1A_1B_1 = 0, \quad x_1A_1^2B_1 = 0, \quad \dots$$

The controllability of (A_1, B_1) yields that $x_1 = 0$. Then $x_2 \neq 0$ and

$$x_2(\lambda I - A_3) = 0$$

which implies that λ is an eigenvalue of A_3 , that is, an uncontrollable mode. \square

Example 4.19 Consider the matrix pair from Example 4.17. Then

$$H(\lambda) = \begin{bmatrix} \lambda - 1 & 0 & 1 \\ -1 & \lambda & 1 \end{bmatrix}$$

which has rank 2 for all $\lambda \neq 0$. However, $\text{rank}(H(0)) = 1$, showing again that 0 is an uncontrollable mode of this system.

As a direct consequence of this, we obtain another characterization of controllable matrix pairs.

Corollary 4.20 (Hautus test for controllability) The following are equivalent:

1. (A, B) is controllable.
2. The matrix $H(\lambda) = \begin{bmatrix} \lambda I - A & B \end{bmatrix}$ has full row rank for all $\lambda \in \mathbb{C}$.

The polynomial matrix $H = \begin{bmatrix} sI - A & B \end{bmatrix} \in \mathbb{R}[s]^{n \times (n+m)}$ is called **Hautus controllability matrix**. Since $\lambda I - A$ is non-singular whenever λ is not an eigenvalue of A , it suffices to check condition 2 from above for $\lambda \in \text{spec}(A)$.

Remark 4.21 Controllability is a **generic** property, that is, if a matrix pair (A, B) is chosen “at random”, then it is very likely to be a controllable one. More precisely, the set of controllable matrix pairs (A, B) is open and dense in the set $\mathbb{R}^{n \times n} \times \mathbb{R}^{n \times m}$. This is due to the fact that (A, B) is uncontrollable if and only if all $n \times n$ subdeterminants of the associated Kalman matrix vanish. This defines a set of polynomial equations to be satisfied by the entries A_{ij} , B_{ij} of A, B . Thus the set of uncontrollable matrix pairs is a proper algebraic variety in $\mathbb{R}^{n \times n} \times \mathbb{R}^{n \times m}$.

4.3 Asymptotic controllability

Sometimes, it is not required that a system should go from one state to another in *finite* time τ . Instead, one is satisfied if this happens asymptotically as $\tau \rightarrow \infty$.

Definition 4.22 We say that a state space system $\dot{x} = Ax + Bu$ or $\sigma x = Ax + Bu$ is **asymptotically controllable (to zero)** if for any $x_0 \in X = \mathbb{R}^n$, there exists an input function $u \in \mathcal{U}$ such that

$$\lim_{t \rightarrow \infty} \varphi(t, 0, x_0, u) = 0.$$

Clearly, controllability implies asymptotic controllability. The proof of the following theorem will be done in the next chapter.

Theorem 4.23 A state space system is asymptotically controllable if and only if its uncontrollable modes λ are asymptotically stable, that is, $\text{Re}(\lambda) < 0$ in continuous time, and $|\lambda| < 1$ in discrete time.

Let us convince ourselves at least of the simple direction of the proof. We may assume, without loss of generality, that

$$A = \begin{bmatrix} A_1 & A_2 \\ 0 & A_3 \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} B_1 \\ 0 \end{bmatrix}.$$

Then $x_2(t) = e^{A_3 t} x_2(0)$ or $x_2(t) = A_3^t x_2(0)$, respectively. If the system is asymptotically controllable, we must have $\lim_{t \rightarrow \infty} x_2(t) = 0$ for all $x_2(0)$, which implies that A_3 must be asymptotically stable.

4.4 Controllable behaviors

The Hautus test gives us an idea about how to generalize the notion of controllability from state space systems $\dot{x} = Ax + Bu$ to general systems $R(\frac{d}{dt})w = 0$ where $R \in \mathbb{R}[s]^{p \times q}$, and $w \in \mathcal{A}^q$. In a state space system,

$$R = \begin{bmatrix} sI - A & -B \end{bmatrix} \in \mathbb{R}[s]^{n \times (n+m)} \quad \text{and} \quad w = \begin{bmatrix} x \\ u \end{bmatrix} \in \mathcal{A}^{n+m}.$$

The polynomial matrix R is recognized as the Hautus controllability matrix (up to the sign of B , which does not influence the rank).

In this section, we restrict to continuous systems, and we return to our original signal spaces, that is, $\mathcal{A} = \mathcal{D}'(T)$, where $T = \mathbb{R}, \mathbb{R}_+$.

Definition 4.24 Let $t_0 \in T$ be fixed. Let $R \in \mathbb{R}[s]^{p \times q}$ and

$$\mathcal{B} = \{w \in \mathcal{A}^q \mid R(\frac{d}{dt})w = 0\}$$

where \mathcal{A} is as described above. Let $w^{(1)}, w^{(2)} \in \mathcal{B}$ be two trajectories. We say that $w^{(2)}$ is reachable from $w^{(1)}$ in time $\tau \geq 0$ (or: $w^{(1)}$ is controllable to $w^{(2)}$ in time τ) if there exists a trajectory $w \in \mathcal{B}$ which coincides with $w^{(1)}$ on $(-\infty, t_0) \cap T$ and with $w^{(2)}$ on $(t_0 + \tau, \infty) \cap T$. If there is any such $\tau \geq 0$, then we say that $w^{(2)}$ is reachable from $w^{(1)}$ (or: $w^{(1)}$ is controllable to $w^{(2)}$). One says that \mathcal{B} is **controllable** if any $w^{(1)} \in \mathcal{B}$ can be controlled to any $w^{(2)} \in \mathcal{B}$.

To say that two distributions coincide on an open set $U \subseteq \mathbb{R}$ means that they assign the same value to each test function whose support lies in U . In the situation described above, one calls w a **connecting trajectory** for $w^{(1)}, w^{(2)}$. For classical functions, this means

$$w(t) = \begin{cases} w^{(1)}(t) & \text{if } t < t_0 \\ w^{(2)}(t) & \text{if } t > t_0 + \tau. \end{cases}$$

The choice of the starting time t_0 makes no difference, because we consider only time-invariant systems. Similarly as with state space systems, the transition time τ can be made arbitrarily small, independently of the choice of the two trajectories to be connected.

Theorem 4.25 (Generalized Hautus test) Let $\mathcal{A} = \mathcal{D}'(T)$ for $T = \mathbb{R}$ or \mathbb{R}_+ and let $R \in \mathbb{R}[s]^{p \times q}$. Without loss of generality, let R have full row rank. Then the behavior

$$\mathcal{B} = \{w \in \mathcal{A}^q \mid R(\frac{d}{dt})w = 0\}$$

is controllable if and only if $\text{rank}(R(\lambda)) = p$ for all $\lambda \in \mathbb{C}$.

Proof: Let R have full row rank and let

$$URV = [D \ 0]$$

be the Smith form of R , with $D = \text{diag}(d_1, \dots, d_p)$. Since U, V are unimodular, we have, for all $\lambda \in \mathbb{C}$,

$$\text{rank}(R(\lambda)) = \text{rank}(U(\lambda)R(\lambda)V(\lambda)) = \text{rank}(D(\lambda)).$$

Thus, $\text{rank}(R(\lambda)) = p$ for all $\lambda \in \mathbb{C}$ if and only if $\det(D(\lambda)) \neq 0$ for all $\lambda \in \mathbb{C}$, that is, if no d_i has a zero in \mathbb{C} . This is true if and only if all d_i are constants. Since we may always assume that the d_i are monic polynomials (i.e., their leading coefficients are equal to one), we have

$$\text{rank}(R(\lambda)) = p \text{ for all } \lambda \in \mathbb{C} \iff D = I.$$

Consider as usual

$$\tilde{\mathcal{B}} = \{ \tilde{w} \in \mathcal{A}^q \mid [D \ 0] \tilde{w} = 0 \}$$

which is related to \mathcal{B} via the isomorphism $\tilde{\mathcal{B}} \rightarrow \mathcal{B}$, $\tilde{w} \mapsto w = V(\frac{d}{dt})\tilde{w}$. This $\tilde{\mathcal{B}}$ is controllable if and only if \mathcal{B} is controllable. However, if $D = I$, then

$$\tilde{\mathcal{B}} = \{0\} \times \mathcal{A}^{q-p} \subset \mathcal{A}^q$$

which is clearly controllable. Conversely, if $D \neq I$, there exists at least one d_i , say d_1 , which is not constant. Then the equation for the first component \tilde{w}_1 of \tilde{w} reads

$$d_1(\frac{d}{dt})\tilde{w}_1 = 0$$

which has precisely the solutions

$$\tilde{w}_1(t) = \sum_{\lambda} a_{\lambda}(t)e^{\lambda t} \tag{4.16}$$

where $\lambda \in \mathbb{C}$ are the zeros of d_1 (since d_1 is not a constant, there exists at least one such λ). Now consider $t_0 \in T$ and two trajectories $\tilde{w}^{(1)}, \tilde{w}^{(2)}$ in $\tilde{\mathcal{B}}$, where the first component of $\tilde{w}^{(1)}$ is not identically zero on $(-\infty, t_0) \cap T$, and $w^{(2)}$ is the zero trajectory. Then $\tilde{w}^{(2)}$ is not reachable from $\tilde{w}^{(1)}$: A connecting trajectory \tilde{w} would have to satisfy $\tilde{w}_1 = 0$ on $(t_0 + \tau, \infty)$, which implies that $\tilde{w}_1 = 0$ everywhere because (4.16) shows that \tilde{w}_1 must be an analytic function. Thus \tilde{w} cannot coincide with $\tilde{w}^{(1)}$ on $(-\infty, t_0) \cap T$. This shows that the system is not controllable. \square

For instance, a scalar input-output representation

$$y^{(n)} + a_{n-1}y^{(n-1)} + \dots + a_1\dot{y} + a_0y = u$$

is controllable, because putting

$$R = \begin{bmatrix} -1 & s^n + a_{n-1}s^{n-1} + \dots + a_1s + a_0 \end{bmatrix} \quad \text{and} \quad w = \begin{bmatrix} u \\ y \end{bmatrix},$$

we see that $\text{rank}(R(\lambda)) = 1$ for all $\lambda \in \mathbb{C}$.

Polynomial matrices that represent controllable behaviors are characterized by the following theorem.

Theorem 4.26 Let $R \in \mathbb{R}[s]^{p \times q}$ be a polynomial matrix with full row rank. The following are equivalent:

1. $\text{rank}(R(\lambda)) = p$ for all $\lambda \in \mathbb{C}$;
2. The Smith form of R is $\begin{bmatrix} I & 0 \end{bmatrix}$;
3. There exists a matrix $T \in \mathbb{R}[s]^{(q-p) \times q}$ such that $\begin{bmatrix} R \\ T \end{bmatrix}$ is unimodular;
4. There exists a matrix $S \in \mathbb{R}[s]^{q \times p}$ such that $RS = I$;
5. If $R = UR_1$ for some $U \in \mathbb{R}[s]^{p \times p}$, $R_1 \in \mathbb{R}[s]^{p \times q}$, then U must be unimodular.

If the equivalent conditions are satisfied, we say that R is **left prime** (or: left irreducible).

Proof: We have already seen in the previous proof that “1 \Rightarrow 2”. The converse is obvious.

For “2 \Rightarrow 3”, assume that

$$R = U \begin{bmatrix} I & 0 \end{bmatrix} V = U \begin{bmatrix} I & 0 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = UV_1$$

where U and V are unimodular. Define $T := V_2$, then

$$\begin{bmatrix} R \\ T \end{bmatrix} = \begin{bmatrix} UV_1 \\ V_2 \end{bmatrix} = \begin{bmatrix} U & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix}$$

which shows that the matrix is unimodular.

For “3 \Rightarrow 4”, let T be a matrix according to assertion 3. Then there exist matrices S_1, S_2 such that

$$\begin{bmatrix} R \\ T \end{bmatrix} \begin{bmatrix} S_1 & S_2 \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix}$$

In particular, $RS_1 = I$.

For “4 \Rightarrow 5”, let S be such that $RS = I$ and $R = UR_1$. Then $UR_1S = I$, which shows that U is unimodular.

Finally, we show “5 \Rightarrow 2” by negation. Assume that the Smith form is $\begin{bmatrix} D & 0 \end{bmatrix}$ with $D = \text{diag}(d_1, \dots, d_p)$ and at least one of the d_i is not a constant, say $d_1 \notin \mathbb{R}$. Then

$$R = U \begin{bmatrix} D & 0 \end{bmatrix} V = U_1 R_1$$

with

$$U_1 = U \begin{bmatrix} d_1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \quad \text{and} \quad R_1 = \begin{bmatrix} 1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_p & 0 \end{bmatrix} V.$$

We have $\det(U_1) = \det(U)d_1$ and thus we have found a factorization $R = U_1 R_1$ in which U_1 is not unimodular. \square

A polynomial matrix $R \in \mathbb{R}[s]^{p \times q}$ with full row rank is left prime if and only if its $p \times p$ subdeterminants have no common zeros in \mathbb{C} (since we have $\text{rank}(R(\lambda)) < p$ if and only if $\lambda \in \mathbb{C}$ is a common zero of all $p \times p$ minors). By the fundamental theorem of algebra, it is also equivalent to say that the $p \times p$ subdeterminants of R are devoid of common factors, i.e., they are coprime polynomials in $\mathbb{R}[s]$.

4.5 Non-linear systems and accessibility

Consider

$$\dot{x}(t) = F(x(t), u(t)), \tag{4.17}$$

where $t \in T$ for some open interval $T \subseteq \mathbb{R}$, and $x(t) \in \mathbb{R}^n$, $u(t) \in \mathbb{R}^m$ denote the state and the input at time t , respectively. The map $F : X \times U \rightarrow \mathbb{R}^n$, where $X \subseteq \mathbb{R}^n$, $U \subseteq \mathbb{R}^m$ are open sets, is supposed to be continuously differentiable. Let \mathcal{U} denote the set of all piecewise continuous functions from T to $U \subseteq \mathbb{R}^m$. It follows from classical ODE theory that for all $u \in \mathcal{U}$, $t_0 \in T$, and $x_0 \in X$ there exists a unique solution

$$x : J \rightarrow X, t \mapsto x(t)$$

to (4.17) with $x(t_0) = x_0$, which is continuous and piecewise continuously differentiable. Here, J is the so-called maximal existence interval, which is an open subinterval of T with $t_0 \in J$. In general, J depends on t_0, x_0 and u , that is, we should actually write $J = J(t_0, x_0, u)$ to be precise. We will use the notation $\varphi(t, t_0, x_0, u)$ for the solution of the initial value problem

$$\begin{aligned}\dot{x}(t) &= F(x(t), u(t)) \\ x(t_0) &= x_0\end{aligned}$$

at time t . When we write $\varphi(t, t_0, x_0, u)$, we tacitly assume that $t \in J(t_0, x_0, u)$, otherwise the notation would not make sense. In the following, we put $t_0 = 0$ without loss of generality (since the right hand side of the differential equation does not explicitly depend on t , we have time-invariance as in (4.7)).

Just as in the linear case, we call $x_1 \in X$ **reachable from $x_0 \in X$ in time $\tau \geq 0$** if $x_1 = \varphi(\tau, 0, x_0, u)$ for some $u \in \mathcal{U}$. Equivalently, we say that x_0 is **controllable to x_1 in time τ** . Define

$$\begin{aligned}\mathcal{R}(\tau, x_0) &:= \{x \in X \mid x \text{ is reachable from } x_0 \text{ in time } \tau\} \\ \mathcal{C}(\tau, x_1) &:= \{x \in X \mid x \text{ is controllable to } x_1 \text{ in time } \tau\}.\end{aligned}$$

Unlike the linear case, these sets will not have an affine structure, in general. Also, there is no reason why $x_0 = 0$ should play a special role here. Nevertheless, we will also use

$$\mathcal{R}(x_0) = \bigcup_{\tau \geq 0} \mathcal{R}(\tau, x_0) \quad \text{and} \quad \mathcal{C}(x_1) = \bigcup_{\tau \geq 0} \mathcal{C}(\tau, x_1),$$

which are the sets of states that are reachable from x_0 or controllable to x_1 , respectively (in any finite time).

The system (4.17) is called **globally reachable from x_0** if $\mathcal{R}(x_0) = X$. If this holds for any x_0 , we call the system **globally reachable**. It turns out that global reachability is usually a too strong requirement for non-linear systems. Therefore, one studies the following weaker version: The system is said to be **locally reachable from x_0** if $\mathcal{R}(x_0)$ contains an open neighborhood of x_0 . Intuitively, this means that starting from x_0 , the system can be steered into “any direction”, or to any desired terminal state that is close enough to x_0 . The following examples show some reachability phenomena that can occur in the non-linear setting.

Example 4.27 1. Consider $\dot{x} = xu$, $x(0) = x_0 \in X = \mathbb{R}$. The unique solution of this initial value problem is given by

$$x(t) = x_0 e^{\int_0^t u(\tau) d\tau}.$$

If $x_0 = 0$, then $x \equiv 0$ and thus $\mathcal{R}(t_1, 0) = \{0\}$ for all t_1 . In particular, the system is not locally reachable from $x_0 = 0$. If $x_0 > 0$, then $x(t) > 0$ for all $t \in T$. Thus we cannot steer the system to any $x_1 \leq 0$. On the other hand, if $x_1 > 0$, then the constant function $u \equiv \frac{1}{t_1} \ln\left(\frac{x_1}{x_0}\right)$ yields $x(t_1) = x_1$. We conclude that $\mathcal{R}(t_1, x_0) = (0, \infty)$ for all $x_0 > 0$ and all $t_1 > 0$. Similarly, one can show that $\mathcal{R}(t_1, x_0) = (-\infty, 0)$ for all $x_0 < 0$ and all $t_1 > 0$. Thus the system is locally reachable from any $x_0 \neq 0$.

2. Consider $\dot{x}_1 = x_2$, $\dot{x}_2 = \sin(x_1) + u$, $x(0) = x_0 \in X = \mathbb{R}^2$. Setting $v := \sin(x_1) + u$, the problem becomes $\dot{x}_1 = x_2$, $\dot{x}_2 = v$, $x(0) = x_0$ which can be solved in closed form. One can show that for any $t_1 > 0$ and any $x_1 \in \mathbb{R}^2$, there exists $a, b \in \mathbb{R}$ such that $v(t) = a + bt$ steers the transformed system to $x(t_1) = x_1$. From this v , one can compute a suitable input u for the original system. Thus $\mathcal{R}(t_1, x_0) = X$ for all $t_1 > 0$ and all $x_0 \in X$. In particular, the system is globally reachable.
3. Consider $\dot{x}_1 = x_2^2$, $\dot{x}_2 = u$, $x(0) = x_0 \in X = \mathbb{R}^2$. Then $x_1(t) \geq x_{01}$ for all t , which implies that $\mathcal{R}(x_0)$ is contained in the half-plane $\{x \in \mathbb{R}^2 \mid x_1 \geq x_{01}\}$. Thus the system is not locally reachable from any x_0 .

Similarly to the notion of local reachability, (4.17) will be called **locally controllable to x_1** if $\mathcal{C}(x_1)$ contains an open neighborhood of x_1 , which says, roughly speaking, that any state that is sufficiently close to x_1 can be controlled to x_1 . The latter notion is particularly interesting when x_1 is an equilibrium of the system, because it means that small perturbations of the equilibrium can be corrected, i.e., there exists an input function such that the system returns to the equilibrium in finite time.

A state $x_0 \in X$ is called an equilibrium of (4.17) for the zero input function if $F(x_0, 0) = 0$. This is true if and only if $x \equiv x_0$ is a solution to the initial value problem from above, or equivalently, $\varphi(t, 0, x_0, 0) = x_0$ for all $t \in T$. The linear approximation of F at $x = x_0$, $u = 0$ is given by

$$F(x, u) \approx \underbrace{F(x_0, 0)}_{=0} + \underbrace{\frac{\partial F}{\partial x}(x_0, 0)}_{=: A \in \mathbb{R}^{n \times n}}(x - x_0) + \underbrace{\frac{\partial F}{\partial u}(x_0, 0)}_{=: B \in \mathbb{R}^{n \times m}} u.$$

Setting $\tilde{x} := x - x_0$, we obtain the linearization of (4.17) at x_0 , which is the linear system

$$\dot{\tilde{x}} = A\tilde{x} + Bu.$$

Our main goal in this section is the following result giving a sufficient condition for local reachability from x_0 in terms of the linearization at x_0 .

Theorem 4.28 Consider $\dot{x}(t) = F(x(t), u(t))$ with $F(x_0, 0) = 0$. Set

$$A = \frac{\partial F}{\partial x}(x_0, 0) \quad \text{and} \quad B = \frac{\partial F}{\partial u}(x_0, 0).$$

If (A, B) is a controllable matrix pair, then the reachability set $\mathcal{R}(t_1, x_0)$ of the non-linear system contains an open neighborhood of x_0 for all $t_1 > 0$.

In other words, if the linearization of (4.17) at an equilibrium point x_0 is completely reachable, then the underlying non-linear system is locally reachable from x_0 . Moreover, $\mathcal{R}(t, x_0)$ contains an open neighborhood for any arbitrarily small $t > 0$ (compare this with Remark 4.4).

For the proof, we need two auxiliary facts.

Fact 1: Parameter-dependent ODE. Consider

$$\begin{aligned} \dot{x}(t) &= f(t, x(t), \xi) \\ x(0) &= x_0, \end{aligned}$$

where $\xi \in \mathbb{R}^n$ is a parameter vector. If f is continuous, and even continuously differentiable with respect to x and ξ , then there exists a unique solution to this initial value problem, which we denote by $\phi(t, x_0, \xi)$. Then ϕ is also continuously differentiable (with respect to all three variables). Let ξ_0 be given and consider

$$G(t) := \frac{\partial \phi}{\partial \xi}(t, x_0, \xi_0).$$

This matrix-valued function is again continuously differentiable and satisfies

$$\dot{G}(t) = \frac{\partial f}{\partial x}(t, \phi(t, x_0, \xi_0), \xi_0)G(t) + \frac{\partial f}{\partial \xi}(t, \phi(t, x_0, \xi_0), \xi_0)$$

and $G(0) = 0$.

This fact will be used in the proof as follows: We will have n continuous input functions $u^{(i)}$ for $1 \leq i \leq n$ and we will consider a linear combination

$$u_\xi(t) = \xi_1 u^{(1)}(t) + \dots + \xi_n u^{(n)}(t),$$

where $\xi \in \mathbb{R}^n$. Plugging this into $\dot{x} = F(x, u)$, we get a parameter-dependent ODE by putting $f(t, x, \xi) := F(x, u_\xi(t))$. Then $\varphi(t, 0, x_0, u_\xi) = \phi(t, x_0, \xi)$. Suppose that x_0 is an equilibrium and set $\xi_0 = 0$. Then

$$G(t) = \frac{\partial \phi}{\partial \xi}(t, x_0, 0)$$

satisfies

$$\dot{G}(t) = \frac{\partial F}{\partial x}(\phi(t, x_0, 0), 0)G(t) + \frac{\partial F}{\partial u}(\phi(t, x_0, 0), 0)\frac{\partial u_\xi}{\partial \xi}(t).$$

Since $\phi(t, x_0, 0) = \varphi(t, 0, x_0, 0) = x_0$ for all t by assumption, we get

$$\begin{aligned}\dot{G}(t) &= AG(t) + B[u^{(1)}(t), \dots, u^{(n)}(t)] \\ G(0) &= 0.\end{aligned}$$

Fact 2: Inverse function theorem. Let $\Psi : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $\xi \mapsto \Psi(\xi)$ be a continuously differentiable function, and let $\xi_0 \in \mathbb{R}^n$ be such that $\det(\frac{\partial \Psi}{\partial \xi}(\xi_0)) \neq 0$. Then there exists an open neighborhood U of ξ_0 and an open neighborhood V of $\Psi(\xi_0)$ such that $U \rightarrow V$, $\xi \mapsto \Psi(\xi)$ is a diffeomorphism (i.e., continuously differentiable and bijective, with a continuously differentiable inverse).

