

UNIVERSIDADE FEDERAL DO CEARÁ  
DEPARTAMENTO DE ENGENHARIA DE TELEINFORMÁTICA  
CURSO DE GRADUAÇÃO EM ENGENHARIA DE COMPUTAÇÃO

# OPTIMAL CONTROL: AN APPLICATION TO A NON-ISOTHERMAL CONTINUOUS REACTOR

Otacílio Bezerra Leite Neto

FORTALEZA – CEARÁ  
2019





UNIVERSIDADE FEDERAL DO CEARÁ  
DEPARTAMENTO DE ENGENHARIA DE TELEINFORMÁTICA  
CURSO DE GRADUAÇÃO EM ENGENHARIA DE COMPUTAÇÃO

# OPTIMAL CONTROL: AN APPLICATION TO A NON-ISOTHERMAL CONTINUOUS REACTOR

**Autor**

Otacílio Bezerra Leite Neto

**Orientador**

Prof. Dr. Francesco Corona

**Co-orientador**

Profa. Dra. Michela Mulas

*Trabalho de Conclusão de Curso submetido à  
Coordenação do Programa de Graduação em  
Engenharia de Computação da Universidade  
Federal do Ceará como parte dos requisitos  
para a obtenção do grau de **Bacharel em  
Engenharia de Computação**.*

FORTALEZA – CEARÁ  
2019

# Contents

<b>List of Figures</b>	<b>iii</b>
<b>List of Tables</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Control Systems Engineering . . . . .	1
1.2 Chemical reactor Systems . . . . .	3
1.3 Motivation . . . . .	4
1.4 Objectives . . . . .	5
1.5 Chapters Organization . . . . .	5
<b>2 Dynamical System Analysis</b>	<b>7</b>
2.1 First Principles Models . . . . .	7
2.2 Canonical Representations for Dynamical Models . . . . .	11
2.3 Response Analysis in the Time Domain . . . . .	18
2.4 Similarity Transformations . . . . .	25
2.5 Stability, Controlability and Observability . . . . .	28
2.6 Response Analysis in the Frequency Domain . . . . .	31
<b>3 Controller Synthesis</b>	<b>33</b>
3.1 State Feedback Controllers . . . . .	33
3.2 Regulation and Reference Tracking . . . . .	36
3.3 Deterministic State Observers . . . . .	41
3.4 Properties of State-Feedback Controllers . . . . .	45
<b>4 Optimal Control and Estimation</b>	<b>49</b>
4.1 Formulation . . . . .	49
4.2 Linear Quadratic (LQ) Controllers . . . . .	52
4.3 Optimal State Estimators . . . . .	56
4.4 Linear Quadratic Gaussian (LQG) . . . . .	61
4.5 Stability and Robustness Analysis . . . . .	64
<b>5 Methodology</b>	<b>67</b>
5.1 Non-Isothermal Continuous Stirred Tank . . . . .	67
5.2 Simulation . . . . .	68
<b>6 Results and Discussion</b>	<b>71</b>
6.1 Dynamical Analysis . . . . .	71
6.2 Optimal Control . . . . .	76
<b>7 Conclusion</b>	<b>85</b>

<b>Bibliography</b>	<b>86</b>
<b>A Properties of Gaussian Distributions</b>	<b>91</b>

# List of Figures

1.1	The schematic of a simple Mass-Spring-Damper system and its dynamical response to being forced by an unitary constant force. . . . .	2
1.2	The two most popular configurations for connecting controllers to dynamical systems. . . . .	2
1.3	General graphical representations of common chemical reactor systems. . . . .	4
2.1	Schematic representations of industrial reactor tanks for (a) a simple isothermal process and (b) a non-isothermal process with heating/cooling system. . . . .	9
2.2	Graphical interpretation of (a) Input-Output models and (b) State-Spaces models. . . . .	12
2.3	Comparison between the nonlinear and linear response of the reactor system. . . . .	17
2.4	Simulation of the total response decomposition for the two first states of the model from (2.39). Each state response is shown separately (left) and together into a trajectory visualization (right), in which the steady-state point is indicated by a circle. . . . .	20
2.5	Simulation of the state-transition matrix $e^{At}$ from Example 2.4 for $t \in [0, 1.5]$ . . . . .	22
2.6	Simulation of the state-transition matrix $e^{At}$ from Example 2.5 for $t \in [0, 50]$ . . . . .	23
2.7	One-state forced response of the two systems from (2.60). . . . .	24
2.8	Eigenvalues of a system and the forced response associated with their modes. . . . .	25
2.9	Stability of forced responses given the positions of the system poles. . . . .	29
2.10	Illustration of the steady-state response of LTI systems to sinusoidal inputs. . . . .	31
2.11	Visualization of the two-states system in (2.75) in Bode plots (left) and Nyquist diagrams (right). The blue and orange lines represents the states $x_1$ and $x_2$ , respectively. . . . .	32
3.1	Block diagram of a state-feedback closed-loop system. . . . .	34
3.2	Simulation of three closed-loop systems from Example 3.1 showing the input signal (left) and state responses (right). . . . .	37
3.3	Block diagram of a state-feedback closed-loop system with integral action. . . . .	39
3.4	Simulation of three closed-loop systems showing the input signal (left) and state response (right) for reference tracking. The reference signal is indicated by the black dashed line. . . . .	41
3.5	Block diagram of a State-Space system with Luenberger observer. . . . .	42
3.6	Block diagram of a state-feedback closed-loop system with integral action and Luenberger observer. . . . .	44
3.7	Simplified block diagram of a perturbed state-feedback closed-loop. . . . .	45
3.8	Stability margins visualizations given Bode plots (left) and Nyquist diagrams (right) of closed-loop dynamical systems. . . . .	47

4.1	Illustration of the Bellman's Principle of Optimality for a system with discrete set of states for a discrete time evolution. Each column represents a time instance and each node represents a possible state. The straight lines represents state transitions given an action with associated costs, whereas the curves represents the trajectory from that state to the terminal state with associated optimal cost. The optimal trajectory between the initial and terminal state is shown in red. . .	52
4.2	Block diagram of a stochastic discrete State-Space system. . . . .	57
4.3	A Controlled Hidden Markov Model for a discrete-time system, where the shaded boxes indicates that the inputs are not random variables. . . . .	59
4.4	Block diagram of a Linear Quadratic Gaussian control configuration. . . . .	62
4.5	Simulation of the LQG controllers showing the input signal (left) and correspondent actual state responses (right). The dots denotes represents the measurements from the real system. . . . .	63
4.6	Nyquist diagram of a single state system with an open-loop unstable pole but which is closed-loop stable. The gain calculated by a LQR stabilizing controller never crosses the unit circle centered at $s = -1 + j0$ . . . . .	65
5.1	The non-isothermal continuous stirred tank reactor proposed. . . . .	67
6.1	Non-linear simulation showing the manipulated variables (left) and correspondent response of the system variables (right) for a pulse change in the input flow-rate. . . . .	71
6.2	Non-linear simulation showing the manipulated variables (left) and correspondent response of the system variables (right) for a pulse change in the cooling capacity. . . . .	72
6.3	Non-linear simulation showing the manipulated variables (left) and correspondent response of the system variables (right) for a change in the input flow-rate and cooling capacity. . . . .	73
6.4	Visualization of the steady-state values for the linearized model. . . . .	74
6.5	Simulation of both the nonlinear and linearized State-Space models. The nonlinear response is represented by the dashed lines. . . . .	75
6.6	Simulation of the $LQR_1$ controller for "starting" the tank reactor to the steady-state point. The reference is indicated by the black dashed lines, and the manipulated variables restrictions are indicated by the red dotted lines. . . . .	76
6.7	Simulation of the $LQR_2$ controller for "starting" the tank reactor to the steady-state point, but subjected to two disturbances pulses in the system variables. . . . .	77
6.8	Simulation of the $LQG_1$ controller for "starting" the tank reactor to the steady-state point. The dots represents the measurements of the output variables, and the colored dashed lines represents the state estimations. . . . .	78
6.9	Simulation of the $LQG_2$ controller for "starting" the tank reactor to the steady-state point, but subjected to two disturbances pulses in the system variables. . .	79
6.10	Simulation of the $LQRI_1$ for reference tracking. . . . .	80
6.11	Simulation of the $LQRI_2$ controller for reference tracking, but subjected to three pulse of disturbances in the system variables. . . . .	81
6.12	Simulation of the $LQGI_1$ controller for reference tracking. . . . .	82
6.13	Simulation of the $LQGI_2$ controller for reference tracking, but subjected to three pulse of disturbances in the system variables. . . . .	83

# List of Tables

2.1	Necessary and sufficient conditions for different classes of models. . . . .	13
5.1	Model parameters and variables nomenclature. . . . .	69
6.1	Comparison between the LQR and LQG controllers obtained. The pair $(LQR_2, LQG_2)$ comprises the simulations in which the system was subjected to disturbances. . .	78
6.2	Comparison between the LQRI and LQGI controllers. The pair $(LQRI_2, LQGI_2)$ comprises the simulations in which the system was subjected to disturbances. . .	82





# Chapter 1

## Introduction

This chapter presents the main problem in discussion and the basic concepts concerning its formulation and solutions, which are detailed further in the next chapters. This is a work on Control Theory and its application to Chemical Reactors, therefore the discussion will follow the notation common to the literature of this field, and the modern approach to this theory is explored.

The sections are organized as follows: Section 1.1 provides a general view on control systems engineering, Section 1.2 discuss chemical reactor systems and its importance in both industry and academia, Sections 1.3 and 1.4 describes the motivation and justification of this work, respectively, and Section 1.5 details the subsequent chapters in this document.

### 1.1 Control Systems Engineering

The discipline of Control Systems Engineering deals with the design of devices, named *controllers*, that are integrated to a physical system (a *dynamical system*, in most cases) in order to impose a desired behavior to this system. To achieve this goal, the discipline covers topics ranging from applied mathematics, such as dynamical systems theory and signal processing, to a more engineering discussion, regarding instrumentation and implementation of these controllers in a real-life plant or individual system.

A system, in a broad physical sense, is defined as an ensemble of interacting components that responds to external stimuli producing a determined dynamical response, and whose individual parts are not able to produce the same functionality by their own. Thus, the first essential element in Control Theory is a mathematical model of the system of interest. One such model is the *Input-Output Representation*, in which an input stimuli, a signal  $u(t)$ , acts on the system producing an output response, a signal  $y(t)$ , described by the following differential equation:

$$\alpha_n \frac{d^n y(t)}{dt^n} + \alpha_{n-1} \frac{d^{n-1} y(t)}{dt^{n-1}} + \cdots + \alpha_1 \frac{dy(t)}{dt} + \alpha_0 y(t) = \beta_m \frac{d^m u(t)}{dt^m} + \beta_{m-1} \frac{d^{m-1} u(t)}{dt^{m-1}} + \cdots + \beta_1 \frac{du(t)}{dt} + \beta_0 u(t) \quad (1.1)$$

In this representation, the input  $u(t)$  is called the *manipulated variable*, since it represents an arbitrary stimuli that can be given directly by human action or a by an automatic controller, while the output  $y(t)$  is called the *controlled variable*, since it can only be modified indirectly through  $u(t)$ . This also leads to a *cause-and-effect* interpretation of the system.

A model can provide a quantitative understanding of the system that is useful both to access some response specifications and to design controllers to modify them based on some requirements. In the case of the model in (1.1), it is possible to calculate the response  $y(t)$ , and its derivatives, resulting from any specific action  $u(t)$ . Besides, a model can be used to perform computer simulations, in order to visualize the dynamical behavior of the system without actually

manipulating it, since real experiments could be expensive or even damage the system. Consider, for instance, a schematic and a simulation for a model representing a mass-spring-damper system, shown at Fig. 1.1.

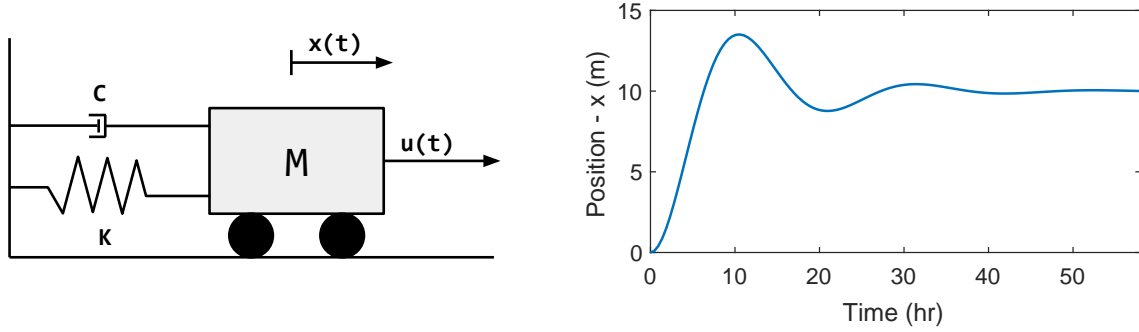


Figure 1.1: The schematic of a simple Mass-Spring-Damper system and its dynamical response to being forced by an unitary constant force.

In this simulation, quantities as the rise time, peak time, overshoot ratio and steady-state value of the visualized signal are examples of response specifications that can be defined to describe the system behavior to an external stimuli (in this case, a constant force of unit magnitude). These specified parameters are characteristic to responses of a class of systems known as *underdamped second-order systems*, that will be discussed further in the document.

A controller is used to calculate, for a time  $t \in [t_0, t_N]$ , the necessary input  $u(t)$  to produce an output  $y(t)$  as close as possible to a desired reference signal  $r(t)$ . There are two common configurations, shown in Fig. 1.2, of how to connect the controller to the system.

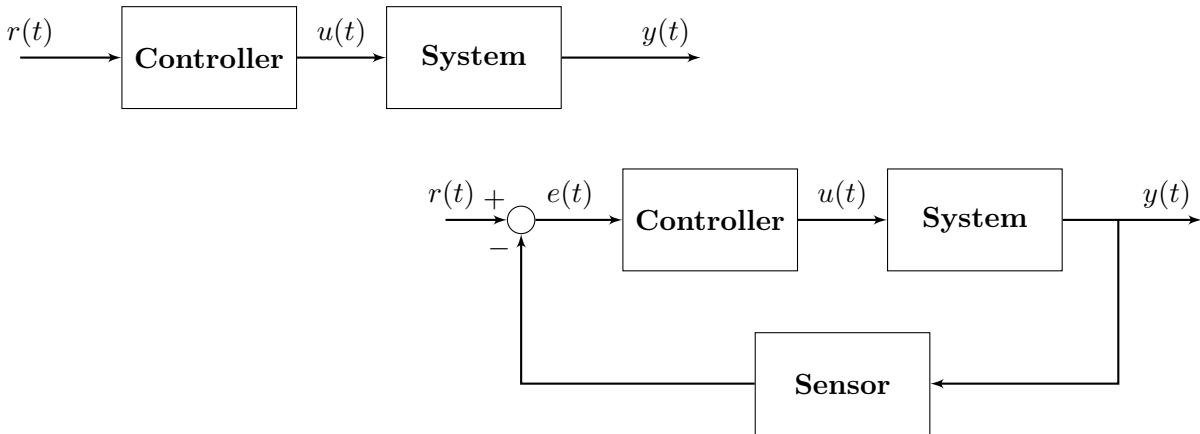


Figure 1.2: The two most popular configurations for connecting controllers to dynamical systems.

The configuration shown in the upper left of Fig. 1.2, known as the *Open-Loop Controller*, calculates the action as a function  $u(t) = \pi(r, t)$ , given an initial condition  $y(t_0) = y_0$ . In this case, the controller does not observe the output  $y(t)$ , and relies on the model to guarantee that the system is driven to the reference. Of course, if there are any external disturbances acting on this configuration, or if the model is not reliable enough, it is not possible to guarantee that the requirements are met. Thus, these type of controllers are not suitable for critical applications, and its use is restricted to systems where deviance from the desired reference can be tolerated.

In contrast, the configuration in the lower right of Fig. 1.2, known as the *Closed-Loop Controller* or *Feedback Controller*, calculates the action as a function  $u(t) = \pi(e, t)$ , where  $e(t) = r(t) - y(t)$  is the error between the reference and the actual response. Now, the controller

will observe the system output, through some sensor device, and compares it to the desired reference in order to calculate a *corrective action*. This feedback property can make the system reject disturbances while still driving it to the desired reference. Thus, the Feedback Controller became the most popular choice of controller configuration in industry for a wide range of applications, even for critical ones. [Syrmos et al., 1997].

## 1.2 Chemical reactor Systems

A chemical reaction, the transformation of a chemical substance into another, is a process central to chemistry and to nature itself. A reaction equation is an intuitive representation of such transformations. For instance, consider the following equation representing a *synthesis* reaction:



In this equation, the compounds  $A$  and  $2B$  forms the set of reactants,  $\mathcal{R}$ , while  $3C$  and  $D$  forms the set of products,  $\mathcal{P}$ . The coefficients in such equations are the stoichiometric numbers, providing an information about proportionality between the quantity of each substances in this reaction.

Usually the products can be directly used as reactants in another reaction, in which case they can also be referred as an intermediate product (or byproduct), and the equations can be appended in a “series” representation. In this case, each  $k$ -th intermediate product forms a set  $\mathcal{I}_k$ . In addition to a chain of series reactions, there is also the possibility of different reactions to occur in parallel, in the same system. The combination of these sets of reactants, byproducts, products and reactions are often referred as a *chemical reaction network*, and the associated equation can be represented in general form as:

$$\left\{ \begin{array}{ccccccc} \mathbb{R}^{(1)} & \longrightarrow & \mathcal{I}_1^{(1)} & \longrightarrow & \cdots & \longrightarrow & \mathcal{I}_{M_1}^{(1)} \longrightarrow \mathcal{P}^{(1)} \\ \mathbb{R}^{(2)} & \longrightarrow & \mathcal{I}_1^{(2)} & \longrightarrow & \cdots & \longrightarrow & \mathcal{I}_{M_2}^{(2)} \longrightarrow \mathcal{P}^{(2)} \\ \vdots & & \vdots & & \vdots & & \vdots \\ \mathbb{R}^{(N)} & \longrightarrow & \mathcal{I}_1^{(N)} & \longrightarrow & \cdots & \longrightarrow & \mathcal{I}_{M_N}^{(N)} \longrightarrow \mathcal{P}^{(N)} \end{array} \right. \quad (1.3)$$

Moreover, chemical reactions displays a dynamical behavior concerning the speed at which a reaction occurs. This rate of reaction, its *kinetics*, are dependent on the conditions in the environment, such as temperature and pressure, and on some properties of the reaction itself. In the case of a *isothermal process*, i.e., when the temperature in the environment remains constant, this rate can be calculated as a constant  $K$ , leading to a representation on the form:

$$\mathcal{R} \xrightarrow{K} \mathcal{P}. \quad (1.4)$$

When the temperature in the environment is not constant, the process is said to be endothermic or exothermic if, respectively, it consumes or produces energy. The kinetics of the reactions in such processes are usually functions of the temperature  $T$  which are assumed to follow the following form, known as the Arrhenius equation

$$K(T) = K_0 e^{\frac{-E/R}{T}}, \quad (1.5)$$

where  $K_0$  is the nominal kinetic rate,  $E$  is an activation energy associated with the reaction,  $R$  is the universal gas constant and  $T$  is the absolute temperature (in kelvins).

In practice, these chemical reactions are produced by mixing the reactants in some environment with adequate conditions. In order to control the quantities of these substances, actual processes consists in a manipulation of the concentrations of reactants in some container, usually by providing a mass flow of these substances through some fluid. A major interest is

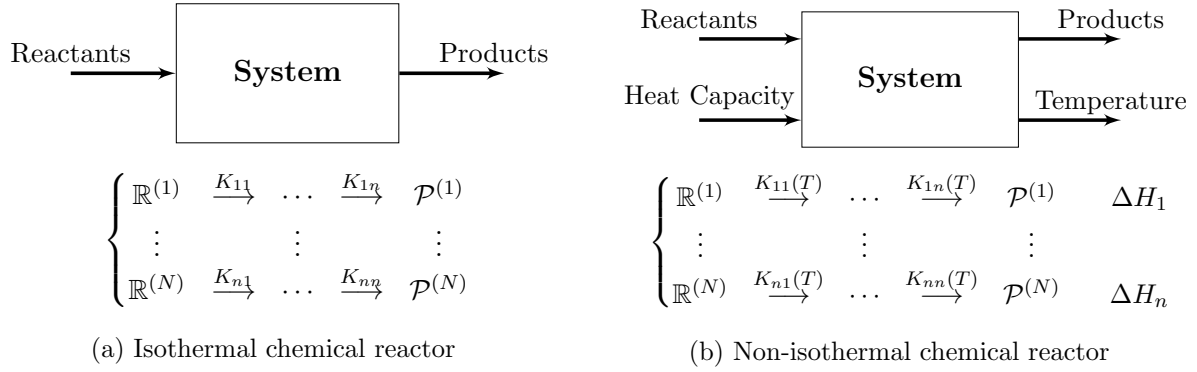


Figure 1.3: General graphical representations of common chemical reactor systems.

to manipulate the reactants in some way to produce a desired concentration of one or more products in the chemical reaction network, allowing this problem to be addressed by a control engineering perspective. A *chemical reactor system*, depicted in Fig. Fig. 1.3a, is a system where the concentration of some reactants is manipulated to produce a desired concentration of some products, given a chemical reaction network describing the reaction between these entities.

When the process is non-isothermal, the occurrence of a reaction contributes to the entropy of the environment, and consequently affects the kinetics of the subsequent ones. To compensate for this, practical applications also try to control the conditions in the environment using instruments external to the reactions themselves. Because of the use of the Arrhenius equation to model these reaction rates, this control is usually implemented through a cooling or heating system coupled to the original reactor system, resulting in the schematic on Fig. 1.3b. Moreover, note that the equation associated with non-isothermal reactions suffers a change of notation to include the value of entropy change, denoted  $\Delta H$ , resulting in the following general equation known as the *reaction heat equation*:

$$\mathcal{R} \xrightarrow{K(T)} \mathcal{P} \quad \Delta H. \quad (1.6)$$

### 1.3 Motivation

The use of automatic controllers to impose a desired behaviour to physical systems is a practice ubiquitous in many engineering fields. In the last years, the price of digital computers have been dropping while their performance have been growing. Consequently, digital controllers have become the central key in developments in important and innovative fields such as aeronautics [Eren et al., 2017, Goerzen et al., 2009, Wang and Su, 2015], autonomous driving [González et al., 2016, Falcone et al., 2007, Cairano et al., 2014] and several industrial applications [Qin and Badgwell, 2003, Tang et al., 2016, Holanda et al., 2008, Lucia et al., 2013]. In parallel, this theory is also useful to understand and bring inspiration from nature itself since, for instance, feedback controller characteristics were observed in the mechanisms for temperature regulation observed in vertebrate animals [Heller et al., 1978] and in [Todorov and Jordan, 2002].

Most recent developments in Control Theory focus on using Feedback Controllers to achieve *Robust and Optimal Control*. This theory accounts for the design of controllers that deals with uncertainty, either from the model or from the observation of the system, and are able to achieve the control objectives in a *optimal manner*. Despite being a few decades old, these fields have gained a lot of interest in the last years thanks to recent results in *Machine Learning* theory, particularly in *Reinforcement Learning*, which have been successful in applying optimization techniques for controlling artificial agents in environments loaded with uncertainty [Sutton and Barto, 2018, Mnih et al., 2015].