Proof: Let $t_1 > 0$ be given. Since (A, B) is a controllable matrix pair, there exists, for all $1 \leq i \leq n$, a continuous input function $u^{(i)}$ such that the solution of

$$\begin{aligned}\dot{\tilde{x}}(t) &= A\tilde{x}(t) + Bu^{(i)}(t) \\ \tilde{x}(0) &= 0\end{aligned}$$

satisfies $\tilde{x}(t_1) = e_i$, where e_i is the i -th standard basis vector of \mathbb{R}^n . Set

$$u_\xi(t) = \xi_1 u^{(1)}(t) + \dots + \xi_n u^{(n)}(t)$$

for $\xi \in \mathbb{R}^n$. Consider

$$\begin{aligned}\dot{x}(t) &= F(x(t), u_\xi(t)) \\ x(0) &= x_0\end{aligned}$$

and let $\phi(t, x_0, \xi)$ denote its solution at time t . Then the function

$$G(t) = \frac{\partial \phi}{\partial \xi}(t, x_0, 0)$$

satisfies $G(0) = 0$ and

$$\dot{G}(t) = AG(t) + B[u^{(1)}(t), \dots, u^{(n)}(t)].$$

By the construction of the input functions $u^{(i)}$, we have

$$G(t_1) = [e_1, \dots, e_n] = I_n.$$

On the other hand,

$$G(t_1) = \frac{\partial \phi}{\partial \xi}(t_1, x_0, 0).$$

Define

$$\Psi : \mathbb{R}^n \rightarrow \mathbb{R}^n, \xi \mapsto \phi(t_1, x_0, \xi).$$

Then $\Psi(0) = x_0$, since x_0 is an equilibrium, and

$$\frac{\partial \Psi}{\partial \xi}(0) = G(t_1) = I_n.$$

The inverse function theorem implies that there exists an open neighborhood U of $\xi_0 = 0$ and an open neighborhood V of $\Psi(0) = x_0$ such that $U \rightarrow V$, $\xi \mapsto \Psi(\xi)$ is a diffeomorphism. Thus

$$V = \Psi(U) = \{\phi(t_1, x_0, \xi) \mid \xi \in U\} \subseteq \mathcal{R}(t_1, x_0),$$

that is, $\mathcal{R}(t_1, x_0)$ contains an open neighborhood of x_0 . \square

Example 4.29 1. For $\dot{x} = xu$, any x_0 is an equilibrium for $u = 0$. The linearization at x_0 is given by $A = 0$, $B = x_0$. Thus the linearization is controllable if and only if $x_0 \neq 0$. Therefore, the system is locally reachable from any $x_0 \neq 0$.

2. For $\dot{x}_1 = x_2$, $\dot{x}_2 = \sin(x_1) + u$, the equilibria are $x_0 = [k\pi, 0]^T$ for $k \in \mathbb{Z}$. The linearization is given by

$$A = \begin{bmatrix} 0 & 1 \\ (-1)^k & 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

whose Kalman matrix is

$$K = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

showing that the linearization is controllable for all x_0 . Thus the theorem says that the original system is locally reachable from the equilibria. In fact, we know that the system is even globally reachable from any x_0 .

3. For $\dot{x}_1 = x_2^2$, $\dot{x}_2 = u$, any $x_0 = [a, 0]^T$ for $a \in \mathbb{R}$ is an equilibrium. The linearization is

$$A = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

which is uncontrollable. Thus the theorem makes no statement for this example.

We have an analogous statement for controllability instead of reachability.

Corollary 4.30 In the situation of Theorem 4.28, if (A, B) is a controllable matrix pair, then $\mathcal{C}(t_1, x_0)$ contains an open neighborhood of x_0 for all $t_1 > 0$.

Proof: Let $t_1 > 0$ be given, consider

$$\dot{x}(t) = F(x(t), u(t)) \quad (4.18)$$

and set $\tilde{x}(t) := x(t_1 - t)$. Then

$$\dot{\tilde{x}}(t) = -\dot{x}(t_1 - t) = -F(x(t_1 - t), u(t_1 - t)) =: -F(\tilde{x}(t), \tilde{u}(t)),$$

where $\tilde{u}(t) := u(t_1 - t)$. Set $\tilde{F}(x, u) = -F(x, u)$ and consider

$$\dot{\tilde{x}}(t) = \tilde{F}(\tilde{x}(t), \tilde{u}(t)) \quad (4.19)$$

Let x_0 be an equilibrium of (4.18), that is, $F(x_0, 0) = 0$. Then x_0 is also an equilibrium of (4.19). Moreover, x_1 can be controlled to x_0 in time t_1 in (4.18) if and only if x_1 is reachable in time t_1 from x_0 in (4.19). Thus $\mathcal{C}(t_1, x_0) = \tilde{\mathcal{R}}(t_1, x_0)$. If (A, B) is the linearization of (4.18) at x_0 , then $(-A, -B)$ is the linearization of (4.19) at x_0 . If (A, B) is a controllable matrix pair, then so is $(-A, -B)$. Then $\tilde{\mathcal{R}}(t_1, x_0)$ and hence $\mathcal{C}(t_1, x_0)$ contains an open neighborhood of x_0 . \square

Remark 4.31 The sufficient condition for local reachability is often not satisfying, in particular, with systems of the form $\dot{x}(t) = g(x(t))u(t)$, where any x_0 is an equilibrium, but the linearization is $A = 0$ and $B = g(x_0)$ which is usually not controllable (unless $\text{rank}(g(x_0)) = n$). For instance, it can be shown that

$$\begin{aligned} \dot{x}_1 &= \cos(x_3)u_1 \\ \dot{x}_2 &= \sin(x_3)u_1 \\ \dot{x}_3 &= u_2 \end{aligned}$$

which is a simplistic model of driving a car ((x_1, x_2) represents the position, x_3 the angle into which the car is heading, u_1 is the speed, and u_2 represents the steering), is globally controllable (as would be expected from our everyday experience with driving a car), but we cannot derive this using the linearization criterion.

We state the following theorem without proof.

Theorem 4.32 Consider

$$\dot{x}(t) = f(x(t)) + g(x(t))u(t),$$

where $f : X \rightarrow \mathbb{R}^n$ and $g : X \rightarrow \mathbb{R}^{n \times m}$ are smooth for some open set $X \subseteq \mathbb{R}^n$. The accessibility algebra \mathfrak{A} is the smallest Lie algebra (with respect to the Lie bracket

$$[h, f] = \frac{\partial f}{\partial x}h - \frac{\partial h}{\partial x}f$$

for smooth $f, h : X \rightarrow \mathbb{R}^n$) that contains the columns of g and is invariant under f . If $\dim(\mathfrak{A}(x_0)) = n$, then for all $t_1 > 0$, the set $\mathcal{R}(t_1, x_0)$ has non-empty interior.

Local accessibility from x_0 (that is, $\mathcal{R}(x_0)$ has non-empty interior) is weaker than local reachability from x_0 . It means that $\mathcal{R}(x_0)$ contains some non-empty open set (but not necessarily a neighborhood of x_0). Thus, the notion of local accessibility includes the case where x_0 lies on the boundary of $\mathcal{R}(x_0)$.

Example 4.33 1. For $\dot{x} = xu$, we have $f \equiv 0$ and $g(x) = x$. Therefore, the accessibility algebra at x_0 is generated by x_0 and its dimension is $n = 1$ if and only if $x_0 \neq 0$.

2. For $\dot{x}_1 = x_2$, $\dot{x}_2 = \sin(x_1) + u$, we have

$$f(x_1, x_2) = \begin{bmatrix} x_2 \\ \sin(x_1) \end{bmatrix} \quad \text{and} \quad g(x_1, x_2) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

The accessibility algebra is generated by g and

$$[g, f] = \frac{\partial f}{\partial x}g - \frac{\partial g}{\partial x}f = \begin{bmatrix} 0 & 1 \\ \cos(x_1) & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

Thus its dimension is 2 at any point.

3. For $\dot{x}_1 = x_2^2$, $\dot{x}_2 = u$, one can show similarly that the dimension is $n = 2$ at every point.
4. For the simplistic model of driving a car from above, we obtain that the dimension of the accessibility algebra equals $n = 3$ for any x_0 .
5. For $\dot{x} = Ax + Bu$, the smallest Lie algebra that contains the columns of $g(x) = B$ and is invariant under $f(x) = Ax$ is nothing but the reachability space \mathcal{R} , and its dimension does not depend on the specific choice of x_0 .

For driftless systems, i.e., systems with $f \equiv 0$, the theorem can be strengthened as follows: If $\dim(\mathfrak{A}(x_0)) = n$, then $\mathcal{R}(t, x_0)$ contains an open neighborhood of x_0 for all $t > 0$, in particular, the system $\dot{x} = g(x)u$ is locally reachable from x_0 . Finally, if $\dim(\mathfrak{A}(x_0)) = n$ holds for all $x_0 \in X$, where X is connected, then $\mathcal{R}(x_0) = X$ for all $x_0 \in X$. For instance, the simple model of driving a car is globally reachable.

We conclude this chapter by studying what the unproven theorem from above implies for linear but time-varying systems

$$\dot{x}(t) = A(t)x(t) + B(t)u(t),$$

where $A : T \rightarrow \mathbb{R}^{n \times n}$ and $B : T \rightarrow \mathbb{R}^{n \times m}$ are smooth. Introducing $\xi = [x^T, t]^T$, such a system can be rewritten as

$$\dot{\xi} = \begin{bmatrix} A(t)x \\ 1 \end{bmatrix} + \begin{bmatrix} B(t) \\ 0 \end{bmatrix} u$$

and thus we have brought the system into the form $\dot{\xi} = f(\xi) + g(\xi)u$. It turns out that the accessibility algebra condition amounts to checking whether the infinite matrix (there is no Hamilton-Cayley like stopping criterion here)

$$K(t) = [B, (A - \frac{d}{dt})B, (A - \frac{d}{dt})^2 B, \dots](t),$$

has rank n at t_0 (the matrix does not depend on x_0). Of course, we recognize this as the time-varying generalization of the Kalman matrix (indeed, if A, B are constant, then we obtain the usual Kalman matrix from the time-invariant case, taking into account that it suffices to consider the powers of A up to A^{n-1}).

Theorem 4.34 Consider

$$\dot{x}(t) = A(t)x(t) + B(t)u(t)$$

where $A : T \rightarrow \mathbb{R}^{n \times n}$ and $B : T \rightarrow \mathbb{R}^{n \times m}$ are smooth. If the time-varying Kalman matrix has rank n for some $t_0 \in T$, then the system is completely reachable at time t_0 , that is, for all $x_0, x_1 \in \mathbb{R}^n$ and all $\varepsilon > 0$ there exists u such that

$$\varphi(\varepsilon + t_0, t_0, x_0, u) = x_1.$$

If A, B are even analytic functions, then the rank condition becomes necessary as well as sufficient. Moreover in that case, there exists an open and dense subset T_0 of T such that for all $t_0 \in T_0$, it suffices to check the rank of the finite matrix

$$[B, (A - \frac{d}{dt})B, \dots, (A - \frac{d}{dt})^{n-1}B](t_0).$$

Example 4.35 1. Consider $\dot{x}(t) = t^2 u(t)$, which can easily be seen to be completely reachable at any t_0 . We have $A(t) = 0$ and $B(t) = t^2$ and hence

$$K(t) = [t^2, -2t, 2, 0, \dots]$$

whose rank equals 1 for all t_0 . Note that for $t_0 = 0$, one has to compute $K(t_0)$ up to the third column (unlike the time-invariant case), but for almost all t_0 , it suffices to consider only the first column (à la Hamilton-Cayley).

2. Let

$$A(t) = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \quad \text{and} \quad B(t) = \begin{bmatrix} \cos(t) \\ -\sin(t) \end{bmatrix}.$$

Here, we find that the system is not reachable at any time t_0 . Nevertheless, the “snap shot” matrix pair $(A(t_0), B(t_0))$ is controllable for all t_0 .

Chapter 5

Feedback control

5.1 Static state feedback

The problem with pre-computed control (“open loop control”) is that it may be sensitive with respect to noisy data. For instance, let u be such that

$$\varphi(t, 0, x, u) = 0,$$

that is, the control function u steers the system from x to 0 in time t . If we start in a slightly perturbed initial state, say $x(0) = x + \xi$ with $\|\xi\|$ small but non-zero, we obtain, due to the linearity of φ ,

$$\varphi(t, 0, x(0), u) = \varphi(t, 0, \xi, 0)$$

where $\varphi(t, 0, \xi, 0) = e^{At}\xi$ in continuous time, and $\varphi(t, 0, \xi, 0) = A^t\xi$ in discrete time. If A is not stable, this deviation from 0 (the desired state) can become arbitrarily large. Thus pre-computed control will usually not work with unstable systems, in practice.

Feedback control is an alternative approach. Its basic assumption is that we can measure the state x , and that we can use this information for control. For a state space system

$$\dot{x} = Ax + Bu \quad \text{or} \quad \sigma x = Ax + Bu \tag{5.1}$$

where $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times m}$, a **static state feedback law** reads

$$u = Fx + v$$

where $F \in \mathbb{R}^{m \times n}$. Combining the given system with the feedback law yields the so-called **closed loop system**

$$\dot{x} = (A + BF)x + Bv \quad \text{or} \quad \sigma x = (A + BF)x + Bv. \quad (5.2)$$

We say that the matrix pair $(A + BF, B)$ results from (A, B) via static state feedback with feedback matrix F . The goal of this chapter is to answer questions like the **stabilization problem**: Given A, B , can we find F such that $A + BF$ is asymptotically stable? If yes, then all the trajectories of the closed loop system with $v = 0$ will tend to zero as $t \rightarrow \infty$. Thus the system is asymptotically controllable to zero, and this will work even if the initial state is subject to disturbance.

5.2 Feedback and controllability

Lemma 5.1 The reachable spaces of (5.1) and (5.2) coincide. In particular, $(A + BF, B)$ is controllable if and only if (A, B) is controllable.

Proof: The reachable space of $(A + BF, B)$ is

$$\mathcal{R} = \text{im}(B) + (A + BF)\text{im}(B) + \dots + (A + BF)^{n-1}\text{im}(B).$$

Let

$$\mathcal{R}_1 = \text{im}(B) \quad \text{and} \quad \mathcal{R}_{i+1} = (A + BF)\mathcal{R}_i + \text{im}(B) \quad \text{for } i = 1, \dots, n-1.$$

Then $\mathcal{R}_n = \mathcal{R}$. Since

$$\text{im}(B) + (A + BF)\mathcal{V} = \text{im}(B) + A\mathcal{V}$$

holds for any vector space $\mathcal{V} \subseteq \mathbb{R}^n$, we conclude that

$$\mathcal{R}_1 = \text{im}(B) \quad \text{and} \quad \mathcal{R}_{i+1} = A\mathcal{R}_i + \text{im}(B).$$

This shows that

$$\mathcal{R} = \mathcal{R}_n = \text{im}(B) + A\text{im}(B) + \dots + A^{n-1}\text{im}(B)$$

which is the reachable space of (A, B) . □

Lemma 5.2 (Heymann) Let (A, B) be controllable, and let $0 \neq b \in \mathbb{R}^m$. Then there exists a matrix F such that $(A + BF, b)$ is controllable.

Proof: First, we claim that there exist vectors $u^{(1)}, \dots, u^{(n-1)} \in \mathbb{R}^m$ such that the vectors $v^{(i)}$ defined by

$$v^{(1)} := b \quad \text{and} \quad v^{(i+1)} := Av^{(i)} + Bu^{(i)} \quad \text{for } i = 1, \dots, n-1$$

are a basis of \mathbb{R}^n . We prove this by induction. Since $v^{(1)} = b \neq 0$, we have a one-dimensional vector space

$$\mathcal{V}^{(1)} := \text{span}\{v^{(1)}\}.$$

Now assume that $v^{(1)}, \dots, v^{(k)}$ have already been constructed and that

$$\mathcal{V}^{(k)} = \text{span}\{v^{(1)}, \dots, v^{(k)}\}$$

has dimension k . We wish to choose $u^{(k)}$ such that

$$v^{(k+1)} = Av^{(k)} + Bu^{(k)} \notin \mathcal{V}^{(k)}.$$

Then $\dim(\mathcal{V}^{(k+1)}) = k + 1$ as desired. We need to show that such a choice is always possible for $k < n$. Assume conversely that

$$Av^{(k)} + Bu \in \mathcal{V}^{(k)} \quad \text{for all } u \in \mathbb{R}^m.$$

In particular,

$$Av^{(k)} \in \mathcal{V}^{(k)} \tag{5.3}$$

and thus

$$\text{im}(B) \subseteq \mathcal{V}^{(k)}. \tag{5.4}$$

On the other hand, $\mathcal{V}^{(k)}$ is A -invariant, because

$$A\mathcal{V}^{(k)} = \text{span}\{Av^{(1)}, \dots, Av^{(k)}\}.$$

To see that $Av^{(i)} \in \mathcal{V}^{(k)}$, we use (5.3) for the case $i = k$, and for $1 \leq i \leq k-1$, we have $Av^{(i)} = v^{(i+1)} - Bu^{(i)} \in \mathcal{V}^{(k)}$ due to (5.4). Thus $\mathcal{V}^{(k)}$ is an A -invariant subspace of \mathbb{R}^n that contains the image of B . Thus $\mathcal{V}^{(k)}$ contains the reachable space of (A, B) which is the *smallest* such space. Since (A, B) is controllable, this reachable space is all of \mathbb{R}^n , and hence $\mathcal{V}^{(k)} = \mathbb{R}^n$ which implies that $k = n$. Thus the construction works whenever $k < n$, and $\mathcal{V}^{(n)} = \text{span}\{v^{(1)}, \dots, v^{(n)}\} = \mathbb{R}^n$.

Now let $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be such that $Fv^{(i)} = u^{(i)}$ for $i = 1, \dots, n-1$. Then

$$(A + BF)v^{(i)} = Av^{(i)} + Bu^{(i)} = v^{(i+1)}.$$

Since $v^{(1)} = b$, this implies

$$v^{(i)} = (A + BF)^{i-1}b \quad \text{for } i = 1, \dots, n.$$

Thus the Kalman matrix of $(A + BF, b)$ is $K = [v^{(1)}, \dots, v^{(n)}]$ which has rank n by construction. \square

5.3 Pole placement

In this section, we investigate the characteristic polynomial and the spectrum of $A + BF$, where A, B are given, and F may be chosen. The goal is to place the eigenvalues of $A + BF$ in some desirable region of the complex plane (e.g., for obtaining asymptotic stability). This is known as pole placement (or: pole shifting, pole assignment).

Definition 5.3 Let $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times m}$ be given. Let $p \in \mathbb{R}[s]$ be a monic polynomial of degree n . We say that p is **assignable** to (A, B) if there exists a matrix $F \in \mathbb{R}^{m \times n}$ such that

$$\chi_{A+BF} = p,$$

that is, the characteristic polynomial of $A + BF$ equals p .

Theorem 5.4 Let $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times m}$. The matrix pair (A, B) is controllable if and only if every monic polynomial of degree n can be assigned to (A, B) .

Remark 5.5 This means that if (A, B) is controllable, then the eigenvalues of $A + BF$ can be shifted, by choice of the feedback matrix F , to any desired location in the complex plane. More precisely, let $\Lambda \subset \mathbb{C}$ be a non-empty set, with

$$\lambda \in \Lambda \quad \Rightarrow \quad \bar{\lambda} \in \Lambda.$$

For $\lambda \in \Lambda$, let $\mu(\lambda)$ be a positive integer with $\sum_{\lambda \in \Lambda} \mu(\lambda) = n$ and $\mu(\lambda) = \mu(\bar{\lambda})$. Then there exists a matrix F such that

$$\chi_{A+BF}(s) = \prod_{\lambda \in \Lambda} (s - \lambda)^{\mu(\lambda)}.$$

In particular, $\text{spec}(A + BF) = \Lambda$.

Recall that the Kalman controllability decomposition yields a factorization $\chi_A = \chi_{A_1} \cdot \chi_{A_3}$, where $\chi_u := \chi_{A_3}$ is the uncontrollable part of χ_A with respect to B . If (A, B) is controllable, one puts $A_1 := A$ and $\chi_u := 1$. Thus (A, B) is uncontrollable if and only if the degree of χ_u is at least one.

Proof: Since a similarity transform does not change the set of assignable polynomials, we may assume, without loss of generality, that a Kalman controllability decomposition has already been performed. Then

$$A+BF = \begin{bmatrix} A_1 & A_2 \\ 0 & A_3 \end{bmatrix} + \begin{bmatrix} B_1 \\ 0 \end{bmatrix} \begin{bmatrix} F_1 & F_2 \end{bmatrix} = \begin{bmatrix} A_1+B_1F_1 & A_2+B_1F_2 \\ 0 & A_3 \end{bmatrix} \quad (5.5)$$

and thus $\chi_{A+BF} = \chi_{A_1+B_1F_1} \cdot \chi_{A_3}$. This shows that χ_{A+BF} will always be a multiple of the polynomial $\chi_u = \chi_{A_3}$. If (A, B) is not controllable, then χ_u is not a constant, and thus not every polynomial can be assigned to (A, B) .

Conversely, let (A, B) be controllable. Let $p = s^n + p_{n-1}s^{n-1} + \dots + p_1s + p_0$ be given. We wish to construct F such that $\chi_{A+BF} = p$. If we choose any $0 \neq b \in \text{im}(B)$, then there exists \tilde{F} such that $(A+B\tilde{F}, b)$ is controllable, according to Heymann's Lemma. Thus we can transform this matrix pair into controller form, that is, there exists a non-singular matrix T such that

$$T^{-1}(A+B\tilde{F})T = \begin{bmatrix} 0 & 1 & & \\ \vdots & & \ddots & \\ 0 & & & 1 \\ -a_0 & -a_1 & \cdots & -a_{n-1} \end{bmatrix} \quad \text{and} \quad T^{-1}b = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

where a_i are the coefficients of $\chi_{A+B\tilde{F}}$. Now let

$$f = \begin{bmatrix} a_0 - p_0 & \dots & a_{n-1} - p_{n-1} \end{bmatrix} \in \mathbb{R}^{1 \times n}.$$

Then

$$T^{-1}(A+B\tilde{F})T + T^{-1}bf = \begin{bmatrix} 0 & 1 & & \\ \vdots & & \ddots & \\ 0 & & & 1 \\ -p_0 & -p_1 & \cdots & -p_{n-1} \end{bmatrix}$$

and therefore

$$\chi_{A+B\tilde{F}+bfT^{-1}} = s^n + p_{n-1}s^{n-1} + \dots + p_1s + p_0 = p.$$

Finally, we use that $b \in \text{im}(B)$, that is, $b = Bv$ for some $v \in \mathbb{R}^m$. Thus

$$A+B\tilde{F}+bfT^{-1} = A+B(\tilde{F}+vfT^{-1})$$

which yields the desired result putting $F := \tilde{F} + vfT^{-1}$. \square

As a direct consequence, we obtain the main theorem of this chapter.

Theorem 5.6 (Pole shifting theorem) The polynomials that can be assigned to (A, B) are precisely the ones of the form

$$p = p_1 \cdot \chi_u$$

where χ_u is the uncontrollable part of χ_A with respect to B , and p_1 is an arbitrary monic polynomial of degree $r = n - \deg(\chi_u)$.

Proof: It follows from (5.5) that

$$\chi_{A+BF} = \chi_{A_1+B_1F_1} \cdot \chi_u,$$

where $A_1 \in \mathbb{R}^{r \times r}$ and $B_1 \in \mathbb{R}^{r \times m}$ are the usual matrices from the Kalman controllability decomposition, and $F = [F_1, F_2]$. Thus any assignable polynomial must be a multiple of χ_u . The first factor can be any monic polynomial of degree r , because (A_1, B_1) is a controllable matrix pair. \square

Remark 5.7 The pole shifting theorem says that the uncontrollable modes of (A, B) cannot be influenced by static state feedback, whereas the other eigenvalues of the system can be moved to any desired location in the complex plane.

5.4 Stabilization

For stabilization, we don't require that the spectrum of $A + BF$ should coincide with some specific set of eigenvalues; we just want the eigenvalues to be contained in some given region of the complex plane. We put

$$\mathbb{C}_g = \{\lambda \in \mathbb{C} \mid \operatorname{Re}(\lambda) < 0\} \quad \text{or} \quad \mathbb{C}_g = \{\lambda \in \mathbb{C} \mid |\lambda| < 1\}$$

in the continuous or discrete case, respectively.