Furthermore, the specific application of controlling chemical reactor systems brings benefits from the fact that this class of dynamical system are present in many important biological and industrial processes, such as wastewater treatment [Mulas et al., 2015, Gupta et al., 2012]. In general, controllers can be used to guarantee safety constraints, maximize productivity and minimize the use of resources, in such a way that could not be achieved without automatic and high performing machinery. The inherent complexity and uncertainty existent in these systems also poses this control problem as a challenging research field still open to be explored.

## 1.4 Objectives

This work aims to provide a self-contained discussion of modeling and control of chemical reactor systems in the perspective of modern control, with focus on optimal control theory. Therefore, the results are focused on *state-feedback controllers* modeled in continuous time. Several properties of the dynamical models of such systems, both in the open-loop and feedback regime, are summarized in the document and the intention is to have a generalized mathematical framework to understand, evaluate and design those control systems. Finally, a specific problem concerning the control of a non-isothermal continuous reactor is explored to demonstrate an application of the theory developed throughout the document.

## 1.5 Chapters Organization

The chapters of this document are mainly organized in two parts. The first part, comprised by the chapters 2, 3 and 4, builds the necessary theoretical background and provides the mathematical framework for modeling and control of the discussed systems. The second part, comprised by the chapters 5 and 6, describes the experiments and results of applying these methods in real-world applications.

Individually, the chapters are organized as follows. Chapter 2 introduces the concept of dynamical models and their several properties with respect to a real system behavior. Chapter 3 discusses general results of state-feedback controllers and state observers. Chapter 4 presents the theory of optimal estimation and optimal control, in a broad perspective. Chapter 5 describes practical experiments in respect to a real-world continuous reactor system. Chapter 6 summarizes and discusses the results of these experiments, emphasizing the theory from the previous sections. Finally, chapter 7 concludes the document and highlight the most important results.



## Chapter 2

# Dynamical System Analysis

This chapter discusses the mathematical models for dynamical systems. The sections starts by introducing a procedure to build models from physical principles and presenting equivalent common representations. Next, the time evolution of dynamical systems is analyzed in relation to the mathematical structure of such models. Finally, some important properties of the systems are defined and discussed. A discussion of the system evolution in the frequency domain is also presented.

### 2.1 First Principles Models

A dynamical system is a physical system whose state evolves with time. One can represent a dynamical system using the *first principles* from physics, and formulate the evolution in time by calculating the rate of change of its state in respect to time. Thus, dynamical models are represented using differential equations with time derivatives.

A straightforward procedure to model a system consists of identifying the variables of interest and relate them using conservation laws, such as conservation of mass, conservation of energy or conservation of momentum. The resulting models are in the form:

$$\begin{pmatrix} \text{Rate of} \\ \text{accumulation of} \\ \text{Mass/Energy/Momentum} \end{pmatrix} = \begin{pmatrix} \text{Mass/Energy/Momentum} \\ \text{entering} \\ \text{the system} \end{pmatrix} - \begin{pmatrix} \text{Mass/Energy/Momentum} \\ \text{leaving} \\ \text{the system} \end{pmatrix}. \quad (2.1)$$

The choice of which conservation law to use depends on the system itself. Usually, conservation of mass is used to relate dynamics of concentrations and volumes, or other material variables, while conservation of momentum is used to relate dynamics of motion. Since energy can be converted on form, the conservation laws of this quantity can be used to model several dynamics, such as the rate of change in heat, electrical charges or velocity of a system.

In the case of a chemical reactor system, the variables of interest are the concentrations of the chemical substances in the system. Hence, the rate of accumulation of a substance can be represented using the mass conservation law, or mass balance:

$$\begin{aligned} \begin{pmatrix} \text{Accumulation} \\ \text{of mass} \\ \text{in the system} \end{pmatrix} &= \begin{pmatrix} \text{Mass} \\ \text{entering} \\ \text{the system} \end{pmatrix} - \begin{pmatrix} \text{Mass} \\ \text{leaving} \\ \text{the system} \end{pmatrix} \\ &= \left[ \begin{pmatrix} \text{Mass flow} \\ \text{entering} \\ \text{system} \end{pmatrix} + \begin{pmatrix} \text{Mass} \\ \text{produced} \\ \text{by reactions} \end{pmatrix} \right] - \left[ \begin{pmatrix} \text{Mass flow} \\ \text{leaving} \\ \text{system} \end{pmatrix} + \begin{pmatrix} \text{Mass} \\ \text{consumed} \\ \text{by reactions} \end{pmatrix} \right]. \end{aligned} \quad (2.2)$$



**Example 2.1.** (*Mass Balance of Reactors*) For the sake of illustration, consider the process of obtaining a first principles model for the concentration of a general substance in a reactor system, as follows. Consider a system such as the one depicted in 1.3a. First of all, one has to evaluate the mass flow through the system, as denoted in (2.2). The mass of any substance  $A$  entering and leaving the system, denoted respectively by  $M_{in}$  and  $M_{out}$ , are assumed to be provided by a flow of fluid carrying the substance. Given a fluid inflow  $F_{in}$ , with density  $\rho_{in}^{(A)}$ , and an outflow  $F_{out}$ , with density  $\rho_{out}^{(A)}$ , the mass flow is calculated by equations:

$$M_{in} = \rho_{in}^{(A)} F_{in} \quad ; \quad M_{out} = \rho_{out}^{(A)} F_{out}. \quad (2.3)$$

Now, one has to consider the quantities of mass that are consumed and produced by the reactions of the chemical reaction network associated to the system. In this case, consider the following reaction between two chemical compounds  $X$  and  $Y$ , with stoichiometric numbers  $\alpha$  and  $\beta$ :



Under the assumption that the reactant is in a dilute solution, the rate of this equation obeys the *law of mass action* [Horn and Jackson, 1972]. Given a constant kinetic rate  $K_{XY}$  and the volume of the solution as  $V$ , the mass of  $X$  consumed,  $M_{cons}^{(X)}$ , and the mass of  $Y$  produced,  $M_{prod}^{(Y)}$ , are given by the power-laws:

$$M_{cons}^{(X)} = \frac{V}{\beta} K_{XY} (\rho_X)^\alpha \quad ; \quad M_{prod}^{(Y)} = \frac{V}{\alpha} K_{XY} (\rho_X)^\alpha, \quad (2.5)$$

where  $\rho_X$  and  $\rho_Y$  are the respective densities of these compounds. Assuming that the network represents a set of reactions occurring within an chemical solution of volume  $V$ , the mass of a substance  $A$  that is consumed and produced by the reactions, named respectively  $M_{cons}$  and  $M_{prod}$ , are given by summing over the contribution of each reaction on the network where  $A$  is either a reactant or a product to any other compound  $X$ :

$$M_{cons} = V \sum_{\alpha X \rightarrow \beta X} \frac{1}{\beta} K_{AX} (\rho_A)^\alpha \quad ; \quad M_{prod} = V \sum_{\alpha X \rightarrow \beta A} \frac{1}{\beta} K_{XA} (\rho_X)^\alpha. \quad (2.6)$$

Finally, directly substituting these values in (2.2), the mass balance of any substance  $A$  in an isothermal chemical reactor system can be represented by the general dynamical model:

$$\begin{aligned} \left( \begin{array}{c} \text{Accumulation} \\ \text{of mass} \\ \text{in the system} \end{array} \right) &= \left[ \left( \begin{array}{c} \text{Mass flow} \\ \text{entering} \\ \text{System} \end{array} \right) + \left( \begin{array}{c} \text{Mass} \\ \text{produced} \\ \text{by reactions} \end{array} \right) \right] - \left[ \left( \begin{array}{c} \text{Mass flow} \\ \text{leaving} \\ \text{System} \end{array} \right) + \left( \begin{array}{c} \text{Mass} \\ \text{consumed} \\ \text{by reactions} \end{array} \right) \right] \\ \frac{d(\rho_A V)}{dt} &= \left[ \rho_{in}^{(A)} F_{in} + V \sum_{\alpha X \rightarrow \beta A} \frac{1}{\beta} K_{XA} (\rho_X)^\alpha \right] - \left[ \rho_{out}^{(A)} F_{out} + V \sum_{\alpha A \rightarrow \beta X} \frac{1}{\beta} K_{AX} (\rho_A)^\alpha \right]. \end{aligned} \quad (2.7)$$

Since the system is closed, i.e., there are no leaks or unknown sources of fluids, the assumptions of a constant volume implies that  $F_{in} = F_{out} = F$ . Normalizing each term by the volume and substituting a new variable  $q = F/V$  results in:

$$\frac{d(\rho_A)}{dt} = q(\rho_{in}^{(A)} - \rho_{out}^{(A)}) + \left( \sum_{\alpha X \rightarrow \beta A} \frac{1}{\beta} K_{XA} (\rho_X)^\alpha \right) - \left( \sum_{\alpha A \rightarrow \beta X} \frac{1}{\beta} K_{AX} (\rho_A)^\alpha \right). \quad (2.8)$$

Notice some important restrictions to the use of the model just presented. First of all, to calculate the mass contribution of an individual reaction was necessary to use a model which assumes that the reactor system is actually comprised of a dilute solution in some closed container. In industry, this means that the reactor system is actually a tank containing the solution. The inflow and outflow of fluid can be represented by flows through pipes which can be manipulated by some valve or pump. An illustration of such physical system is exhibited at Fig. 2.1a.

Furthermore, this model accounts for a single substance  $A$ , but the system is actually a solution of several compounds, each one with a specific concentration. From the model presented, it is visible that it is necessary to compute each concentration  $\rho_X$  before actually computing the rate of change in  $\rho_A$ . However, from the same model, the computation of the rate of change of any  $\rho_X$  may depend on  $\rho_A$  itself. Therefore, the change of concentration inside the whole system is actually the result of a system of differential equations:

$$\begin{cases} \frac{d(\rho_{X_1})}{dt} = f(\rho_{X_1}, \rho_{X_2}, \dots, \rho_{X_n}, \rho_{in}^{(X_1)}, \rho_{out}^{(X_1)}, t) \\ \frac{d(\rho_{X_2})}{dt} = f(\rho_{X_1}, \rho_{X_2}, \dots, \rho_{X_n}, \rho_{in}^{(X_2)}, \rho_{out}^{(X_2)}, t) \\ \vdots \\ \frac{d(\rho_{X_n})}{dt} = f(\rho_{X_1}, \rho_{X_2}, \dots, \rho_{X_n}, \rho_{in}^{(X_n)}, \rho_{out}^{(X_n)}, t) \end{cases}, \quad (2.9)$$

where  $X_1, X_2, \dots, X_n$  are the chemical compounds inside the reactor and  $f(\cdot)$  is the dynamical model presented in Example 2.1.

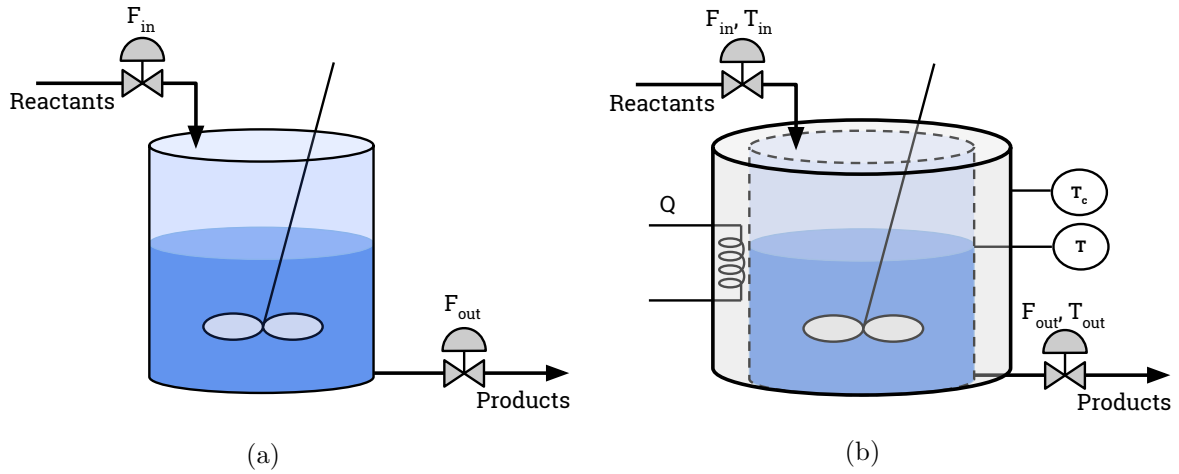


Figure 2.1: Schematic representations of industrial reactor tanks for (a) a simple isothermal process and (b) a non-isothermal process with heating/cooling system.

The model discussed so far is fairly simple. To account for more complex systems, the same modeling procedure can still be applied. For instance, it is possible to extend the description to account for exothermic and endothermic processes, when the temperature inside the system has a dynamical evolution and the dynamics of the reactions starts to depend on it.

When discussing non-isothermal processes, it is also common to discuss heating or cooling systems that tries to impose certain operational conditions to the reactions, as illustrated in Fig. 2.1b. In exothermic processes, for instance, the heat accumulated in the system tends to grows as the reactions occurs, which can be very dangerous. One approach to regulate the temperature consists in involving the chemical solution, or the container containing it, with a material whose temperature can be manipulated, transferring or absorbing heat by conductance. The temperature of this material can be manipulated by, for instance, running a heated fluid or converting electrical energy to heat energy.

**Example 2.2.** (Non-isothermal Reactor System) Consider the process of obtaining a first principle model for a non-isothermal reactor, such as the one depicted in 1.3b. In this case, the kinetic rates will follow the Arrhenius equation from (1.5). Following the same procedure from Example 2.1, a dynamical model for the concentrations of a substance  $A$ , inside this reactor, can be formulated as:

$$\frac{d(\rho_A)}{dt} = q(\rho_{in}^{(A)} - \rho_{out}^{(A)}) + \left( \sum_{\alpha X \rightarrow \beta A} \frac{1}{\beta} K_{XA} e^{-\frac{E_{XA}}{T}} (\rho_X)^\alpha \right) - \left( \sum_{\alpha A \rightarrow \beta X} \frac{1}{\beta} K_{AX} e^{-\frac{E_{AX}}{T}} (\rho_A)^\alpha \right), \quad (2.10)$$

where  $E_{XA}$  and  $E_{AX}$  are the activation energy needed for each reaction, and the rest of the parameters are still the same as defined in Example 2.1.

Furthermore, since the process is non-isothermal, the reactions will cause a change of temperature inside the reactor system. Consider the temperature inside the reactor, denoted as  $T$ . In this case, it is possible to relate its dynamics to the conservation law:

$$\begin{aligned} \left( \begin{array}{c} \text{Accumulation} \\ \text{of thermal energy} \\ \text{in the system} \end{array} \right) &= \left( \begin{array}{c} \text{Thermal energy} \\ \text{entering} \\ \text{the system} \end{array} \right) - \left( \begin{array}{c} \text{Thermal energy} \\ \text{leaving} \\ \text{the system} \end{array} \right) \\ &= \left( \begin{array}{c} \text{Heat flow} \\ \text{entering} \\ \text{the system} \end{array} \right) - \left( \begin{array}{c} \text{Heat flow} \\ \text{leaving} \\ \text{the system} \end{array} \right) + \left( \begin{array}{c} \text{Entropy} \\ \text{contribution} \\ \text{from reactions} \end{array} \right). \end{aligned} \quad (2.11)$$

It is well known, from Fourier's law [Bergman et al., 2017], that the transfer of heat by conduction between the reactor and the coolant/heater, in their respective contact interfaces, is proportional to  $(T_C - T)$ , where  $T_C$  represents the temperature of the latter. Moreover, another identifiable source of heat transfer consists in the fluid flow entering and leaving the system. Similarly for the conduction, this heat contribution should also be proportional to  $(T_{in} - T_{out})$ , where  $T_{in}$  and  $T_{out}$  are respectively the temperature of the fluid entering and leaving the reactor.

The entropy contribution from a reaction  $\alpha X \rightarrow \beta Y$ , denoted as  $S_{XY}$ , is proportional to the concentration of  $\rho_X$  consumed, or produced, by the reactions multiplied by the energy that it liberates or absorbs, as stated by the Hess' Law [Atkins and de Paula, 2011]:

$$S_{XY} \sim K_{XY} e^{-\frac{E_{XY}}{T}} (\rho_X)^\alpha \Delta H_{XY}. \quad (2.12)$$

All the proportionality can be turned into equalities by imposing real constant factors that are calculated independent of the dynamical variables. In the case of the heat transfer from the fluid, this constant factor is the flow-rate  $q$  itself. Plugging all together, and summing the entropy contribution from each reaction, the accumulation of heat can be modeled:

$$\frac{d(T)}{dt} = q(T_{in} - T_{out}) + \eta(T_C - T) + \delta \sum_{\alpha A \rightarrow \beta X} K_{AX} e^{-\frac{E_{AX}}{T}} (\rho_A)^\alpha \Delta H_{AX}. \quad (2.13)$$

Consider now the heating/cooling system involving the reactor system. From the choice of design of this apparatus, it is possible to manipulate the temperature of its material using a heat capacity  $Q$ , up to a real factor of  $\gamma$ . Similar to the temperature inside the container, there is the conduction of heat between the heating/cooling system and the walls of the container for the reactor. Therefore, the model of heat accumulation for this quantity is simply given by:

$$\frac{d(T_C)}{dt} = \gamma Q + \beta(T - T_C). \quad (2.14)$$

The non-isothermal reactor system is a more general model that accounts for the fact that the temperature of the environment is usually not constant. From this assumption, the flow entering and leaving the system are also not assumed to have the same temperature that the fluid inside the reactor. In practical applications, the temperature of the fluid inflow can either be manipulated or measured, where the temperature of the fluid outflow is actually assumed to be equal to the temperature inside the reactor. In addition, the proportionality constants are not functions of any dynamical variable, so they can be calculated before the operation of the system by using the properties of the materials and containers.

In the case of this model, the rate of changes in the chemical concentrations depends on the temperature through the Arrhenius equation. However, the temperature of the reactor itself depends on those concentrations. So, as noted in (2.9), the dynamical model of the entire reactor is a system of differential equations relating all those quantities.

## 2.2 Canonical Representations for Dynamical Models

The last section presented the foundation for modeling a dynamical system using first principles from physics. Although it is a well-defined formulation, the resulting models are not guaranteed to be practical in a mathematical sense. This motivates a discussion about possible canonical formats for such models, which would allow for systematic analysis approaches. In the perspective of control theory, there are two main formats for the model of a system: the *Input-Output* (IO) and the *State-Space* (SS) representations.

An input-output representation is a simple model that describes the entire system using only two types of variables, and their derivatives. Therefore, the dimension of these variables and the order of derivatives at each differential equation provides the information about the structure of the model. For instance, in the case where  $p = r = 1$  the system can be classified as a *Single-Input Single-Output* (SISO) configuration, whereas it is classified as a *Multiple-Input Multiple-Output* (MIMO) configuration if  $p, r > 1$ . A formal definition is given below.

**Definition 2.1.** (Input-Output Representation) An Input-Output (IO) representation of a dynamical system with  $p \geq 1$  output variables, represented by  $\mathbf{y} : \mathbb{R} \rightarrow \mathbb{R}^p$ , and  $r \geq 1$  input variables, represented by  $\mathbf{u} : \mathbb{R} \rightarrow \mathbb{R}^r$ , is the system of differential equations:

$$\begin{cases} h_1 \left( y_1, \dot{y}_1, \dots, y_1^{(n_1)}, u_1, \dot{u}_1, \dots, u_1^{(m_{11})}, u_2, \dot{u}_2, \dots, u_2^{(m_{12})}, \dots, u_r, \dot{u}_r, \dots, u_r^{(m_{1r})}, t \right) = 0 \\ h_2 \left( y_2, \dot{y}_2, \dots, y_2^{(n_2)}, u_1, \dot{u}_1, \dots, u_1^{(m_{21})}, u_2, \dot{u}_2, \dots, u_2^{(m_{22})}, \dots, u_r, \dot{u}_r, \dots, u_r^{(m_{2r})}, t \right) = 0 \\ \vdots \\ h_p \left( y_p, \dot{y}_p, \dots, y_p^{(n_p)}, u_1, \dot{u}_1, \dots, u_1^{(m_{p1})}, u_2, \dot{u}_2, \dots, u_2^{(m_{p2})}, \dots, u_r, \dot{u}_r, \dots, u_r^{(m_{pr})}, t \right) = 0 \end{cases}, \quad (2.15)$$

where:

$$\dot{y}(t) = \frac{dy(t)}{dt}, \quad \ddot{y}(t) = \frac{d^2y(t)}{dt^2}, \quad \dots, \quad y^{(n)}(t) = \frac{d^n y(t)}{dt^n}$$

and

$$\dot{u}(t) = \frac{du(t)}{dt}, \quad \ddot{u}(t) = \frac{d^2u(t)}{dt^2}, \quad \dots, \quad u^{(n)}(t) = \frac{d^n u(t)}{dt^n}.$$

This model presents a cause-and-effect interpretation of the system where the direct relationship between the input and output signal, and its derivatives, are equated as if the system is a processing unit. In practice, the input signals  $\mathbf{u}(t)$  are the manipulated variables of the system, where the output signals  $\mathbf{y}(t)$  are the observations of the controlled variables. This

representation brings an easy visualization on how a desired system behavior can be achieved by applying a specific input signal, posing as a practical framework for designing controllers.

The State-Space representation is another formulation for a dynamical model, centered in the concept of *state variables*. In a formal definition, the set of state variables is the smallest set of linearly independent variables that can unequivocally determine the value of all the states variables given an initial state  $\mathbf{x}(t_0)$ , at time  $t_0 \in \mathbb{R}$ , and a forcing function  $\mathbf{u}(t)$ , for any time  $t \geq t_0$ . From a physical perspective, however, these variables may account for quantities that can describe the dynamics of the system, such as position or velocities. In comparison to the Input-Output representation, this formulation poses a simpler mathematical model, since it is composed by a system of first order ordinary differential equations and a system of algebraic equations. However, this model presents a semantical improvement over the latter since the inclusion of the state variables expands the internal description of the system. A formal definition is given below.

**Definition 2.2.** (State-Space Representation) A State-Space (SS) representation of a system with  $n \geq 1$  states variables, represented by  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^n$ , for  $p \geq 1$  output variables, represented by  $\mathbf{y} : \mathbb{R} \rightarrow \mathbb{R}^p$ , and  $r \geq 1$  input variables, represented by  $\mathbf{u} : \mathbb{R} \rightarrow \mathbb{R}^r$ , is given by the systems of state and output equations:

$$\begin{array}{ll} \text{State Equations:} & \text{Output Equations:} \\ \left\{ \begin{array}{l} \dot{x}_1(t) = f_1(x_1, \dots, x_n, u_1, \dots, u_r, t) \\ \dot{x}_2(t) = f_2(x_1, \dots, x_n, u_1, \dots, u_r, t) \\ \vdots \\ \dot{x}_n(t) = f_n(x_1, \dots, x_n, u_1, \dots, u_r, t) \end{array} \right. & ; \left\{ \begin{array}{l} y_1(t) = g_1(x_1, \dots, x_n, u_1, \dots, u_r, t) \\ y_2(t) = g_2(x_1, \dots, x_n, u_1, \dots, u_r, t) \\ \vdots \\ y_p(t) = g_p(x_1, \dots, x_n, u_1, \dots, u_r, t) \end{array} \right. , \end{array} \quad (2.16)$$

or, in the matrix form:

$$\begin{cases} \dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t), t) \\ \mathbf{y}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{u}(t), t) \end{cases} . \quad (2.17)$$

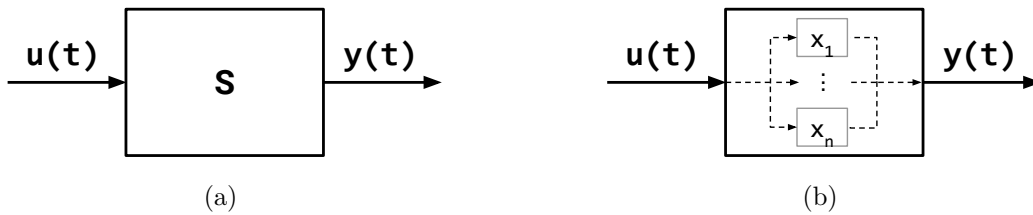


Figure 2.2: Graphical interpretation of (a) Input-Output models and (b) State-Spaces models.

A graphical illustration of both representations is shown at Fig. 2.2. In the light of these formulations, a system can be classified with respect to the model mathematical structure. There are five main properties used for this classification: whether the system is causal or non-causal, linear or nonlinear, dynamical or instantaneous, time-invariant or time-varying and with or without delay. The necessary and sufficient conditions for each one of these properties are summarized in Table 2.1.

This work focus on dynamical linear systems, since its models are the most well studied in the control theory community. In reality, a physical system is always causal, nonlinear and

	Input-Output	State-Space
<b>Causal</b>	$m_{ij} \leq n_k$ $i \in [1, \dots, p], j \in [1, \dots, r]$	Always causal
<b>Linear</b>	$h_i(\cdot) = \sum_{j=0}^{n_i} y^{(j)} + \dots$ $\dots + \sum_{k=1}^r \sum_{l=0}^{m_{ik}} u_k^{(l)}$ $i \in [1, 2, \dots, p]$	$f_i = \mathbf{a}_i(t)\mathbf{x}(t) + \mathbf{b}_i(t)\mathbf{u}(t), i = 1, 2, \dots, n$ $g_j = \mathbf{c}_j(t)\mathbf{x}(t) + \mathbf{d}_j(t)\mathbf{u}(t), j = 1, 2, \dots, p$ $\mathbf{a}_i, \mathbf{c}_j \in \mathbb{R}^{1 \times n}$ and $\mathbf{b}_i, \mathbf{d}_j \in \mathbb{R}^{1 \times r}$
<b>Dynamical</b>	$n_i > 0$ or $m_{jk} > 0$ $i, j \in [1, \dots, p], k \in [1, \dots, r]$	$n > 0$
<b>Time-Invariant</b>	$h_i(y_i(t), \dots, u_1(t), \dots, u_r(t)) = 0$ $i \in [1, 2, \dots, p]$	$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t))$ $\mathbf{y}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{u}(t))$
<b>Without-Delay</b>	All the signals share the same time arguments	All the signals share the same time arguments

Table 2.1: Necessary and sufficient conditions for different classes of models.

time-varying [Vidyasagar, 2002], but the models can be assumed differently with fairly accuracy. The benefit of linear systems is that it obeys the superposition principle, and a linear combination of the inputs directly causes the exact same linear combination of the individual outputs. Under the assumption of a linear system, a nice result is that the vectorial functions  $\mathbf{f}(\cdot)$  and  $\mathbf{g}(\cdot)$  of the State-Space representation in (2.27) reduces to simple matrix forms.

**Definition 2.3.** (Linear State-Space Representation) A State-Space representation describing a linear system with state vector  $\mathbf{x}(t) : \mathbb{R} \rightarrow \mathbb{R}^n$ , output vector  $\mathbf{y}(t) : \mathbb{R} \rightarrow \mathbb{R}^p$  and input vector  $\mathbf{u}(t) : \mathbb{R} \rightarrow \mathbb{R}^r$  is given by the system of equations:

$$\begin{cases} \dot{\mathbf{x}}(t) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{B}(t)\mathbf{u}(t) \\ \mathbf{y}(t) = \mathbf{C}(t)\mathbf{x}(t) + \mathbf{D}(t)\mathbf{u}(t) \end{cases}, \quad (2.18)$$

where  $\mathbf{A}(t) : \mathbb{R} \rightarrow \mathbb{R}^{n \times n}$ ,  $\mathbf{B}(t) : \mathbb{R} \rightarrow \mathbb{R}^{n \times r}$ ,  $\mathbf{C}(t) : \mathbb{R} \rightarrow \mathbb{R}^{p \times n}$  and  $\mathbf{D}(t) : \mathbb{R} \rightarrow \mathbb{R}^{p \times r}$ . In the case of a time-invariant linear system, these matrices becomes constants.

This formulation has the advantage that the time response of the system can be easily calculated and that the analysis of the dynamics follows well-established results from linear algebra applied to the matrices  $\mathbf{A}(t)$ ,  $\mathbf{B}(t)$ ,  $\mathbf{C}(t)$  or  $\mathbf{D}(t)$ , as well as for the vectors  $\mathbf{x}(t)$  and  $\mathbf{u}(t)$ . Furthermore, the physical interpretation of the system through the state variables becomes straightforward in this model.

In addition to the State-Space representation, the linear assumption also benefits Input-Output representations. One major analytical tool that can be used in these cases is to transform this model to a frequency domain, using a linear transform operator, in order to simplify the solution for the differential equations. The most popular choice of transformation is the *Laplace transform*,  $\mathcal{L}\{h(t)\}$ , which converts functions in time to functions in complex frequencies. Using the properties of this operator, differential equations are converted to simple algebraic equations.

**Definition 2.4.** (Transfer Function) Given a linear model for a SISO system, with initial conditions  $\mathbf{y}(0^-) = \mathbf{u}(0^-) = \mathbf{0}$ , in the Input-Output formulation:

$$\alpha_n \frac{d^n y(t)}{dt^n} + \dots + \alpha_1 \frac{dy(t)}{dt} + \alpha_0 y(t) = \beta_m \frac{d^m u(t)}{dt^m} + \dots + \beta_1 \frac{du(t)}{dt} + \beta_0 u(t), \quad (2.19)$$

its correspondent transfer function, calculated in the Laplace domain, is defined as:

$$G(s) = \frac{Y(s)}{U(s)} = \frac{\beta_m s^m + \beta_{m-1} s^{m-1} + \dots + \beta_1 s + \beta_0}{\alpha_n s^n + \alpha_{n-1} s^{n-1} + \dots + \alpha_1 s + \alpha_0}. \quad (2.20)$$

An indirect result of this is that the SS representation can be converted to the IO representation using the Laplace transform operator, leading to a notion of equivalence between the two representations. Notice that the extension to the MIMO case is straightforward: just compute the transfer function between each pair of input and output, leading to the matrix  $\mathbf{G} \in \mathbb{C}^{p \times r}$ , for a system with  $p$  outputs and  $r$  inputs. For the general linear case, this matrix has the form:

$$\mathbf{G}(s) = \begin{bmatrix} \frac{Y_1(s)}{U_1(s)} & \frac{Y_1(s)}{U_2(s)} & \dots & \frac{Y_1(s)}{U_r(s)} \\ \frac{Y_2(s)}{U_1(s)} & \frac{Y_2(s)}{U_2(s)} & \dots & \frac{Y_2(s)}{U_r(s)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{Y_p(s)}{U_1(s)} & \frac{Y_p(s)}{U_2(s)} & \dots & \frac{Y_p(s)}{U_r(s)} \end{bmatrix} = \begin{bmatrix} \frac{\beta_{m_{11}}^{(11)} s^{m_{11}} + \dots + \beta_0^{(11)}}{\alpha_{n_1}^{(1)} s^{n_1} + \dots + \alpha_0^{(1)}} & \dots & \frac{\beta_{m_{1r}}^{(1r)} s^{m_{1r}} + \dots + \beta_0^{(1r)}}{\alpha_{n_1}^{(1)} s^{n_1} + \dots + \alpha_0^{(1)}} \\ \frac{\beta_{m_{21}}^{(21)} s^{m_{21}} + \dots + \beta_0^{(21)}}{\alpha_{n_2}^{(2)} s^{n_2} + \dots + \alpha_0^{(2)}} & \dots & \frac{\beta_{m_{2r}}^{(2r)} s^{m_{2r}} + \dots + \beta_0^{(2r)}}{\alpha_{n_2}^{(2)} s^{n_2} + \dots + \alpha_0^{(2)}} \\ \vdots & \ddots & \vdots \\ \frac{\beta_{m_{p1}}^{(p1)} s^{m_{p1}} + \dots + \beta_0^{(p1)}}{\alpha_{n_p}^{(p)} s^{n_p} + \dots + \alpha_0^{(p)}} & \dots & \frac{\beta_{m_{pr}}^{(pr)} s^{m_{pr}} + \dots + \beta_0^{(pr)}}{\alpha_{n_p}^{(p)} s^{n_p} + \dots + \alpha_0^{(p)}} \end{bmatrix}. \quad (2.21)$$

**Theorem 2.1.** (Passage from SS to IO) Consider a linear and time-invariant system in State-Space form with initial states  $\mathbf{x}(0^-) = \mathbf{0}$  and represented as:

$$\begin{cases} \dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) \\ \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t) \end{cases}. \quad (2.22)$$

The equivalent system in Input-Output representation is given by the transfer function:

$$\mathbf{G}(s) = \mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}. \quad (2.23)$$

*Proof.* Applying the Laplace formation to both sides of the equation, using the property of derivative in time, results in:

$$\begin{aligned} \mathcal{L}\{\dot{\mathbf{x}}(t)\} &= \mathbf{A}\mathcal{L}\{\mathbf{x}(t)\} + \mathbf{B}\mathcal{L}\{\mathbf{u}(t)\} \\ s\mathbf{X}(s) - \dot{\mathbf{x}}(0^-) &= \mathbf{A}\mathbf{X}(s) + \mathbf{B}\mathbf{U}(s) \end{aligned}. \quad (2.24)$$

Substituting the derivative  $\dot{\mathbf{x}}(0^-) = \mathbf{0}$  and doing some manipulations results in:

$$\begin{aligned} (s\mathbf{I} - \mathbf{A})\mathbf{X}(s) &= \mathbf{B}\mathbf{U}(s) \\ \mathbf{X}(s) &= (s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B}\mathbf{U}(s). \end{aligned} \quad (2.25)$$

Finally, applying the same procedure to the output equations, and using the previous result, the Laplace transform of the output is:

$$\begin{aligned} \mathcal{L}\{\mathbf{y}(t)\} &= \mathbf{C}\mathcal{L}\{\mathbf{x}(t)\} + \mathbf{D}\mathcal{L}\{\mathbf{u}(t)\} \\ \mathbf{Y}(s) &= \mathbf{C}\mathbf{X}(s) + \mathbf{D}\mathbf{U}(s) \\ &= \mathbf{C}((s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B}\mathbf{U}(s)) + \mathbf{D}\mathbf{U}(s) \\ &= (\mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D})\mathbf{U}(s) \end{aligned} \quad (2.26)$$



Since, by definition,  $\mathbf{Y}(s) = \mathbf{G}(s)\mathbf{U}(s)$ , it is possible to obtain the transfer function matrix  $\mathbf{G}(s) = \mathbf{C}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$  as an equivalent representation of the system.  $\square$

Despite of the discussion about the benefits of linear models, it is necessary to account for the fact that physical systems will present, in most situations, nonlinear behavior. For this reason, some effort must be done to develop a linear model that can describe the nonlinear behavior with certain accuracy, even if over some small region of the space. With this motivation, a technique for *linearization* of a nonlinear model is detailed below.

**Theorem 2.2.** (*Linearization by Taylor Expansion*) Consider a nonlinear time-invariant system:

$$\begin{cases} \dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)) \\ \mathbf{y}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{u}(t)) \end{cases} \quad (2.27)$$

Given steady-state operating points  $\mathbf{x}_o$ ,  $\mathbf{y}_o$  and  $\mathbf{u}_o$ , the dynamics of the system in the neighborhood of these points can be represented by the linear model:

$$\begin{cases} \Delta\dot{\mathbf{x}}(t) = \mathbf{A}\Delta\mathbf{x}(t) + \mathbf{B}\Delta\mathbf{u}(t) \\ \mathbf{y}(t) = \mathbf{C}\Delta\mathbf{x}(t) + \mathbf{D}\Delta\mathbf{u}(t) \end{cases}, \quad (2.28)$$

where

$$\mathbf{A} = \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\mathbf{x}_o, \mathbf{u}_o}; \quad \mathbf{B} = \left. \frac{\partial \mathbf{f}}{\partial \mathbf{u}} \right|_{\mathbf{x}_o, \mathbf{u}_o}; \quad \mathbf{C} = \left. \frac{\partial \mathbf{g}}{\partial \mathbf{x}} \right|_{\mathbf{x}_o, \mathbf{u}_o}; \quad \mathbf{D} = \left. \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right|_{\mathbf{x}_o, \mathbf{u}_o} \quad (2.29)$$

and

$$\Delta\mathbf{x}(t) = \mathbf{x}(t) - \mathbf{x}_o; \quad \Delta\mathbf{u}(t) = \mathbf{u}(t) - \mathbf{u}_o. \quad (2.30)$$

*Proof.* Consider a system represented by state equations  $\mathbf{f}(\cdot)$  and output equations  $\mathbf{g}(\cdot)$ , with steady-state points  $\mathbf{x}_o$ ,  $\mathbf{y}_o$  and  $\mathbf{u}_o$ . Now, consider a very small perturbation  $\Delta\mathbf{u}(t)$  in the input signal around these operation points. This perturbation will result in small changes in the state and output variables:

$$\mathbf{x}(t) = \mathbf{x}_o + \Delta\mathbf{x}(t); \quad \mathbf{u}(t) = \mathbf{u}_o + \Delta\mathbf{u}(t); \quad \mathbf{y}(t) = \mathbf{y}_o + \Delta\mathbf{y}(t). \quad (2.31)$$

This results in the following configuration on the State-Space:

$$\begin{cases} \frac{d(\mathbf{x}_o + \Delta\mathbf{x}(t))}{dt} = \mathbf{f}(\mathbf{x}_o + \Delta\mathbf{x}(t), \mathbf{u}_o + \Delta\mathbf{u}(t)) \\ \mathbf{y}_o + \Delta\mathbf{y}(t) = \mathbf{g}(\mathbf{x}_o + \Delta\mathbf{x}(t), \mathbf{u}_o + \Delta\mathbf{u}(t)) \end{cases}, \quad (2.32)$$

where  $d(\mathbf{x}_o + \Delta\mathbf{x}(t))/dt = d(\Delta\mathbf{x}(t))/dt$ , since  $\mathbf{x}_o$  is constant. The perturbed variables are very close to the steady-state points, hence the functions  $\mathbf{f}(\cdot)$  and  $\mathbf{g}(\cdot)$  can be approximated by a Taylor series expansion, yielding:

$$\begin{cases} \frac{d(\Delta\mathbf{x}(t))}{dt} = \mathbf{f}(\mathbf{x}_o, \mathbf{u}_o) + \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\mathbf{x}_o, \mathbf{u}_o} \Delta\mathbf{x}(t) + \left. \frac{\partial \mathbf{f}}{\partial \mathbf{u}} \right|_{\mathbf{x}_o, \mathbf{u}_o} \Delta\mathbf{u}(t) + \mathcal{O}(\Delta\mathbf{x}^2, \Delta\mathbf{u}^2) \\ \mathbf{y}_o + \Delta\mathbf{y}(t) = \mathbf{g}(\mathbf{x}_o, \mathbf{u}_o) + \left. \frac{\partial \mathbf{g}}{\partial \mathbf{x}} \right|_{\mathbf{x}_o, \mathbf{u}_o} \Delta\mathbf{x}(t) + \left. \frac{\partial \mathbf{g}}{\partial \mathbf{u}} \right|_{\mathbf{x}_o, \mathbf{u}_o} \Delta\mathbf{u}(t) + \mathcal{O}(\Delta\mathbf{x}^2, \Delta\mathbf{u}^2) \end{cases}. \quad (2.33)$$



Since the steady-state condition implies zero variation, it is possible to assume  $\mathbf{f}(\mathbf{x}_o, \mathbf{u}_o) = \mathbf{0}$  and  $g(\mathbf{x}_o, \mathbf{u}_o) = \mathbf{0}$ , since they are ordinary differential equations. Truncating in the first order terms and substituting  $\mathbf{y}(t) = \mathbf{y}_o + \Delta\mathbf{y}(t)$  results in:

$$\begin{cases} \frac{d(\Delta\mathbf{x}(t))}{dt} = \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\mathbf{x}_o, \mathbf{u}_o} \Delta\mathbf{x}(t) + \left. \frac{\partial \mathbf{f}}{\partial \mathbf{u}} \right|_{\mathbf{x}_o, \mathbf{u}_o} \Delta\mathbf{u}(t) \\ \mathbf{y}(t) = \left. \frac{\partial g}{\partial \mathbf{x}} \right|_{\mathbf{x}_o, \mathbf{u}_o} \Delta\mathbf{x}(t) + \left. \frac{\partial g}{\partial \mathbf{u}} \right|_{\mathbf{x}_o, \mathbf{u}_o} \Delta\mathbf{u}(t) \end{cases} \quad (2.34)$$

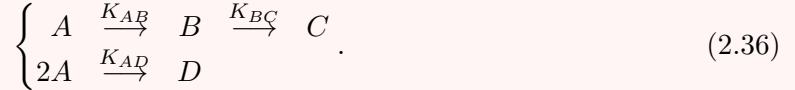
Finally, since all the Jacobians involved are actually matrices of appropriate dimensions, the final linear approximation of the system is the SS model given by:

$$\begin{cases} \Delta\dot{\mathbf{x}}(t) = \mathbf{A}\Delta\mathbf{x}(t) + \mathbf{B}\Delta\mathbf{u}(t) \\ \mathbf{y}(t) = \mathbf{C}\Delta\mathbf{x}(t) + \mathbf{D}\Delta\mathbf{u}(t) \end{cases} \quad (2.35)$$

□

Note that this method can be used to convert non-linear models obtained from the first-principle procedure to a desirable State-Space form, since those models are usually in the non-linear State-Space representation.

**Example 2.3.** (Van de Vusse Reactor) Consider, for instance, a reactor system describing an isothermal process that follows the Van de Vusse reaction scheme [Van de Vusse, 1964]:



Consider, also, that the concentration of substances  $A$ ,  $B$ ,  $C$  and  $D$  inside the fluid inflow and outflow are described by the following conditions:

$$\begin{cases} \rho_{in}^{(B)}(t) = \rho_{in}^{(C)}(t) = \rho_{in}^{(D)}(t) = 0 \\ \rho_{out}^{(k)}(t) = \rho_k(t), \forall k \in \{A, B, C, D\} \end{cases} \quad (2.37)$$

Now, let  $\mathbf{x} = [\rho_A, \rho_B, \rho_C, \rho_D]^T$  and  $u = q$ , and assume that the states are perfectly observed, so that  $\mathbf{y} = \mathbf{x}$ . Using Example 2.1, the nonlinear State-Space representation of this reactor is can be modeled by the state-equations:

$$\begin{cases} \dot{x}_1(t) = u(t)(\rho_{in}^{(A)} - x_1(t)) - (K_{AB}x_1(t) + K_{AD}(x_1(t))^2) \\ \dot{x}_2(t) = -u(t)x_2(t) + K_{AB}x_1(t) - K_{BC}x_2(t) \\ \dot{x}_3(t) = -u(t)x_3(t) + K_{BC}x_2(t) \\ \dot{x}_4(t) = -u(t)x_4(t) + \frac{1}{2}K_{AD}(x_1(t))^2 \end{cases} \quad (2.38)$$

To select steady-state points  $\mathbf{x}_o$  and  $u_o$ , it is necessary to find values for  $\mathbf{x}(t)$  and  $u(t)$  that satisfies the system  $\dot{\mathbf{x}}(t) = \mathbf{0}$ . This can be done analytically, in some cases, or by using numerical methods. Finally, the linear State-Space matrices  $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$  are obtained by taking the

derivative each one of these equations in respect to each state and input, resulting in:

$$\left\{ \begin{array}{l} \mathbf{A} = \begin{bmatrix} -u_o - K_{AB} - 2K_{AC}x_{1o} & 0 & 0 & 0 \\ K_{AB} & -u_o - K_{BC} & 0 & 0 \\ 0 & K_{BC} & -u_o & 0 \\ K_{AD}x_{1o} & 0 & 0 & -u_o \end{bmatrix}; \quad \mathbf{B} = \begin{bmatrix} \rho_{in}^{(A)} - x_{1o} \\ -x_{2o} \\ -x_{3o} \\ -x_{4o} \end{bmatrix}; \\ \mathbf{C} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}; \quad \mathbf{D} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \end{array} \right. . \quad (2.39)$$

Notice that the response of a linearized model represents the deviations around the given steady-state points. If the errors in approximating the non-linearities of the actual system are somehow tolerable, however, this model can still be used as a nominal model to represent the system in any point in space. For the sake of illustration, the simulation of both the nonlinear and the linear systems are displayed in Fig. 2.3, for a realization of the system accounted in Example 2.3. In this realization, the kinetic rates were set as  $K_{AB} = 5/6$ ,  $K_{BC} = 5/3$  and  $K_{AC} = 1/6$ , and the inflow concentration was consider constant as  $\rho_{in}^{(A)} = 10$  (mol/l). In this simulation, the dashed line represents the linear response of the linearized model given  $\mathbf{x}_o = [6.19, 1.09, 0.60, 1.05]^T$  and  $u_0 = 3.03$ . Notice how the linearized system is capable of approximating the behavior of the non-linear system when closer to the steady-state point, but starts to display an error when the system leaves that state.

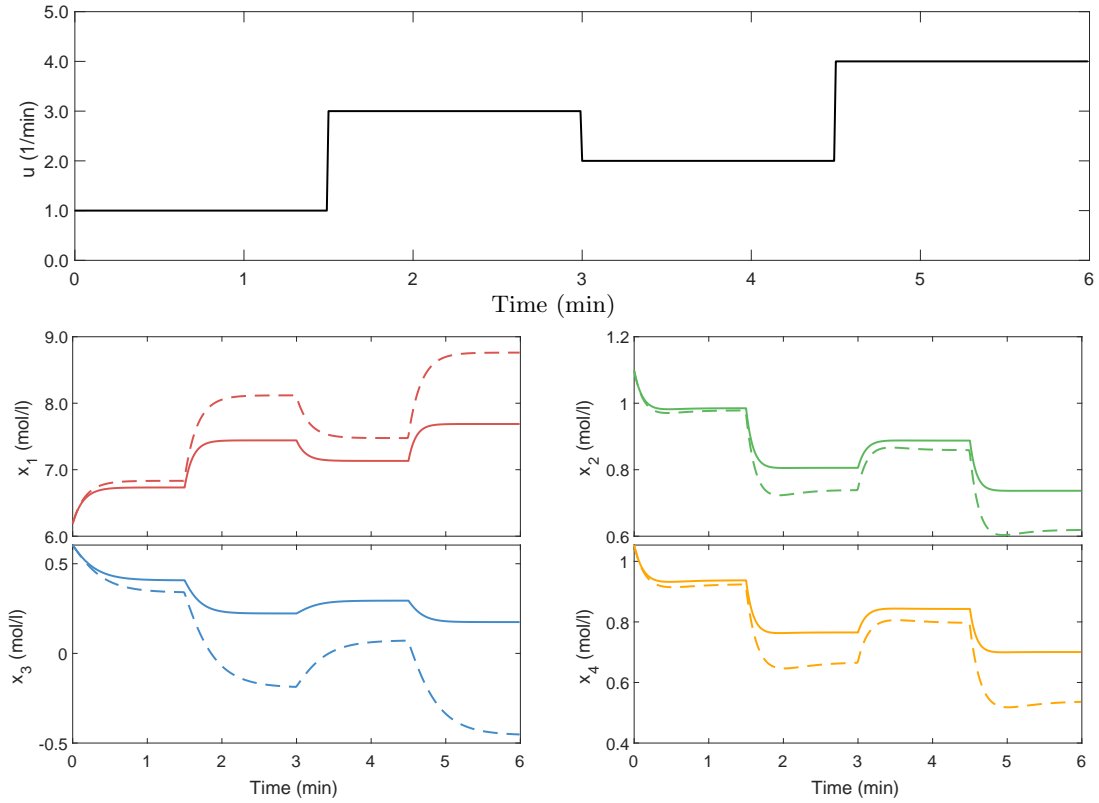


Figure 2.3: Comparison between the nonlinear and linear response of the reactor system.