Definition 5.8 A matrix pair (A, B) is called **stabilizable** if there exists F such that $A + BF$ is asymptotically stable, that is, $\operatorname{spec}(A + BF) \subset \mathbb{C}_g$.

Theorem 5.9 The following are equivalent:

1. (A, B) is stabilizable;

2. (A, B) is asymptotically controllable to zero;
3. All uncontrollable modes of (A, B) lie in \mathbb{C}_g ;
4. Any eigenvalue λ of A which is not in \mathbb{C}_g satisfies $\text{rank} \begin{bmatrix} \lambda I - A & B \end{bmatrix} = n$.

Remark 5.10 Part of this theorem was already stated (without proof) in Section 4.3.

Proof: Conditions 3 and 4 are logically equivalent.

For “1 \Rightarrow 2”, let F be such that $A + BF$ is asymptotically stable. Put $u = Fx$, then the closed loop system reads

$$\dot{x} = (A + BF)x \quad \text{or} \quad \sigma x = (A + BF)x$$

respectively, and thus $\lim_{t \rightarrow \infty} x(t) = 0$ for all $x(0) = x_0$. An explicit formula for u is

$$u(t) = F e^{(A+BF)t} x_0 \quad \text{or} \quad u(t) = F(A + BF)^t x_0.$$

For “2 \Rightarrow 3”, we assume, without loss of generality, that a Kalman controllability decomposition has already been performed. Then the system law reads

$$\begin{aligned} \dot{x}_1 &= A_1 x_1 + A_2 x_2 + B_1 u & \text{or} & & \sigma x_1 &= A_1 x_1 + A_2 x_2 + B_1 u \\ \dot{x}_2 &= A_3 x_2 & & & \sigma x_2 &= A_3 x_2. \end{aligned}$$

By assumption, there exists, for any $x(0) = x_0$, an input function u such that $\lim_{t \rightarrow \infty} x(t) = 0$. Since $x_2(t) = e^{A_3 t} x_{02}$ or $x_2(t) = A_3^t x_{02}$, this can only be true if A_3 is asymptotically stable, which means that all the uncontrollable modes of (A, B) are asymptotically stable.

For “3 \Rightarrow 1”, choose a monic polynomial p_1 of degree $r = n - \deg(\chi_u)$ whose zeros are all in \mathbb{C}_g . By the pole shifting theorem, there exists F such that

$$\chi_{A+BF} = p_1 \cdot \chi_u.$$

By assumption, all zeros of χ_u lie in \mathbb{C}_g . Hence, all eigenvalues of $A + BF$ are contained in \mathbb{C}_g , that is, $A + BF$ is asymptotically stable. \square

5.5 Feedback equivalence

Consider $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times m}$, where (A, B) is a controllable matrix pair and $\text{rank}(B) = m \leq n$. Define the reachability spaces

$$\begin{aligned}\mathcal{R}_1 &= \text{im}(B) \\ \mathcal{R}_2 &= \text{im}[B, AB] = \text{im}(B) + A\text{im}(B) \\ &\vdots \\ \mathcal{R}_k &= \text{im}[B, AB, \dots, A^{k-1}B] = \text{im}(B) + A\text{im}(B) + \dots + A^{k-1}\text{im}(B).\end{aligned}$$

Clearly, we have $\mathcal{R}_1 \subseteq \mathcal{R}_2 \subseteq \dots \subseteq \mathbb{R}^n$ and, due to Hamilton-Cayley, $\mathcal{R}_n = \mathcal{R}_{n+l}$ for all $l \geq 0$. Moreover, because of the assumption of controllability, we have $\mathcal{R}_n = \mathbb{R}^n$. Define

$$d_i := \dim(\mathcal{R}_i)$$

for $1 \leq i \leq n$. Then $d_1 \leq d_2 \leq \dots \leq d_n$ and, according to our assumptions,

$$d_1 = m \quad \text{and} \quad d_n = n.$$

Next, we define a sequence of non-negative integers α_i for $1 \leq i \leq n$ by setting

$$\alpha_1 := d_1 \quad \text{and} \quad \alpha_i := d_i - d_{i-1} \text{ for } i > 1.$$

Lemma 5.11 We have $m = \alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_n$ and $\sum_{i=1}^n \alpha_i = d_n = n$.

Proof: Set $\mathcal{R}_0 := \{0\}$ and $d_0 := 0$. Then we have

$$\alpha_i = d_i - d_{i-1} = \dim(\mathcal{R}_i) - \dim(\mathcal{R}_{i-1}) = \dim(\mathcal{R}_i/\mathcal{R}_{i-1})$$

for all $1 \leq i \leq n$. Since $\mathcal{R}_{i+1} = A\mathcal{R}_i + \text{im}(B)$, we obtain

$$\alpha_{i+1} = \dim(\mathcal{R}_{i+1}/\mathcal{R}_i) = \dim\left(\frac{A\mathcal{R}_i + \text{im}(B)}{A\mathcal{R}_i + \text{im}(B)}\right) \leq \dim\left(\frac{A\mathcal{R}_i}{A\mathcal{R}_i}\right),$$

where we have used the dimension formula

$$\begin{aligned}\dim\left(\frac{U+W}{V+W}\right) &= \dim(U+W) - \dim(V+W) \\ &= \dim(U) - \dim(U \cap W) - \dim(V) + \dim(V \cap W) \\ &\leq \dim(U) - \dim(V) = \dim(U/V),\end{aligned}$$

which holds for any subspaces $V \subseteq U$ and W of \mathbb{R}^n . Thus

$$\alpha_{i+1} \leq \dim\left(\frac{A\mathcal{R}_i}{A\mathcal{R}_i}\right) \leq \dim(\mathcal{R}_i/\mathcal{R}_{i-1}) = \alpha_i,$$

since

$$\begin{aligned} \dim\left(\frac{AU}{AV}\right) &= \dim(AU) - \dim(AV) \\ &= \dim(U) - \dim(U \cap \ker(A)) - \dim(V) + \dim(V \cap \ker(A)) \\ &\leq \dim(U) - \dim(V) = \dim(U/V) \end{aligned}$$

holds for any subspaces $V \subseteq U \subseteq \mathbb{R}^n$. Finally,

$$\sum_{i=1}^n \alpha_i = \sum_{i=1}^n (d_i - d_{i-1}) = d_n - d_0 = d_n = n.$$

□

Definition 5.12 The **controllability indices** $\kappa_1, \dots, \kappa_m$ of (A, B) are defined by

$$\kappa_j := |\{i \mid \alpha_i \geq j\}|.$$

Remark 5.13 The controllability index κ_1 is characterized by

$$\kappa_1 = |\{i \mid \alpha_i \geq 1\}|.$$

Thus we have $\alpha_i = 0$ for all $i > \kappa_1$. In other words, we have $d_i = d_{i-1}$ for all $i > \kappa_1$ and thus $d_{\kappa_1} = d_{\kappa_1+l}$ for all $l \geq 0$. This shows that κ_1 is the index at which

$$\{0\} = \mathcal{R}_0 \subseteq \mathcal{R}_1 \subseteq \mathcal{R}_2 \subseteq \dots \subseteq \mathbb{R}^n$$

becomes stationary, that is, $\mathcal{R}_{\kappa_1} = \mathcal{R}_{\kappa_1+l}$ for all $l \geq 0$. Due to the assumption of controllability, κ_1 is therefore the smallest integer with $\mathcal{R}_{\kappa_1} = \mathbb{R}^n$. Still equivalently, we have

$$\kappa_1 = \min\{k \mid \text{rank}[B, \dots, A^{k-1}B] = n\}.$$

Sometimes, κ_1 is called “the” controllability index of (A, B) .

Example 5.14 Let $n = 5$ and $m = 3$. Then there are two possibilities for $\alpha = (\alpha_1, \dots, \alpha_5)$, $d = (d_1, \dots, d_5)$, and $\kappa = (\kappa_1, \kappa_2, \kappa_3)$.

Case 1: $\alpha = (3, 2, 0, 0, 0)$, that is, $d = (3, 5, 5, 5, 5)$. Then $\kappa = (2, 2, 1)$.

Case 2: $\alpha = (3, 1, 1, 0, 0)$, that is, $d = (3, 4, 5, 5, 5)$. Then $\kappa = (3, 1, 1)$.

We note that the two outcomes for κ correspond to the two ways in which 5 can be written as a sum of 3 positive integers (disregarding the order). This is due to the following fact.

Lemma 5.15 We have $\kappa_1 \geq \dots \geq \kappa_m \geq 1$ and $\sum_{j=1}^m \kappa_j = n$.

Proof: Since $\alpha_i \geq j$ implies $\alpha_i \geq j - 1$, we clearly have $\kappa_j \leq \kappa_{j-1}$. More precisely,

$$\kappa_{j-1} = \kappa_j + |\{i \mid \alpha_i = j - 1\}|.$$

By assumption, $\alpha_1 = m$. Hence $\kappa_m \geq 1$. We have

$$\sum_{j=1}^m \kappa_j = \sum_{j=1}^m |\{i \mid \alpha_i \geq j\}| = \sum_{j=1}^m j \cdot |\{i \mid \alpha_i = j\}| = \sum_{i=1}^n \alpha_i = n.$$

□

The importance of the controllability indices lies in the fact that they form a complete set of invariants under an equivalence relation of state space systems with the properties mentioned at the beginning of this section. This equivalence relation will be introduced next. For this, let

$$\mathcal{S} = \{(A, B) \in \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times m} \mid (A, B) \text{ controllable and } \text{rank}(B) = m\}.$$

Definition 5.16 Two elements (A, B) , (\tilde{A}, \tilde{B}) of \mathcal{S} are called **feedback equivalent** if there exist non-singular matrices $T \in \mathbb{R}^{n \times n}$ and $G \in \mathbb{R}^{m \times m}$, and a matrix $F \in \mathbb{R}^{m \times n}$ such that

$$(\tilde{A}, \tilde{B}) = (T^{-1}(A + BF)T, T^{-1}BG).$$

The obvious interpretation of feedback equivalence is that a state space system $\dot{x} = Ax + Bu$ is combined with a feedback law $u = Fx + Gv$, yielding $\dot{x} = (A + BF)x + BGv$, and transformed into

$$\dot{\tilde{x}} = T^{-1}(A + BF)T\tilde{x} + T^{-1}BGv = \tilde{A}\tilde{x} + \tilde{B}v$$

via a similarity transform (corresponding to the coordinate change $\tilde{x} = T^{-1}x$ in the state space). It is easy to check that feedback equivalence is indeed an equivalence relation on \mathcal{S} .

Now consider the columns of the Kalman controllability matrix from left to right:

$$K = [B, AB, \dots, A^{n-1}B].$$

When adding the block $A^{i-1}B$, the dimension of the column space of K increases by α_i . The number κ_m tells us how many times the dimension goes up by m . Thus, the columns

$$b_1, \dots, b_m, \dots, A^{\kappa_m-1}b_1, \dots, A^{\kappa_m-1}b_m$$

are all linearly independent. When adding the next $\kappa_{m-1} - \kappa_m$ blocks, the dimension increases by $m - 1$. Thus, in $A^{\kappa_m} B$, one column is superfluous (for spanning the column space of K). Without loss of generality, assume that it is the column $A^{\kappa_m} b_m$ (otherwise, permute the columns of B correspondingly). This means that also this holds also in the remaining blocks, that is, the columns

$$\begin{aligned} & b_1, \dots, b_m, \dots, A^{\kappa_{m-1}} b_1, \dots, A^{\kappa_{m-1}} b_m, \\ & A^{\kappa_m} b_1, \dots, A^{\kappa_m} b_{m-1}, \dots, A^{\kappa_{m-1}-1} b_1, \dots, A^{\kappa_{m-1}-1} b_{m-1} \end{aligned}$$

are linearly independent. Proceeding like this, we find that the following columns of K must be linearly independent:

$$\begin{aligned} & b_1, \dots, b_m, \dots, A^{\kappa_{m-1}} b_1, \dots, A^{\kappa_{m-1}} b_m, \\ & A^{\kappa_m} b_1, \dots, A^{\kappa_m} b_{m-1}, \dots, A^{\kappa_{m-1}-1} b_1, \dots, A^{\kappa_{m-1}-1} b_{m-1}, \\ & \quad \vdots \\ & A^{\kappa_3} b_1, A^{\kappa_3} b_2, \dots, A^{\kappa_2-1} b_1, A^{\kappa_2-1} b_2, \\ & A^{\kappa_2} b_1, \dots, A^{\kappa_1-1} b_1. \end{aligned} \tag{5.6}$$

These are

$$m\kappa_m + (m-1)(\kappa_{m-1} - \kappa_m) + \dots + 2(\kappa_2 - \kappa_3) + (\kappa_1 - \kappa_2) = \sum_{j=1}^m j \cdot |\{i \mid \alpha_i = j\}| = n$$

columns, in accordance with the assumption that K has rank n .

We conclude that we have (after a suitable permutation of the columns of B) for all $1 \leq j \leq m$

$$\kappa_j = \min\{k \mid A^k b_j \in \text{im}[B, \dots, A^{k-1} B, A^k b_1, \dots, A^k b_{j-1}]\}.$$

In particular, there exist coefficients $\lambda_{jk} \in \mathbb{R}$ such that

$$A^{\kappa_j} (b_j - \sum_{k=1}^{j-1} \lambda_{jk} b_k) \in \text{im}[B, \dots, A^{\kappa_j-1} B].$$

Define $\tilde{b}_j := b_j - \sum_{k=1}^{j-1} \lambda_{jk} b_k$. Due to the structure of this transformation of the basis of $\text{im}(B)$, the column vectors listed in (5.6) remain linearly independent if each b_j is replaced by \tilde{b}_j . Thus we can drop the tildes, without loss of generality.

Lemma 5.17 Let $b \in \text{im}(B)$ and let $k \geq 1$ be the smallest integer such that

$$A^k b \in \text{im}[B, AB, \dots, A^{k-1} B].$$

Then there exists F such that $(A+BF)^i b$ for $0 \leq i \leq k-1$ are linearly independent and $(A+BF)^k b = 0$.

The proof uses an argument similar to Heymann's Lemma and will be done as an exercise.

Definition 5.18 A subspace $\mathcal{V} \subseteq \mathbb{R}^n$ is called a **controllability subspace** of (A, B) if there exist F and G (not necessarily invertible) such that

$$\mathcal{V} = \mathcal{R}(A + BF, BG),$$

where $\mathcal{R}(A, B) := \text{im}[B, AB, \dots, A^{n-1}B]$.

Theorem 5.19 Let $(A, B) \in \mathcal{S}$ and let κ_j be the controllability indices of (A, B) . Then

$$\mathbb{R}^n = \bigoplus_{j=1}^m \mathcal{V}_j$$

where each \mathcal{V}_j is a controllability subspace of (A, B) of dimension κ_j .

Proof: As we have seen, there exist $b_j \in \text{im}(B)$ such that

$$\kappa_j = \min\{k \mid A^k b_j \in \text{im}[B, AB, \dots, A^{k-1}B]\}. \quad (5.7)$$

Thus there exist F_j such that $(A + BF_j)^i b_j$ for $0 \leq i \leq \kappa_j - 1$ are linearly independent and $(A + BF_j)^{\kappa_j} b_j = 0$. Define

$$\mathcal{V}_j := \mathcal{R}(A + BF_j, b_j).$$

Then \mathcal{V}_j is a controllability subspace of (A, B) of dimension κ_j . We have

$$\begin{aligned} \sum_{j=1}^m \mathcal{V}_j &= \text{span}\{b_1, \dots, (A + BF_1)^{\kappa_1-1} b_1, \dots, b_m, \dots, (A + BF_m)^{\kappa_m-1} b_m\} \\ &\supseteq \text{span}\{b_1, \dots, A^{\kappa_1-1} b_1, \dots, b_m, \dots, A^{\kappa_m-1} b_m\}. \end{aligned}$$

However, the space on the right hand side is all of \mathbb{R}^n due to (5.6). Thus a dimensional argument shows that the sum must be direct. \square

Corollary 5.20 Let $(A, B) \in \mathcal{S}$ and let κ_j be its controllability indices. Then (A, B) is feedback equivalent to a system of the form $\tilde{A} = \text{diag}(A_1, \dots, A_m)$, $\tilde{B} = \text{diag}(B_1, \dots, B_m)$ with

$$A_j = \begin{bmatrix} 0 & \dots & \dots & 0 \\ 1 & \ddots & & \vdots \\ & \ddots & \ddots & \vdots \\ & & 1 & 0 \end{bmatrix} \in \mathbb{R}^{\kappa_j \times \kappa_j} \quad \text{and} \quad B_j = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{\kappa_j}.$$

This is called the **Brunovsky form** of (A, B) . Two matrix pairs (A, B) , (\tilde{A}, \tilde{B}) are feedback equivalent if and only if they have the same Brunovsky form, i.e., the same controllability indices.

Remark 5.21 The Brunovsky form generalizes the controllability form from single- to multi-input systems, incorporating the effect of a feedback. There exists also an alternative version with

$$A_j = \begin{bmatrix} 0 & 1 & & \\ \vdots & \ddots & \ddots & \\ \vdots & & \ddots & 1 \\ 0 & \dots & \dots & 0 \end{bmatrix} \in \mathbb{R}^{\kappa_j \times \kappa_j} \quad \text{and} \quad B_j = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \in \mathbb{R}^{\kappa_j}$$

which corresponds to the controller form.

Proof: There exists an invertible matrix G such that the columns b_j of BG satisfy (5.7). Thus there exist F_j such that

$$\mathbb{R}^n = \bigoplus_{j=1}^m \mathcal{V}_j \quad \text{with} \quad \mathcal{V}_j = \mathcal{R}(A + BF_j, b_j),$$

where $(A + BF_j)^i b_j$ for $0 \leq i \leq \kappa_j - 1$ are linearly independent and $(A + BF_j)^{\kappa_j} b_j = 0$. Thus there exists $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $F|_{\mathcal{V}_j} = F_j|_{\mathcal{V}_j}$. Then

$$\mathcal{V}_j = \mathcal{R}(A + BF_j, b_j) = \mathcal{R}(A + BF, b_j).$$

Let $x_1^{(j)}, \dots, x_{\kappa_j}^{(j)}$ be a basis of \mathcal{V}_j , set $T_j := [x_1^{(j)}, \dots, x_{\kappa_j}^{(j)}] \in \mathbb{R}^{n \times \kappa_j}$, and $T := [T_1, \dots, T_m] \in \mathbb{R}^{n \times n}$, which is clearly non-singular. Then

$$(A + BF)T = (A + BF)[T_1, \dots, T_m] = [(A + BF_1)T_1, \dots, (A + BF_m)T_m]$$

Since each \mathcal{V}_j is $(A + BF_j)$ -invariant, there exist matrices A_j such that

$$(A + BF)T = [T_1 A_1, \dots, T_m A_m] = T \text{diag}(A_1, \dots, A_m).$$

Since $b_j \in \mathcal{V}_j$, we have $b_j = T_j B_j$ for some B_j and thus

$$BG = T \text{diag}(B_1, \dots, B_m).$$

Finally, using the special basis $x_1^{(j)} = b_j$, $x_2^{(j)} = (A + BF_j)b_j$ etc. of \mathcal{V}_j , we find that each b_j is the first column of the corresponding T_j , that is,

$$b_j = T_j \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

and that $(A + BF)x_k^{(j)} = (A + BF_j)x_k^{(j)} = x_{k+1}^{(j)}$ for $1 \leq k \leq \kappa_j - 1$, whereas $(A + BF)x_{\kappa_j}^{(j)} = 0$, that is,

$$A_j = \begin{bmatrix} 0 & \dots & \dots & 0 \\ 1 & \ddots & & \vdots \\ & \ddots & \ddots & \vdots \\ & & & 1 & 0 \end{bmatrix}.$$

□

5.6 Stabilization and feedback equivalence for non-linear systems

Consider

$$\dot{x}(t) = F(x(t), u(t))$$

as in Section 4.5. Let x_0 be an equilibrium of the system for the zero input function, that is, $F(x_0, 0) = 0$. One calls x_0 a **stable** equilibrium (compare with the end of Section 3.1) if

$$\forall \varepsilon > 0 \exists \delta > 0 : \|x_0 - \tilde{x}_0\| \leq \delta \quad \Rightarrow \quad \|x(t) - \tilde{x}(t)\| \leq \varepsilon \quad \text{for all } t > 0,$$

where $x(t) = \varphi(t, 0, x_0, 0)$ is the solution at time t of $\dot{x} = F(x, u)$ with $x(0) = x_0$ and $u \equiv 0$ (and thus $x(t) = x_0$ for all t , since x_0 is an equilibrium), and $\tilde{x}(t) = \varphi(t, 0, \tilde{x}_0, 0)$. One calls x_0 **asymptotically stable** if additionally, we have

$$\exists \gamma > 0 : \lim_{t \rightarrow \infty} \|x(t) - \tilde{x}(t)\| = 0 \quad \text{for all } \tilde{x}_0 \text{ with } \|x_0 - \tilde{x}_0\| \leq \gamma.$$

Let

$$A := \frac{\partial F}{\partial x}(x_0, 0) \in \mathbb{R}^{n \times n}.$$

An important stability criterion from classical ODE theory says that: If A is asymptotically stable (i.e., all its eigenvalues have a negative real part), then x_0 is an asymptotically stable equilibrium (of the underlying non-linear system $\dot{x} = F(x, 0)$). If A has an eigenvalue with positive real part, then x_0 is an unstable equilibrium.

The non-linear stabilization problem can be posed as follows: Find a \mathcal{C}^1 -function $\alpha : X \rightarrow U$ such that with the feedback law $u(t) = \alpha(x(t))$, the resulting closed loop system

$$\dot{x}(t) = F(x(t), \alpha(x(t)))$$

has x_0 as an asymptotically stable equilibrium.

Remark 5.22 To guarantee that x_0 , which is an equilibrium of $\dot{x} = F(x, 0)$ by assumption, becomes an equilibrium of $\dot{x} = F(x, \alpha(x))$ as well, we admit only feedback functions α with $\alpha(x_0) = 0$. Then $F(x_0, \alpha(x_0)) = F(x_0, 0) = 0$.

Definition 5.23 The non-linear system $\dot{x} = F(x, u)$ is called **stabilizable at x_0** if the stabilization problem is solvable.

Theorem 5.24 Consider $\dot{x} = F(x, u)$ with $F(x_0, 0) = 0$. Set

$$A = \frac{\partial F}{\partial x}(x_0, 0) \quad \text{and} \quad B = \frac{\partial F}{\partial u}(x_0, 0).$$

If (A, B) is a stabilizable matrix pair, then $\dot{x} = F(x, u)$ is stabilizable at x_0 . If (A, B) has an uncontrollable mode with positive real part, then $\dot{x} = F(x, u)$ is not stabilizable at x_0 .

Proof: Consider the closed loop system

$$\dot{x} = F(x, \alpha(x))$$

with its equilibrium x_0 . Computing its linearization at x_0 , we get

$$\tilde{A} = \frac{\partial F}{\partial x}(x_0, \alpha(x_0)) + \frac{\partial F}{\partial u}(x_0, \alpha(x_0)) \frac{\partial \alpha}{\partial x}(x_0) = A + B \frac{\partial \alpha}{\partial x}(x_0).$$

If (A, B) is stabilizable, then there exists F_1 such that $A + BF_1$ is asymptotically stable. Let α be such that

$$\frac{\partial \alpha}{\partial x}(x_0) = F_1$$

and $\alpha(x_0) = 0$. Such an α exists always, e.g., $\alpha(x) = F_1(x - x_0)$. Then α solves the stabilization problem. On the other hand, if (A, B) has an uncontrollable mode λ with positive real part, then this λ belongs to the spectrum of every matrix of the form $A + BF_1$, that is, \tilde{A} has an eigenvalue with positive real part for every choice of $F_1 = \frac{\partial \alpha}{\partial x}(x_0)$. This shows that x_0 is an unstable equilibrium of $\dot{x} = F(x, \alpha(x))$ for all admissible α , and thus the stabilization problem is not solvable. \square

Remark 5.25 We have just seen that if the linearization at x_0 is stabilizable, then the non-linear system can be stabilized at x_0 . Moreover, the proof shows that this can be done using an affine feedback law $u(t) = F_1(x(t) - x_0)$, where F_1 is chosen such that it stabilizes the linearization.