## 2.3 Response Analysis in the Time Domain

Once that a model is well-established, it is possible both to simulate the system and to analyze its response, given an initial state and input signal. This section focus on developing a quantitative understanding of a system behavior through a linear model. The results are focused on continuous-time response of linear and time-invariant systems in the form:

$$\begin{cases} \dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) \\ \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t) \end{cases} \quad (2.40)$$

First of all, it is necessary to access some properties of the matrix  $\mathbf{A}$ , which describes the influence of the states in the evolution of the system. Consider, in this case, the following matrix defined by computing the exponential of the matrix  $\mathbf{A}$ .

**Definition 2.5.** (State-Transition Matrix) Consider a system in State-Space representation with matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Its *State-Transition Matrix*,  $e^{\mathbf{A}t} \in \mathbb{R}^{n \times n}$ , is the converging series:

$$e^{\mathbf{A}t} = \mathbf{I} + \mathbf{A}t + \frac{\mathbf{A}^2 t^2}{2!} + \frac{\mathbf{A}^3 t^3}{3!} + \cdots = \sum_{k=0}^{\infty} \frac{\mathbf{A}^k t^k}{k!}. \quad (2.41)$$

This matrix is central to the computation of the system time response, since it can directly relate the initial state of a system to any other state forward in time. Moreover, the exponential of a matrix is not a trivial operation. However, the fact that the State-Space matrix  $\mathbf{A}$  is always a square matrix results in the following properties for the State-Transition matrix, stated below.

**Theorem 2.3.** Consider a matrix exponential as in Definition 2.5, for a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ . Then, the following properties holds:

$$\text{I) } \frac{d(e^{\mathbf{A}t})}{dt} = \mathbf{A}e^{\mathbf{A}t}; \quad \text{II) } e^{\mathbf{A}t}e^{\mathbf{A}\tau} = e^{\mathbf{A}(t+\tau)}; \quad \text{III) } e^{-\mathbf{A}t}e^{\mathbf{A}t} = e^{\mathbf{A}t}e^{-\mathbf{A}t} = \mathbf{I}. \quad (2.42)$$

The proof of these properties can be obtained directly from using the Definition 2.5 and doing some algebra. Given these results, the calculation of the time response of a system in SS representation becomes straightforward.

**Theorem 2.4.** (Lagrange Formula) Consider a LTI system in State-Space representation. Its response for any time  $t \geq t_0$ , initial state  $\mathbf{x}(t_0)$  and input signal  $\mathbf{u}(t)$  is given by the solutions of the state and output equations:

$$\begin{cases} \mathbf{x}(t) = e^{\mathbf{A}(t-t_0)}\mathbf{x}(t_0) + \int_{t_0}^t e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)d\tau \\ \mathbf{y}(t) = \mathbf{C}e^{\mathbf{A}(t-t_0)}\mathbf{x}(t_0) + \mathbf{C} \int_{t_0}^t e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)d\tau + \mathbf{D}\mathbf{u}(t) \end{cases} \quad (2.43)$$

*Proof.* First of all, consider a system in State-Space representation with state equation as defined in (2.27). Multiplying both sides by  $e^{-\mathbf{A}t}$ :

$$\begin{aligned} e^{-\mathbf{A}t}\dot{\mathbf{x}}(t) &= e^{-\mathbf{A}t}(\mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)) \\ e^{-\mathbf{A}t}\dot{\mathbf{x}}(t) - \mathbf{A}e^{-\mathbf{A}t}\mathbf{x}(t) &= e^{-\mathbf{A}t}\mathbf{B}\mathbf{u}(t) \end{aligned} \quad (2.44)$$

Using the multiplication rule and the first property from Theorem 2.3, it is easy to see that  $d[e^{-\mathbf{A}t}\mathbf{x}(t)]/dt = e^{-\mathbf{A}t}\dot{\mathbf{x}}(t) - \mathbf{A}e^{-\mathbf{A}t}\mathbf{x}(t)$ . Substituting this result in (2.44) and integrating both sides from  $t_0$  to  $t$ :

$$\begin{aligned}\frac{d(e^{-\mathbf{A}t}\mathbf{x}(t))}{dt} &= e^{-\mathbf{A}t}\mathbf{B}\mathbf{u}(t) \\ e^{-\mathbf{A}t}\mathbf{x}(t)|_{t_0}^t &= \int_{t_0}^t e^{-\mathbf{A}\tau}\mathbf{B}\mathbf{u}(\tau)d\tau. \\ e^{-\mathbf{A}t}\mathbf{x}(t) - e^{-\mathbf{A}t_0}\mathbf{x}(t_0) &= \int_{t_0}^t e^{-\mathbf{A}\tau}\mathbf{B}\mathbf{u}(\tau)d\tau\end{aligned}\tag{2.45}$$

Multiplying both sides of (2.45) by  $e^{\mathbf{A}t}$  and using the second and third property from Theorem 2.3, the state response can be calculated as:

$$\begin{aligned}e^{\mathbf{A}t}(e^{-\mathbf{A}t}\mathbf{x}(t) - e^{-\mathbf{A}t_0}\mathbf{x}(t_0)) &= e^{\mathbf{A}t} \int_{t_0}^t e^{-\mathbf{A}\tau}\mathbf{B}\mathbf{u}(\tau)d\tau \\ \mathbf{I}\mathbf{x}(t) - e^{\mathbf{A}(t-t_0)}\mathbf{x}(t_0) &= \int_{t_0}^t e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)d\tau. \\ \mathbf{x}(t) &= e^{\mathbf{A}(t-t_0)}\mathbf{x}(t_0) + \int_{t_0}^t e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)d\tau\end{aligned}\tag{2.46}$$

Finally, substituting (2.46) into the output equation leads to:

$$\begin{aligned}\mathbf{y}(t) &= \mathbf{C} \left( e^{\mathbf{A}(t-t_0)}\mathbf{x}(t_0) + \int_{t_0}^t e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)d\tau \right) + \mathbf{D}\mathbf{u}(t) \\ &= \mathbf{C}e^{\mathbf{A}(t-t_0)}\mathbf{x}(t_0) + \mathbf{C} \int_{t_0}^t e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau)d\tau + \mathbf{D}\mathbf{u}(t)\end{aligned}\tag{2.47}$$

□

When discussing the response of a system, the focus is directed to the state equation describing its dynamics, since the output equation represents an observation of the system through the states. In this sense, the Lagrange formula shows the nice characteristic of linear systems that the total response is a composition of two separated actions:

$$\mathbf{x}(t) = \mathbf{x}_n(t) + \mathbf{x}_f(t),\tag{2.48}$$

where the natural response,  $\mathbf{x}_n(t)$ , corresponds to the state-transition matrix multiplication term and the forced response,  $\mathbf{x}_f(t)$ , corresponds to the integral term. This concept is visualized in Fig. 2.4 for the total response of the first two states of the model derived in (2.39), given the operation points  $\mathbf{x}_o = [6.19, 1.09]$  and  $u_o = 3.03$ , excited with a step input  $u(t) = 3$ ,  $t \in [0, 1.2]$ . In this plot, the solid line represents the total response, whereas the dashed and dotted lines represents the natural and forced response, respectively. It is easy to verify that the decomposition of the total response is equal to the sum of those independent components.

In order to characterize the dynamical response of a system, there is still the need to compute the state-transition matrix. This can be done in several ways [Moler and Loan, 2003]. A specific method, known as the *Sylvester expansion*, is an analytical solution that brings an understanding of the system behavior through the eigenvalues of the state-state matrix  $\mathbf{A}$ .

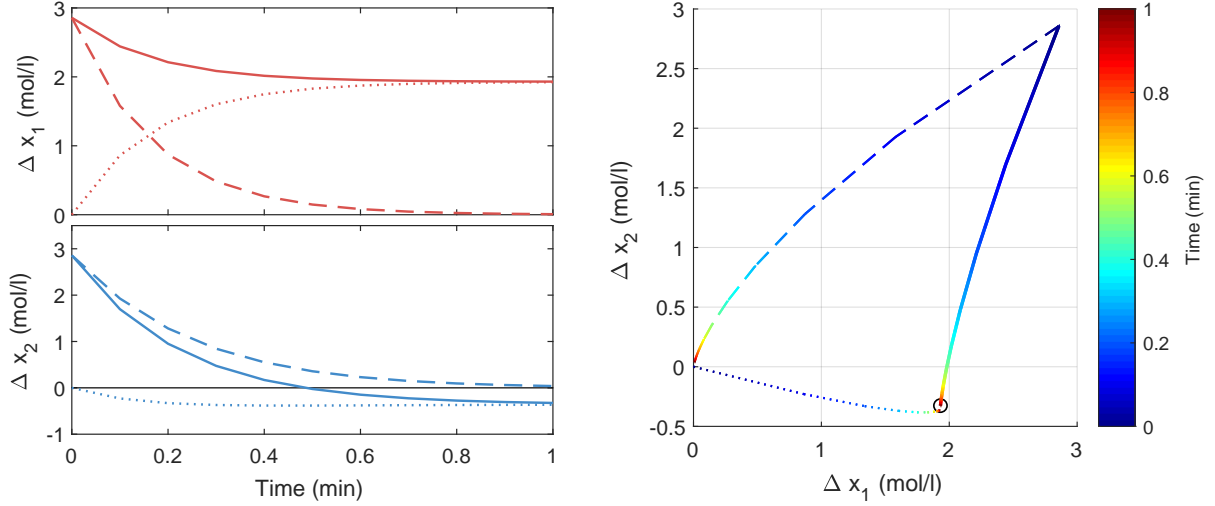


Figure 2.4: Simulation of the total response decomposition for the two first states of the model from (2.39). Each state response is shown separately (left) and together into a trajectory visualization (right), in which the steady-state point is indicated by a circle.

**Theorem 2.5.** (*Sylvester expansion*) Consider a matrix exponential function  $e^{At}$  for any square matrix  $A \in \mathbb{R}^{n \times n}$ , whose distinct eigenvalues are  $\lambda \in \mathbb{R}^m$ ,  $m \leq n$ , with associated multiplicity vector  $\nu \in \mathbb{R}^m$  such that  $\|\nu\|_1 = n$ . The result of this function can be expanded as

$$e^{At} = \beta_0(t)I + \beta_1(t)A + \beta_2(t)A^2 + \dots + \beta_{n-1}(t)A^{n-1} = \sum_{i=0}^{n-1} \beta_i(t)A^i, \quad (2.49)$$

where  $\beta_i(t) : \mathbb{R} \rightarrow \mathbb{R}$ ,  $i \in [1, 2, \dots, n-1]$ , are scalar functions of time. These functions are obtained by solving the linear system

$$V\beta = \eta, \quad (2.50)$$

for the parameter vector  $\beta = [\beta_0(t), \beta_1(t), \dots, \beta_{n-1}(t)]^T$ , given vector of modes  $\eta = [\eta_1, \eta_2, \dots, \eta_m]^T$  such that

$$\eta_i = [e^{\lambda_i t} \quad t e^{\lambda_i t} \quad t^2 e^{\lambda_i t} \quad \dots \quad t^{\nu_i-1} e^{\lambda_i t}]^T$$

and the confluent Vandermonde matrix  $V = [V_1, V_2, \dots, V_m]^T$  such that

$$V_j = \begin{bmatrix} 1 & \lambda_j & \lambda_j^2 & \dots & \lambda_j^{(\nu_j-1)} & \dots & \lambda_j^{n-1} \\ 0 & 1 & 2\lambda_j & \dots & (\nu_j-1)\lambda_j^{(\nu_j-1)} & \dots & (n-1)\lambda_j^{n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & (\nu_j-1)! & \dots & \frac{(n-1)!}{(n-\nu_j)!} \lambda_j^{n-\nu_j} \end{bmatrix}$$

The proof of this expansion is somewhat extensive, but a detailed discussion can be found in [Chen, 1998]. Basically, the expansion in (2.49) is a direct application of the Cayley-Hamilton theorem [Atiyah, 2018] and the linear system in (2.50) is a result of the Sylvester's matrix theorem [Horn and Johnson, 2012].

From the expansion presented in Theorem 2.5, it is possible to understand the relationship

between the state-transition matrix  $e^{\mathbf{A}t}$  and each eigenvalue  $\lambda$  of the matrix  $\mathbf{A}$ , also known as the *poles* of the system. First of all, notice that the formulation of the linear system that defines the parameters  $\beta_0(t), \beta_1(t), \dots, \beta_{n-1}(t)$  implies that each one of these functions are linear combinations of the exponentials  $e^{\lambda_i t}$  for each eigenvalue  $\lambda_i$ ,  $i = 1, 2, \dots, m$ . These exponentials are known as the *modes* of the matrix  $\mathbf{A}$ . Since the Sylvester expansion is linear in those parameters, it is possible to conclude that the state-transition matrix, and consequently the response of a system, is a linear combination of the modes.

**Example 2.4.** Consider, for the sake of illustration, the same State-Space model as the one simulated in Fig. 2.3, but considering only the first two states (since they are independent of the others). Its matrix  $\mathbf{A}$  is shown below:

$$\mathbf{A} = \begin{bmatrix} -5.93 & 0 \\ 0.83 & -4.70 \end{bmatrix}. \quad (2.51)$$

It is easy to see that  $\boldsymbol{\lambda} = [-5.93, -4.70]$ , since  $\mathbf{A}$  is diagonal. The Sylvester expansion for the exponential of this matrix has the parameter vector  $\boldsymbol{\beta} = [\beta_0(t), \beta_1(t)]$  which solves the system

$$\begin{bmatrix} 1 & -5.93 \\ 1 & -4.70 \end{bmatrix} \begin{bmatrix} \beta_0(t) \\ \beta_1(t) \end{bmatrix} = \begin{bmatrix} e^{-5.93t} \\ e^{-4.70t} \end{bmatrix}. \quad (2.52)$$

The solution to this system is  $\boldsymbol{\beta} = [4.82e^{-4.7t} - 3.82e^{-5.93t}, 0.813e^{-4.7t} - 0.813e^{-5.93t}]$ . Hence, state-transition matrix is calculated through the expansion:

$$\begin{aligned} e^{\mathbf{A}t} &= \beta_0(t)\mathbf{I} + \beta_1(t)\mathbf{A} \\ &= 4.82e^{-4.7t} - 3.82e^{-5.93t} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + 0.813e^{-4.7t} - 0.813e^{-5.93t} \begin{bmatrix} -5.93 & 0 \\ 0.83 & -4.70 \end{bmatrix}. \\ &= \begin{bmatrix} e^{-5.93t} & 0 \\ 0.675e^{-4.7t} - 0.675e^{-5.93t} & e^{-4.7t} \end{bmatrix} \end{aligned} \quad (2.53)$$

The simulation in Fig 2.5 shows the evolution of each element of the state-transition matrix from Example 2.4, for a given time interval. Considering only the natural response  $\mathbf{x}_n(t)$ , i.e., setting  $\mathbf{u}(t) = \mathbf{0}$ ,  $t \in [t_0, t]$ , it is evident from the Lagrange formula that the time evolution of each state is given by a row-wise weighted sum of these elements. The weights, in this case, are given by the initial state  $\mathbf{x}(t_0)$ . Therefore, the  $(i, j)$  element of this matrix describes how the  $j$ -th state affects the response of the  $i$ -th state.

Now, attention must be drawn to the case where the eigenvalues are not real, but complex and conjugate. Despite that the Sylvester expansion is still defined as in Theorem 2.5, this case introduces a slightly different interpretation of the modes contributions to the natural response.

**Theorem 2.6.** Consider the same expansion defined in Theorem 2.5. Consider, now, that the matrix  $\mathbf{A}$  has two distinct eigenvalues  $\lambda_c, \lambda'_c \in \mathbb{C}$  in the form  $\lambda_c, \lambda'_c = \alpha \pm j\omega$ . In this case, the linear system solved by the parameters  $\beta_0(t), \beta_1(t), \dots, \beta_{n-1}(t)$  will have two equations:

$$\begin{cases} \beta_0 + \alpha\beta_1 + \alpha^2\beta_2 + \dots + \alpha^{n-1}\beta_{n-1} = e^{\alpha t} \cos(\omega t) \\ 0 + \omega\beta_1 + \omega^2\beta_2 + \dots + \omega^{n-1}\beta_{n-1} = e^{\alpha t} \sin(\omega t) \end{cases}. \quad (2.54)$$

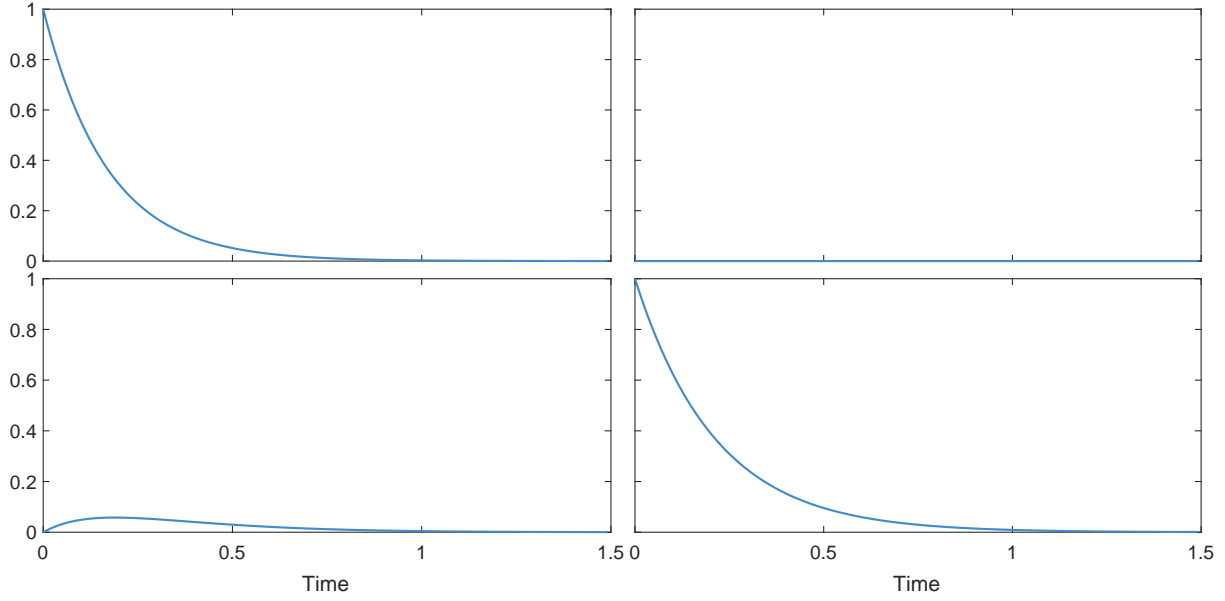


Figure 2.5: Simulation of the state-transition matrix  $e^{At}$  from Example 2.4 for  $t \in [0, 1.5]$ .

*Proof.* Consider the matrix  $\mathbf{A}$  with eigenvalues as specified and the Sylvester expansion as presented. In this case, there will be two equations in the system:

$$\begin{cases} \beta_0 + \lambda_c \beta_1 + \lambda_c^2 \beta_2 + \dots + \lambda_c^{n-1} \beta_{n-1} = e^{\lambda_c t} \\ \beta_0 + \lambda'_c \beta_1 + (\lambda'_c)^2 \beta_2 + \dots + (\lambda'_c)^{n-1} \beta_{n-1} = e^{\lambda'_c t} \end{cases} \quad (2.55)$$

Since the eigenvalues are complex and conjugate, it has that  $\lambda_c + \lambda'_c = 2\text{Re}[\lambda_c] = 2\alpha$  and  $\lambda_c - \lambda'_c = 2j\text{Im}[\lambda_c] = 2j\omega$ . Moreover, the Euler identity  $e^{\alpha \pm j\omega} = e^{\alpha t}(\cos(\omega t) \pm j\sin(\omega t))$  shows that  $e^{\lambda} + e^{\lambda'} = 2e^{\alpha t}\cos(\omega t)$  and  $e^{\lambda} - e^{\lambda'} = 2je^{\alpha t}\sin(\omega t)$ . In this case, summing the two rows and subtracting the first row by the second one results in:

$$\begin{cases} \beta_0 + 2\alpha\beta_1 + 2\alpha^2\beta_2 + \dots + 2\alpha^{n-1}\beta_{n-1} = 2e^{\alpha t}\cos(\omega t) \\ 0 + 2j\omega\beta_1 + 2j\omega^2\beta_2 + \dots + 2j\omega^{n-1}\beta_{n-1} = 2je^{\alpha t}\sin(\omega t) \end{cases} \quad (2.56)$$

Finally, dividing the first row by 2 and the second row by  $2j$  results in (2.54).  $\square$

From the same reasons stated before, this result implies that the actual response of the system will have sinusoidal components that produces oscillations in the response. The modes associated with complex and conjugate eigenvalues are classified as *pseudo-periodic*, since they are composed by an exponential growth (or decay) enveloping a sinusoidal function.

**Example 2.5.** Consider, again for the sake of illustration, the following example. A toy system is described by the matrix

$$\mathbf{A} = \begin{bmatrix} -0.1 & 0.5 \\ -0.5 & -0.1 \end{bmatrix}. \quad (2.57)$$

The eigenvalues of this matrix are given as  $\lambda, \lambda' = -0.1 \pm j0.5$ . The Sylvester expansion for the exponential of this matrix has the parameter vector  $\boldsymbol{\beta} = [\beta_0(t), \beta_1(t)]$  which solves the system

$$\begin{bmatrix} 1 & -0.1 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} \beta_0(t) \\ \beta_1(t) \end{bmatrix} = \begin{bmatrix} e^{-0.1t}\cos(0.5t) \\ e^{-0.1t}\sin(0.5t) \end{bmatrix}. \quad (2.58)$$

The solution to this system is  $\beta = [e^{-0.1t}\cos(0.5t) + 0.2e^{-0.1t}\sin(0.5t), 2e^{-0.1t}\sin(0.5t)]$ . Hence, state-transition matrix is calculated through the expansion:

$$\begin{aligned} e^{\mathbf{A}t} &= \beta_0(t)\mathbf{I} + \beta_1(t)\mathbf{A} \\ &= e^{-0.1t}\cos(0.5t) + 0.2e^{-0.1t}\sin(0.5t) \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + 2e^{-0.1t}\sin(0.5t) \begin{bmatrix} -0.1 & 0.5 \\ -0.5 & -0.1 \end{bmatrix}. \quad (2.59) \\ &= \begin{bmatrix} e^{-0.1t}\cos(0.5t) & e^{-0.1t}\sin(0.5t) \\ -e^{-0.1t}\sin(0.5t) & e^{-0.1t}\cos(0.5t) \end{bmatrix} \end{aligned}$$

The elements of the resulting state-transition matrix are shown in Fig. 2.6. It is possible to see, in this case, the pseudo-periodic behavior of the complex conjugate modes, where the dashed lines represents the exponential envelope.

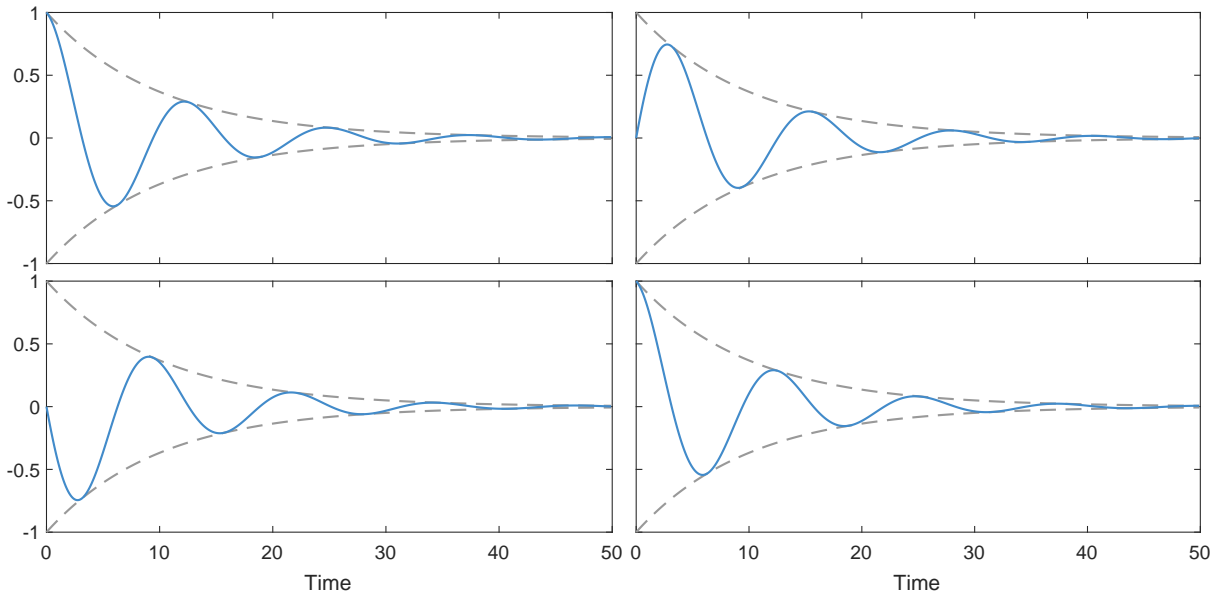


Figure 2.6: Simulation of the state-transition matrix  $e^{\mathbf{A}t}$  from Example 2.5 for  $t \in [0, 50]$ .

The previous discussion on the modes of matrix  $\mathbf{A}$  introduced the importance of the eigenvalues of this matrix in analyzing the system response. The analysis, however, focused on the natural response  $\mathbf{x}_n(t)$ , whereas the forced response  $\mathbf{x}_f(t)$  was deliberately omitted. Now, the analysis focus on the opposite case. Consider the two systems described by the matrices  $(\mathbf{A}^{(1)}, b^{(1)})$  and  $(\mathbf{A}^{(2)}, b^{(2)})$  given below. These matrices are more complete descriptions of the systems from Examples 2.4 and 2.5.

$$\mathbf{A}^{(1)} = \begin{bmatrix} -5.93 & 0 \\ 0.83 & -4.70 \end{bmatrix} \quad b^{(1)} = \begin{bmatrix} 3.81 \\ -1.09 \end{bmatrix}; \quad \mathbf{A}^{(2)} = \begin{bmatrix} -0.1 & 0.5 \\ -0.5 & -0.1 \end{bmatrix} \quad b^{(2)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}. \quad (2.60)$$

Considering initial states  $\mathbf{x}^{(1)}(0) = \mathbf{x}^{(2)}(0) = \mathbf{0}$  and the input signals  $\mathbf{u}^{(1)}(t) = \mathbf{u}^{(2)}(t) = \mathbf{1}$ , for time  $t \in [0, 60]$ , the evolution of the first state of both systems are shown in Fig. 2.7. Since  $\mathbf{A}^{(1)}$  is diagonal and  $\mathbf{A}^{(2)}$  has a pair of complex conjugate eigenvalues, it is possible to consider that those are characteristic responses of aperiodic and pseudo-periodic modes, respectively, to an unitary step input signal.

The unitary step is the default input signal used to analyze the forced response of dynamical systems. This is due to the fact that any continuous signal can be approximated by a sequence of step signals. When discussing the time evolution of systems given forcing input signals,



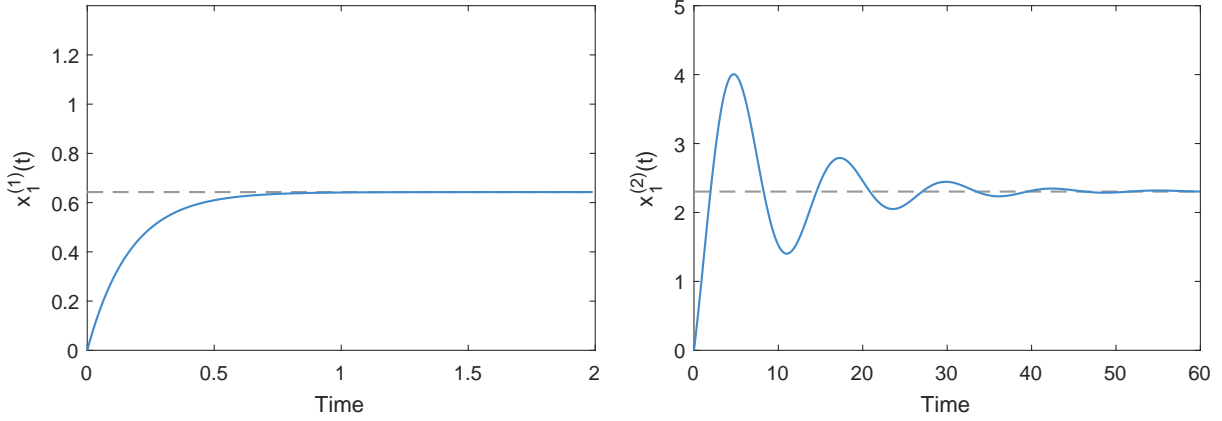


Figure 2.7: One-state forced response of the two systems from (2.60).

the response can be separated into a transient regime and a subsequent steady-state regime. The transient part of the response corresponds to the dynamical evolution from the previous state equilibrium to the next steady-state point. The steady-state point, as already mentioned, corresponds to the point in which the derivatives of the states are zero, which represents the discharge of the initial energy given by the forcing input. These regimes can be characterized in respect to the specifications given below.

**Definition 2.6.** (Unit-Step Response Specifications) Given a mode  $e^{\lambda t}$  associated with the eigenvalue  $\lambda \in \mathbb{R}$ , its contribution to the response has a time constant ( $\tau$ ) defined as

$$\tau = -\frac{1}{\lambda}. \quad (2.61)$$

Furthermore, given pseudo-periodic modes  $e^{\lambda t}$  and  $e^{\lambda' t}$  associated with the eigenvalues  $\lambda, \lambda' = \alpha \pm j\omega$  of a matrix  $\mathbf{A}$ , their contribution to the response has a time constant ( $\tau$ ), a natural frequency ( $\omega_n$ ) and a damping coefficient ( $\zeta$ ) defined as:

$$\tau = -\frac{1}{\alpha} \quad ; \quad \omega_n = \sqrt{\alpha^2 + \omega^2} \quad ; \quad \zeta = -\frac{\alpha}{\omega_n}. \quad (2.62)$$

The specifications just defined provides a quantitative way to describe the response of a system in terms of time and frequencies. The time constant, for instance, is a quantity that represents the time needed for the mode to lost 63% of its initial value, since  $e^{\lambda\tau} = e^{-1} = 0.37$ . A greater value of a time constant indicates that the system is able to “discharge” energy faster. The damping coefficient, in turn, provides an information about the intensity of the peak in the pseudo-periodic responses, which is known as *overshoot* (or *undershoot* in the case of a negative peak) and the natural frequency represents the oscillation of the response before reaching steady-state. From the perspective of control theory, these are some of the specifications used to define desirable transient responses to a controlled system.

Notice that the response specifications are always functions of the real and imaginary parts of the eigenvalues. This brings the possibility of a visualization in the complex plane to interpret how each eigenvalue contributes to the total response. A straightforward notion is that the closer an eigenvalue is to the imaginary axis, the faster is its contribution. Similarly, the furthest an eigenvalue is to the real axis, the more oscillatory is its contribution. Finally, a vector from the origin of the plane to a complex eigenvalue has a norm equal to the natural frequency ( $\omega_n$ )

and the cosine of the angle formed with the imaginary axis is equal to the damping factor ( $\zeta$ ). A simulation of the contributions from different eigenvalues are shown in Fig. 2.8.

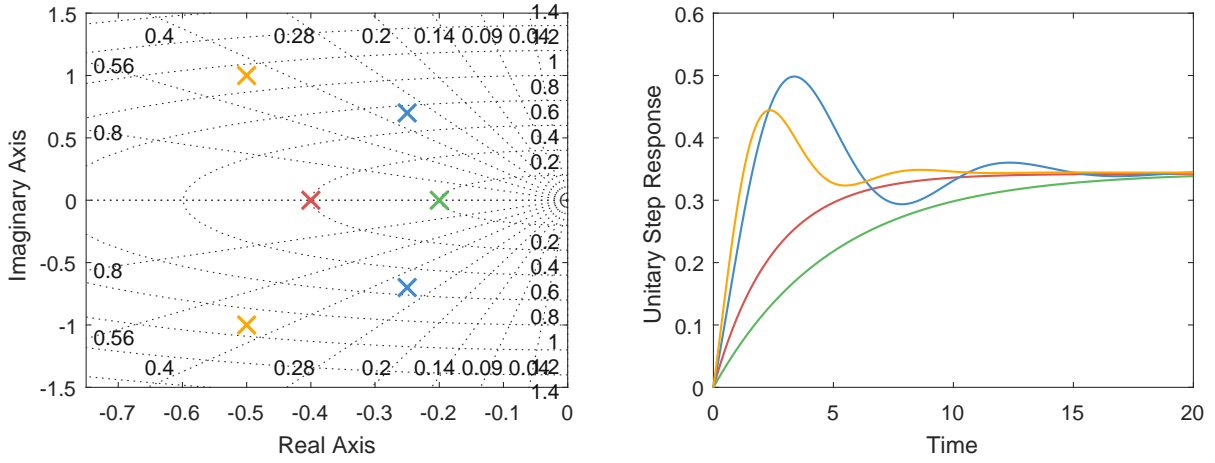


Figure 2.8: Eigenvalues of a system and the forced response associated with their modes.

## 2.4 Similarity Transformations

A State-Space representation can be interpreted through a system of coordinates. A state, in this context, represents a point as visualized through this reference. Under the assumption of a linear time-invariant system, there is an intuition that is possible to change the representation of the states by changing this system of coordinates through some linear transformation, obtaining a different representation to the same system. This is the motivation for this section.

First of all, consider a brief reflection about the geometrical interpretation of a State-Space model. Let any state-vector at an arbitrary time  $t$  be  $\mathbf{x}(t)$ , defined in the  $\mathbb{R}^2$  space. Naturally, it is possible to associate to this vector an orthonormal basis  $\mathbf{I} \in \mathbb{R}^2$ , the 2-dimensional identity matrix. There is, however, the possibility to associate any other arbitrary basis to represent a state-vector and visualize both states through this new perspective. When the basis is not orthogonal, a change in the state-vector in a direction parallel to a component of the basis, i.e., a change in only one element of this vector, also produces a change in other directions if observed through the original orthonormal basis. This is an interesting result, since it allows for the use of a single direction to represent the simultaneous evolution of both states.

In the State-Space formulation, the matrix  $\mathbf{A}$  represents a linear function that maps state-vectors from  $\mathbb{R}^n$  to itself. When applying a new basis to represent the state-vectors, it is intuitive that the mapping performed by this function also changes so that it still represents the same linear combination of the states.

**Theorem 2.7.** (*Similarity Transformation*) Consider a system in SS representation described by the matrices  $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$  and a nonsingular transformation matrix  $\mathbf{P} \in \mathbb{R}^{n \times n}$ . An equivalent representation for the transformation  $\mathbf{z}(t) = \mathbf{P}\mathbf{x}(t)$  is:

$$\begin{cases} \dot{\mathbf{z}}(t) = \tilde{\mathbf{A}}\mathbf{z}(t) + \tilde{\mathbf{B}}\mathbf{u}(t) \\ \mathbf{y}(t) = \tilde{\mathbf{C}}\mathbf{z}(t) + \tilde{\mathbf{D}}\mathbf{u}(t) \end{cases}, \quad (2.63)$$

where:

$$\tilde{\mathbf{A}} = \mathbf{P}\mathbf{A}\mathbf{P}^{-1} \quad ; \quad \tilde{\mathbf{B}} = \mathbf{P}\mathbf{B} \quad ; \quad \tilde{\mathbf{C}} = \mathbf{C}\mathbf{P}^{-1} \quad ; \quad \tilde{\mathbf{D}} = \mathbf{D}. \quad (2.64)$$

*Proof.* Consider a SS representation and any nonsingular matrix  $\mathbf{P} \in \mathbb{R}^{n \times n}$ . Making  $\mathbf{z}(t) = \mathbf{P}\mathbf{x}(t)$  leads to  $\mathbf{x}(t) = \mathbf{P}^{-1}\mathbf{z}(t)$ . Substituting this in the state equation results in:

$$\begin{aligned}\dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) \\ \mathbf{P}^{-1}\dot{\mathbf{z}}(t) &= \mathbf{A}\mathbf{P}^{-1}\mathbf{z}(t) + \mathbf{B}\mathbf{u}(t) \\ \mathbf{I}\dot{\mathbf{z}}(t) &= \mathbf{P}\mathbf{A}\mathbf{P}^{-1}\mathbf{z}(t) + \mathbf{P}\mathbf{B}\mathbf{u}(t) \\ \dot{\mathbf{z}}(t) &= \tilde{\mathbf{A}}\mathbf{z}(t) + \tilde{\mathbf{B}}\mathbf{u}(t)\end{aligned}\tag{2.65}$$

Thus, substituting  $\mathbf{x}(t) = \mathbf{P}^{-1}\mathbf{z}(t)$  in the output equation results in:

$$\begin{aligned}\mathbf{y}(t) &= \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t) \\ &= \mathbf{C}\mathbf{P}^{-1}\mathbf{z}(t) + \mathbf{D}\mathbf{u}(t). \\ &= \tilde{\mathbf{C}}\mathbf{z}(t) + \tilde{\mathbf{D}}\mathbf{u}(t)\end{aligned}\tag{2.66}$$

□

By this theorem it is clear that, after transforming the state-vector, the entire dynamical model  $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$  changes. This would imply that the analysis of the original system, for its properties and time response, is not valid for the similar transformed system. However, as discussed before, these transformations accounts for the same system when observed through a different reference basis. Therefore, it is expected that the model presents the same analysis results, as demonstrated below.

**Theorem 2.8.** *Consider a system in State-Space form with matrix  $\mathbf{A}$ . Consider also a similarity transformation  $\mathbf{z}(t) = \mathbf{P}\mathbf{x}(t)$  that results in a matrix  $\tilde{\mathbf{A}} = \mathbf{P}\mathbf{A}\mathbf{P}^{-1}$ . In this case,  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$  shares the same set of eigenvalues.*

*Proof.* The eigendecomposition problem is defined as  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ . From the similarity transformation,  $\tilde{\mathbf{A}} = \mathbf{P}\mathbf{A}\mathbf{P}^{-1}$  leads to  $\mathbf{A} = \mathbf{P}^{-1}\tilde{\mathbf{A}}\mathbf{P}$ . Substituting this in the eigendecomposition:

$$\begin{aligned}\mathbf{P}^{-1}\tilde{\mathbf{A}}\mathbf{P}\mathbf{v} &= \lambda\mathbf{v} \\ \tilde{\mathbf{A}}\mathbf{P}\mathbf{v} &= \lambda\mathbf{P}\mathbf{v}\end{aligned}\tag{2.67}$$

Considering the transformed eigenvector  $\tilde{\mathbf{v}} = \mathbf{P}\mathbf{v}$ , it is clear that  $\tilde{\mathbf{A}}\tilde{\mathbf{v}} = \lambda\tilde{\mathbf{v}}$ . This implies that matrices  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$  shares the same set of eigenvalues  $\lambda$ . □

It is clear from previous results that, since both the matrices shares the same set of eigenvalues, the original and transformed model provides the same dynamical responses and, as shown later, the same general dynamic properties. Therefore, the similarity transformation consists in a method to produce new State-Spaces representations that emphasizes some geometrical perspective of the model, hopefully helping to analyze a specific property, without actually changing the relationship between the original model and the physical system.

The use of similarity transformations can also benefits the computation of functions of the matrices of the State-Space representation, given that they impose a desirable structure to this matrix. With this motivation, a popular transformation that provides a new representation with computational advantages is the Similarity transformation.

**Theorem 2.9.** (*Diagonalization*) Consider a  $n$ -dimensional system in State-Space form represented by the matrices  $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ . Let matrix  $\mathbf{A}$  have  $n$  distinct real eigenvalues, i.e.,  $\lambda \in \mathbb{R}^n$ . Performing the transformation  $\mathbf{z}(t) = \mathbf{V}\mathbf{x}(t)$ , where  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n]$  is the modal matrix composed by the eigenvectors  $\mathbf{v}_i \in \mathbb{R}^n$ ,  $i \in [1, \dots, n]$ , of matrix  $\mathbf{A}$ , will result in a transformed matrix  $\mathbf{\Lambda} = \mathbf{V}\mathbf{A}\mathbf{V}^{-1}$  is diagonal.

*Proof.* Since the eigenvalues are real and distinct, the eigenvectors must be linearly independent, proving that the inverse  $\mathbf{V}^{-1}$  always exist and that it is a feasible transformation matrix. Using the identity for the eigendecomposition of matrix  $\mathbf{A}$ :

$$\begin{aligned}
 \lambda \mathbf{v} &= \mathbf{A} \mathbf{v} \\
 [\lambda_1 \mathbf{v}_1 \quad \lambda_2 \mathbf{v}_2 \quad \dots \quad \lambda_n \mathbf{v}_n] &= [\mathbf{A} \mathbf{v}_1 \quad \mathbf{A} \mathbf{v}_2 \quad \dots \quad \mathbf{A} \mathbf{v}_n] \\
 [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_n] \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix} &= \mathbf{A} [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_n] \\
 &= [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_n]^{-1} \mathbf{A} [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_n] \\
 \mathbf{\Lambda} &= \mathbf{V}^{-1} \mathbf{A} \mathbf{V}
 \end{aligned} \tag{2.68}$$

Thus,  $\mathbf{\Lambda}$  is a diagonal matrix whose elements are the eigenvalues of matrix  $\mathbf{A}$ .  $\square$

A model with a diagonal matrix  $\mathbf{A}$  has the nice property that the evolution of the states are decoupled, in the sense that each state evolution is a linear function of itself. A geometrical interpretation of this transformation is that the eigenvectors of this matrix produces a basis that encode information about the interaction between those states, while the interaction in the original formulation was the linear combination of the modes. This result can be easily extended to the case where the eigenvalues are distinct conjugate complex pairs.

In this diagonalization procedure, the elements of the resulting matrix  $\mathbf{\Lambda}$  are the very own eigenvalues of the original matrix  $\mathbf{A}$ . Furthermore, it is easy to verify that the state-transition matrix for the transformed matrix,  $e^{\mathbf{\Lambda}t}$ , is also a diagonal matrix whose elements are the modes of the system, and can be easily computed as:

$$\begin{aligned}
 e^{\mathbf{\Lambda}t} &= \sum_{k=0}^{\infty} \frac{\mathbf{A}^k t^k}{k!} = \sum_{k=0}^{\infty} \frac{t^k}{k!} \begin{bmatrix} \lambda_1^k & 0 & \dots & 0 \\ 0 & \lambda_2^k & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n^k \end{bmatrix} = \begin{bmatrix} \sum_{k=0}^{\infty} \frac{t^k \lambda_1^k}{k!} & 0 & \dots & 0 \\ 0 & \sum_{k=0}^{\infty} \frac{t^k \lambda_2^k}{k!} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sum_{k=0}^{\infty} \frac{t^k \lambda_n^k}{k!} \end{bmatrix} \\
 &= \begin{bmatrix} e^{\lambda_1 t} & 0 & \dots & 0 \\ 0 & e^{\lambda_2 t} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & e^{\lambda_n t} \end{bmatrix}.
 \end{aligned} \tag{2.69}$$

In the case that the eigenvalues are not all distinct, it may be not possible to design a modal matrix in the same way, since the eigenvectors could not be all linearly independent, and the matrix would not form a basis. In those cases, however, it is still possible to design a generalized modal matrix that transform the original matrix  $\mathbf{A}$  to a quasi-diagonal matrix  $\mathbf{J}$ , in which there are decoupled states between blocks. This similarity transformation is known as the Jordan form [Strang, 2016], and it generalizes the notion of diagonalization for any arbitrary matrix.

## 2.5 Stability, Controlability and Observability

When moving from the dynamical system analysis to a practical perspective of control theory, it is necessary to define and discuss some important properties system and their models. These are definitions concerning the stability, controllability and observability of dynamical systems. Although those properties are established from the dynamical model of the system, they can be directly related to questions about control and measurement instrumentation applied to the control system. .

### Stability

The first property to be discussed consists in the stability of a system. An instable system, as the name suggests, is a system whose response does not converge to a specific value and rather oscillates or grows unbounded. In physical scenarios, unstable systems are problematic, since their response to external stimuli can result in dangerous situations to itself and, maybe, to the environment around it. Because of this, determining the stability of a system is a crucial procedure into analyzing a system that will be controlled. Under the several quantitative methods to determine if a system is indeed stable, given a mathematical model, a popular and practical one is the Bounded-Input Bounded-Output (BIBO) stability criteria.

**Definition 2.7.** (BIBO Stability) A dynamical system is defined as BIBO stable if every bounded input stimuli  $|\mathbf{u}(t)| \leq \epsilon < \infty$  produces in it a bounded output response  $|\mathbf{y}(t)| \leq \delta < \infty$ .

The main result behind this criteria is that the natural response of a system should vanish as time evolves, i.e.,  $\mathbf{x}(t) \rightarrow \mathbf{0}$  as  $t \rightarrow \infty$ . This result is very intuitive from Theorem 2.4, since the vanishing of the natural response implies that  $e^{\mathbf{A}t} \rightarrow \mathbf{0}$  as  $t \rightarrow \infty$  and the forced response is expected to be bounded, if  $\mathbf{u}(t)$  is bounded. Since the state-transition matrix is a linear combination of the modes, and the modes are exponential functions of the eigenvalues of the matrix  $\mathbf{A}$ , it is possible to determine a condition for stability in the light of these quantities, as shown in the following theorem.

**Theorem 2.10.** (BIBO Stability in SS) A system in State-Space form, represented by a matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  with eigenvalues  $\lambda \in \mathbb{R}^n$ , is BIBO stable if and only if  $\text{Re}[\lambda_i] < 0$ ,  $\forall i \in [1, 2, \dots, n]$ .