Note that the theorem makes no statement in the case where all uncontrollable modes of (A, B) have non-positive real part, but some of them lie on the imaginary axis. This is due to the fact that the mentioned criterion from ODE theory does also not work in this case.

Example 5.26 Consider the pendulum equations $\dot{x}_1 = x_2$, $\dot{x}_2 = \sin(x_1) + u$ with its equilibrium $x_0 = [0, 0]^T$. The linearization is given by

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

Since $\text{spec}(A) = \{\pm 1\}$, the equilibrium is unstable. However, the linearization is controllable and thus stabilizable. Any $F_1 = [f_1, f_2]$ with $f_1 < -1$ and $f_2 < 0$ stabilizes (A, B) . Thus any such $u = F_1 x$ stabilizes also the non-linear pendulum.

Now consider a non-linear system of the form

$$\dot{x}(t) = f(x(t)) + g(x(t))u(t)$$

as in Theorem 4.32. Consider a feedback law

$$u = \alpha(x) + \beta(x)v,$$

where $\alpha : X \rightarrow \mathbb{R}^m$ and $\beta : X \rightarrow \mathbb{R}^{m \times m}$ are smooth. Then the closed loop system is

$$\dot{x}(t) = (f + g\alpha)(x(t)) + (g\beta)(x(t))v(t).$$

Of course, one would want the new input v to be recomputable from u , at least locally, near some point x_0 of interest, that is,

$$\det(\beta(x_0)) \neq 0.$$

(This corresponds to “ G invertible” with linear feedback equivalence.) The role of the linear transformation matrix T^{-1} is taken over by a non-linear state transformation $\Phi : X \rightarrow \mathbb{R}^n$ that is smooth and locally invertible near x_0 , that is, using the inverse function theorem,

$$\det\left(\frac{\partial \Phi}{\partial x}(x_0)\right) \neq 0.$$

The systems

$$\dot{x} = f(x) + g(x)u$$

and

$$\dot{\xi} = \tilde{f}(\xi) + \tilde{g}(\xi)v$$

are called **feedback equivalent at x_0** if there exist maps α, β, Φ as described above such that

$$\tilde{f}(\xi) = \frac{\partial \Phi}{\partial x}(\Phi^{-1}(\xi))(f + g\alpha)(\Phi^{-1}(\xi))$$

and

$$\tilde{g}(\xi) = \frac{\partial \Phi}{\partial x}(\Phi^{-1}(\xi))(g\beta)(\Phi^{-1}(\xi)).$$

Example 5.27 Consider

$$\begin{aligned}\dot{x}_1 &= \cos(x_3)x_4 \\ \dot{x}_2 &= \sin(x_3)x_4 \\ \dot{x}_3 &= u_2 \\ \dot{x}_4 &= u_1\end{aligned}$$

which is a variant of the model for driving a car discussed in Remark 4.31, with the modification that u_2 is the acceleration here, and not the velocity. Putting $\xi_1 = x_1$, $\xi_2 = \cos(x_3)x_4$, $\xi_3 = x_2$, $\xi_4 = \sin(x_3)x_4$ (which is an admissible state transformation at every $x_0 \in \mathbb{R}^4$ with $x_{04} \neq 0$) and $v_1 = \cos(x_3)u_1 - x_4 \sin(x_3)u_2$, $v_2 = \sin(x_3)u_1 + x_4 \cos(x_3)u_2$ (which is also admissible at every $x_0 \in \mathbb{R}^4$ with $x_{04} \neq 0$), this system is feedback equivalent to

$$\begin{aligned}\dot{\xi}_1 &= \xi_2 \\ \dot{\xi}_2 &= v_1 \\ \dot{\xi}_3 &= \xi_4 \\ \dot{\xi}_4 &= v_2\end{aligned}$$

at every $x_0 \in \mathbb{R}^4$ with $x_{04} \neq 0$. We notice that the new system is linear and in Brunovsky form with controllability indices $(2, 2)$.

Remark 5.28 The question under which conditions a non-linear control system is (locally) feedback equivalent to a linear system is of fundamental importance in non-linear control theory. It is well studied, and there is also a non-linear theory of controllability indices etc.

5.7 Control as interconnection

Feedback control is based on the interconnection of systems: Given a dynamical system (the “plant”), the goal of feedback control is to design another system (a

“controller”), in a way that the interconnection of the two systems has certain desired properties. As an example, consider a plant given in classical state space form

$$\dot{x} = Ax + Bu$$

and let the controller be specified by the feedback law $u = Fx + v$. Then the interconnection (the “closed loop” system) is

$$\dot{x} = (A + BF)x + Bv. \quad (5.8)$$

A typical aim of the controller design in this setting is spectral assignment, that is, a condition is given on the desired location of the eigenvalues of $A + BF$. Note that interconnection means nothing but combining the equations that determine plant and controller, respectively, and to look at their common solutions. For instance, combining the plant given by

$$\begin{bmatrix} \frac{d}{dt}I - A & -B & 0 \end{bmatrix} \begin{bmatrix} x \\ u \\ v \end{bmatrix} = 0$$

with the controller given by

$$\begin{bmatrix} F & -I & I \end{bmatrix} \begin{bmatrix} x \\ u \\ v \end{bmatrix} = 0$$

yields the interconnected system

$$\begin{bmatrix} \frac{d}{dt}I - A & -B & 0 \\ F & -I & I \end{bmatrix} \begin{bmatrix} x \\ u \\ v \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

from which we may eliminate u to get (5.8). It is worth noting that this interconnection is regular in a sense to be defined below.

The behavioral approach to systems theory provides an elegant framework for dealing with such, and more general, interconnection problems. Suppose that the plant is given by $R_1(\frac{d}{dt})w = 0$. Similarly, let the controller be determined by $R_2(\frac{d}{dt})w = 0$. Then the interconnection is determined by the system law

$$\begin{bmatrix} R_1 \\ R_2 \end{bmatrix} (\frac{d}{dt})w = 0.$$

Let \mathcal{B}_i be the plant and the controller,

$$\mathcal{B}_i = \{w \in \mathcal{A}^q \mid R_i(\frac{d}{dt})w = 0\}$$

respectively ($i = 1, 2$). Then their **interconnection** is defined as

$$\mathcal{B} = \mathcal{B}_1 \cap \mathcal{B}_2 = \{w \in \mathcal{A}^q \mid R(\frac{d}{dt})w = 0\},$$

where

$$R = \begin{bmatrix} R_1 \\ R_2 \end{bmatrix}.$$

The interconnection is said to be **regular** if

$$\text{rank}(R_1) + \text{rank}(R_2) = \text{rank}(R).$$

This means that the system laws of plant and controller are independent of each other. Equivalently, the interconnection is regular if $\mathcal{B}_1 + \mathcal{B}_2 = \mathcal{A}^q$.

Let \mathcal{B}_1 be a given plant, and let \mathcal{B} be some desired behavior. The goal is to find a controller \mathcal{B}_2 such that the interconnection of plant and controller equals this desired behavior, i.e.,

$$\mathcal{B}_1 \cap \mathcal{B}_2 = \mathcal{B}.$$

Obviously, $\mathcal{B} \subseteq \mathcal{B}_1$ is a necessary condition. If it is satisfied, we say that \mathcal{B} is a **subsystem** of \mathcal{B}_1 . Then the problem is always solvable, because there is the trivial solution $\mathcal{B}_2 = \mathcal{B}$. The question becomes more interesting if we require the interconnection to be regular. Then we say that **\mathcal{B} can be achieved from \mathcal{B}_1 by regular interconnection**.

Theorem 5.29 Let \mathcal{B}_1 be controllable. Then every subsystem of \mathcal{B}_1 can be achieved from \mathcal{B}_1 by regular interconnection.

Proof: We have seen in the proof of Theorem 4.25 that a controllable behavior \mathcal{B}_1 is isomorphic to $\{0\} \times \mathcal{A}^m$, where m is the number of inputs of \mathcal{B}_1 . Therefore, there exists a behavior \mathcal{B}'_1 such that $\mathcal{B}_1 \oplus \mathcal{B}'_1 = \mathcal{A}^q$. Now let $\mathcal{B} \subseteq \mathcal{B}_1$ be given. We set $\mathcal{B}_2 := \mathcal{B} + \mathcal{B}'_1$. Then

$$\mathcal{B}_1 \cap \mathcal{B}_2 = \mathcal{B}_1 \cap (\mathcal{B} + \mathcal{B}'_1) = \mathcal{B}.$$

To see the last equality, note that the inclusion “ \supseteq ” is obvious, since \mathcal{B} is contained both in \mathcal{B}_1 and in $\mathcal{B} + \mathcal{B}'_1$. For the converse inclusion, let $w_1 \in \mathcal{B}_1$ have a decomposition $w_1 = w + w'_1$ with $w \in \mathcal{B} \subseteq \mathcal{B}_1$ and $w'_1 \in \mathcal{B}'_1$. Then $w'_1 = w_1 - w \in \mathcal{B}_1 \cap \mathcal{B}'_1 = \{0\}$. Thus $w'_1 = 0$ and $w_1 = w \in \mathcal{B}$ as desired. Finally,

$$\mathcal{B}_1 + \mathcal{B}_2 = \mathcal{B}_1 + \mathcal{B} + \mathcal{B}'_1 = \mathcal{A}^q$$

from which it follows that the interconnection is regular. \square

Series connection: Let $P_i(\frac{d}{dt})y_i = Q_i(\frac{d}{dt})u_i$, where $i = 1, 2$, be two input-output representations. The series connection is defined by taking the output of the first system as the input of the second system. Of course, this is only possible if the dimensions match, that is, $p_1 = m_2$, which we assume. We set $u = u_1$, $y_1 = u_2$, and $y = y_2$. The interconnection is therefore described by

$$\begin{bmatrix} -Q_1 & P_1 & 0 \\ 0 & -Q_2 & P_2 \end{bmatrix} \begin{bmatrix} u \\ y_1 \\ y \end{bmatrix} = 0.$$

Since P_i are both square and non-singular, this is a regular interconnection.

Parallel connection: Let $P_i(\frac{d}{dt})y_i = Q_i(\frac{d}{dt})u_i$, where $i = 1, 2$, be two input-output representations. The parallel connection is defined by giving the same input to both systems, and by summing the outputs. This is only possible if $m := m_1 = m_2$ and $p := p_1 = p_2$. We put $u = u_1 = u_2$ and $y = y_1 + y_2$. The interconnection is given by

$$\begin{bmatrix} -Q_1 & P_1 & 0 & 0 \\ -Q_2 & 0 & P_2 & 0 \\ 0 & I & I & -I \end{bmatrix} \begin{bmatrix} u \\ y_1 \\ y_2 \\ y \end{bmatrix} = 0,$$

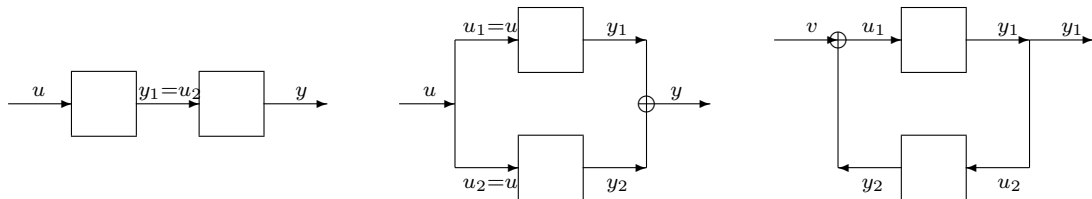
which is actually a regular interconnection of three systems (the two given systems and the summing system).

Feedback connection: Let $P_i(\frac{d}{dt})y_i = Q_i(\frac{d}{dt})u_i$, where $i = 1, 2$, be two input-output representations. The feedback connection is defined by taking the output of the first system as the input of the second, and vice versa, that is, $y_1 = u_2$ and $y_2 = u_1$. The connection is given by

$$\begin{bmatrix} P_1 & -Q_1 \\ -Q_2 & P_2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = 0.$$

This interconnection is not necessarily regular. It becomes regular if we modify the equation $u_1 = y_2$, using a new input v , say $u_1 = y_2 + v$.

The three standard interconnections are illustrated in the following figure.



Chapter 6

Observability

6.1 Basic notions for state space systems

Consider the state space equations

$$\begin{aligned} \dot{x} &= Ax + Bu & \text{or} & & \sigma x &= Ax + Bu \\ y &= Cx + Du & & & y &= Cx + Du, \end{aligned}$$

where $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{p \times n}$ and $D \in \mathbb{R}^{p \times m}$. Let $X = \mathbb{R}^n$ and $Y = \mathbb{R}^p$.

The state transition map

$$\varphi : \{(t, t_0) \in T^2 \mid t \geq t_0\} \times X \times \mathcal{U} \rightarrow X, \quad (t, t_0, x_0, u) \mapsto \varphi(t, t_0, x_0, u)$$

gives the state at time t if the state at time t_0 was x_0 and the control function u was applied. We define the **state-to-output map**

$$\eta : \{(t, t_0) \in T^2 \mid t \geq t_0\} \times X \times \mathcal{U} \rightarrow Y, \quad (t, t_0, x_0, u) \mapsto \eta(t, t_0, x_0, u)$$

as the map which gives the output at time t if the state at time t_0 was x_0 and the control function was u . For state space systems, we have $\eta(t, t_0, x_0, u) = C\varphi(t, t_0, x_0, u) + Du(t)$, and thus

$$\eta(t, t_0, x_0, u) = Ce^{A(t-t_0)}x_0 + \int_{t_0}^t Ce^{A(t-\tau)}Bu(\tau)d\tau + Du(t) \quad (6.1)$$

and

$$\eta(t, t_0, x_0, u) = CA^{t-t_0}x_0 + \sum_{i=t_0}^{t-1} CA^{t-i-1}Bu(i) + Du(t). \quad (6.2)$$

The state-to-output map is **causal**, i.e., if u_1 and u_2 coincide for all $t_0 \leq t \leq t_1$, then

$$\eta(t_1, t_0, x, u_1) = \eta(t_1, t_0, x, u_2)$$

for all x . We say that it is **strictly causal** if this is already implied by $u_1(t) = u_2(t)$ for $t_0 \leq t < t_1$. In a state space system, strict causality holds if and only if $D = 0$. Due to **linearity**, we have

$$\eta(t, t_0, \lambda_1 x_1 + \lambda_2 x_2, \lambda_1 u_1 + \lambda_2 u_2) = \lambda_1 \eta(t, t_0, x_1, u_1) + \lambda_2 \eta(t, t_0, x_2, u_2)$$

for all $t \geq t_0$, $\lambda_i \in \mathbb{R}$, $x_i \in X$, $u_i \in \mathcal{U}$.

Observability is concerned with the following problem. Usually, we only know the system's manifest variables (in a state space system, the manifest variables are input and output): The input is free and can be chosen by the control engineer, and the output is the system's response which can be measured. The latent variables (in a state space system, the latent variable is the state) are not directly measurable, in general. Remember that they are usually auxiliary variables introduced during modelling, or simply for mathematical convenience, e.g., for reducing a system to first order. Hence the physical meaning of the latent variables may be obscure. The following question arises: If we know the manifest variables of a system, what can we conclude about the latent variables? An observable system is one in which the latent variables can be reconstructed from the manifest variables. Due to the property of state, we only need to reconstruct the state at a specific time t_0 , then we know it everywhere in the "future", i.e., for $t \geq t_0$.

Definition 6.1 Let $t_0 \in T$ be fixed. One says that the state $x \in X$

can be distinguished from state $x' \in X$ in time $\tau \in T$ ($\tau \geq 0$) if there exists $u \in \mathcal{U}$ and $t_0 \leq t \leq t_0 + \tau$ such that

$$\eta(t, t_0, x, u) \neq \eta(t, t_0, x', u);$$

then we also say that **u distinguishes between x and x' in time τ** .

can be distinguished from $x' \in X$ if this holds for at least one $\tau \geq 0$.

We say that the system is

observable if for any $x \neq x' \in X$, the state x can be distinguished from the state x' .

As usual, the starting time t_0 is not important as long as we deal only with time-invariant systems, and thus we put $t_0 = 0$. Define the set

$$\mathcal{J}(\tau, x) := \{x' \in X \mid x' \text{ is indistinguishable from } x \text{ in time } \tau\}$$

and let $\mathcal{J}(\tau) := \mathcal{J}(\tau, 0)$ denote the set of states that are indistinguishable from state 0 in time τ . We have

$$x \in \mathcal{J}(\tau) \iff \eta(t, 0, x, u) = \eta(t, 0, 0, u) \text{ for all } 0 \leq t \leq \tau, u \in \mathcal{U}.$$

Because of linearity,

$$\eta(t, 0, x, u) = \eta(t, 0, x, 0) + \eta(t, 0, 0, u).$$

Therefore we have

$$x \in \mathcal{J}(\tau) \iff \eta(t, 0, x, 0) = 0 \text{ for all } 0 \leq t \leq \tau.$$

Finally,

$$\mathcal{J} := \bigcap_{\tau \geq 0} \mathcal{J}(\tau)$$

is the set of states that are indistinguishable from state 0. We have

$$x \in \mathcal{J} \iff \eta(t, 0, x, 0) = 0 \text{ for all } t \geq 0.$$

Since $\eta(t, 0, 0, 0) = 0$ for all $t \geq 0$, this means that x cannot be distinguished from zero if and only if the zero input function does not distinguish between x and 0. In other words: If x can be distinguished from 0 at all, then it can also be distinguished from 0 by the zero input function. This shows that the choice of the input function plays no role for the question of observability, i.e., we may put $u = 0$ without loss of generality in most of this chapter.

Theorem 6.2 Let $s, t \in T$, $0 \leq s \leq t$. We have

1. $\mathcal{J}(t) \subseteq \mathcal{J}(s)$;
2. $\mathcal{J}(t), \mathcal{J}$ are subspaces of $X = \mathbb{R}^n$;
3. There exists $\tau^* \in T$, $\tau^* \geq 0$ such that

$$\mathcal{J} = \mathcal{J}(\tau) \text{ for all } \tau \geq \tau^*.$$

Corollary 6.3 In discrete time,

$$\mathcal{J}(n-1) = \mathcal{J},$$

where n is the dimension of the state space. In continuous time,

$$\mathcal{J}(\varepsilon) = \mathcal{J}$$

for every $\varepsilon > 0$.

Corollary 6.4 The following are equivalent:

1. The system is observable;
2. Any non-zero state can be distinguished from zero, that is, $\mathcal{J} = \{0\}$.

We define the **observability Gramians**

$$W(t) = \int_0^t e^{A^T \tau} C^T C e^{A \tau} d\tau \quad \text{or} \quad W(t) = \sum_{i=0}^t (A^T)^i C^T C A^i \in \mathbb{R}^{n \times n}$$

and the **Kalman observability matrix**

$$O = \begin{bmatrix} C \\ CA \\ CA^2 \\ \vdots \\ CA^{n-1} \end{bmatrix} \in \mathbb{R}^{np \times n}.$$

Theorem 6.5 Consider $\dot{x} = Ax$ or $\sigma x = Ax$, $y = Cx$. In continuous time, let $\tau^* = \varepsilon > 0$ be arbitrary. In discrete time, let $\tau^* = n-1$, where n is the dimension of the state space. We have

$$\mathcal{J} = \mathcal{J}(\tau^*) = \ker(W(\tau^*)) = \ker(O).$$

Therefore, the following are equivalent:

1. $\dot{x} = Ax$ or $\sigma x = Ax$, $y = Cx$ is observable;
2. $W(\tau^*)$ is non-singular;
3. O has full column rank.

Moreover in that case, we have the reconstruction formulas

$$x(0) = W(\tau^*)^{-1} \int_0^{\tau^*} e^{A^T t} C^T y(t) dt \quad \text{or} \quad x(0) = W(\tau^*)^{-1} \sum_{i=0}^{\tau^*} (A^T)^i C^T y(i).$$

Since $W(\tau^*)$ is always positive semi-definite due to its form, condition 2 from above is also equivalent to: $W(\tau^*) > 0$.

6.2 Observable matrix pairs

Let $A \in \mathbb{R}^{n \times n}$ and $C \in \mathbb{R}^{p \times n}$. We say that the matrix pair (A, C) is **observable** if the associated Kalman observability matrix

$$O = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{bmatrix}$$

has full column rank, that is, $\text{rank}(O) = n$.

If a state space system $\dot{x} = Ax + Bu$, $y = Cx + Du$ is subject to a coordinate transform $x = Tz$, where $T \in \mathbb{R}^{n \times n}$ is invertible, then we get

$$\begin{aligned} \dot{z} &= T^{-1}ATz + T^{-1}Bu \\ y &= CTz + Du. \end{aligned}$$

Discrete systems behave analogously. We say that the matrix pair $(T^{-1}AT, CT)$ is **similar** to the matrix pair (A, C) . A coordinate transform does not change structural system properties such as stability and observability.

The following result is limited to the **single-output** case, that is, $p = 1$. Then C is a single row vector. In that case, we simply write c instead of C . The associated Kalman observability matrix is then a square matrix.

Theorem 6.6 Let $A \in \mathbb{R}^{n \times n}$ and $c \in \mathbb{R}^{1 \times n}$, and let (A, c) be an observable matrix pair. Then there exists an invertible matrix $T \in \mathbb{R}^{n \times n}$ such that

$$\tilde{A} := T^{-1}AT = \begin{bmatrix} 0 & \cdots & 0 & -a_0 \\ 1 & & & -a_1 \\ & \ddots & & \vdots \\ & & 1 & -a_{n-1} \end{bmatrix} \quad \text{and} \quad \tilde{c} := cT = [0 \quad \cdots \quad 0 \quad 1].$$

The numbers a_i are precisely the coefficients of the characteristic polynomial, that is,

$$\chi_A(s) = \chi_{T^{-1}AT}(s) = s^n + a_{n-1}s^{n-1} + \dots + a_1s + a_0.$$

This is called the **observer form** of (A, c) . Moreover, there exists an invertible matrix $T_1 \in \mathbb{R}^{n \times n}$ such that

$$T_1^{-1}AT_1 = \begin{bmatrix} 0 & 1 & & \\ \vdots & & \ddots & \\ 0 & & & 1 \\ -a_0 & -a_1 & \cdots & -a_{n-1} \end{bmatrix} \quad \text{and} \quad cT_1 = [1 \ 0 \ \cdots \ 0].$$

This is called the **observability form** of (A, c) . The coefficients a_i are the same as with the observer form.

If a scalar input-output representation

$$y^{(n)} + a_{n-1}y^{(n-1)} + \dots + a_1\dot{y} + a_0y = u$$

is reduced to first order in the usual way, i.e., via putting $x = [y, \dot{y}, \dots, y^{(n-1)}]^T$ and $y = [1, 0, \dots, 0]x$, then the resulting state space system is precisely in observability form.

We return to the general **multi-output** case, and we give another result about transforming a given matrix pair into some special form via a similarity transform.

Theorem 6.7 (Kalman observability decomposition) Let $A \in \mathbb{R}^{n \times n}$ and $C \in \mathbb{R}^{p \times n}$. Let O be the associated Kalman observability matrix and let $r := \text{rank}(O)$. Then there exists an invertible matrix $T \in \mathbb{R}^{n \times n}$ such that

$$\tilde{A} := T^{-1}AT = \begin{bmatrix} A_1 & 0 \\ A_2 & A_3 \end{bmatrix} \quad \text{and} \quad \tilde{C} := CT = [C_1 \ 0]$$

where $A_1 \in \mathbb{R}^{r \times r}$, $C_1 \in \mathbb{R}^{p \times r}$ is an observable matrix pair.