A detailed proof of this theorem can be found in Appendix A. First of all, note that this criteria depends only on the eigenvalues of matrix  $\mathbf{A}$ , so the stability property of a system is invariant to any similarity transformation, since  $\mathbf{A}$  and any transformed matrix  $\tilde{\mathbf{A}} = \mathbf{PAP}^{-1}$  shares the same eigenvalues. From the previous results it is also known that the real part of the eigenvalues, independent of their multiplicity or domain, appears as the arguments of the exponential functions that are the system modes. Therefore, an eigenvalue with a negative

real part will produce a mode that decays exponentially, as this theorem shows. By the same argument, if the matrix  $\mathbf{A}$  has at least one eigenvalue  $\lambda_j = 0$  such that  $\Re[\lambda_i] \leq 0, \forall i \in [1, 2, \dots, n]$ , then the mode associated with this eigenvalue is a constant and the natural response becomes bounded. This configuration is known as a marginally stable condition, in the BIBO perspective. The time response of a 2-nd order system is shown Fig. 2.9 for three different poles configurations, given an unitary step.

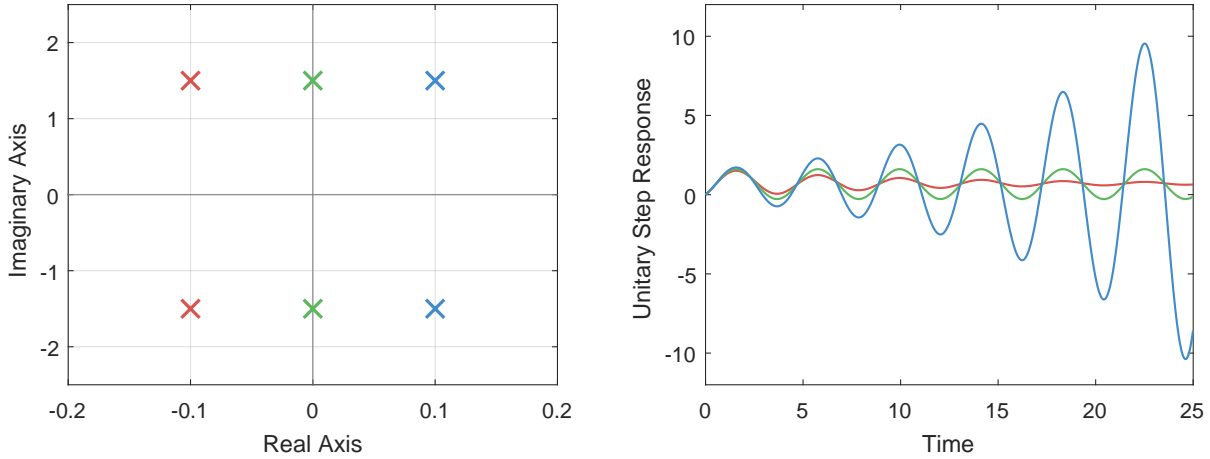


Figure 2.9: Stability of forced responses given the positions of the system poles.

### Controllability

Later chapters will discuss the possibility of stabilize an unstable system through a controller. There are, however, some restrictions to the possibility of controlling or not the states of a system, which includes the necessity of discussing *controllability*. The controllability of a system states whether it is possible to calculate an input signal that drives the system to any arbitrary point in the space, in a finite time interval. This property accounts exclusively for this possibility, in the sense that it does not account for the operational feasibility of actually applying this input signal into a physical system, since it may need more energy than an actuator can produce.

**Definition 2.8.** (Controllability) A system in State-Space form with matrices  $(\mathbf{A}, \mathbf{B})$  is said to be controllable if, for any initial state  $\mathbf{x}(t_0) = \mathbf{x}_0$  and terminal state  $\mathbf{x}(T) = \mathbf{x}_T, T < \infty$ , there exists an input signal  $\mathbf{u}(t), t \in [t_0, T]$ , that can transfer  $\mathbf{x}(t_0)$  to  $\mathbf{x}(T)$ . Otherwise, the system is said to be uncontrollable.

There are several methods to analyze the controllability of a system given a mathematical model and the definition above. A popular criteria introduces the concept of a controllability matrix and has a nice geometrical interpretation.

**Theorem 2.11.** (Controllability in SS) Consider a system in linear State-Space form with matrices  $(\mathbf{A}, \mathbf{B})$  and the controllability matrix  $\mathbf{C} \in \mathbb{R}^{n \times nr}$  defined as:

$$\mathbf{C} = [\mathbf{B} \quad \mathbf{AB} \quad \mathbf{A}^2\mathbf{B} \quad \dots \quad \mathbf{A}^{n-1}\mathbf{B}]. \quad (2.70)$$

The system is controllable if and only if  $\mathbf{C}$  has full row rank.

An intuition behind this theorem is that the full row rank condition implies that  $\mathcal{C}$  has  $n$  linearly independent columns, therefore these columns can be used as a basis for what is known as the *controllable subspace*. To better understand that, consider the following forced solution  $\mathbf{x}_f(t)$  given by the Lagrange formula:

$$\mathbf{x}_f(t) = \int_0^t e^{\mathbf{A}(t-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau. \quad (2.71)$$

From the Cayley-Hamilton theorem, it is possible to represent  $e^{\mathbf{A}(t-\tau)}$  as a linear combination of scalars  $\beta_i(t-\tau)$  and powers of the matrix  $\mathbf{A}^i$ ,  $i \in [0, 1, \dots, n-1]$ . Using this theorem to substitute the matrix exponential at (2.71), and letting  $\tau_2 = t - \tau$  for easier manipulation, the forced response can be represented as:

$$\begin{aligned} \mathbf{x}_f(t) &= \int_0^t \left( \sum_{i=0}^{n-1} \beta_i(\tau_2) \mathbf{A}^i \right) \mathbf{B} \mathbf{u}(t - \tau_2) d\tau_2 = \sum_{i=0}^{n-1} (\mathbf{A}^i \mathbf{B}) \int_0^t \beta_i(\tau_2) \mathbf{u}(t - \tau_2) d\tau_2 \\ &= \sum_{i=0}^{n-1} (\mathbf{A}^i \mathbf{B}) \tilde{\beta}_i(\mathbf{u}, t), \end{aligned} \quad (2.72)$$

where  $\mathbf{A}^i \mathbf{B}$  are the columns of the matrix  $\mathcal{C}$  and  $\tilde{\beta}_i(\mathbf{u}, t)$ , comprising the integral term, is a function that depends only on the input signal  $\mathbf{u}(t)$  and time  $t$ . This result implies that the forced response  $\mathbf{x}_f(t)$  is a linear combination given by the columns of  $\mathcal{C}$ . If  $\mathcal{C}$  has  $n$  linearly independent columns, then it spans the entire  $n$ -dimensional space, i.e., the entire state space, and thus any desirable state vector can be reached. If the column rank of  $\mathcal{C}$  is less than  $n$ , then only a subspace of smaller dimension can be reached.

### Observability

While the discussion on controllability concerns the possibility of driving a system to a desirable state through an actuator signal, there is also the necessity to discuss the possibility of determining the internal state of a system given the output signal  $\mathbf{y}(t)$ . This property, known as *observability*, comes from the fact that any output of a system, related to the states through the matrix  $\mathcal{C}$ , may be a combination of states, and that some states may not even be present in the output signal. Thus, it is necessary to know if it is possible to reconstruct  $\mathbf{x}(t)$  through  $\mathbf{y}(t)$ .

**Definition 2.9.** Observability A system in State-Space form with matrices  $(\mathbf{A}, \mathcal{C})$  is said to be observable if, given an input signal  $\mathbf{u}(t)$  and output signal  $\mathbf{y}(t)$ , over the interval  $t \in [t_0, T]$ , it is possible to uniquely determine the value of the initial state  $\mathbf{x}(t_0)$ . Otherwise, the system is said to be unobservable.

Similarly to the controllability property, there are several ways to analyze the observability of a system, given a mathematical model. A popular criteria introduces the concept of an observability matrix.

**Theorem 2.12.** (Observability in SS) Consider a system in linear State-Space form with matrices  $(\mathbf{A}, \mathcal{C})$  and the observability matrix  $\mathcal{O} \in \mathbb{R}^{nq \times n}$  defined as:

$$\mathcal{O} = [\mathcal{C} \quad \mathcal{C}\mathbf{A} \quad \mathcal{C}\mathbf{A}^2 \quad \dots \quad \mathcal{C}\mathbf{A}^{n-1}]^T. \quad (2.73)$$

The system is observable if and only if  $\mathcal{O}$  has full column rank.



The interpretation of this theorem follows the same intuition of before: if the matrix  $\mathbf{O}$  has full column rank, then it can be used as a basis to span a subspace with the same dimension as the State-Space. In fact, the proof of both theorems follows the same procedures and there is a direct relationship between controllability and observability, known as the Theorem of Duality.

The concepts of controllability, observability and their duality, together with the conditions to fulfill these properties, were first introduced by [Kalman, 1960b]. The geometrical interpretations of both theorems may suggest a practical solution for the cases where a system is uncontrollable or unobservable. In the first case, the solution would be to add specific actuators to the system, in the condition that they are linearly independent between themselves and the ones actually in operation. The same procedure can be done to solve the unobservable problem, but adding more sensors instead. Those procedures would change the matrices  $\mathbf{B}$  and  $\mathbf{C}$  and could ensure the necessary conditions in matrices  $\mathbf{C}$  and  $\mathbf{O}$ . Of course, the implementation of such instrumentation may not be practical, due to technical or economical constraints. In the case which these additional devices are not feasible, one is restricted to control and/or observe only a subspace of the state variables.

## 2.6 Response Analysis in the Frequency Domain

Although the response of dynamical systems are naturally perceived in time, there are advantages of analyzing the models in a frequency domain perspective. This analysis differs from a simple time domain analysis from the fact that, in a steady-state regime, the response of a linear system for a sinusoidal input is itself sinusoidal, with the same frequency but different amplitude and phase, as illustrated at Fig. 2.10.

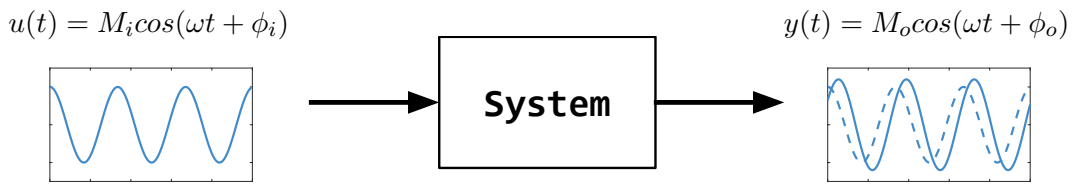


Figure 2.10: Illustration of the steady-state response of LTI systems to sinusoidal inputs.

A common representation is the phasor representation, where  $M\angle\phi = M\cos(\omega t + \phi)$ . In the context of SISO Input-Output models, the response of the system can be summarized by a transfer function which is also a phasor:

$$g(t) = \frac{y(t)}{u(t)} = \frac{M_o\angle\phi_o}{M_i\angle\phi_i} = M_g\angle\phi_g, \quad (2.74)$$

where  $M_g = M_o/M_i$  and  $\phi_g = \phi_o - \phi_i$ . Notice that this formulation makes the time-dependance implicit in the system response, since it is periodic in this case. For this reason, the system response can be visualized as a function of frequency rather than a function of time, and the properties of the system can be accessed in this way. The two most popular techniques for frequency response analysis are Bode plots [Bode, 1945] and Nyquist diagrams [Nyquist, 1932]. The first is a direct plot of  $M_g(\omega)$  and  $\phi_g(\omega)$  for several values of  $\omega$ , making  $M_g(\omega) = 20\log|G(j\omega)|$  and  $\phi_g(\omega) = \angle G(j\omega)$ , where  $G(j\omega)$  is a transfer function of a system evaluated for an input signal with exclusively oscillatory components. The Nyquist diagram, in the other hand, is a direct phasor visualization obtained by applying the Argument Principle to a contour containing the entire right-hand side of the complex plane. Both visualizations are depicted in Fig. 2.11, for a transfer function obtained by applying the transformation of Theorem 2.1 to the system from (2.51):

$$\mathbf{G}(s) = \left[ \frac{3.81}{s + 5.92} \quad \frac{-1.09s - 3.33}{s^2 + 10.62s + 27.84} \right]. \quad (2.75)$$



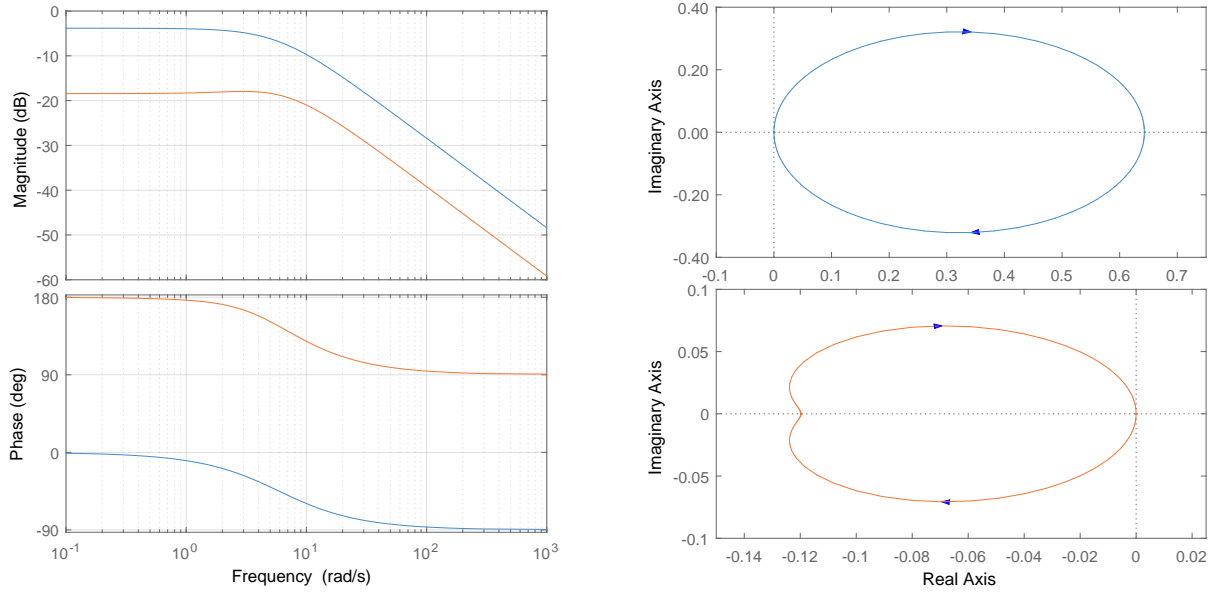


Figure 2.11: Visualization of the two-states system in (2.75) in Bode plots (left) and Nyquist diagrams (right). The blue and orange lines represents the states  $x_1$  and  $x_2$ , respectively.

The immediate advantage of these visualizations is that is not necessary to compute an entire simulation of this system, until it reaches steady-state, to be able to perform analysis, which is really critic for high-dimensional systems with slow time-constants. In addition to that, these visualizations (specially the Bode plot) can be easily sketched by hand with fairly accuracy, allowing for some understanding of the system without need for solving differential equations or the inverse Laplace transforms. For these reasons, frequency responses methods for analyzing systems were ubiquitous for many years in industry applications, and some properties assessments, such as closed-loop stability, are still better understood under this formulation, as it will be shown in later chapters.

## Chapter 3

# Controller Synthesis

This chapter discusses the general results and properties for the design of controllers, focusing on feedback architectures. The devices are motivated and formulated using the State-Space model for dynamical systems, so the feedback is performed on the state response rather than on the system output. For this reason, the results in this section focus on this formulation explicitly through the state-equations, where the output-equations are made implicit.

### 3.1 State Feedback Controllers

In modern control theory, advances in computer performance and the maturing of State-Space model analysis have made the design of controllers using State-Space models feasible for real-world applications. This is usually desirable since these kind of models provides a practical solution to understand dynamical system response and properties, so it is natural to want a controller design technique that accounts for that representation. The most basic, yet very popular, feedback controller used in those settings is the *Full-State Feedback Controller*, defined below.

**Definition 3.1.** Full-State Linear Feedback Given a linear system in State-Space representation, an input action  $\mathbf{u}(t)$  is calculated by the linear control law  $\pi(\cdot)$  through state-feedback as:

$$\mathbf{u}(t) = \pi(\mathbf{r}(t), \mathbf{x}(t)) = \mathbf{r}(t) - \mathbf{K}\mathbf{x}(t), \quad (3.1)$$

where  $\mathbf{r} : \mathbb{R} \rightarrow \mathbb{R}^n$  is a state reference signal that the system must follows and  $\mathbf{K} \in \mathbb{R}^{r \times n}$  is the *feedback gain matrix*.

A block diagram of the closed-loop system for the state-feedback controller is shown at Fig. 3.1. In this schematic each blocks represents a function or operation that is applied to the denoted signals, which themselves are represented by the arrows. The direction of the arrows indicates whether the signal variable is an operand or a result of an operation.

Notice that the control law of the full-state feedback is both linear and time-invariant, which makes the analysis of the closed-loop system similar to the one used in open-loop configurations. Of course, this is a specific choice of control law, and feedback controllers can also be defined using nonlinear or time-dependent functions. Notice that this new definition for the calculation of  $\mathbf{u}(t)$  allows for the following closed-loop representation of the system:

$$\begin{aligned} \dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t) + \mathbf{B}(\mathbf{r}(t) - \mathbf{K}\mathbf{x}(t)) \\ &= (\mathbf{A} - \mathbf{BK})\mathbf{x}(t) + \mathbf{B}\mathbf{r}(t) \end{aligned} \quad (3.2)$$

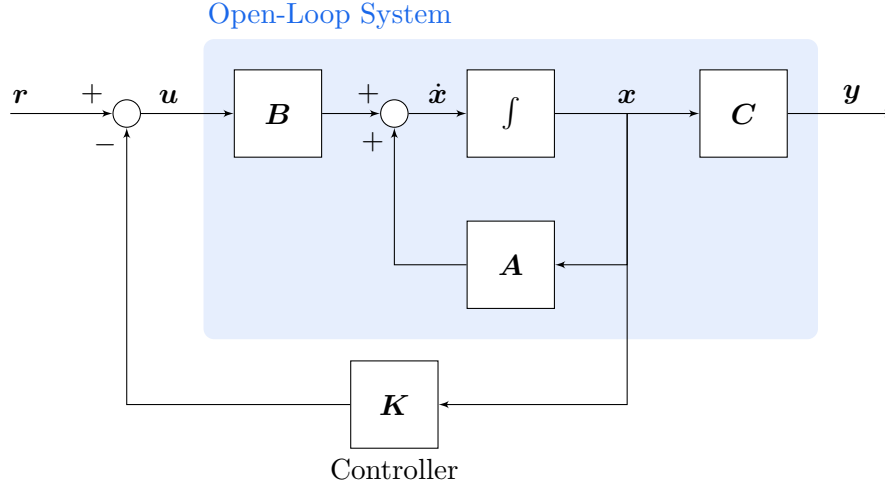


Figure 3.1: Block diagram of a state-feedback closed-loop system.

Clearly, the inclusion of a feedback controller in the loop is equivalent to transform an open-loop system into a new system  $(\mathbf{A}_{cl}, \mathbf{B})$ , with  $\mathbf{A}_{cl} = \mathbf{A} - \mathbf{B}\mathbf{K}$ , whose manipulated variable is now the reference signal  $\mathbf{r}(t)$ . Therefore, an appropriate choice of  $\mathbf{K}$  allows to modify the behavior of the closed-loop system such that a desirable behavior is obtained. To understand better the capabilities of the state feedback, consider the following theorems.

**Theorem 3.1.** (*Controller Canonical Form*) If a SISO system in State-Space representation is controllable, then by applying the transformation  $\mathbf{z}(t) = \mathbf{P}\mathbf{x}(t)$ , for a matrix  $\mathbf{P}$  calculated as the inverse of:

$$\mathbf{P}^{-1} = \mathbf{C} \begin{bmatrix} 1 & \alpha_1 & \cdots & \alpha_{n-1} \\ 0 & 1 & \cdots & \alpha_{n-2} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}, \quad (3.3)$$

where  $[\alpha_1, \alpha_2, \dots, \alpha_{n-1}]$  are the  $n - 1$  first coefficients of the characteristic polynomial  $\Delta(s) = \det(s\mathbf{I} - \mathbf{A})$ , the resulting representation is in the controller canonical form given as:

$$\begin{cases} \dot{\mathbf{z}}(t) = \begin{bmatrix} -\alpha_1 & -\alpha_2 & \cdots & -\alpha_{n-1} & -\alpha_n \\ 1 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix} \mathbf{z}(t) + \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} u(t) \\ y(t) = [\beta_1 \ \beta_2 \ \cdots \ \beta_n] \mathbf{z}(t) \end{cases} \quad (3.4)$$

Details for the proof of this theorem can be found in [Chen, 1998]. An equivalent result can be provided for MIMO systems, but the result is cumbersome and does not highlight the main results for the state-feedback. Using the canonical form just presented, it is possible to discuss an important result for the state feedback controller.

**Theorem 3.2.** (*Pole-Placement Method*) If a system in State-Space representation is controllable, then by state feedback using a gain matrix  $\mathbf{K} \in \mathbb{R}^{r \times n}$  the eigenvalues of  $\mathbf{A}_{cl} = \mathbf{A} - \mathbf{B}\mathbf{K}$ , the

poles of the closed-loop system, can be placed arbitrarily in the complex plane, as long as complex conjugate eigenvalues are assigned in pairs.

*Proof.* Suppose that the system is controllable. In this case, it can be converted to the controller canonical form of Theorem 3.1. Substituting  $\mathbf{z}(t) = \mathbf{P}\mathbf{x}(t)$  results in the following control law for the state feedback:

$$u(t) = r(t) - \mathbf{K} (\mathbf{P}^{-1}\mathbf{z}(t)) = r(t) - \tilde{\mathbf{K}}\mathbf{z}(t). \quad (3.5)$$

Applying the state feedback, the transformed closed-loop is given by:

$$\tilde{\mathbf{A}}_{cl} = \mathbf{P}(\mathbf{A} - \mathbf{BK})\mathbf{P}^{-1} = \mathbf{PAP}^{-1} - \mathbf{PBK}\mathbf{P}^{-1} = \tilde{\mathbf{A}} - \tilde{\mathbf{B}}\tilde{\mathbf{K}}. \quad (3.6)$$

From Theorem 2.8 it is known that  $\mathbf{A}_{cl}$  and  $\tilde{\mathbf{A}}_{cl}$  share the same set of eigenvalues, and, consequently, the same characteristic equations. Consider this characteristic equation as:

$$\Delta(s) = \det(s\mathbf{I} - \mathbf{A}) = s^n + \alpha_1 s^{n-1} + \alpha_2 s^{n-2} + \cdots + \alpha_{n-1}s + \alpha_n. \quad (3.7)$$

Given a set of coefficients  $[\tilde{\alpha}_1, \tilde{\alpha}_2, \dots, \tilde{\alpha}_n]$  for a characteristic polynomial whose roots are the closed-loop eigenvalues, define the transformed feedback gain matrix as:

$$\tilde{\mathbf{K}} = [\tilde{\alpha}_1 - \alpha_1 \quad \tilde{\alpha}_2 - \alpha_2 \quad \cdots \quad \tilde{\alpha}_n - \alpha_n]. \quad (3.8)$$

Substituting this in (3.6), it is easy to see that the resulting representation is:

$$\begin{cases} \dot{\mathbf{z}}(t) = \begin{bmatrix} -\tilde{\alpha}_1 & -\tilde{\alpha}_2 & \cdots & -\tilde{\alpha}_{n-1} & -\tilde{\alpha}_n \\ 1 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix} \mathbf{z}(t) + \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} u(t) \\ y(t) = [\beta_1 \quad \beta_2 \quad \cdots \quad \beta_n] \mathbf{z}(t) \end{cases}, \quad (3.9)$$

whose characteristic polynomial is now described by the designed coefficients that yield the desirable eigenvalues. Since  $\tilde{\mathbf{A}}_{cl}$  and  $\mathbf{A}_{cl}$  shares the same set of eigenvalues, it can be concluded that it is possible to assign the system poles directly through matrix  $\tilde{\mathbf{K}}$ .  $\square$

Notice that an “original” feedback gain matrix can be obtained as  $\mathbf{K} = \tilde{\mathbf{K}}\mathbf{P}$  and still yield the same eigenvalues assignment directly in  $\mathbf{A}_{cl}$  (since  $\mathbf{P}$  is just a linear transformation). The theorem has the direct result that, under full-state feedback, the transient response of a linear system can be completely determined by including a controller, which is described by this matrix  $\mathbf{K}$ . The result is preserved for MIMO systems, although the design of  $\mathbf{K}$  is not as straightforward, since it is not unique for a desired set of eigenvalues in this case [Moore, 1975].