Remark 6.8 The theorem says that via a suitable coordinate transform, namely $x = Tz$, the given system $\dot{x} = Ax$, $y = Cx$ can be put into the form

$$\begin{aligned} \dot{z}_1 &= A_1 z_1 \\ \dot{z}_2 &= A_2 z_1 + A_3 z_2 \\ y &= C_1 z_1 \end{aligned}$$

where (A_1, C_1) is observable. The state z_2 is certainly not observable (the output does not depend on z_2 at all, hence there is no chance of reconstructing z_2 by measuring y). On the other hand, since (A_1, C_1) is observable, $z_1(0)$ is reconstructible. Thus

$$\{z \in \mathbb{R}^n \mid z \text{ cannot be distinguished from } 0\} = \{0\} \times \mathbb{R}^{n-r}.$$

This can be used to determine the indistinguishable space of the original system, because

$$\{x \in \mathbb{R}^n \mid x \text{ cannot be distinguished from } 0\} = T(\{0\} \times \mathbb{R}^{n-r}).$$

Note that if the original (A, C) is observable, then $r = n$, and the Kalman observability decomposition becomes trivial. Thus the interesting case arises when (A, C) itself is not observable.

Remark 6.9 In a Kalman observability decomposition, we clearly have

$$\chi_A = \chi_{A_1} \cdot \chi_{A_3}$$

and thus

$$\text{spec}(A) = \text{spec}(A_1) \cup \text{spec}(A_3).$$

One calls χ_{A_3} the **unobservable part** of the characteristic polynomial of A with respect to C , and $\lambda \in \text{spec}(A_3)$ an **unobservable mode** of (A, C) .

There is also a direct way to characterize the unobservable modes of a matrix pair.

Theorem 6.10 Let $A \in \mathbb{R}^{n \times n}$, $C \in \mathbb{R}^{p \times n}$, and $\lambda \in \mathbb{C}$. The following are equivalent:

1. λ is an unobservable mode of (A, C) ;
2. $\text{rank} \begin{bmatrix} \lambda I - A \\ C \end{bmatrix} < n$.

As a direct consequence of this, we obtain another characterization of observable matrix pairs.

Corollary 6.11 (Hautus test for observability) The following are equivalent:

1. (A, C) is observable.
2. The matrix $H(\lambda) = \begin{bmatrix} \lambda I - A \\ C \end{bmatrix}$ has full column rank for all $\lambda \in \mathbb{C}$.

The polynomial matrix $H = \begin{bmatrix} sI - A \\ C \end{bmatrix} \in \mathbb{R}[s]^{(n+p) \times n}$ is called **Hautus observability matrix**. It suffices to check condition 2 from above for $\lambda \in \text{spec}(A)$.

Remark 6.12 Observability is a **generic** property, that is, if a matrix pair (A, C) is chosen “at random”, then it is very likely to be an observable one. More precisely, the set of observable matrix pairs (A, C) is open and dense in the set $\mathbb{R}^{n \times n} \times \mathbb{R}^{p \times n}$.

6.3 Asymptotic observability

Observability means

$$x \in \mathcal{J} \quad \Rightarrow \quad x = 0.$$

An alternative formulation is

$$\eta(t, 0, x, 0) = 0 \text{ for all } t \geq 0 \quad \Rightarrow \quad \varphi(t, 0, x, 0) = 0 \text{ for all } t \geq 0.$$

For asymptotic observability, one is satisfied if this is true in the limit as $t \rightarrow \infty$.

Definition 6.13 We say that a state space system $\dot{x} = Ax$ or $\sigma x = Ax$, $y = Cx$, is **asymptotically observable** if $\eta(t, 0, x, 0) = 0$ for all $t \geq 0$ implies that

$$\lim_{t \rightarrow \infty} \varphi(t, 0, x, 0) = 0.$$

Clearly, observability implies asymptotic observability.

Theorem 6.14 A state space system is asymptotically observable if and only if its unobservable modes λ are asymptotically stable, that is, $\text{Re}(\lambda) < 0$ in continuous time, and $|\lambda| < 1$ in discrete time.

6.4 Observable latent variable descriptions

The Hautus test gives us an idea about how to generalize the notion of observability from state space systems $\dot{x} = Ax + Bu$, $y = Cx + Du$ to general systems $R(\frac{d}{dt})w = M(\frac{d}{dt})l$ where $R \in \mathbb{R}[s]^{p \times q}$, $M \in \mathbb{R}[s]^{p \times r}$ and $w \in \mathcal{A}^q$, $l \in \mathcal{A}^r$. In a state space system,

$$R = \begin{bmatrix} B & 0 \\ -D & I \end{bmatrix}, \quad M = \begin{bmatrix} sI - A \\ C \end{bmatrix} \quad \text{and} \quad w = \begin{bmatrix} u \\ y \end{bmatrix}, \quad l = x.$$

The polynomial matrix M is recognized as the Hautus observability matrix.

In this section, we restrict to continuous systems, and we return to our original signal spaces, that is, $\mathcal{A} = \mathcal{D}'(T)$, where $T = \mathbb{R}, \mathbb{R}_+$.

Definition 6.15 We say that the latent variables l can be observed from the manifest variables w if l is uniquely determined by w , which means that $w_1 = w_2$ implies that $l_1 = l_2$.

Theorem 6.16 (Generalized Hautus test) Let $\mathcal{A} = \mathcal{D}'(T)$ for $T = \mathbb{R}$ or \mathbb{R}_+ and let $M \in \mathbb{R}[s]^{p \times r}$. Without loss of generality, let M have full column rank. Then the latent variable description of

$$\mathcal{B} = \{w \in \mathcal{A}^q \mid \exists l \in \mathcal{A}^r : R(\frac{d}{dt})w = M(\frac{d}{dt})l\}$$

is observable if and only if $\text{rank}(M(\lambda)) = r$ for all $\lambda \in \mathbb{C}$. In this case, one calls the matrix M **right prime** (or: right irreducible).

Remark 6.17 In Theorem 4.26, the notion of left primeness was characterized. The connection with right primeness is easy: A matrix is right prime if and only if its transpose is left prime.

If M does not have full column rank, we can write (up to a permutation of columns) $M = [M_1, M_2]$, where M_1 has full column rank, and $M_2 = M_1 X$ for some rational matrix X . Then $X = \frac{N}{d}$ for some polynomial matrix N and some $0 \neq d \in \mathbb{R}[s]$. By the elimination of latent variables, we have $\forall l_2 \exists \tilde{l}_2 : d(\frac{d}{dt})\tilde{l}_2 = l_2$. Thus

$$M(\frac{d}{dt})l = M_1(\frac{d}{dt})l_1 + M_2(\frac{d}{dt})l_2 = M_1(\frac{d}{dt})l_1 + M_2(\frac{d}{dt})d(\frac{d}{dt})\tilde{l}_2 = M_1(\frac{d}{dt})(l_1 + N(\frac{d}{dt})\tilde{l}_2).$$

This shows that we may assume w.l.o.g. that M has full column rank.

Proof: We need to show that the right primeness of M is equivalent to

$$\mathcal{B}_{\text{unobs}} := \{l \in \mathcal{A}^r \mid M(\frac{d}{dt})l = 0\} = \{0\}.$$

If M is right prime, there exists a polynomial matrix S such that $SM = I$. Then

$$M(\frac{d}{dt})l = 0 \quad \Rightarrow \quad S(\frac{d}{dt})M(\frac{d}{dt})l = 0 \quad \Rightarrow \quad l = 0$$

and thus $\mathcal{B}_{\text{unobs}} = \{0\}$ as desired. Conversely, if M is not right prime, then there exists $\lambda \in \mathbb{C}$ and $0 \neq z \in \mathbb{C}^r$ such that $M(\lambda)z = 0$. Set $l(t) := \text{Re}(e^{\lambda t}z)$. Then $(M(\frac{d}{dt})l)(t) = \text{Re}(M(\lambda)e^{\lambda t}z) = 0$. We have $l \in \mathcal{C}^\infty(T)^r \subset \mathcal{A}^r$, and thus $0 \neq l \in \mathcal{B}_{\text{unobs}}$. \square

6.5 Non-linear systems and zero-input observability

Consider a non-linear system

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

where $F : X \times U \rightarrow \mathbb{R}^n$ and $H : X \rightarrow \mathbb{R}^p$ are continuously differentiable, and X, U are open subsets of \mathbb{R}^n and \mathbb{R}^m , respectively. As usual, let $\varphi(t, t_0, x_0, u)$ denote the state transition map. Additionally, let $\eta(t, t_0, x_0, u)$ denote the output at time t when the state at time t_0 was x_0 and the input function was u . Let t_0 be fixed. Just like in the linear case, we say that **$x \in X$ can be distinguished from $x' \in X$ in time $\tau \geq 0$** if there exists an input function u and some $t_0 \leq t \leq t_0 + \tau$ such that

$$\eta(t, t_0, x, u) \neq \eta(t, t_0, x', u).$$

Then we say that **u distinguishes between x and x' in time τ** . Also, we say that **x can be distinguished from x'** if it can be distinguished in some finite time. The definitions

$$\mathcal{J}(\tau, x) := \{x' \in X \mid x' \text{ is indistinguishable from } x \text{ in time } \tau\}$$

and

$$\mathcal{J}(x) := \bigcap_{\tau \geq 0} \mathcal{J}(\tau, x) = \{x' \in X \mid x' \text{ is indistinguishable from } x\}$$

are just like in the linear setting (since there is no vector space structure, the state $x = 0$ does not play any special role).

The system is said to be **globally observable** if any $x \neq x' \in X$ can be distinguished, that is, $\mathcal{J}(x) = \{x\}$ for all $x \in X$. Just as with reachability, global notions are typically too strong requirements for non-linear system. Therefore, we say that the system is **locally observable at x_0** if there exists an open neighborhood X_0 of x_0 such that any $x, x' \in X_0$ with $x \neq x'$ can be distinguished.

Unlike the linear case, it is restrictive to consider only the zero-input case. However, it suffices for our purposes, and thus we consider

$$\begin{aligned}\dot{x} &= F(x) \\ y &= H(x)\end{aligned}$$

in the following. We will write $\varphi(t, x)$ for the state at time t when starting in $x(0) = x$ and $\eta(t, x) = H(\varphi(t, x))$ for the output at time t when starting in $x(0) = x$.

Let x_0 be an equilibrium of the system, that is, $F(x_0) = 0$. Setting

$$A := \frac{\partial F}{\partial x}(x_0) \quad \text{and} \quad C := \frac{\partial H}{\partial x}(x_0),$$

we get the linearization at x_0 , which is given by

$$\begin{aligned}\dot{\tilde{x}} &= A\tilde{x} \\ \tilde{y} &= C\tilde{x},\end{aligned}$$

where $\tilde{x} = x - x_0$ and $\tilde{y} = y - y_0$ for $y_0 := H(x_0)$.

Theorem 6.18 Consider $\dot{x}(t) = F(x(t))$, $y(t) = H(x(t))$ with $F(x_0) = 0$. Define A and C as above. If (A, C) is an observable matrix pair, then for all $t_1 > 0$, there exists an open neighborhood X_0 of x_0 such that $\mathcal{J}(x, t_1) \cap X_0 = \{x\}$ for all $x \in X_0$.

In words: If the linearization at x_0 is observable, then the non-linear system is locally observable at x_0 . Moreover, we have local observability in arbitrarily small time.

Outline of proof: Let $t_1 > 0$ be given. Consider

$$L : X \rightarrow \mathcal{Y} := \mathcal{C}^0([0, t_1], \mathbb{R}^p), \xi \mapsto \eta(\cdot, \xi),$$

that is, $L(\xi)(t) = \eta(t, \xi) = H(\varphi(t, \xi))$. Then

$$\frac{\partial L}{\partial \xi}(\xi)(t) = \frac{\partial H}{\partial x}(\varphi(t, \xi)) \frac{\partial \varphi}{\partial \xi}(t, \xi).$$

Fact from ODE theory: the matrix-valued function

$$G(t) := \frac{\partial \varphi}{\partial \xi}(t, \xi_0)$$

satisfies

$$\dot{G}(t) = \frac{\partial F}{\partial x}(\varphi(t, \xi_0))G(t)$$

and $G(0) = I$. If $\xi_0 = x_0$, this simplifies to (since x_0 is an equilibrium)

$$\dot{G}(t) = AG(t)$$

and $G(0) = I$, which clearly implies that $G(t) = e^{At}$.

Thus (plugging in $\xi_0 = x_0$)

$$\frac{\partial L}{\partial \xi}(x_0) = Ce^{A \cdot} \in \mathcal{Y}^{1 \times n}.$$

We can interpret this as the linearization of L at x_0 ,

$$\frac{\partial L}{\partial \xi}(x_0) : X \rightarrow \mathcal{Y}, x \mapsto Ce^{A \cdot} x.$$

By assumption, this linear map is injective (observability of (A, C) says that $Ce^{At}x = 0$ for all $0 \leq t \leq t_1$ implies that $x = 0$). The application of an appropriate infinite-dimensional variant of the inverse function theorem (\mathcal{Y} is a Banach space) yields that L is injective on some neighborhood X_0 of x_0 . \square

Another approach to non-linear observability assumes that F and H are even smooth. Then both $x(\cdot)$ and $y(\cdot)$ are also smooth. Set $t_0 = 0$. When we observe the output $y|_{[0, t_1]}$ for some $t_1 > 0$ in order to reconstruct the initial state $x(0) = x_0$, we get also information about the derivatives of y at zero, for instance

$$y(0) = H(x_0), \quad \dot{y}(0) = \frac{\partial H}{\partial x}(x_0)F(x_0)$$

etc. So we may ask ourselves: Can x_0 be reconstructed from the knowledge of $y(0), \dot{y}(0), \ddot{y}(0), \dots$? To express the higher order derivatives of y along the system in terms of F and H , the concept of **Lie derivative** is useful:

$$L_F H := \frac{\partial H}{\partial x} F = \sum_{i=1}^n \frac{\partial H}{\partial x_i} F_i$$

is called the Lie derivative of H along F . This corresponds to differentiating $y = H(x)$ according to $\dot{y} = \frac{\partial H}{\partial x} \dot{x}$ and to replace \dot{x} by $F(x)$ according to our differential equation. Then $\dot{y} = L_F H$ and iteratively,

$$y^{(k)} = L_F^k H.$$

By the inverse function theorem, we have the following result.

Theorem 6.19 If there exists $k \in \mathbb{N}$ such that

$$\frac{\partial \Psi}{\partial x}(x_0)$$

has full column rank, where Ψ is defined by

$$\Psi : X \rightarrow \mathbb{R}^{kp}, x \mapsto \begin{bmatrix} H \\ L_F H \\ \vdots \\ L_F^{k-1} H \end{bmatrix} (x),$$

then for all $t_1 > 0$, the non-linear system is locally observable at x_0 in time t_1 .

In particular, for $\dot{x} = Ax$, $y = Cx$ we get $L_F^k H(x) = CA^k x$ and thus for $k = n$, $\frac{\partial \Psi}{\partial x}(x_0)$ is precisely the Kalman observability matrix (independently of x_0).

The second approach can also be applied to linear but time-varying systems

$$\begin{aligned} \dot{x}(t) &= A(t)x(t) \\ y(t) &= C(t)x(t), \end{aligned}$$

where $A : T \rightarrow \mathbb{R}^{n \times n}$ and $C : T \rightarrow \mathbb{R}^{p \times n}$ are smooth. Considering $y(0), \dot{y}(0), \dots$, we find the time-varying Kalman matrix

$$O^T(t) = [C^T, (A^T + \frac{d}{dt})C^T, (A^T + \frac{d}{dt})^2 C^T, \dots](t),$$

which coincides with

$$O^T = [C^T, A^T C^T, \dots]$$

in the case where A, C are constant (except for the fact that $O(t)$ is an infinite matrix, whereas the time-invariant O is finite due to Hamilton-Cayley).

Note that we use the transposed notation, since differential operators are traditionally written on the left of the function they are acting on.

Theorem 6.20 Consider

$$\begin{aligned}\dot{x}(t) &= A(t)x(t) \\ y(t) &= C(t)x(t)\end{aligned}$$

where A, C are smooth. If the time-varying Kalman observability matrix has rank n for some $t_0 \in T$, then the system is observable at time t_0 , that is, for all $x, x' \in \mathbb{R}^n$ and all $\varepsilon > 0$, we have

$$\eta(t, t_0, x, 0) = \eta(t, t_0, x', 0) \text{ for all } t_0 \leq t \leq t_0 + \varepsilon \quad \Rightarrow \quad x = x'.$$

If A, C are even analytic, then the rank condition becomes necessary as well as sufficient. Moreover in that case, there exists an open and dense subset T_0 of T such that for all $t_0 \in T_0$, it suffices to check the rank of the finite matrix

$$[C^T, (A^T + \frac{d}{dt})C^T, \dots, (A^T + \frac{d}{dt})^{n-1}C^T](t_0).$$

Of course, this theorem does also apply to

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

because we only need linearity to reduce the observability problem (w.l.o.g.) to the case where $u = 0$.

Chapter 7

Observers

7.1 State observers

Feedback control requires the knowledge of the state of the system. However, the state is usually not directly measurable. In an observable system, it can – in principle – be reconstructed from the inputs and outputs. However, the reconstruction procedure is sensitive with respect to noisy data.

Observer design is an alternative approach. Its basic idea is to build another system whose state converges to the state of the given system, independently of the initial conditions. For a state space system (in this chapter, we put $D = 0$ for simplicity),

$$\begin{aligned} \dot{x} &= Ax + Bu & \text{or} & & \sigma x &= Ax + Bu \\ y &= Cx & & & y &= Cx \end{aligned} \tag{7.1}$$

where $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{p \times n}$, a **state observer** has the form (in continuous time, the discrete case is analogous)

$$\begin{aligned} \dot{z} &= Az + Bu + Ld \\ \hat{y} &= Cz \end{aligned}$$

where $L \in \mathbb{R}^{n \times p}$. The signal $d := \hat{y} - y$ is the difference between the observer output \hat{y} and the output y of the original system. Thus, the observer is a state space system with state z , inputs u and d , and output \hat{y} . Since A, B, C are the same matrices as with the given system, the observer can be seen as a copy of the original system, with the additional input d . The observer equations can be

rewritten as

$$\begin{aligned}\dot{z} &= (A + LC)z + Bu - Ly \\ \hat{y} &= Cz.\end{aligned}\tag{7.2}$$

We say that the matrix pair $(A + LC, C)$ results from (A, C) via the observer gain matrix L . In this form, the observer has precisely the input and output of the original system as its inputs. Thus it requires only the knowledge of the manifest variables of the given system. One should think of an observer as a signal processing algorithm rather than a physical system. Its goal is to produce an estimate z of the unknown state x .

Consider the difference between the true state x and the observer state z . For this, put $e := z - x$. Then (7.1) and (7.2) together imply that

$$\dot{e} = (A + LC)e$$

and thus $e(t) = e^{(A+LC)t}e(0)$. The goal of this chapter is to answer questions like the **detection problem**: Given A, C , can we find L such that $A + LC$ is asymptotically stable? If yes, then the error signal e will tend to zero as $t \rightarrow \infty$, for all $e(0)$. Thus we have

$$\lim_{t \rightarrow \infty} \|z(t) - x(t)\| = 0$$

for all $z(0), x(0)$. This shows that the observer state z will asymptotically approach the true state x , and this will work even if the initial states are subject to disturbance.

Lemma 7.1 The indistinguishable spaces of (7.1) and (7.2) coincide. In particular, $(A + LC, C)$ is observable if and only if (A, C) is observable.

7.2 Pole placement

In this section, we investigate the characteristic polynomial and the spectrum of $A + LC$, where A, C are given, and L may be chosen. The goal is to place the eigenvalues of $A + LC$ in some desirable region of the complex plane.

Definition 7.2 Let $A \in \mathbb{R}^{n \times n}$ and $C \in \mathbb{R}^{p \times n}$ be given. Let p be a monic polynomial of degree n . We say that p is **assignable** to (A, C) if there exists a matrix $L \in \mathbb{R}^{n \times p}$ such that $\chi_{A+LC} = p$.

Theorem 7.3 Let $A \in \mathbb{R}^{n \times n}$ and $C \in \mathbb{R}^{p \times n}$. The matrix pair (A, C) is observable if and only if every monic polynomial of degree n can be assigned to (A, C) .

Theorem 7.4 (Pole shifting theorem) The polynomials that can be assigned to (A, C) are precisely the ones of the form

$$p = p_1 \cdot \chi_u$$

where χ_u is the unobservable part of χ_A with respect to C , and p_1 is an arbitrary monic polynomial of degree $r = n - \deg(\chi_u)$.

7.3 Detection

For detection, we don't require that the spectrum of $A + LC$ should coincide with some specific set of eigenvalues; we just want the eigenvalues to be contained in some given region of the complex plane. We put

$$\mathbb{C}_g = \{\lambda \in \mathbb{C} \mid \operatorname{Re}(\lambda) < 0\} \quad \text{or} \quad \mathbb{C}_g = \{\lambda \in \mathbb{C} \mid |\lambda| < 1\}$$

in the continuous or discrete case, respectively.

Definition 7.5 A matrix pair (A, C) is called **detectable** if there exists L such that $A + LC$ is asymptotically stable, that is, $\operatorname{spec}(A + LC) \subset \mathbb{C}_g$.

Theorem 7.6 The following are equivalent:

1. (A, C) is detectable;
2. (A, C) is asymptotically observable;
3. All unobservable modes of (A, C) lie in \mathbb{C}_g ;
4. Any eigenvalue λ of A which is not in \mathbb{C}_g satisfies $\operatorname{rank} \begin{bmatrix} \lambda I - A \\ C \end{bmatrix} = n$.

Remark 7.7 Part of this theorem was already stated in Section 6.3.

7.4 Compensators

Finally, we combine feedback control with state observation. Consider the state space system

$$\begin{aligned}\dot{x} &= Ax + Bu \\ y &= Cx.\end{aligned}$$

We wish to stabilize the system via a feedback $u = Fx + v$. However, since we do not know x , we build an observer

$$\dot{z} = (A + LC)z + Bu - Ly.$$

We take the observer state z instead of x in the feedback law and put

$$u = Fz + v.$$

Then we obtain the closed loop system

$$\frac{d}{dt} \begin{bmatrix} x \\ z \end{bmatrix} = \begin{bmatrix} A & BF \\ -LC & A + LC + BF \end{bmatrix} \begin{bmatrix} x \\ z \end{bmatrix} + \begin{bmatrix} B \\ B \end{bmatrix} v.$$

If we can make this system asymptotically stable, the combination of feedback and observer outlined above will work nicely.

Definition 7.8 We call (F, L) a **compensator** for (A, B, C) if

$$K := \begin{bmatrix} A & BF \\ -LC & A + LC + BF \end{bmatrix}$$

is asymptotically stable.

Theorem 7.9 The matrix pair (F, L) is a compensator for (A, B, C) if and only if $A + BF$ and $A + LC$ are both asymptotically stable. Therefore, (A, B, C) possesses a compensator if and only if (A, B) is stabilizable and (A, C) is detectable. If we even have that (A, B) is controllable and (A, C) is observable, then for any monic polynomial p of degree $2n$, there exists (F, L) such that $\chi_K = p$.

Proof: Let

$$T = \begin{bmatrix} I & 0 \\ I & I \end{bmatrix}.$$

Then

$$T^{-1}KT = \begin{bmatrix} A + BF & BF \\ 0 & A + LC \end{bmatrix}$$

and thus $\chi_K = \chi_{A+BF} \cdot \chi_{A+LC}$ and $\text{spec}(K) = \text{spec}(A + BF) \cup \text{spec}(A + LC)$. \square

Chapter 8

Transfer matrices

An input-output representation has the form

$$P\left(\frac{d}{dt}\right)y = Q\left(\frac{d}{dt}\right)u \quad \text{or} \quad P(\sigma)y = Q(\sigma)u$$

where $P \in \mathbb{R}[s]^{p \times p}$ is non-singular, and $Q \in \mathbb{R}[s]^{p \times m}$. The rational matrix

$$H := P^{-1}Q \in \mathbb{R}(s)^{p \times m}$$

is called **transfer matrix** (or: transfer function) of the input-output representation.