Using the parameters from Definition 2.6, the positions of the closed-loop system poles can be designed given desirable operations, and the matrix  $\mathbf{K}$  can be hand-designed to meet these requirements. This method is known as the *Pole-Placement method* for control synthesis. Algorithm 1 summarizes a simple procedure of designing an appropriate feedback gain matrix, given desirable pole positions.

Albeit being a simple formula, Pole-Placement can be used in several applications to yield controllers capable of achieve performance requirements. Of course, the designer must take into account that, although the eigenvalue assignment allows for the whole complex plane, a careless choice of eigenvalues could result in controllers requesting very aggressive or oscillatory input signals. For this reason, it is necessary some knowledge of the instruments limits before designing the matrix  $\mathbf{K}$ .

**Algorithm 1:** Pole-Placement Method for SISO Systems**Input** : State-space model  $(\mathbf{A}, \mathbf{B})$  and a set of  $n$  desired eigenvalues  $\boldsymbol{\lambda}^*$ .**Output** : Feedback gain matrix  $\mathbf{K}$ .

- 1 Calculate  $[\alpha_1, \alpha_2, \dots, \alpha_n]$  as the coefficients of the polynomial  $\Delta(s) = \det(s\mathbf{I} - \mathbf{A})$ ;
- 2 Let  $\mathbf{P}^{-1} = [\mathbf{B} \quad \mathbf{AB} \quad \mathbf{A}^2\mathbf{B} \quad \dots \quad \mathbf{A}^{n-1}\mathbf{B}] \begin{bmatrix} 1 & \alpha_1 & \dots & \alpha_{n-1} \\ 0 & 1 & \dots & \alpha_{n-2} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$ ;
- 3 Let  $\mathbf{P} = (\mathbf{P}^{-1})^{-1}$ ;
- 4 Calculate  $[\tilde{\alpha}_1, \tilde{\alpha}_2, \dots, \tilde{\alpha}_n]$  as the coefficients of the polynomial  $\Delta_{cl}(s) = \prod_{i=1}^n (s - \lambda_i^*)$ ;
- 5 Let  $\tilde{\mathbf{K}} = [\alpha_1 - \tilde{\alpha}_1 \quad \alpha_2 - \tilde{\alpha}_2 \quad \dots \quad \alpha_n - \tilde{\alpha}_n]$ ;
- 6 Return  $\mathbf{K} = \tilde{\mathbf{K}}\mathbf{P}$ ;

## 3.2 Regulation and Reference Tracking

When discussing controller synthesis it is also necessary to account for the objective that the device is expected to fulfill. There exists two modes of operation that the controller may require for the system: regulation and reference tracking. In the case of state-feedback, these two classes of controllers differs only by what type of reference signal  $\mathbf{r}(t)$ , for an operation in a time interval  $t \in [t_0, t_f]$ , the system is expected to follow.

### Regulator Problem

The following statements gives a formal definition of a controller for regulation:

**Definition 3.2.** (Regulator) If a state-feedback controller has to make a system follows the reference  $\mathbf{r}(t) = \mathbf{0}$ , as  $t \rightarrow \infty$ , it is said to be a *regulator*. In this case, the closed-loop state equation and equivalent solutions reduces to the following:

$$\begin{array}{ll} \textbf{State Equation:} & \textbf{Lagrange solution:} \\ \dot{\mathbf{x}}(t) = (\mathbf{A} - \mathbf{BK}) \mathbf{x}(t) & \mathbf{x}(t) = e^{(\mathbf{A} - \mathbf{BK})t} \mathbf{x}(t_0). \end{array} \quad (3.10)$$

Of course, if the feedback gain matrix impose that all poles of the system are in the left-half plane, the closed-loop is stable and the natural response will eventually converge to zero (Theorem 2.10). Therefore, all stable feedback controllers are able to impose regulation to a system, and the characteristics of the transient response can be fully determined by the matrix  $\mathbf{K}$ . These type of controllers are used to make systems goes from nonzero initial states to the zero-state  $\mathbf{x}(t) = \mathbf{0}$  and stays there, meaning that the the controller can also be used to reject disturbances. Now, one may wonder if this operation is too restrictive in the sense that the zero-state is not the desirable state in many control objectives. However, it is common to have linear systems that are linearized versions of nonlinear models, using the approximation from Theorem 2.2. In these cases, the regulator actually has to impose  $\Delta \mathbf{x}(t) = \mathbf{x}(t) - \mathbf{x}_o = \mathbf{0}$ , i.e., drives the system to the steady-state point  $\mathbf{x}_o$  and reject disturbances.

**Example 3.1.** For the sake of illustration, considers the first system of (2.60), repeated here:

$$\mathbf{A} = \begin{bmatrix} -5.93 & 0 \\ 0.83 & -4.70 \end{bmatrix} ; \quad b = \begin{bmatrix} 3.81 \\ -1.09 \end{bmatrix}. \quad (3.11)$$

It is easy to verify that this system is controllable, since its controllability matrix is full-rank. In this case, consider the gains  $K_1 = [-0.993, 1.221]$ ,  $K_2 = [0.782, 0.549]$  and  $K_3 = [0.81, 5.20]$  to yield regulator controllers. The correspondent closed-loop systems  $\mathbf{A}_{cl}^{(i)}$  for each  $i$ -th controller are given as:

$$\mathbf{A}_{cl}^{(1)} = \begin{bmatrix} -2.144 & -4.651 \\ -0.257 & -3.355 \end{bmatrix}; \quad \mathbf{A}_{cl}^{(2)} = \begin{bmatrix} -8.906 & -2.092 \\ 1.692 & -4.093 \end{bmatrix}; \quad \mathbf{A}_{cl}^{(3)} = \begin{bmatrix} -9.014 & -19.812 \\ 1.723 & 1.014 \end{bmatrix}. \quad (3.12)$$

Since the open-loop system was obtained by the linearized model in (2.39) for  $\mathbf{x}_o = [6.19, 1.09]^T$  and  $u_o = 3.03$ , the regulator will impose the zero-state  $\Delta \mathbf{x} = \mathbf{0}$  which actually means imposing  $\mathbf{x}(t) = [6.19, 1.09]^T$  as  $t \rightarrow \infty$ .

The responses of the three closed-loop systems calculated in Example 3.1 are shown in Fig. 3.2. The simulations assumed initial state  $\Delta \mathbf{x}(t) = \mathbf{0} - \mathbf{x}_o$ , and the response and input signal were vertically corrected by  $\mathbf{x}_o$  and  $u_o$ , respectively. Notice that, as expected, the feedback gains  $K_i$  can effectively change the system behavior, even causing pseudo-periodic responses.

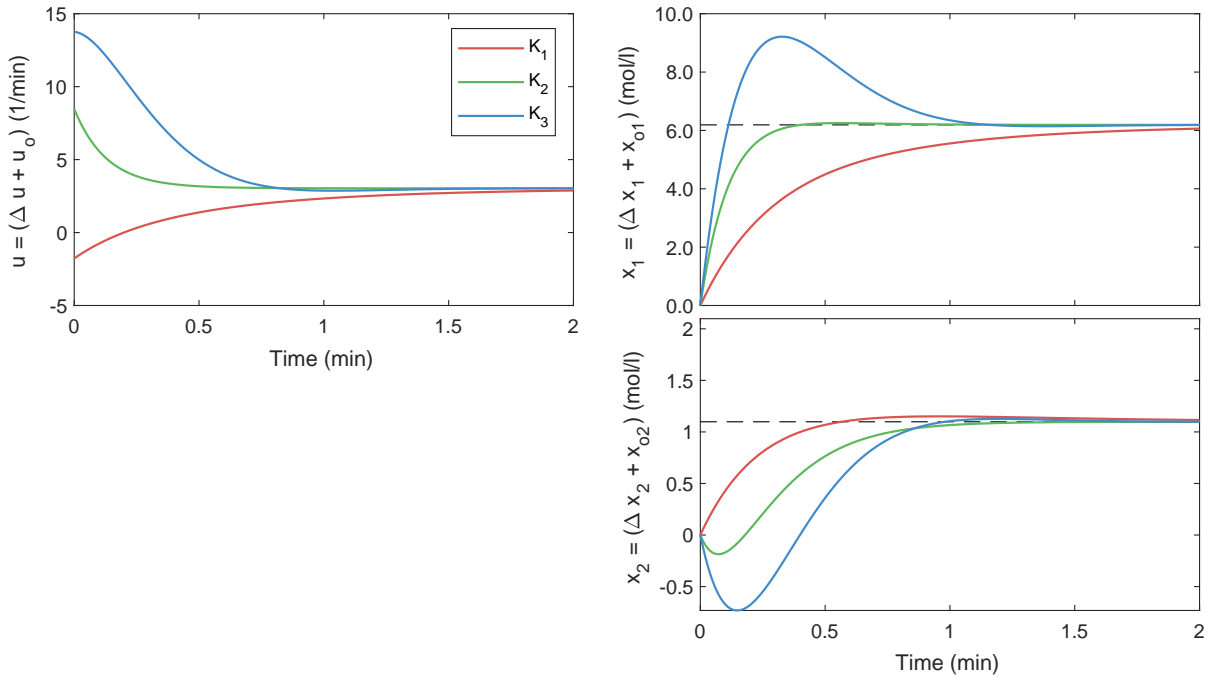


Figure 3.2: Simulation of three closed-loop systems from Example 3.1 showing the input signal (left) and state responses (right).

### Reference Tracking Problem

In contrast, the control objective could be to follow a non-constant signal  $\mathbf{r}(t)$ , or to follow a constant signal different from the zero-state, giving rise to the tracking or servomechanism controllers. A more complete discussion is needed in such cases, since there is a possibility that

state-feedback is not capable of actually perform the tracking. To understand this question, consider, for simplicity, that the system must follow a constant reference  $\mathbf{r}(t) = a$ . Consider, now, an Input-Output conversion of a closed-loop SISO State-Space model, which from (3.9) directly results in the transfer function:

$$G(s) = \frac{Y(s)}{R(s)} = \frac{\beta_1 s^{n-1} + \beta_2 s^{n-2} + \cdots + \beta_{n-1} s + \beta_n}{s^n + \tilde{\alpha}_1 s^{n-1} + \tilde{\alpha}_2 s^{n-2} + \cdots + \tilde{\alpha}_{n-1} s + \tilde{\alpha}_n}. \quad (3.13)$$

From that formulation it is clear that the response  $Y(s) = G(s)R(s)$  will yield a perfect tracking if  $G(s) = 1$ . Moreover, if the system has to track asymptotically track this reference, this operation can be evaluated as time  $t \rightarrow \infty$  or, equivalently, as the frequency  $s \rightarrow 0$ . Plugging this limit in the transfer function implies that a perfect tracking is always possible if  $G(0) = \beta_n/\tilde{\alpha}_n = 1$ , which is not guaranteed a priori. A possible solution is to transform the reference with as  $\tilde{r}(t) = Fr(t)$ , so that  $Y(s) = G(s)\tilde{R}(s) = G(s)FR(s)$ , resulting that:

$$G(0)F = 1 \Rightarrow F = \frac{\tilde{\alpha}_n}{\beta_n}, \quad (3.14)$$

which allows for perfect asymptotically tracking in all cases but when  $\beta_n = 0$ . This same reasoning can easily be extended to MIMO systems (the gain  $F$  turns into a matrix). In the case of non-constant references, the same intuition could still be used, but the design of  $F$  becomes more complex since it would be clearly time-varying. This, however, allows for the definition of tracking controllers.

**Definition 3.3.** (Tracking Controllers) If a state-feedback controller has to make a system track any step reference  $\mathbf{r}(t) \neq \mathbf{0}$ , as  $t \rightarrow \infty$ , it is said to be a *tracking controller*. In this case, one has to apply the *feedforward gain*  $\mathbf{F}$  to adapt the reference as  $\tilde{\mathbf{r}}(t) = \mathbf{F}\mathbf{r}(t)$ , resulting in the following closed-loop state equation and equivalent solution:

<p><b>State Equation:</b></p> $\dot{\mathbf{x}}(t) = (\mathbf{A} - \mathbf{BK})\mathbf{x}(t) + \mathbf{BF}\mathbf{r}(t)$	<p><b>Lagrange solution:</b></p> $\mathbf{x}(t) = e^{(\mathbf{A}-\mathbf{BK})t}\mathbf{x}(t_0) + \int_{t_0}^t e^{(\mathbf{A}-\mathbf{BK})(t-\tau)}\mathbf{BF}\mathbf{r}(\tau)d\tau$
--	---

(3.15)

Despite being a feasible solution, there are still problems with this definition of tracking controllers. For instance, if the system is subject to a *constant additive disturbance*, which has not been included in the model, the resulting operation will not yield a perfect tracking.

The problem of the previous formulation for a tracking controller is that it is not robust to actions that happens outside the model. A direct cause of this is the fact that the feedforward gain  $\mathbf{F}$  does not benefit from the real-time corrective action of the state-feedback, but rather is calculated *off-line* using the model properties. Therefore, a way to ensure a more robust operation could be to insert real-time information about the tracking error directly to the feedback corrective action. With this motivation, a new formulation of the tracking controller is given below.

**Definition 3.4.** (Robust Tracking Controllers) Consider a State-space system and augmented state  $\mathbf{x}_a : \mathbb{R} \rightarrow \mathbb{R}^p$  defined by:

$$\mathbf{x}_a(t) = \int_0^t \mathbf{r}(\tau) - \mathbf{y}(\tau)d\tau \implies \dot{\mathbf{x}}_a(t) = \mathbf{r}(t) - \mathbf{y}(t). \quad (3.16)$$

A robust tracking (or servo) controller, defined by the gain  $\tilde{\mathbf{K}} = [\mathbf{K} \ \mathbf{K}_a]$ , for  $\mathbf{K} \in \mathbb{R}^{r \times n}$  and  $\mathbf{K}_a \in \mathbb{R}^{r \times p}$ , operates on the following augmented version of the original system:

$$\begin{cases} \begin{bmatrix} \dot{\mathbf{x}}(t) \\ \dot{\mathbf{x}}_a(t) \end{bmatrix} = \begin{bmatrix} \mathbf{A} - \mathbf{BK} & -\mathbf{BK}_a \\ -\mathbf{C} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}(t) \\ \mathbf{x}_a(t) \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_p \end{bmatrix} \mathbf{r}(t) \\ \mathbf{y}(t) = \begin{bmatrix} \mathbf{C} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}(t) \\ \mathbf{x}_a(t) \end{bmatrix} \end{cases} \quad (3.17)$$

or, equivalently:

$$\begin{cases} \dot{\tilde{\mathbf{x}}}(t) = \tilde{\mathbf{A}}_{cl} \tilde{\mathbf{x}}(t) + \tilde{\mathbf{B}} \mathbf{r}(t) \\ \mathbf{y}(t) = \tilde{\mathbf{C}} \tilde{\mathbf{x}}(t) \end{cases} \quad (3.18)$$

Since the augmented state  $\mathbf{x}_a(t)$  represents an integral of the tracking error until a time  $t$ , this formulation is usually characterized as imposing “integral action” to the controller. The schematic in Fig. 3.3 illustrates how an integrator can be included to the block diagram of the control loop.

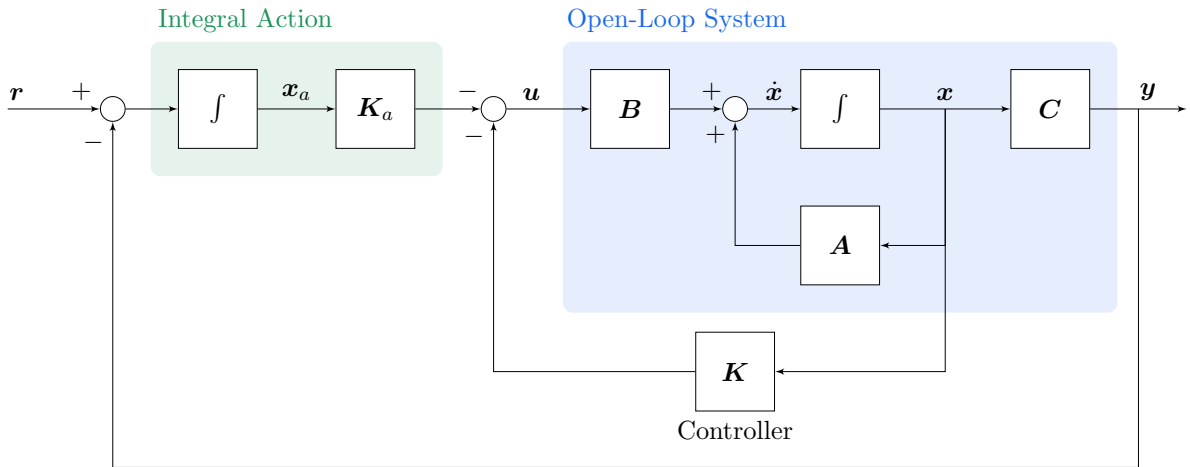


Figure 3.3: Block diagram of a state-feedback closed-loop system with integral action.

Notice that state-equation for closed-loop of this new augmented system is resulting from

$$\tilde{\mathbf{A}}_{cl} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ -\mathbf{C} & \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{BK} & \mathbf{BK}_a \\ \mathbf{0} & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ -\mathbf{C} & \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{B} \\ \mathbf{0} \end{bmatrix} [\mathbf{K} \ \mathbf{K}_a] = \tilde{\mathbf{A}} - \tilde{\mathbf{B}} \tilde{\mathbf{K}}. \quad (3.19)$$

In this case, it must be discussed whether the new gain  $\tilde{\mathbf{K}}$  still preserves the eigenvalue assignment property of a state-feedback gain defined on the original State-Space model.

**Theorem 3.3.** *If the SISO system described by  $(\mathbf{A}, \mathbf{B})$  is controllable and its transfer function  $G(s)$  has no zero at  $s = 0$ , then the eigenvalues of the augmented matrix  $\tilde{\mathbf{A}}$  can be assigned arbitrary by the feedback gain  $\tilde{\mathbf{K}}$ .*

*Proof.* Consider a SISO controllable system. After the augmentation, the new controllability



matrix  $\tilde{\mathcal{C}} \in \mathbb{R}^{(n+p) \times nr}$  is calculated as:

$$\begin{aligned} \tilde{\mathcal{C}} &= \begin{bmatrix} B & AB & A^2B & A^3B & \cdots & A^{n-1}B \\ 0 & -CB & -CA^2B & -CA^2B & \cdots & -CA^{n-2}B \end{bmatrix} \\ &= \begin{bmatrix} 1 & -\alpha_1 & -\alpha_1^2 - \alpha_2 & -\alpha_1(\alpha_1^2 - \alpha_2) + \alpha_2\alpha_1 - \alpha_3 & \cdots & \Delta_1(\alpha_1, \dots, \alpha_n) \\ 0 & 1 & -\alpha_1 & -\alpha_1^2 - \alpha_2 & \cdots & \Delta_2(\alpha_1, \dots, \alpha_n) \\ 0 & 0 & 1 & -\alpha_1 & \cdots & \Delta_1(\alpha_1, \dots, \alpha_n) \\ 0 & 0 & 0 & 1 & \cdots & \Delta_3(\alpha_1, \dots, \alpha_n) \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \Delta_n(\alpha_1, \dots, \alpha_n) \\ 0 & -\beta_1 & -\beta_1\alpha_1 - \beta_2 & -\beta_1(\alpha_1^2 - \alpha_2) + \beta_2\alpha_1 - \beta_3 & \cdots & \Delta_{n+1}(\alpha_1, \dots, \alpha_n) \end{bmatrix}, \end{aligned} \quad (3.20)$$

where  $\Delta_i(\alpha_1, \dots, \alpha_n)$  is a polynomial used to save space in the equations. By inspection of this matrix, it is possible to discover a pattern between the rows. Since elementary operations between the rows  $r_1, r_2, \dots, r_n$  do not change the matrix row rank, the last row of the matrix can be transformed as  $r_n = r_n + r_{n-1}\beta_{n-2} + r_{n-2}\beta_{n-3} + \cdots + r_2\beta_1$ . The result is the triangular matrix in the form:

$$\tilde{\mathcal{C}} = \begin{bmatrix} 1 & -\alpha_1 & -\alpha_1^2 - \alpha_2 & -\alpha_1(\alpha_1^2 - \alpha_2) + \alpha_2\alpha_1 - \alpha_3 & \cdots & \Delta_1(\alpha_1, \dots, \alpha_n) \\ 0 & 1 & -\alpha_1 & -\alpha_1^2 - \alpha_2 & \cdots & \Delta_2(\alpha_1, \dots, \alpha_n) \\ 0 & 0 & 1 & -\alpha_1 & \cdots & \Delta_3(\alpha_1, \dots, \alpha_n) \\ 0 & 0 & 0 & 1 & \cdots & \Delta_4(\alpha_1, \dots, \alpha_n) \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \Delta_n(\alpha_1, \dots, \alpha_n) \\ 0 & 0 & 0 & 0 & \cdots & \beta_n \end{bmatrix}. \quad (3.21)$$

Since  $G(s)$  has no zeros at  $s = 0$ , then  $\beta_n \neq 0$ , meaning that  $\tilde{\mathcal{C}}$  is nonsingular and, therefore has full row rank. From this it can be concluded that the augmented system  $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$  is controllable and, from Theorem 3.2, the eigenvalues of  $\tilde{\mathbf{A}}$  can be assigned anywhere in the complex plane.  $\square$

Although the result was only shown for the SISO case, the main intuition still holds in the MIMO case, but the equations become cumbersome. Basically, if the augmented matrices are controllable, the robust tracking can be achieved by the same formulation.

A possible intuition on how this system performs robust tracking and disturbance rejection can be taken from the fact that the first row of (3.19) has the same form of a basic regulator, except from the term  $\mathbf{B}\mathbf{K}_a$  which is a linear function of  $\mathbf{x}_a(t)$ . Because of this, a control action will always be requested whenever  $\mathbf{x}_a(t) \neq \mathbf{0}$ , i.e., when there is an error between the reference and the output signal. When  $\mathbf{x}_a(t) = \mathbf{0}$ , the equation reduces to a simple regulator. A more quantitative analysis on tracking State-Space controllers can be found in [Franklin et al., 2018].