Lemma 8.1 The transfer matrix of a state space system is

$$H = C(sI - A)^{-1}B + D. \tag{8.1}$$

Proof: We need to eliminate the state from

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

where $A \in \mathbb{R}^{n \times n}, \dots, D \in \mathbb{R}^{p \times m}$ (the discrete case is analogous). Let U be a unimodular matrix such that

$$U \begin{bmatrix} sI - A \\ C \end{bmatrix} = \begin{bmatrix} R_1 \\ 0 \end{bmatrix} \tag{8.2}$$

where $R_1 \in \mathbb{R}[s]^{n \times n}$ is non-singular (see Corollary 2.4). Then

$$\exists x : \begin{cases} \dot{x} = Ax + Bu \\ y = Cx + Du \end{cases}$$

is equivalent to

$$\exists x : U\left(\frac{d}{dt}\right) \begin{bmatrix} sI - A \\ C \end{bmatrix} \left(\frac{d}{dt}\right)x = \begin{bmatrix} R_1 \\ 0 \end{bmatrix} \left(\frac{d}{dt}\right)x = U\left(\frac{d}{dt}\right) \begin{bmatrix} B & 0 \\ -D & I \end{bmatrix} \begin{bmatrix} u \\ y \end{bmatrix}.$$

According to the fundamental principle, this is also equivalent to

$$(U_3B - U_4D)\left(\frac{d}{dt}\right)u + U_4\left(\frac{d}{dt}\right)y = 0 \quad \text{where} \quad U = \begin{bmatrix} U_1 & U_2 \\ U_3 & U_4 \end{bmatrix}. \quad (8.3)$$

The matrix $[U_3, U_4]$ has full row rank p , and it follows from (8.2) that

$$U_3(sI - A) + U_4C = 0,$$

that is, $U_3 = -U_4C(sI - A)^{-1}$ which shows that the columns of U_3 are (rational) linear combinations of the columns of U_4 . Therefore $p = \text{rank}[U_3, U_4] = \text{rank}(U_4)$. We conclude that $U_4 \in \mathbb{R}[s]^{p \times p}$ is non-singular. Thus (8.3) is an input-output representation with transfer function

$$H = -U_4^{-1}(U_3B - U_4D) = -U_4^{-1}U_3B + D.$$

Using $U_4^{-1}U_3 = -C(sI - A)^{-1}$, we have the desired result. \square

8.1 Realization theory

By the previous lemma, it is easy to compute H if A, B, C, D are known. However, one often faces the inverse problem: Given H , find matrices A, B, C, D such that (8.1) holds. This is known as the **realization problem**. If (8.1) is satisfied, the matrix quadruple (A, B, C, D) is called a **realization** of H , and H is called **realizable** if it possesses a realization.

We first observe that any H according to (8.1) will be a **proper** rational matrix, that is, if we write $H = \frac{N}{d}$, where $N \in \mathbb{R}[s]^{p \times m}$ is a polynomial matrix, and $0 \neq d \in \mathbb{R}[s]$ is a scalar polynomial, then

$$\deg(N_{ij}) \leq \deg(d) \quad \text{for all } i, j. \quad (8.4)$$

This follows from $(sI - A)^{-1} = \frac{\text{adj}(sI - A)}{\det(sI - A)}$ using Cramer's rule. In (8.4), strict inequality holds for all i, j if and only if $D = 0$. In this case, one says that H is **strictly proper**. It turns out that properness is not only necessary but also sufficient for realizability.

Theorem 8.2 A rational matrix is realizable if and only if it is proper.

Proof: Let $H \in \mathbb{R}(s)^{p \times m}$ be proper, then we can write $H = D + H_1$ with $D \in \mathbb{R}^{p \times m}$ and H_1 strictly proper. Thus $H_1 = \frac{N}{d}$ with

$$d = s^\nu + d_{\nu-1}s^{\nu-1} + \dots + d_1s + d_0 \quad \text{and} \quad N = N_{\nu-1}s^{\nu-1} + \dots + N_1s + N_0$$

for some $\nu = \deg(d)$, $d_i \in \mathbb{R}$, $N_i \in \mathbb{R}^{p \times m}$. Put $n = \nu m$ and

$$A = \begin{bmatrix} 0 & I & & \\ \vdots & & \ddots & \\ 0 & & & I \\ -d_0I & -d_1I & \dots & -d_{\nu-1}I \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ I \end{bmatrix} \quad C = [N_0 \quad \dots \quad N_{\nu-1}]$$

then (A, B, C, D) is a realization of H . To see this, note that

$$(sI - A) \begin{bmatrix} I \\ sI \\ \vdots \\ s^{\nu-1}I \end{bmatrix} = \begin{bmatrix} sI & -I & & & \\ 0 & sI & -I & & \\ \vdots & & \ddots & \ddots & \\ 0 & & & sI & -I \\ d_0I & d_1I & \dots & \dots & sI + d_{\nu-1}I \end{bmatrix} \begin{bmatrix} I \\ sI \\ \vdots \\ s^{\nu-1}I \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ dI \end{bmatrix}$$

and hence

$$\begin{bmatrix} I \\ sI \\ \vdots \\ s^{\nu-1}I \end{bmatrix} = (sI - A)^{-1} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ dI \end{bmatrix} = (sI - A)^{-1} B d.$$

Pre-multiplying this by C , we obtain

$$C \begin{bmatrix} I \\ sI \\ \vdots \\ s^{\nu-1}I \end{bmatrix} = N = C(sI - A)^{-1} B d$$

which yields the desired result, after division by d . \square

Thus, any proper rational matrix H is realizable. Let (A, B, C, D) be a realization of H , with $A \in \mathbb{R}^{n \times n}$. We call the number n the **size** of the realization. Of course, it is desirable to have small realizations. We say that a realization of H is **minimal** if there exists no realization of H with a smaller size. The subsequent lemma gives an important relation between two realizations of a transfer function.

Lemma 8.3 If (A, B, C, D) and $(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D})$ are two realizations of the same transfer matrix, then $D = \tilde{D}$ and

$$CA^i B = \tilde{C} \tilde{A}^i \tilde{B}$$

for all $i \in \mathbb{N}$.

Proof: If $H = C(sI - A)^{-1}B + D$, then

$$\lim_{s \rightarrow \infty} H(s) = D.$$

This shows that we must have $D = \tilde{D}$. Moreover, we can expand $H - D$ into a Laurent series

$$H - D = C(sI - A)^{-1}B = \sum_{i=0}^{\infty} CA^i B s^{-i-1}$$

and this is convergent on $|s| > \rho(A)$, where $\rho(A)$ is the spectral radius of A . Therefore, by comparing coefficients, $CA^i B = \tilde{C}\tilde{A}^i\tilde{B}$ for all i . \square

Now we can give a sufficient condition for minimality which will soon turn out to be also necessary.

Lemma 8.4 Let (A, B, C, D) be such that (A, B) is controllable and (A, C) is observable. Then (A, B, C, D) is a minimal realization of $H = C(sI - A)^{-1}B + D$.

Remark 8.5 The proof of Lemma 8.4 uses **Sylvester's inequality**: If O is a real matrix with n columns, and K is a real matrix with n rows, then

$$\text{rank}(OK) \geq \text{rank}(O) + \text{rank}(K) - n.$$

Proof: Suppose that (A, B, C, D) is a realization of H , with size n , in which (A, B) is controllable and (A, C) is observable. Let $(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D})$ be another realization of H , with size \tilde{n} . We need to prove that $n \leq \tilde{n}$. Define

$$O = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{bmatrix} \quad K = [B \quad AB \quad \dots \quad A^{n-1}B]$$

and

$$\tilde{O} = \begin{bmatrix} \tilde{C} \\ \tilde{C}\tilde{A} \\ \vdots \\ \tilde{C}\tilde{A}^{n-1} \end{bmatrix} \quad \tilde{K} = [\tilde{B} \quad \tilde{A}\tilde{B} \quad \dots \quad \tilde{A}^{n-1}\tilde{B}].$$

Note that whereas O, K are precisely the Kalman matrices associated to (A, B) and (A, C) , this is not true for \tilde{O}, \tilde{K} (we have n instead of \tilde{n} in the highest power of A). Then, because $CA^i B = \tilde{C}\tilde{A}^i\tilde{B}$ for all i ,

$$OK = \tilde{O}\tilde{K}.$$

We have

$$\text{rank}(K) \geq \text{rank}(OK) \geq \text{rank}(O) + \text{rank}(K) - n.$$

By assumption, K and O both have rank n . Therefore,

$$n = \text{rank}(OK) = \text{rank}(\tilde{O}\tilde{K}) \leq \text{rank}(\tilde{O}) \leq \tilde{n}$$

as desired. \square

Theorem 8.6 (Reduction to minimality) Let (A, B, C, D) be a realization of H .

1. Consider a Kalman controllability decomposition

$$T^{-1}AT = \begin{bmatrix} A_1 & A_2 \\ 0 & A_3 \end{bmatrix} \quad T^{-1}B = \begin{bmatrix} B_1 \\ 0 \end{bmatrix} \quad CT = [C_1 \ C_2].$$

Then $H = C(sI - A)^{-1}B + D = C_1(sI - A_1)^{-1}B_1 + D$, that is, (A_1, B_1, C_1, D) is another realization of H , with size $r = \text{rank}(K)$, where K is the Kalman controllability matrix.

2. Consider a Kalman observability decomposition

$$T^{-1}AT = \begin{bmatrix} A_1 & 0 \\ A_2 & A_3 \end{bmatrix} \quad T^{-1}B = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} \quad CT = [C_1 \ 0].$$

Then $H = C(sI - A)^{-1}B + D = C_1(sI - A_1)^{-1}B_1 + D$, that is, (A_1, B_1, C_1, D) is another realization of H , with size $r = \text{rank}(O)$, where O is the Kalman observability matrix.

3. If the two reduction steps are done successively, one ends up with a minimal realization of H .

Proof: One can easily check that a similarity transform does not change the transfer function. Therefore, assume that a Kalman controllability decomposition has already been performed. Then

$$\begin{aligned} H &= [C_1 \ C_2] \begin{bmatrix} sI - A_1 & -A_2 \\ 0 & sI - A_3 \end{bmatrix}^{-1} \begin{bmatrix} B_1 \\ 0 \end{bmatrix} + D \\ &= [C_1 \ C_2] \begin{bmatrix} (sI - A_1)^{-1} & * \\ 0 & (sI - A_3)^{-1} \end{bmatrix} \begin{bmatrix} B_1 \\ 0 \end{bmatrix} + D \\ &= [C_1 \ C_2] \begin{bmatrix} (sI - A_1)^{-1}B_1 \\ 0 \end{bmatrix} + D = C_1(sI - A_1)^{-1}B_1 + D. \end{aligned}$$

The second statement is analogous. Recall that after a Kalman controllability decomposition, the matrix pair (A_1, B_1) is controllable. Now if one performs a Kalman observability decomposition with the already reduced system (A_1, B_1, C_1, D) , then one obtains

$$T_1^{-1}A_1T_1 = \begin{bmatrix} A_{11} & 0 \\ A_{12} & A_{13} \end{bmatrix} \quad T_1^{-1}B_1 = \begin{bmatrix} B_{11} \\ B_{12} \end{bmatrix} \quad C_1T_1 = [C_{11} \quad 0]$$

in which (A_{11}, C_{11}) is observable. We only need to convince ourselves that the controllability of (A_1, B_1) implies the controllability of (A_{11}, B_{11}) (i.e., the controllability established in the first reduction step is not destroyed by the second reduction step in which we achieve observability). Therefore, the resulting realization $(A_{11}, B_{11}, C_{11}, D)$ is both controllable and observable, and hence minimal according to Lemma 8.4. \square

In particular, this theorem shows that if in a realization (A, B) is not controllable or (A, C) is not observable, then the realization can be reduced in size. Moreover, this can be done constructively, using a Kalman decomposition. In other words, a minimal realization will always be both controllable and observable. Combining this result with Lemma 8.4, we obtain the following theorem as a summary.

Theorem 8.7 The matrix quadruple (A, B, C, D) is a minimal realization of $H = C(sI - A)^{-1}B + D$ if and only if (A, B) is controllable and (A, C) is observable.

The next theorem says that minimal realizations are essentially unique (up to similarity transforms).

Theorem 8.8 Any two minimal realizations of a transfer matrix are similar, that is, if (A, B, C, D) and $(\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D})$ are two minimal realizations of H , then there exists a non-singular matrix T such that

$$\tilde{A} = T^{-1}AT, \quad \tilde{B} = T^{-1}B, \quad \tilde{C} = CT, \quad \tilde{D} = D.$$

Proof: Since both realizations are minimal, their size must be the same, that is, $n = \tilde{n}$. Moreover, $D = \tilde{D}$, and $CA^iB = \tilde{C}\tilde{A}^i\tilde{B}$ for all i . Let $O, K, \tilde{O}, \tilde{K}$ be the observability and controllability matrices of the two realizations. Then $OK = \tilde{O}\tilde{K}$ and

$$OAK = \tilde{O}\tilde{A}\tilde{K}, \quad OB = \tilde{O}\tilde{B}, \quad CK = \tilde{C}\tilde{K}.$$

By assumption, \tilde{K} has full row rank, and \tilde{O} has full column rank. Therefore, there exist matrices L, N such that

$$\tilde{K}L = I \quad \text{and} \quad N\tilde{O} = I.$$

Put $T := KL \in \mathbb{R}^{n \times n}$, then $T^{-1} = NO$, because

$$NOKL = N\tilde{O}\tilde{K}L = I.$$

We have

$$\begin{aligned} T^{-1}AT &= NOAKL = N\tilde{O}\tilde{A}\tilde{K}L = \tilde{A} \\ T^{-1}B &= NOB = N\tilde{O}\tilde{B} = \tilde{B} \\ CT &= CKL = \tilde{C}\tilde{K}L = \tilde{C} \end{aligned}$$

which completes the proof. \square

Starting with an arbitrary realization of H , one can determine the size of a minimal realization of H by successively computing two Kalman decompositions, as outlined in Theorem 8.6. However, there is also a direct way to determine the size of a minimal realization. This will be discussed in the next section.

8.2 Matrix fraction descriptions

Let $H \in \mathbb{R}(s)^{p \times m}$ be given. If $H = P^{-1}Q$ for some $Q \in \mathbb{R}[s]^{p \times m}$, $P \in \mathbb{R}[s]^{p \times p}$ with $\det(P) \neq 0$, we call (P, Q) a **left factorization** (or: left matrix fraction description) of H .

Similarly, if $H = QP^{-1}$ for some polynomial matrices $Q \in \mathbb{R}[s]^{p \times m}$, $P \in \mathbb{R}[s]^{m \times m}$ with $\det(P) \neq 0$, we call (Q, P) a **right factorization** of H .

For example, we have already used the representation $H = \frac{N}{d}$ several times. In other words, (dI_p, N) is a left and (N, dI_m) is a right factorization of H .

In the scalar case, it is desirable to write a rational function $h \in \mathbb{R}(s)$ as the ratio (“fraction”) of two coprime polynomials. We wish to do the same with polynomial matrices.

We say that a left factorization (P, Q) is **left coprime** if the matrix $\begin{bmatrix} P & Q \end{bmatrix}$ is left irreducible (see Theorem 4.26). Similarly, a right factorization (Q, P) is called **right coprime** if the matrix $\begin{bmatrix} Q \\ P \end{bmatrix}$ is right irreducible, which means, by definition, that its transpose is left irreducible.

Lemma 8.9 1. Let (P, Q) be a left coprime factorization of H . The degree of the determinant of P is independent of the specific choice of the coprime factorization and therefore

$$d(H) := \deg \det(P) \quad (8.5)$$

is well-defined. If (P, Q) is an arbitrary (not necessarily coprime) left factorization of H , then

$$d(H) \leq \deg \det(P).$$

2. Let (Q, P) be a right coprime factorization of H . The degree of the determinant of P is independent of the specific choice of the coprime factorization and moreover,

$$d(H) = \deg \det(P)$$

where $d(H)$ is the number defined in (8.5). If (Q, P) is an arbitrary (not necessarily coprime) right factorization of H , then

$$d(H) \leq \deg \det(P).$$

Remark 8.10 The degree of the determinant of the “denominator” matrix P is as small as possible if the factorization is coprime. This generalizes the well known fact that the degree of d is minimal if we write a scalar rational function $h = \frac{n}{d} \in \mathbb{R}(s)$ as the ratio of two coprime polynomials. The choice of the coprime factorization does not influence this minimal degree, it does not even matter whether we take a right or left matrix fraction description.

In the proof, we use the formula

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

which holds provided that A is invertible. Thus

$$\det \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \det(A) \det(D - CA^{-1}B).$$

This shows that the block matrix is invertible if and only if $D - CA^{-1}B$ is invertible, and we have

$$\begin{aligned} \begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} &= \begin{bmatrix} I & -A^{-1}B \\ 0 & I \end{bmatrix} \begin{bmatrix} A^{-1} & 0 \\ 0 & (D - CA^{-1}B)^{-1} \end{bmatrix} \begin{bmatrix} I & 0 \\ -CA^{-1} & I \end{bmatrix} \\ &= \begin{bmatrix} * & * \\ * & (D - CA^{-1}B)^{-1} \end{bmatrix}. \end{aligned}$$

The matrix $D - CA^{-1}B$ is called **Schur complement** of A .

Proof: 1. Let (P, Q) and (P_1, Q_1) be two left factorizations of H , that is,

$$H = P^{-1}Q = P_1^{-1}Q_1$$

which implies that $Q_1 = UQ$ and $P_1 = UP$, where $U := P_1P^{-1}$. We show that U is polynomial if (P, Q) is left coprime; and even unimodular if both factorizations are coprime. If (P, Q) is left coprime, then there exist, according to Theorem 4.26, polynomial matrices R, S such that

$$PR + QS = I.$$

Then

$$P^{-1} = R + P^{-1}QS = R + HS = R + P_1^{-1}Q_1S$$

and thus $U = P_1P^{-1} = P_1R + Q_1S$. This shows that U is polynomial. Then

$$\det(P_1) = \det(UP) = \det(U) \det(P)$$

means that $\det(P)$ divides $\det(P_1)$, in particular,

$$\deg \det(P_1) \geq \deg \det(P).$$

Similarly, if also (P_1, Q_1) is coprime, there exist polynomial matrices R_1, S_1 such that $P_1R_1 + Q_1S_1 = I$ and then

$$P_1^{-1} = R_1 + P_1^{-1}Q_1S_1 = R_1 + HS_1 = R_1 + P^{-1}QS_1$$

and thus $U^{-1} = PP_1^{-1} = PR_1 + QS_1$ which shows that U^{-1} is polynomial and hence, U is unimodular. Thus we have

$$\det(P_1) = \det(U) \det(P)$$

where $\det(U)$ is a non-zero constant, and hence

$$\deg \det(P_1) = \deg \det(P).$$

2. In view of part 1, it suffices to show that if (P, Q) is a left coprime, and (Q_1, P_1) is a right coprime factorization of H , then $\deg \det(P) = \deg \det(P_1)$. We have $H = P^{-1}Q = Q_1P_1^{-1}$ and thus

$$QP_1 = PQ_1.$$

There exist polynomial matrices R, S, R_1, S_1 such that

$$PR + QS = I \quad \text{and} \quad R_1P_1 + S_1Q_1 = I.$$

Thus

$$\begin{bmatrix} P & -Q \\ S_1 & R_1 \end{bmatrix} \begin{bmatrix} R & Q_1 \\ -S & P_1 \end{bmatrix} = \begin{bmatrix} I & 0 \\ X & I \end{bmatrix}$$

where $X = S_1R - R_1S$. Post-multiplying this by

$$\begin{bmatrix} I & 0 \\ -X & I \end{bmatrix}$$

yields

$$\begin{bmatrix} P & -Q \\ S_1 & R_1 \end{bmatrix} \begin{bmatrix} \tilde{R} & Q_1 \\ -\tilde{S} & P_1 \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix}.$$

This shows that the matrices are unimodular, and inverse to each other,

$$\begin{bmatrix} P & -Q \\ S_1 & R_1 \end{bmatrix}^{-1} = \begin{bmatrix} \tilde{R} & Q_1 \\ -\tilde{S} & P_1 \end{bmatrix}.$$

Thus P_1^{-1} is the Schur complement of P and

$$\det(P) \det(P_1^{-1}) = \det \begin{bmatrix} P & -Q \\ S_1 & R_1 \end{bmatrix} = c$$

where $0 \neq c \in \mathbb{R}$. Thus $\det(P) = c \det(P_1)$, in particular, $\deg \det(P) = \deg \det(P_1)$. \square

It turns out below that the integer $d(H)$ from (8.5) equals the size of a minimal realization of a proper rational matrix H . For the proof, we need the concept of a row-reduced polynomial matrix.

Definition 8.11 The degree of a non-zero polynomial row vector is defined to be the highest power of s appearing in it with a non-zero coefficient, and the degree of the zero row is set to $-\infty$. For a non-singular polynomial matrix $P \in \mathbb{R}[s]^{p \times p}$, let $\delta_i(P)$ be the degree of the i -th row of P , for $i = 1, \dots, p$. Then P has a unique representation

$$P = SP_{hr} + L$$

where $S = \text{diag}(s^{\delta_1(P)}, \dots, s^{\delta_p(P)})$, $P_{hr} \in \mathbb{R}^{p \times p}$, and $L \in \mathbb{R}[s]^{p \times p}$ is such that $\delta_i(L) < \delta_i(P)$. One calls P_{hr} the **highest row coefficient matrix**. If P_{hr} is non-singular, we say that P is **row-reduced** (or row-proper). For $R \in \mathbb{R}[s]^{p \times q}$ with full row rank, the highest row coefficient matrix $R_{hr} \in \mathbb{R}^{p \times q}$ is defined analogously, and R is called row-reduced if R_{hr} has full row rank.

Lemma 8.12 Let $P \in \mathbb{R}[s]^{p \times p}$ be non-singular.

1. We have

$$\sum_{i=1}^p \delta_i(P) \geq \deg \det(P).$$

The matrix P is row-reduced if and only if

$$\sum_{i=1}^p \delta_i(P) = \deg \det(P).$$

2. If $H = P^{-1}Q$ is strictly proper, then $\delta_i(Q) < \delta_i(P)$ for $i = 1, \dots, p$; if H is proper, then $\delta_i(Q) \leq \delta_i(P)$ for all i . The converse is also true, provided that P is row-proper.
3. For every non-singular $P \in \mathbb{R}[s]^{p \times p}$ there exists a unimodular matrix $U \in \mathbb{R}[s]^{p \times p}$ such that UP is row-proper.
4. For every $R \in \mathbb{R}[s]^{p \times q}$ with full row rank, there exists a unimodular matrix $U \in \mathbb{R}[s]^{p \times p}$ such that UR is row-proper.

Proof:

1. Rewrite $P = SP_{hr} + L$ as $P_{hr} = S^{-1}P - S^{-1}L$. Consider the limit as $s \rightarrow \infty$. Then we have, since $\lim_{s \rightarrow \infty} S^{-1}L = 0$,

$$P_{hr} = \lim_{s \rightarrow \infty} S^{-1}P$$

and, putting $\delta(P) := \sum_{i=1}^p \delta_i(P)$,

$$\det(P_{hr}) = \lim_{s \rightarrow \infty} \frac{\det(P)}{s^{\delta(P)}}.$$

This shows that $\deg \det(P) \leq \delta(P)$ and

$$\det(P_{hr}) = 0 \quad \Leftrightarrow \quad \deg \det(P) < \delta(P).$$

2. Let $H = P^{-1}Q$, then

$$Q_{ij} = \sum_{k=1}^p P_{ik} H_{kj}$$

and

$$\frac{Q_{ij}}{s^{\delta_i(P)}} = \sum_{k=1}^p \frac{P_{ik}}{s^{\delta_i(P)}} H_{kj}.$$

Consider again the limit as $s \rightarrow \infty$. If H is strictly proper, the right hand side tends to zero, and hence all the powers of s appearing in Q_{ij} must be strictly less than $\delta_i(P)$. Since this holds for all j , we obtain $\delta_i(Q) < \delta_i(P)$.

For the converse, write $P = SP_{hr} + L$, and assume that P_{hr} is invertible. Then

$$P^{-1}Q = (SP_{hr} + L)^{-1}Q = (I + P_{hr}^{-1}S^{-1}L)^{-1}P_{hr}^{-1}S^{-1}Q$$

and thus, since $\lim_{s \rightarrow \infty} S^{-1}L = 0$ and $\lim_{s \rightarrow \infty} S^{-1}Q = 0$, we have

$$\lim_{s \rightarrow \infty} P^{-1}Q = 0,$$

that is, $P^{-1}Q$ is strictly proper. The argument for “proper” is similar.

3. If P is row-proper, we are finished. Therefore, assume otherwise, that is, let $\det(P_{hr}) = 0$. We show that there exists a unimodular matrix U such that

$$\deg \det(P) = \deg \det(UP) \leq \sum_{j=1}^p \delta_j(UP) < \sum_{j=1}^p \delta_j(P).$$

Iteratively, this yields the result. We write $\delta_j := \delta_j(P)$ for simplicity. Since $\det(P_{hr}) = 0$, there exists $0 \neq \alpha = (\alpha_1, \dots, \alpha_p) \in \mathbb{R}^{1 \times p}$ with

$$\sum_{j=1}^p \alpha_j P_{hr}^{(j)} = 0$$

where $P_{hr}^{(j)}$ denotes the j -th row of P_{hr} . Among the j with $\alpha_j \neq 0$, select j^* with

$$\delta_{j^*} \geq \delta_j \quad \text{for all } j \text{ with } \alpha_j \neq 0.$$

Without loss of generality, let $\alpha_{j^*} = 1$. Then we have

$$P_{hr}^{(j^*)} + \sum_{j \neq j^*} \alpha_j P_{hr}^{(j)} = 0.$$

Now perform the elementary operation: $P^{(j^*)}$ plus $\sum_{j \neq j^*} \alpha_j s^{\delta_{j^*} - \delta_j} P^{(j)}$. The new matrix $P' = UP$ satisfies $\delta_{j^*}(P') < \delta_{j^*}(P)$ and $\delta_j(P') = \delta_j(P)$ for all $j \neq j^*$. This establishes the claim.

4. The final statement is analogous to the previous one. □

Given a full-row-rank representation R of \mathcal{B} , we can thus assume without loss of generality that R is row-proper. Then R_{hr} has full row rank, and thus it contains a sub-matrix that is square and invertible. If we choose an input-output decomposition of \mathcal{B} such that the output corresponds to such a choice of the columns of R , then one can show that in the resulting representation (2.6), the matrix N_2 is invertible (recall that this assumption was needed to transform (2.6) into state space form). This shows that any $\mathcal{B} = \{w \in \mathcal{A}^q \mid R(\frac{d}{dt})w = 0\}$ admits a partition of its signal components into inputs and outputs such that it has a state space representation, in particular, the resulting transfer matrix is proper.

Theorem 8.13 Let H be a proper rational matrix. The size of a minimal realization of H is given by the integer $d(H)$ from (8.5).

Proof: There is no loss of generality in assuming that H is strictly proper. Let $(A, B, C, 0)$ be a realization of H , that is,

$$H = C(sI_n - A)^{-1}B.$$

Define $G := (sI_n - A)^{-1}B$. Then $(sI_n - A, B)$ is a left factorization of G . Thus

$$\deg \det(sI_n - A) = n \geq d(G).$$

On the other hand, if (Q, P) is a right factorization of G , that is, $G = QP^{-1}$, then $H = CG = CQP^{-1}$, that is, (CQ, P) is a right factorization of H . We conclude that

$$d(G) \geq d(H).$$

This shows that $n \geq d(H)$, that is, the size of any realization of H must be at least $d(H)$. Conversely, we show that a strictly proper H possesses a realization of size $d(H)$. Let $H = P^{-1}Q$ be a left coprime factorization. Without loss of generality, let P be row-proper. Now consider the i -th row of P and Q , and denote them by P_i and Q_i , respectively. According to Theorem 2.22, there exist matrices K_i, L_{ij}, M_i, N_{ij} ($j = 1, 2$) such that

$$P_i\left(\frac{d}{dt}\right)y = Q_i\left(\frac{d}{dt}\right)u \quad \Leftrightarrow \quad \exists x : \begin{cases} \frac{d}{dt}x_i = K_i x_i + L_{i1}u + L_{i2}y \\ 0 = M_i x_i + N_{i1}u + N_{i2}y. \end{cases}$$

Here K_i can be chosen to be a $\delta_i(P) \times \delta_i(P)$ matrix. Combining these representations via

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix} \quad K = \begin{bmatrix} K_1 & & \\ & \ddots & \\ & & K_p \end{bmatrix} \quad L_j = \begin{bmatrix} L_{1j} \\ \vdots \\ L_{pj} \end{bmatrix}$$

$$M = \begin{bmatrix} M_1 & & \\ & \ddots & \\ & & M_p \end{bmatrix} \quad N_j = \begin{bmatrix} N_{1j} \\ \vdots \\ N_{pj} \end{bmatrix}$$

we obtain

$$P\left(\frac{d}{dt}\right)y = Q\left(\frac{d}{dt}\right)u \quad \Leftrightarrow \quad \exists x : \begin{cases} \frac{d}{dt}x = Kx + L_1u + L_2y \\ 0 = Mx + N_1u + N_2y \end{cases}$$

where the size of K is $\sum_{i=1}^p \delta_i(P) = \deg \det(P) = d(H)$. Moreover, $N_1 = 0$ (because $P^{-1}Q$ is strictly proper) and $N_2 = P_{hr}$, which is invertible. Thus we

can find, as in Section 2.7, matrices A, B, C, D , with A of the same size as K , such that

$$P\left(\frac{d}{dt}\right)y = Q\left(\frac{d}{dt}\right)u \quad \Leftrightarrow \quad \exists x : \begin{cases} \frac{d}{dt}x &= Ax + Bu \\ y &= Cx. \end{cases}$$

Then $H = P^{-1}Q = C(sI - A)^{-1}B$ which shows that H has a realization of size $d(H)$. \square

Remark 8.14 In the final step of the proof, we have used the fact that equivalent representations have the same transfer function. This statement has not been proven. Alternatively, one may give a direct proof of

$$H = P^{-1}Q = C(sI - A)^{-1}B$$

using the special form of A, B, C that results from Section 2.7. The details are omitted.

Theorem 8.15 (McMillan form) For each rational matrix $H \in \mathbb{R}(s)^{p \times m}$ there exist unimodular matrices $U \in \mathbb{R}[s]^{p \times p}$ and $V \in \mathbb{R}[s]^{m \times m}$ such that

$$UHV = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} \quad (8.6)$$

where $D \in \mathbb{R}(s)^{r \times r}$ is a diagonal matrix

$$D = \begin{bmatrix} \frac{\gamma_1}{\delta_1} & & \\ & \ddots & \\ & & \frac{\gamma_r}{\delta_r} \end{bmatrix}$$

with polynomials $\gamma_i, \delta_i \neq 0$ such that each pair (γ_i, δ_i) is coprime and $\gamma_1 | \dots | \gamma_r$ and $\delta_r | \dots | \delta_1$. Clearly, $r = \text{rank}(H)$. The matrix on the right hand side of (8.6) is called **McMillan form** of H and the integer

$$n := \sum_{i=1}^r \deg(\delta_i)$$

is called **McMillan-degree** of H .

Theorem 8.16 The size of a minimal realization of H equals the McMillan-degree of H .

Proof: We show that $d(H)$ from (8.5) coincides with the McMillan-degree of H . Let

$$UHV = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix}$$

be the McMillan form of H , with $D = \text{diag}(\frac{\gamma_1}{\delta_1}, \dots, \frac{\gamma_r}{\delta_r})$. Define

$$\Gamma = \begin{bmatrix} \text{diag}(\gamma_1, \dots, \gamma_r) & 0 \\ 0 & 0 \end{bmatrix} \quad \text{and} \quad \Delta = \begin{bmatrix} \text{diag}(\delta_1, \dots, \delta_r) & 0 \\ 0 & I \end{bmatrix}.$$

Then $UHV = \Gamma\Delta^{-1}$, that is, (Γ, Δ) is a right factorization of UHV . Since each pair (γ_i, δ_i) is coprime, it is even a right coprime factorization of UHV . Then $(U^{-1}\Gamma, V\Delta)$ is a right coprime factorization of H . Therefore

$$d(H) = \deg \det(V\Delta) = \deg \det(\Delta) = \deg \prod_{i=1}^r \delta_i = \sum_{i=1}^r \deg(\delta_i)$$

which is precisely the McMillan-degree of H . □

8.3 Poles

A complex number λ is called a **pole** of a rational matrix H if it is a pole of one of its entries. Equivalently, the poles of H are the zeros of the polynomials δ_i in the McMillan form of H . Still equivalently, they are the zeros of $\det(P)$, where P is the denominator matrix in a coprime factorization of H .

Theorem 8.17 Let (A, B, C, D) be a realization of H . Any pole of H is an eigenvalue of A . Conversely, an eigenvalue of A which is not a pole of H must be an uncontrollable mode of (A, B) or an unobservable mode of (A, C) . In particular, if (A, B, C, D) is a minimal realization of H , then the eigenvalues of A are precisely the poles of H .

Remark 8.18 This theorem shows that there is a close relation between the eigenvalues of A and the poles of $H = C(sI - A)^{-1}B + D$. This is the reason why one speaks of “pole shifting” in Section 5.3, although “eigenvalue shifting” would probably be more appropriate.

Proof: Without loss of generality, let H be strictly proper. Since

$$H = C(sI - A)^{-1}B = \frac{C \operatorname{adj}(sI - A)B}{\det(sI - A)}$$

any pole of H must be a zero of $\det(sI - A)$, that is, it must be an eigenvalue of A . Now let λ be an eigenvalue of A which is not a pole of H . Assume that λ is not an unobservable mode, that is,

$$\operatorname{rank} \begin{bmatrix} \lambda I - A \\ C \end{bmatrix} = n. \quad (8.7)$$

We need to show that λ is an uncontrollable mode. Let $G = (sI - A)^{-1}B = QP^{-1}$, where (Q, P) is right coprime. Then we have

$$BP = (sI - A)Q \quad (8.8)$$

and $H = CG = CQP^{-1}$, with

$$\operatorname{rank} \begin{bmatrix} CQ(\lambda) \\ P(\lambda) \end{bmatrix} = m. \quad (8.9)$$

To see this, let v be such that $CQ(\lambda)v = 0$ and $P(\lambda)v = 0$. We need to show that this implies $v = 0$. From (8.8),

$$0 = BP(\lambda)v = (\lambda I - A)Q(\lambda)v$$

and thus

$$\begin{bmatrix} \lambda I - A \\ C \end{bmatrix} Q(\lambda)v = 0$$

which implies that $Q(\lambda)v = 0$ because of (8.7). But then both $P(\lambda)v = 0$ and $Q(\lambda)v = 0$, which implies that $v = 0$ due to the right coprimeness of (Q, P) .

Now we must have $\det(P(\lambda)) \neq 0$. If conversely, $\det(P(\lambda)) = 0$, then there would exist a $v \neq 0$ such that $P(\lambda)v = 0$. Since λ is not a pole of H , the complex matrix $H(\lambda)$ is well-defined, and thus

$$0 = H(\lambda)P(\lambda)v = CQ(\lambda)v$$

and this would be a contradiction to (8.9).

Since λ is an eigenvalue of A , there exists $0 \neq z \in \mathbb{C}^{1 \times n}$ such that $z(\lambda I - A) = 0$. Then (8.8) implies $zBP(\lambda) = 0$ and hence, since $P(\lambda)$ is non-singular, we must have $zB = 0$. Then

$$z \begin{bmatrix} \lambda I - A & B \end{bmatrix} = 0$$

which shows that

$$\operatorname{rank} \begin{bmatrix} \lambda I - A & B \end{bmatrix} < n$$

that is, λ is an uncontrollable mode of (A, B) . \square

8.4 Zeros

A complex number λ is called a **zero** of a rational matrix H if it is not a pole of H and

$$\text{rank}(H(\lambda)) < \text{rank}(H).$$

Equivalently, the zeros of H are the zeros of the polynomials γ_i in the McMillan form of H provided that they are not also zeros of the δ_i . Equivalently, they are the λ with $\text{rank}(Q(\lambda)) < \text{rank}(Q)$ and $\det(P(\lambda)) \neq 0$ in a coprime factorization of H .

Let (A, B, C, D) be a realization of H . A complex number λ is called a **zero** of (A, B, C, D) if

$$\text{rank} \begin{bmatrix} A - \lambda I & B \\ C & D \end{bmatrix} < \text{rank} \begin{bmatrix} A - sI & B \\ C & D \end{bmatrix}.$$

Theorem 8.19 Let (A, B, C, D) be a realization of H . Any zero of H must be a zero of (A, B, C, D) . Conversely, a zero of (A, B, C, D) which is not a zero of H must be a pole of H or an uncontrollable mode of (A, B) or an unobservable mode of (A, C) . In particular, let (A, B, C, D) be a minimal realization of H and let λ be not a pole of H . Then λ is a zero of H if and only if it is a zero of (A, B, C, D) .

Proof: Over $\mathbb{R}(s)$, we have the Schur complement formula

$$\begin{bmatrix} A - sI & B \\ C & D \end{bmatrix} = \begin{bmatrix} I & 0 \\ C(A - sI)^{-1} & I \end{bmatrix} \begin{bmatrix} A - sI & 0 \\ 0 & H \end{bmatrix} \begin{bmatrix} I & (A - sI)^{-1}B \\ 0 & I \end{bmatrix}.$$

Thus we get

$$\text{rank} \begin{bmatrix} A - sI & B \\ C & D \end{bmatrix} = \text{rank}(A - sI) + \text{rank}(H).$$

Without loss of generality, let (A, B, C, D) be such that

$$\begin{bmatrix} A - sI & B \\ C & D \end{bmatrix} = \begin{bmatrix} A_1 - sI & A_2 & B_1 \\ 0 & A_3 - sI & 0 \\ C_1 & C_2 & D \end{bmatrix} \sim \begin{bmatrix} A_1 - sI & B_1 & A_2 \\ C_1 & D & C_2 \\ 0 & 0 & A_3 - sI \end{bmatrix},$$

where the equality comes from a Kalman decomposition, and \sim denotes unimodular equivalence of matrices. By construction, we have $H = C(sI - A)^{-1}B + D =$

$C_1(sI - A_1)^{-1}B_1 + D$, and (A_1, B_1) is controllable. Let $G := (sI - A_1)^{-1}B_1$ and let $G = QP^{-1}$ be a right coprime factorization. Then we have

$$B_1P = (sI - A_1)Q, \quad RP + SQ = I, \quad (A_1 - sI)R_1 + B_1S_1 = I$$

for some polynomial matrices R, S, R_1, S_1 . In matrix notation,

$$\begin{bmatrix} A_1 - sI & B_1 \\ S & R \end{bmatrix} \begin{bmatrix} R_1 & Q \\ S_1 & P \end{bmatrix} = \begin{bmatrix} I & 0 \\ * & I \end{bmatrix}$$

which yields

$$\begin{bmatrix} A_1 - sI & B_1 \\ S & R \end{bmatrix} \begin{bmatrix} \tilde{R}_1 & Q \\ \tilde{S}_1 & P \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix},$$

similarly as in an earlier proof. We may conclude that the matrices on the left hand side are unimodular. Since

$$\begin{bmatrix} A_1 - sI & B_1 \\ C_1 & D \end{bmatrix} \begin{bmatrix} \tilde{R}_1 & Q \\ \tilde{S}_1 & P \end{bmatrix} = \begin{bmatrix} I & 0 \\ * & C_1Q + DP \end{bmatrix},$$

we get

$$\begin{bmatrix} A_1 - sI & B_1 \\ C_1 & D \end{bmatrix} \sim \begin{bmatrix} I & 0 \\ 0 & C_1Q + DP \end{bmatrix}.$$

Coming back to the original state space system, we obtain

$$\begin{bmatrix} A - sI & B \\ C & D \end{bmatrix} \sim \begin{bmatrix} I & 0 & 0 \\ 0 & C_1Q + DP & * \\ 0 & 0 & A_3 - sI \end{bmatrix}.$$

Recalling that $G = QP^{-1}$, we have $H = C_1G + D = (C_1Q + DP)P^{-1}$. Hence the zeros of H are zeros of $C_1Q + DP$, which implies that they are also zeros of (A, B, C, D) .

Now let λ be a zero of (A, B, C, D) .

Case 1: λ is an eigenvalue of A and neither uncontrollable nor unobservable. Then λ is a pole of H .

Case 2: λ is an eigenvalue of A and uncontrollable or unobservable.

Case 3: λ is not an eigenvalue of A . Then it is not a pole of H and we may use the Schur complement formula from above with $s = \lambda$. Noting that $\text{rank}(A - \lambda I) = n = \text{rank}(A - sI)$, we see that $\text{rank}(H(\lambda)) < \text{rank}(H)$. Thus λ is a zero of H . \square

Remark 8.20 Interpretation of zeros: Assume that $\text{rank}(H) = m$. Then

$$\begin{bmatrix} A - sI & B \\ C & D \end{bmatrix}$$

has full column rank. Thus λ is a zero of (A, B, C, D) if and only if there exists $(x_0, u_0) \neq (0, 0)$ such that

$$\begin{bmatrix} A - \lambda I & B \\ C & D \end{bmatrix} \begin{bmatrix} x_0 \\ u_0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

Then the input $u(t) = u_0 e^{\lambda t}$ and the initial condition $x(0) = x_0$ lead to $x(t) = x_0 e^{\lambda t}$, and the corresponding output $y(t) = Cx(t) + Du(t)$ is identically zero. (We admit complex-valued signals, for simplicity; but the same argument holds for the real parts of u, x, y .) Thus the system “annihilates” the input u of frequency λ , and this is why λ can be seen as a zero of the system. If $u_0 = 0$, then $x_0 \neq 0$ is indistinguishable from zero. In particular, any unobservable eigenvalue of A w.r.t. C is a zero of (A, B, C, D) .

A similar interpretation is possible for $\text{rank}(H) = p$. Then any uncontrollable eigenvalue of A w.r.t. B is a zero of (A, B, C, D) .

Example 8.21 Let

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

Then $H = \frac{s}{s-1}$, which has a pole at 1 and a zero at 0. We have $\text{spec}(A) = \{1, 2\}$. The eigenvalue 2 is not a pole, and indeed one can check that it is uncontrollable (but not unobservable). Note that the eigenvalue 1 is both uncontrollable and unobservable, but still a pole of H . The zeros of (A, B, C, D) are 0, 1, 2: 0 is a zero of H , 1 is a pole, and 2 is an uncontrollable eigenvalue. A minimal realization of H is given by $A_1 = B_1 = C_1 = D = 1$. Then 1 is the only eigenvalue of A_1 and 0 is the only zero of (A_1, B_1, C_1, D) .

If we modify this example by taking A, B as above and

$$C' = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}, \quad D' = \begin{bmatrix} 1 \\ 0 \end{bmatrix},$$

then we have

$$H' = \begin{bmatrix} \frac{s}{s-1} \\ 0 \end{bmatrix}$$

and (A, B, C', D') has zeros at 0 and 1, but not at 2. This shows that an uncontrollable mode is not necessarily a zero of the realization if $\text{rank}(H) < p$.

Remark 8.22 In the tutorial, we have seen that the series interconnection of two state space systems is controllable if both systems are controllable and additionally, for all $\lambda \in \text{spec}(A_2)$, we have that

$$\begin{bmatrix} A_1 - \lambda I & B_1 \\ C_1 & D_1 \end{bmatrix}$$

has full row rank. This requirement amounts (more or less) to saying that no pole of H_2 should be a zero of H_1 (to avoid cancellation effects). In terms of the system trajectories, this means that no characteristic frequency of the second system should be blocked by the first system. A similar fact can be shown for $H = H_2H_1$ with $H_1 = P_1^{-1}Q_1$ and $H_2 = P_2^{-1}Q_2$ where the requirement was that for all λ with $\det(P_2(\lambda)) = 0$, the matrix $Q_1(\lambda)$ has full row rank. Again, this can be interpreted (roughly) as meaning that no pole of H_2 should be a zero of H_1 .

Appendix A

Background material from distribution theory

Let \mathcal{D} denote the set of all smooth functions $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ which have compact support. Recall that

$$\text{supp}(\varphi) = \text{cl}\{t \in \mathbb{R} \mid \varphi(t) \neq 0\},$$

where $\text{cl}(\cdot)$ denotes the closure with respect to the standard topology of \mathbb{R} (where the term “compact” coincides with “closed and bounded”). The elements of \mathcal{D} are called **test functions**. A **distribution** D is a linear, continuous (in a sense not to be specified here) functional defined on \mathcal{D} , that is, it assigns to each test function φ a real number $D(\varphi)$ such that $D(\lambda_1\varphi_1 + \lambda_2\varphi_2) = \lambda_1D(\varphi_1) + \lambda_2D(\varphi_2)$ for all $\lambda_1, \lambda_2 \in \mathbb{R}$ and all $\varphi_1, \varphi_2 \in \mathcal{D}$. The set of all distributions is denoted by \mathcal{D}' . (We restrict to real-valued test functions and distributions, although the complex case is analogous.)

A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is called **locally integrable** if

$$\int_I |f(t)| dt$$

exists for any compact interval $I \subset \mathbb{R}$. Each locally integrable function f defines a distribution D_f by

$$D_f(\varphi) := \int_{-\infty}^{\infty} f(t)\varphi(t) dt$$

for all $\varphi \in \mathcal{D}$. We say that f **generates** the distribution D_f . The distributions D_f , where f is locally integrable, are called **regular**, and one often identifies f with D_f , that is, we can interpret L^1_{loc} , the set of locally integrable functions, as a

subset of \mathcal{D}' . Note that every continuous function is locally integrable, and hence, it can be seen as an element of \mathcal{D}' . However, there are also many non-continuous locally integrable functions, such as, for example, the Heaviside function. An important non-regular distribution is the Dirac delta distribution δ defined by

$$\delta(\varphi) := \varphi(0)$$

for all $\varphi \in \mathcal{D}$.

Linear combinations of distributions are defined by

$$(\lambda_1 D_1 + \lambda_2 D_2)(\varphi) := \lambda_1 D_1(\varphi) + \lambda_2 D_2(\varphi)$$

for $\lambda_1, \lambda_2 \in \mathbb{R}$, and thus \mathcal{D}' becomes a real vector space. One can multiply a distribution D by a smooth function $a : \mathbb{R} \rightarrow \mathbb{R}$ via

$$(aD)(\varphi) := D(a\varphi)$$

for all $\varphi \in \mathcal{D}$. This definition is motivated by the requirement $aD_f \stackrel{!}{=} D_{af}$. Note that we need the smoothness of a to guarantee that $a\varphi$ is again a test function. Thus, \mathcal{D}' becomes a module over the ring \mathcal{C}^∞ of smooth functions from \mathbb{R} to \mathbb{R} .

The derivative of a distribution is defined by

$$\dot{D}(\varphi) := -D(\dot{\varphi}) \tag{A.1}$$

for all $\varphi \in \mathcal{D}$. For well-definedness, note that the derivative of a test function is again a test function. This definition is motivated by the law of partial integration. Let D_f be a regular distribution, and assume that f is continuously differentiable. Then the distribution $D_{\dot{f}}$ is well-defined, and we certainly want $\dot{D}_f \stackrel{!}{=} D_{\dot{f}}$ and thus we put

$$\dot{D}_f(\varphi) = D_{\dot{f}}(\varphi) = \int_{-\infty}^{\infty} \dot{f}(t)\varphi(t)dt = f\varphi \Big|_{-\infty}^{\infty} - \int_{-\infty}^{\infty} f(t)\dot{\varphi}(t)dt = -D_f(\dot{\varphi}).$$

For example, the derivative of the distribution D_h generated by the Heaviside function h is the Dirac delta distribution, because

$$\dot{D}_h(\varphi) = -D_h(\dot{\varphi}) = - \int_{-\infty}^{\infty} h(t)\dot{\varphi}(t)dt = - \int_0^{\infty} \dot{\varphi}(t)dt = -\varphi \Big|_0^{\infty} = \varphi(0) = \delta(\varphi)$$

for all $\varphi \in \mathcal{D}$. Note that according to (A.1), distributions can be differentiated arbitrarily often: by repeated application of (A.1), we get $D^{(k)}(\varphi) = (-1)^k D(\varphi^{(k)})$ for all $k \in \mathbb{N}$. The distributional derivative (A.1) provides a generalization of the classical concept of differentiability. For instance, any continuous function can

be differentiated in the distributional sense, but of course, not necessarily in the classical sense.

We would like to use a similar trick for integration. However, a primitive function of a test function $\varphi \in \mathcal{D}$, for instance,

$$\psi(t) := \int_{-\infty}^t \varphi(\tau) d\tau$$

is *not* a test function, in general. This is because

$$\psi(\infty) := \lim_{t \rightarrow \infty} \psi(t) = \int_{-\infty}^{\infty} \varphi(\tau) d\tau$$

is non-zero in general, and thus, ψ may not have compact support.

Therefore, the set \mathcal{D}_0 is introduced as the set of all test functions which have a primitive function that is again a test function. More precisely, let $\mathcal{D}_0 \subset \mathcal{D}$ denote the space of all test functions φ_0 which satisfy the following equivalent conditions:

1. There exists $\psi \in \mathcal{D}$ such that $\varphi_0 = \dot{\psi}$.
2. $\int_{-\infty}^{\infty} \varphi_0(t) dt = 0$.

In that case, a primitive function $\psi \in \mathcal{D}$ of φ_0 is given by $\psi(t) = \int_{-\infty}^t \varphi_0(\tau) d\tau$, in fact, this is the only primitive in \mathcal{D} , because adding a non-zero constant will destroy the compact support property.

Now let α be a fixed test function with $\int_{-\infty}^{\infty} \alpha(t) dt = 1$. Then for any $\varphi \in \mathcal{D}$, let $\lambda = \int_{-\infty}^{\infty} \varphi(t) dt$ and set

$$\varphi_0 := \varphi - \lambda\alpha.$$

Then $\varphi_0 \in \mathcal{D}_0$, because $\int_{-\infty}^{\infty} \varphi_0(t) dt = \int_{-\infty}^{\infty} \varphi(t) dt - \lambda \int_{-\infty}^{\infty} \alpha(t) dt = 0$.

Lemma A.1 The equation $\dot{\xi} = 0$ has no distributional solutions apart from the classical solutions, i.e., the constant functions $\xi(t) = \xi_0$ for all t , where ξ_0 is a real number (more precisely, the distributional solutions are the regular distributions generated by constant functions.)

Proof: Let $\xi \in \mathcal{D}'$ be such that $\dot{\xi} = 0$. Consider first a test function $\varphi_0 \in \mathcal{D}_0$, say, with $\varphi_0 = \dot{\psi}$, then

$$\xi(\varphi_0) = \xi(\dot{\psi}) = -\dot{\xi}(\psi) = 0.$$

Now let $\varphi \in \mathcal{D}$ be arbitrary, and define $\varphi_0 = \varphi - \lambda\alpha$ as above. Then

$$0 = \xi(\varphi_0) = \xi(\varphi) - \lambda\xi(\alpha)$$

and thus

$$\xi(\varphi) = \lambda\xi(\alpha) = \int_{-\infty}^{\infty} c\varphi(t)dt = D_c(\varphi),$$

where the constant c is given by $c := \xi(\alpha)$. Thus $\xi = D_c$, which is the distribution generated by the constant function with value c . \square

Theorem A.2 Let g be a distribution. Then there exists a distribution G with $\dot{G} = g$. We call G a primitive (distribution) of g .

Proof: Let $g \in \mathcal{D}'$ be given. For $\varphi \in \mathcal{D}$, we define $\varphi_0 \in \mathcal{D}_0$ as above and we set

$$\psi(t) = \int_{-\infty}^t \varphi_0(\tau)d\tau.$$

Then $\psi \in \mathcal{D}$. In particular, if $\varphi = \dot{\phi}$ for some $\phi \in \mathcal{D}$, then $\varphi_0 = \varphi$ and $\psi = \phi$. We define

$$G(\varphi) := -g(\psi)$$

for all $\varphi \in \mathcal{D}$. This is linear and continuous, and for any $\phi \in \mathcal{D}$, we have

$$\dot{G}(\phi) = -G(\dot{\phi}) = g(\phi)$$

as desired. \square

According to Lemma A.1, primitives of distributions are unique up to additive constants. If $g = D_f$ is regular, then a primitive of D_f is given by D_F , where

$$F(t) = \int_0^t f(\tau)d\tau$$

is a primitive function of f .

So far, we have only dealt with distributions in $\mathcal{D}' = \mathcal{D}'(\mathbb{R})$, corresponding to the time set $T = \mathbb{R}$. Similarly, one can define $\mathcal{D}'(U)$ for any open set $U \subseteq \mathbb{R}$. This is called the set of distributions on U (which actually means that the relevant test functions are defined on U). For $T = \mathbb{R}_+ = [0, \infty)$, one constructs $\mathcal{D}'(\mathbb{R}_+)$ as the set of all distributions on some (arbitrarily small) open neighborhood of $[0, \infty)$.

Appendix B

Jordan form

For every matrix $A \in \mathbb{R}^{n \times n}$, there exists an invertible matrix $T \in \mathbb{C}^{n \times n}$ such that

$$T^{-1}AT = J = \begin{bmatrix} J_1 & & \\ & \ddots & \\ & & J_k \end{bmatrix} \quad (\text{B.1})$$

where each matrix $J_i \in \mathbb{C}^{n_i \times n_i}$ has the form

$$J_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix}.$$

The matrix on the right hand side of (B.1) is called **Jordan form** of A . The complex numbers $\lambda_1, \dots, \lambda_k$ are the eigenvalues of A . The number of matrices J_i having a particular eigenvalue λ on their main diagonal coincides with the geometric multiplicity of that eigenvalue, and the sum of the sizes of these blocks is precisely the algebraic multiplicity of λ .

We have

$$A^t = TJ^tT^{-1} \quad \text{and} \quad e^{At} = Te^{Jt}T^{-1}.$$

Using

$$J^t = \begin{bmatrix} J_1^t & & \\ & \ddots & \\ & & J_k^t \end{bmatrix} \quad \text{and} \quad e^{Jt} = \begin{bmatrix} e^{J_1t} & & \\ & \ddots & \\ & & e^{J_kt} \end{bmatrix}$$

as well as

$$J_i^t = \begin{bmatrix} \lambda_i^t & \binom{t}{1}\lambda_i^{t-1} & \cdots & \binom{t}{n_i-1}\lambda_i^{t-n_i+1} \\ & \ddots & \ddots & \vdots \\ & & \ddots & \binom{t}{1}\lambda_i^{t-1} \\ & & & \lambda_i^t \end{bmatrix}$$

and

$$e^{J_i t} = \begin{bmatrix} e^{\lambda_i t} & te^{\lambda_i t} & \cdots & \frac{t^{n_i-1}}{(n_i-1)!}e^{\lambda_i t} \\ & \ddots & \ddots & \vdots \\ & & \ddots & te^{\lambda_i t} \\ & & & e^{\lambda_i t} \end{bmatrix}$$

we can see that the entries of A^t and e^{At} have the form

$$\sum_{\lambda} a_{\lambda}(t)\lambda^t \quad \text{and} \quad \sum_{\lambda} a_{\lambda}(t)e^{\lambda t}$$

respectively, where λ are the eigenvalues of A and a_{λ} are polynomials in t .

Strictly speaking, the formula for the discrete case holds only if A is invertible. The reason is that the eigenvalue zero (which is present if and only if A is not invertible) plays a special role in the discrete case. Jordan blocks with zero on the main diagonal are nilpotent, and thus they contribute nothing to A^t for large enough t . Thus, if A is not invertible, the formula is still correct for t that are large enough. However, this is quite sufficient, e.g., for stability analysis, because this is anyhow concerned with the behavior of A^t for large t , and for $t \rightarrow \infty$.

If one restricts to real transformation matrices $T \in \mathbb{R}^{n \times n}$, the real Jordan form can be achieved. Note that since A is real, the non-real eigenvalues come in pairs of complex conjugate numbers. For an eigenvalue pair $a_i \pm ib_i$, where $a_i, b_i \in \mathbb{R}$, the real Jordan blocks take the form

$$J_i = \begin{bmatrix} \Lambda_i & I_2 & & \\ & \ddots & \ddots & \\ & & \ddots & I_2 \\ & & & \Lambda_i \end{bmatrix}, \quad \text{where } \Lambda_i = \begin{bmatrix} a_i & b_i \\ -b_i & a_i \end{bmatrix}.$$

Appendix C

Kronecker-Weierstraß form

Let $K, L \in \mathbb{R}^{n \times n}$ be matrices with $\det(sK - L) \neq 0$. Then there exist non-singular real matrices U, V such that

$$UKV = \begin{bmatrix} I & 0 \\ 0 & N \end{bmatrix} \quad \text{and} \quad ULV = \begin{bmatrix} A & 0 \\ 0 & I \end{bmatrix}$$

where N is a nilpotent matrix.

Proof: Let $\alpha \in \mathbb{R}$ be such that $\det(\alpha K - L) \neq 0$. Define

$$\hat{K} = (\alpha K - L)^{-1}K \quad \text{and} \quad \hat{L} = (\alpha K - L)^{-1}L.$$

Then

$$\alpha \hat{K} - \hat{L} = I. \tag{C.1}$$

Using the real Jordan form, there exists a non-singular matrix T such that

$$T^{-1}\hat{K}T = \begin{bmatrix} E_1 & 0 \\ 0 & E_2 \end{bmatrix}$$

where E_1 is non-singular and E_2 is nilpotent. Then $\alpha E_2 - I$ is non-singular (because zero is the only eigenvalue of a nilpotent matrix). Define

$$U = \begin{bmatrix} E_1^{-1} & 0 \\ 0 & (\alpha E_2 - I)^{-1} \end{bmatrix} T^{-1}(\alpha K - L)^{-1} \quad \text{and} \quad V = T.$$

These matrices are clearly non-singular and

$$UKV = \begin{bmatrix} I & 0 \\ 0 & (\alpha E_2 - I)^{-1}E_2 \end{bmatrix} \quad \text{and} \quad ULV = \begin{bmatrix} E_1^{-1} & 0 \\ 0 & (\alpha E_2 - I)^{-1} \end{bmatrix} T^{-1}\hat{L}T.$$

We set $N := (\alpha E_2 - I)^{-1} E_2$. Since $(\alpha E_2 - I)^{-1}$ and E_2 commute with each other,

$$N^k = (\alpha E_2 - I)^{-k} E_2^k$$

and thus N is nilpotent (since E_2 is nilpotent). On the other hand, (C.1) implies that

$$T^{-1} \hat{L} T = \alpha T^{-1} \hat{K} T - I = \begin{bmatrix} \alpha E_1 - I & 0 \\ 0 & \alpha E_2 - I \end{bmatrix}$$

and thus

$$ULV = \begin{bmatrix} E_1^{-1}(\alpha E_1 - I) & 0 \\ 0 & I \end{bmatrix}.$$

Thus we set $A := E_1^{-1}(\alpha E_1 - I)$ and we are finished. \square

Note that the size of $A \in \mathbb{R}^{\nu \times \nu}$ is uniquely determined by K, L , because

$$\nu = \deg(\det(sK - L)).$$

Similarly, the nilpotency index of N is uniquely determined by K, L , because it is equal to the degree of the polynomial part of the rational matrix $(sK - L)^{-1}$ plus one.

Appendix D

Smith form

The ring $\mathcal{P} = \mathbb{R}[s]$ is a **Euclidean domain**. This means that for any $a, b \in \mathcal{P}$, $b \neq 0$, there exist $c, d \in \mathcal{P}$ with

$$a = bc + d \tag{D.1}$$

where either $d = 0$ or $\deg(d) < \deg(b)$. The representation (D.1) is obtained by **division with remainder**.

By an **elementary operation**, we mean one of the following matrix transformations:

- interchanging two rows/columns of a matrix;
- multiplying a row/column by a unit (that is, a non-zero constant);
- adding a multiple of one row/column to another row/column.

It is easy to see that these operations correspond to multiplication by unimodular matrices from the left/right.

Let $R \in \mathcal{P}^{p \times q}$ be a matrix. Then there exist unimodular matrices $U \in \mathcal{P}^{p \times p}$ and $V \in \mathcal{P}^{q \times q}$ such that

$$URV = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} \tag{D.2}$$

where $D = \text{diag}(d_1, \dots, d_r)$ is a diagonal matrix, and $d_1 | d_2 | \dots | d_r$. The matrix on the right hand side of (D.2) is called the **Smith form** of R . The non-zero polynomials d_1, \dots, d_r are uniquely determined by R (up to multiplication by a non-zero constant).

Proof: Without loss of generality, let $R \neq 0$. It is sufficient to show that by elementary operations, R can be brought into the form

$$R' = \begin{bmatrix} a & 0 & \cdots & 0 \\ 0 & & & \\ \vdots & & Q & \\ 0 & & & \end{bmatrix} \quad (\text{D.3})$$

where a divides all entries of Q . Then one applies the same procedure to Q , and the result follows inductively.

Case 1: There exists i, j such that R_{ij} divides all entries of R . By a suitable interchange of rows and columns, this element can be brought into the (1,1) position of the matrix. Therefore without loss of generality, R_{11} divides all entries of R . Now perform the following elementary operations: for all $i \neq 1$, put i th row minus R_{i1}/R_{11} times 1st row; for all $j \neq 1$, put j th column minus R_{1j}/R_{11} times 1st column. Then we are finished.

Case 2: There is no i, j such that R_{ij} divides all entries of R . Let

$$\delta(R) := \min\{\deg(R_{ij}) \mid R_{ij} \neq 0\}.$$

Without loss of generality, $\deg(R_{11}) = \delta(R)$. We show that by elementary operations, we can transform R into R' with $\delta(R') < \delta(R)$. Then we obtain a strictly decreasing sequence $\delta(R) > \delta(R') > \delta(R'') > \dots \geq 0$. After finitely many steps, we arrive at zero, i.e., we obtain a matrix which has a unit as an entry, and thus we are in Case 1.

Case 2a: R_{11} does not divide all R_{1j}, R_{i1} , say, it does not divide R_{1k} . By the Euclidean algorithm, we can write

$$R_{1k} = R_{11}c + d$$

where $d \neq 0$ and $\deg(d) < \deg(R_{11})$. Perform the elementary operation: k th column minus c times 1st column. Then the new matrix R' has d in the (1, k) position and thus $\delta(R') < \delta(R)$ as desired.

Case 2b: R_{11} divides all R_{1j}, R_{i1} . Similarly as in Case 1, we can transform, by elementary operations, R into the form (D.3). If a divides all entries of Q , then we are finished. If there exists i, j such that a does not divide Q_{ij} , then we perform the elementary operation: 1st row plus $(i+1)$ st row. (Note that the $(i+1)$ st row of R' corresponds to the i th row of Q .) The new matrix has Q_{ij} in the (1, $j+1$) position and therefore we are in Case 2a. \square

Appendix E

McMillan form

Let $H \in \mathbb{R}(s)^{p \times m}$ be a rational matrix. Then there exist unimodular polynomial matrices $U \in \mathbb{R}[s]^{p \times p}$ and $V \in \mathbb{R}[s]^{m \times m}$ such that

$$UHV = \begin{bmatrix} D & 0 \\ 0 & 0 \end{bmatrix} \quad (\text{E.1})$$

where

$$D = \begin{bmatrix} \frac{\gamma_1}{\delta_1} & & \\ & \ddots & \\ & & \frac{\gamma_r}{\delta_r} \end{bmatrix}$$

for some non-zero polynomials γ_i, δ_i , where each pair (γ_i, δ_i) is coprime and $\gamma_1 | \gamma_2 | \dots | \gamma_r$ and $\delta_r | \dots | \delta_2 | \delta_1$. Clearly, $r = \text{rank}(H)$. The right hand side of (E.1) is called **McMillan form** of H . The number

$$n := \sum_{i=1}^r \deg(\delta_i)$$

is called **McMillan-degree** of H . If $\lambda \in \mathbb{C}$ is a zero of δ_{i+1} , then it is also a zero of δ_i , because $\delta_{i+1} | \delta_i$. Therefore, the **poles** of H are precisely the zeros of δ_1 .

Proof: Write $H = \frac{N}{d}$, where $N \in \mathbb{R}[s]^{p \times m}$ and $0 \neq d \in \mathbb{R}[s]$. Compute the Smith form of N , say

$$UNV = \begin{bmatrix} \tilde{D} & 0 \\ 0 & 0 \end{bmatrix} \quad (\text{E.2})$$

where U, V are unimodular and $\tilde{D} = \text{diag}(d_1, \dots, d_r)$ for some non-zero polynomials d_i . We divide (E.2) by d and put $D := \frac{\tilde{D}}{d} = \text{diag}(\frac{d_1}{d}, \dots, \frac{d_r}{d})$. Let γ_i, δ_i be

coprime polynomials with

$$\frac{d_i}{d} = \frac{\gamma_i}{\delta_i}.$$

Since $d_1|d_2|\dots|d_r$, we have $d_{i+1} = d_i e_i$ for some polynomials e_i , where $i = 1, \dots, r-1$. This implies

$$\gamma_{i+1} \delta_i = \gamma_i \delta_{i+1} e_i$$

for all i . Since δ_i and γ_i are coprime, γ_{i+1} must be a multiple of γ_i . Since δ_{i+1} and γ_{i+1} are coprime, δ_i must be a multiple of δ_{i+1} . \square

Appendix F

An optimal control problem

Consider $\dot{x} = Ax + Bu$, where (A, B) is controllable. Let $\varepsilon > 0$ and $\bar{x} \in \mathbb{R}^n$ be given. We wish to steer the system from state 0 to state \bar{x} in time ε . Moreover, we would like to do this with the smallest possible amount of energy, that is,

$$E(u) = \int_0^\varepsilon \|u(\tau)\|^2 d\tau \rightarrow \min!$$

Define

$$V(t, x) := x^T W(t)^{-1} x$$

where

$$W(t) = \int_0^t e^{A\tau} B B^T e^{A^T \tau} d\tau$$

is the controllability Gramian. Let us look at the change of $V(t, x(t))$ along a trajectory x of our system. We have

$$\begin{aligned} \frac{d}{dt} V(t, x(t)) &= \frac{d}{dt} x(t)^T W(t)^{-1} x(t) \\ &= \dot{x}(t)^T W(t)^{-1} x(t) + x(t)^T \left(\frac{d}{dt} W(t)^{-1} \right) x(t) + x(t)^T W(t)^{-1} \dot{x}(t). \end{aligned}$$

Note that for any matrix-valued function W ,

$$\frac{d}{dt} W^{-1} = -W^{-1} \dot{W} W^{-1}.$$

Moreover, we plug in $\dot{x} = Ax + Bu$ and we obtain (omitting the argument t wherever possible)

$$\begin{aligned} \frac{d}{dt} V(t, x) &= (Ax + Bu)^T W^{-1} x - x^T W^{-1} \dot{W} W^{-1} x + x^T W^{-1} (Ax + Bu) \\ &= x^T (A^T W^{-1} + W^{-1} A - W^{-1} \dot{W} W^{-1}) x + 2u^T B^T W^{-1} x. \end{aligned}$$

Consider the matrix

$$X(t) := A^T W(t)^{-1} + W(t)^{-1} A - W(t)^{-1} \dot{W}(t) W(t)^{-1}.$$

We have

$$\begin{aligned} W(t)X(t)W(t) &= W(t)A^T + AW(t) - \dot{W}(t) & (F.1) \\ &= \int_0^t (e^{A\tau} BB^T e^{A^T \tau} A^T + Ae^{A\tau} BB^T e^{A^T \tau}) d\tau - \dot{W}(t) \\ &= e^{A\tau} BB^T e^{A^T \tau} \Big|_0^t - \dot{W}(t) \\ &= e^{At} BB^T e^{A^T t} - BB^T - \dot{W}(t). \end{aligned}$$

Noting that by the definition of W ,

$$\dot{W}(t) = e^{At} BB^T e^{A^T t}$$

we obtain

$$W(t)X(t)W(t) = -BB^T \quad (F.2)$$

and hence

$$X(t) = -W(t)^{-1} BB^T W(t)^{-1}.$$

We use this to rewrite our expression for $\frac{d}{dt}V(t, x)$ and obtain

$$\begin{aligned} \frac{d}{dt}V(t, x) &= -x^T W^{-1} BB^T W^{-1} x + 2u^T B^T W^{-1} x \\ &= -\|B^T W^{-1} x\|^2 + 2\langle u, B^T W^{-1} x \rangle \\ &= \|u\|^2 - \|u - B^T W^{-1} x\|^2. \end{aligned}$$

Let's integrate this from 0 to ε , exploiting that $x(0) = 0$ and $x(\varepsilon) = \bar{x}$. Then

$$V(\varepsilon, \bar{x}) - V(0, 0) = \int_0^\varepsilon \|u(\tau)\|^2 d\tau - \int_0^\varepsilon \|u(\tau) - B^T W(\tau)^{-1} x(\tau)\|^2 d\tau$$

or

$$\bar{x}^T W(\varepsilon)^{-1} \bar{x} = E(u) - \int_0^\varepsilon \|u(\tau) - B^T W(\tau)^{-1} x(\tau)\|^2 d\tau \leq E(u). \quad (F.3)$$

This shows that

$$E(u) \geq E_{\min}(\varepsilon, \bar{x}) := \bar{x}^T W(\varepsilon)^{-1} \bar{x}.$$

Equality is achieved if and only if the integral in (F.3) vanishes, i.e., if

$$u(t) = B^T W(t)^{-1} x(t). \quad (F.4)$$

Plugging that into $\dot{x} = Ax + Bu$, we get

$$\dot{x}(t) = (A + BB^T W(t)^{-1})x(t).$$

Since we know that $x(\varepsilon) = \bar{x}$, the solution of this linear time-varying ordinary differential equation is uniquely determined for all $t > 0$. I claim that this solution is

$$\xi(t) = W(t)e^{A^T(\varepsilon-t)}W(\varepsilon)^{-1}\bar{x}.$$

This can easily be checked: We have $\xi(\varepsilon) = \bar{x}$ and

$$\dot{\xi}(t) = (\dot{W}(t) - W(t)A^T)e^{A^T(\varepsilon-t)}W(\varepsilon)^{-1}\bar{x}.$$

Combining (F.1) with (F.2), we see that

$$\dot{W}(t) = W(t)A^T + AW(t) + BB^T.$$

This implies

$$\begin{aligned}\dot{\xi}(t) &= (AW(t) + BB^T)e^{A^T(\varepsilon-t)}W(\varepsilon)^{-1}\bar{x} \\ &= (A + BB^TW(t)^{-1})W(t)e^{A^T(\varepsilon-t)}W(\varepsilon)^{-1}\bar{x} \\ &= (A + BB^TW(t)^{-1})\xi(t)\end{aligned}$$

as desired. Thus $x = \xi$ is the optimal state trajectory. Then, according to (F.4),

$$u(t) = B^TW(t)^{-1}\xi(t) = B^Te^{A^T(\varepsilon-t)}W(\varepsilon)^{-1}\bar{x}$$

is the minimum energy control function that steers the system from 0 to \bar{x} in time ε .

Now let $0 < \varepsilon < \delta$. Then

$$W(\varepsilon) < W(\delta)$$

(we write $P < Q$ if $Q - P$ is positive definite) which implies that

$$W(\varepsilon)^{-1} > W(\delta)^{-1}$$

and hence

$$E_{\min}(\varepsilon, \bar{x}) = \bar{x}^TW(\varepsilon)^{-1}\bar{x} > \bar{x}^TW(\delta)^{-1}\bar{x} = E_{\min}(\delta, \bar{x})$$

for all $\bar{x} \neq 0$. Thus one needs more energy for doing the transition from 0 to \bar{x} in time ε than in time δ . This explains the trade-off between the speed of control on the one hand and the energy consumption of control on the other.

Bibliography

- [1] D. Hinrichsen, A. J. Pritchard, *Mathematical Systems Theory I*, Springer, 2005.
- [2] T. Kailath, *Linear Systems*, Prentice-Hall, 1980.
- [3] H. W. Knobloch, H. Kwakernaak, *Lineare Kontrolltheorie*, Springer, 1985.
- [4] J. W. Polderman, J. C. Willems, *Introduction to Mathematical Systems Theory*, Texts in Applied Mathematics 26, Springer, 1998.
- [5] H. H. Rosenbrock, *State-space and Multivariable Theory*, Nelson, 1970.
- [6] E. D. Sontag, *Mathematical Control Theory*, Texts in Applied Mathematics 6, Springer, 1990.
- [7] W. M. Wonham, *Linear Multivariable Control*, Applications of Mathematics 10, Springer, 1979.