**Example 3.2.** For the sake of illustration, consider the same system used in (3.11). Unfortunately, the augmented version of this *single-input multiple-output* (SIMO) formulation is not controllable. However, the SISO version obtained by letting  $C = [0, 1]$  obeys Theorem 3.3 and has no zeros at  $s = 0$ . Thus, a robust tracking controller can be designed by state-feedback. Consider the gains  $\tilde{K}_1 = [-0.99, 1.22, -25]$ ,  $\tilde{K}_2 = [0.78, 0.55, -25]$  and  $\tilde{K}_3 = [0.81, 5.20, -50]$ .

The correspondent closed-loop systems  $\tilde{\mathbf{A}}_{cl}^{(i)}$  for each  $i$ -th controller are given as:

$$\begin{aligned}\tilde{\mathbf{A}}_{cl}^{(1)} &= \begin{bmatrix} -9.707 & 4.65 & -95.25 \\ 1.92 & -6.04 & 27.45 \\ 0 & -1.00 & 0 \end{bmatrix}; \quad \tilde{\mathbf{B}} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \\ \tilde{\mathbf{A}}_{cl}^{(2)} &= \begin{bmatrix} -2.94 & 2.09 & -95.25 \\ -0.02 & -5.30 & 27.46 \\ 0 & -1.00 & 0 \end{bmatrix}; \quad \tilde{\mathbf{B}} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}. \\ \tilde{\mathbf{A}}_{cl}^{(3)} &= \begin{bmatrix} -2.84 & 19.81 & -190.50 \\ -0.06 & -10.41 & 54.91 \\ 0 & -1.00 & 0 \end{bmatrix}; \quad \tilde{\mathbf{B}} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.\end{aligned}\tag{3.22}$$

Some simulations of closed-loop systems for this system to track are shown in Fig. 3.4 for a non-constant reference consisting of a sequence of step signals.

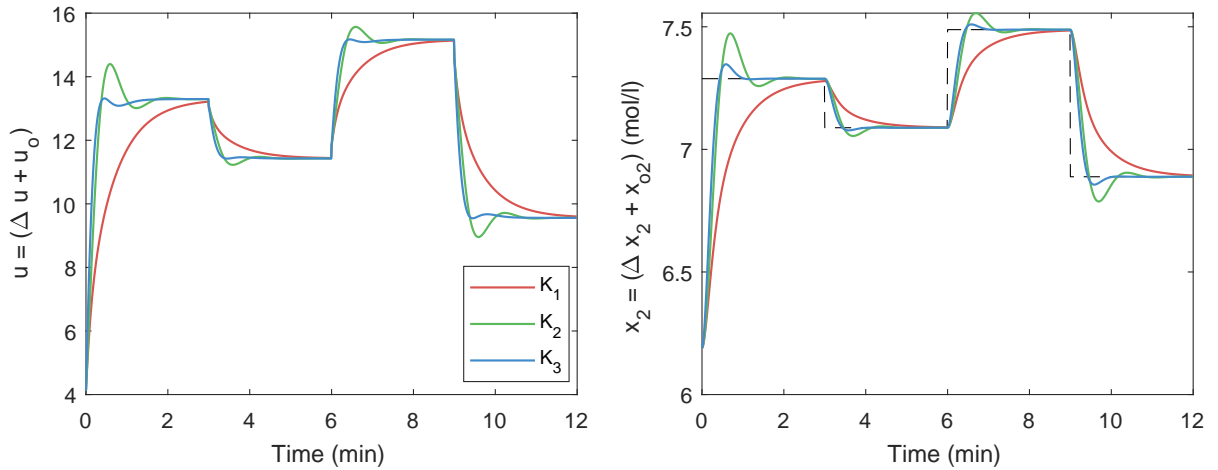


Figure 3.4: Simulation of three closed-loop systems showing the input signal (left) and state response (right) for reference tracking. The reference signal is indicated by the black dashed line.

### 3.3 Deterministic State Observers

Until now, the state feedback was discussed from the perspective that the device has direct access to the real value of all the states of the system. This assumption is unrealistic, since the states are only observed through the output signal  $\mathbf{y}(t)$ , through the matrix  $\mathbf{C}$ , which itself is not assumed to be always equal to the identity matrix. In practice, this means that some states could not be measured, due to technical difficulties or economic reasons, or that the instrumentation available is not perfect, and the observations are prone to deviate from the real value. Since the state-vector is necessary for the state-feedback to compute the input to the system, this section discusses how to develop devices that can reconstruct the states from the observations.

A device that determines a state-vector  $\mathbf{x}(t)$  from the output signal  $\mathbf{y}(t)$  is known as *state observer*, or *state estimator* in some cases. Amongst several possible techniques, a practical and popular one is the *Luenberger observer*, which is defined below.

**Definition 3.5.** (Luenberger Observer) Given a system in State-Space with output signal  $\mathbf{y}(t) : \mathbb{R} \rightarrow \mathbb{R}^p$  and an observer gain  $\mathbf{L} \in \mathbb{R}^{n \times p}$ , the estimated state-vector  $\hat{\mathbf{x}}(t)$  is represented by the observer system:

$$\dot{\hat{\mathbf{x}}}(t) = \mathbf{A}\hat{\mathbf{x}}(t) + \mathbf{B}\mathbf{u}(t) + \mathbf{L}(\mathbf{y}(t) - \mathbf{C}\hat{\mathbf{x}}(t)), \quad (3.23)$$

or, equivalently:

$$\dot{\hat{\mathbf{x}}}(t) = (\mathbf{A} - \mathbf{LC})\hat{\mathbf{x}}(t) + \mathbf{B}\mathbf{u}(t) + \mathbf{L}\mathbf{y}(t). \quad (3.24)$$

The observer system works as a parallel system that is simulated alongside the actual system, as illustrate in Fig. 3.5. The expected result is that the observer yields  $\hat{\mathbf{x}}(t) = \mathbf{x}(t)$ , as time  $t \rightarrow \infty$ . Alternatively, it is possible to create a variable  $\mathbf{e}(t) = \mathbf{x}(t) - \hat{\mathbf{x}}(t)$  such that, using (3.24):

$$\begin{aligned} \dot{\mathbf{e}} &= \dot{\mathbf{x}} - \dot{\hat{\mathbf{x}}} \\ &= (\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}) - ((\mathbf{A} - \mathbf{LC})\hat{\mathbf{x}} + \mathbf{B}\mathbf{u} + \mathbf{LC}\mathbf{x}) \\ &= (\mathbf{A} - \mathbf{LC})\mathbf{x} - (\mathbf{A} - \mathbf{LC})\hat{\mathbf{x}} \\ &= (\mathbf{A} - \mathbf{LC})(\mathbf{x} - \hat{\mathbf{x}}) \\ &= (\mathbf{A} - \mathbf{LC})\mathbf{e} \end{aligned} \quad (3.25)$$

which implies that the observer asymptotically tracks the actual state-vector if  $\mathbf{e}(t) = \mathbf{0}$  as  $t \rightarrow \infty$ . Analyzing the equation above, it is intuitive to notice that this result can be guaranteed if all the eigenvalues of matrix  $\mathbf{A}_{obs} = \mathbf{A} - \mathbf{LC}$  have negative real parts, since it implies that the natural response of  $\dot{\mathbf{e}}$  vanishes and there is no forcing input action to yield forcing response to this variable (Theorem 2.10). Therefore, a specific gain  $\mathbf{L}$  could be designed to ensure this condition and allows the observer to asymptotically track the state-vector with a desired behavior. The following theorem relates this statement with the choice of a gain  $\mathbf{L}$ .

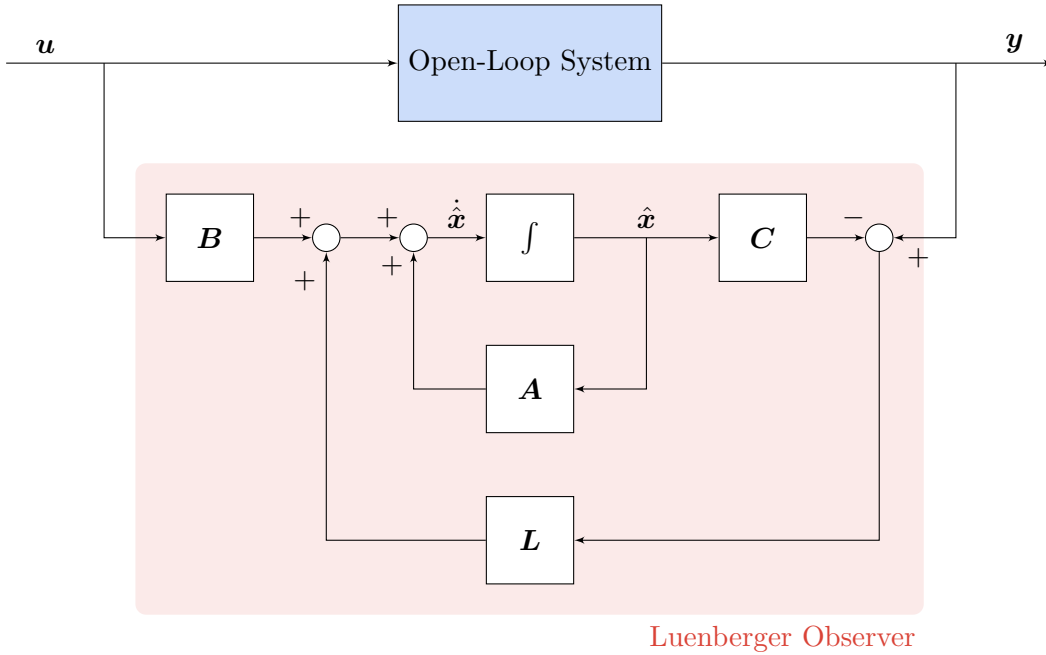


Figure 3.5: Block diagram of a State-Space system with Luenberger observer.

**Theorem 3.4.** *If a system in State-Space representation is observable, then by a Luenberger observer with gain matrix  $\mathbf{L} \in \mathbb{R}^{n \times p}$  the eigenvalues of  $\mathbf{A}_{obs} = \mathbf{A} - \mathbf{LC}$  can arbitrarily be assigned anywhere in the complex plane, as long as that complex conjugate eigenvalues are assigned in pairs.*

*Proof.* Consider that a State-Space with matrices  $(\mathbf{A}, \mathbf{C})$  is observable. From the Duality Theorem [Chen, 1998], if the pair  $(\mathbf{A}, \mathbf{C})$  is observable then the pair  $(\mathbf{A}^T, \mathbf{C}^T)$  is controllable. In this case, it is possible to design a gain matrix  $\mathbf{K}$  to assign the eigenvalues of  $\tilde{\mathbf{A}}_{obs} = \mathbf{A}^T - \mathbf{C}^T \mathbf{K}$  in any desirable points in the complex space. Since the eigenvalues of a matrix are invariant to the transpose operation, the design of  $\mathbf{K}$  can also place the eigenvalues of the matrix  $(\tilde{\mathbf{A}}_{obs})^T = \mathbf{A} - \mathbf{K}^T \mathbf{C}$ . Therefore, making  $\mathbf{L} = \mathbf{K}^T$  proves the theorem.  $\square$

The procedure stated in this proof highlights the similarities between existing closed-loop observers, such as the Luenberger observer, and closed-loop controllers, such as the state-feedback. Basically, the same design considerations that concerns state-feedback are important in the design of the observer gain  $\mathbf{L}$ . For instance, the eigenvalues of  $\tilde{\mathbf{A}}_{obs} = \mathbf{A}^T - \mathbf{C}^T \mathbf{K}$  can be assigned such that the time evolution of the error  $\mathbf{e}(t)$  has a desirable time constant, damping coefficient or natural frequency. In conclusion, the state-vector  $\mathbf{x}(t)$  can be reconstructed by using an observer gain such that each eigenvalue of  $\mathbf{A}_{obs}$  is on the left-half side of the complex plane.

The only reason to develop an observer is to allow for state-feedback controllers to access the values of the state-vector, thus being able to calculate an appropriate action to follow the reference signal. If a controller can only access the estimated state-vector  $\hat{\mathbf{x}}(t)$ , it is possible to define a State-Space formulation for a closed-loop based on feedback from estimated states given data from a, possibly non-linear and time-varying, disturbed system.

**Definition 3.6.** (Feedback from Estimated States) Given a system in State-Space representation whose state-vector is reconstructed from a Luenberger observer of gain  $\mathbf{L} \in \mathbb{R}^{n \times p}$  and input signal is calculated through state-feedback with gain  $\mathbf{K} \in \mathbb{R}^{n \times r}$ , its time evolution can be represented through the model:

$$\begin{cases} \dot{\hat{\mathbf{x}}}(t) = \mathbf{A}\hat{\mathbf{x}}(t) - \mathbf{BK}\hat{\mathbf{x}}(t) + \mathbf{B}\mathbf{r}(t) \\ \dot{\hat{\mathbf{x}}}(t) = (\mathbf{A} - \mathbf{LC} - \mathbf{BK})\hat{\mathbf{x}}(t) + \mathbf{B}\mathbf{r}(t) + \mathbf{LC}\mathbf{x}(t) \end{cases} \quad (3.26)$$

A schematic for this formulation of feedback from estimated states is shown Fig. 3.6, including the possibility for integral action. In fact, this schematic summarizes several different control architectures that includes both feedback action and state estimation, and it will be a reference whenever this work mentions physical control loops and instrumentation.

Notice, now, that the formulation just defined imposes that the dynamics of the estimated state  $\hat{\mathbf{x}}(t)$  is dependent both in  $\mathbf{K}$  and  $\mathbf{L}$ , which are arbitrary matrices chosen by the control designer. This leads to the possible conclusion that the choice of  $\mathbf{K}$  is now restricted by the effect that it will produce in the choice of  $\mathbf{L}$ , which is not true. The following theorem, known as the Separation Principle, states that design of these two gains are independent.

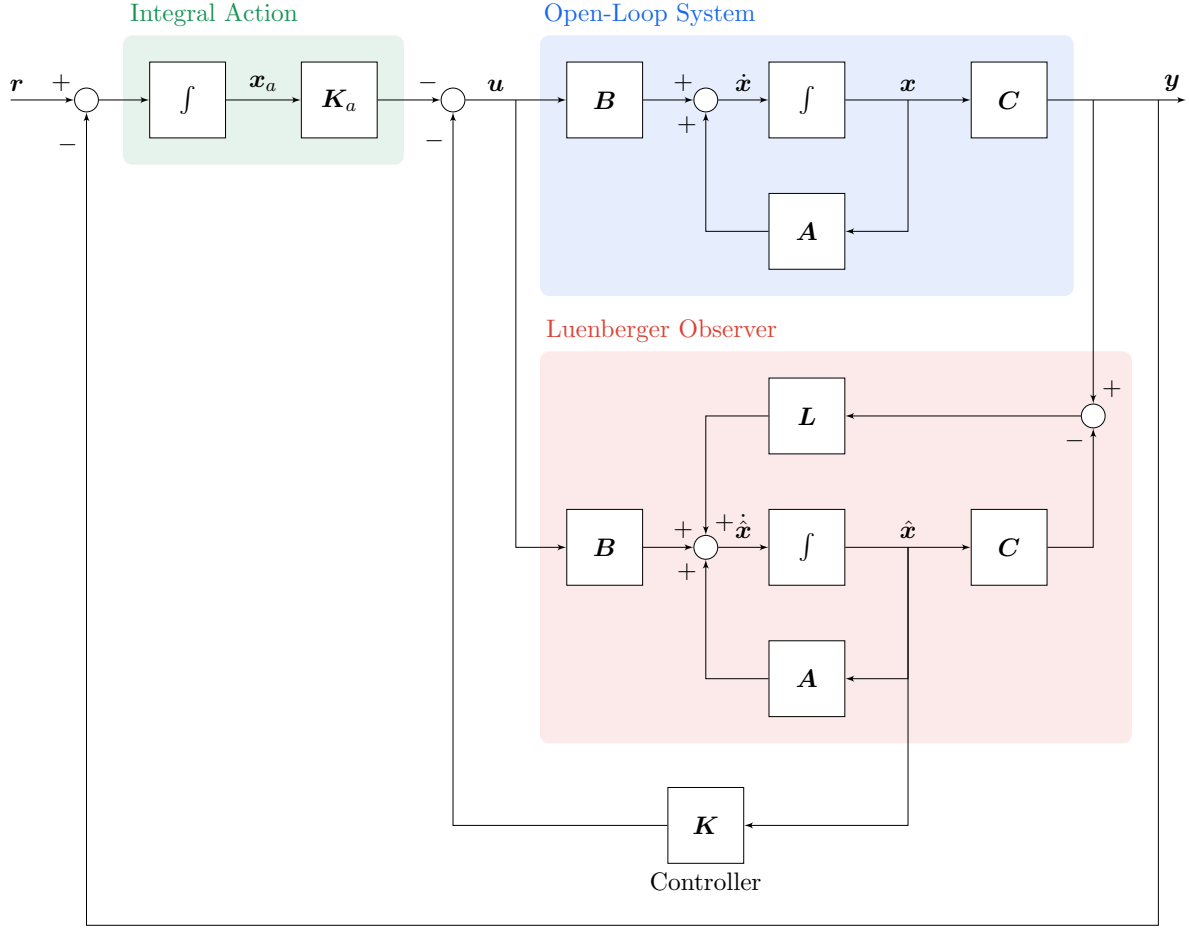


Figure 3.6: Block diagram of a state-feedback closed-loop system with integral action and Luenberger observer.

**Theorem 3.5.** (*Separation Principle*) Given a system in State-Space with a Luenberger observer of gain  $L$  and state-feedback controller of gain  $K$ , the closed-loop eigenvalues contributions of  $(A - BK)$  are independent from those of  $(A - LC)$ .

*Proof.* Consider a feedback from estimated states as defined in (3.26). The controller-estimator system can be rewritten as a single state equation:

$$\begin{bmatrix} \dot{x} \\ \dot{\hat{x}} \end{bmatrix} = \underbrace{\begin{bmatrix} A & -BK \\ LC & A - LC - BK \end{bmatrix}}_{\tilde{A}} \begin{bmatrix} x \\ \hat{x} \end{bmatrix} + \underbrace{\begin{bmatrix} B \\ B \end{bmatrix}}_{\tilde{B}} r. \quad (3.27)$$

Consider, now, the following similarity transformation  $z(t) = Px(t)$ :

$$\underbrace{\begin{bmatrix} I & 0 \\ I & -I \end{bmatrix}}_P \begin{bmatrix} x \\ \hat{x} \end{bmatrix} = \begin{bmatrix} x \\ x - \hat{x} \end{bmatrix} = \begin{bmatrix} x \\ e \end{bmatrix}. \quad (3.28)$$

Since  $P = P^{-1}$ , and this is a valid similarity transformation that does not alter the system eigenvalues, the equivalent system for state  $z(t)$  is obtained as:

$$\begin{bmatrix} \dot{x} \\ \dot{e} \end{bmatrix} = \begin{bmatrix} A - BK & -BK \\ 0 & A - LC \end{bmatrix} \begin{bmatrix} x \\ e \end{bmatrix} + \begin{bmatrix} B \\ 0 \end{bmatrix} r. \quad (3.29)$$

Since the system matrix obtained is block triangular, it is possible to conclude that the system in such configuration has eigenvalues that are separate contributions from the eigenvalues of  $(\mathbf{A} - \mathbf{BK})$  and from  $(\mathbf{A} - \mathbf{LC})$ .  $\square$

The Separation Principle is a nice results that further motivates the topology of Fig. 3.6, since it explicitly states that the design of the controller and the state observer can be done separately. Thus, any structure that obeys the state feedback formulation can be used as a controller and the same is valid for the observer device. In the next chapter, this result will be explored to motivate the use of more advanced control and state estimation techniques without having to redefine the analytical tools and intuitions built for traditional state-feedback controllers from pole-placement methods.

### 3.4 Properties of State-Feedback Controllers

The last section introduces the first considerations into applying state-feedback in real-world systems, given limitations on the instruments and uncertainty on the environment. Basically, a mathematical analysis of such closed-loop systems allows for a full characterization of its behavior, but the real system will exhibit a different response due to these limitations and, essentially, due to unmodelled external disturbances. Therefore, it is desirable to anticipate these variations and discuss the properties of the closed-loop system in the sense of robustness and stability margins. Since this goal requires that the system response is evaluated in as general as possible context, the discussion in this section relies on frequency response analysis, which consider a broader class of input signals: sinusoids of any frequency. Consider the representation of closed-loop systems shown in Fig. 3.7, which consists of a simplified version of Fig. 3.1 in the Laplace frequency domain considering only the state response.

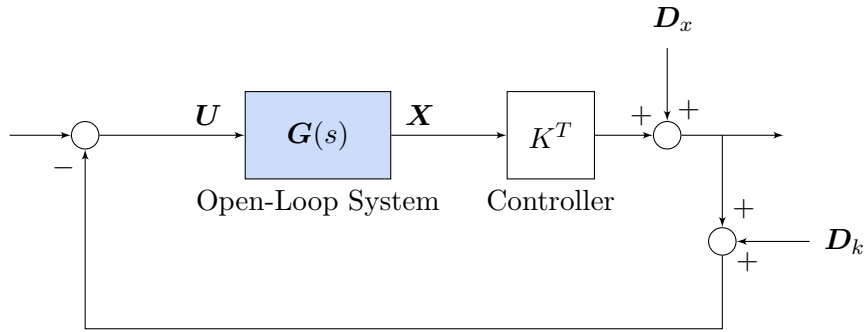


Figure 3.7: Simplified block diagram of a perturbed state-feedback closed-loop.

Using block diagram algebra and Theorem 2.1, it is possible to associate a forward-path transfer function  $\mathbf{G}_f(s)$  given as [Anderson and Moore, 1990]:

$$\mathbf{G}_f(s) = -\mathbf{K}^T \mathbf{G}(s) = -\mathbf{K}^T (s\mathbf{I} - \mathbf{A})^{-1} \mathbf{B}. \quad (3.30)$$

Now, it is possible to characterize the properties of this quantity which relates the state-feedback gain  $\mathbf{K}$  with the system disturbance  $\mathbf{D}_x(s)$  and gain disturbance  $\mathbf{D}_k(s)$ . To facilitate the definitions and discussion, the results are shown for single-state single-input systems, so  $\mathbf{G}_f(s)$  is a scalar function, where the extension for  $n > 1$  states is intuitive in most cases. The first property to be discussed, then, considers the stability of closed-loop feedback controllers. Consider the following closed-loop stability criterion from a Bode plot visualization.

**Theorem 3.6.** (*Bode stability criterion*) Consider a feedback system whose closed-loop transfer function is defined, assuming perfect measuring sensors, as:

$$T(s) = \frac{KG(s)}{1 + KG(s)}. \quad (3.31)$$

The closed-loop system is said to be stable if  $|G(j\omega_{pc})| \leq 0$ , where  $\omega_{pc}$  is the phase crossover frequency obtained such that  $\angle G(j\omega_{pc}) = -180^\circ$ .

Additionally, it is possible to define a stability criterion through a Nyquist diagram of the closed-loop system.

**Theorem 3.7.** (*Nyquist criterion*) Consider a system with feedforward transfer function as defined in (3.30). Now, let  $P$  and  $Z$  be respectively the number of poles of  $G_f(s)$  and zeros of  $1 + G_f(s)$  that are in the right-half plane. In this case, the Nyquist contour shall clockwise encircle the point  $s = -1$  a number of times  $N$  such that  $N = Z - P$ .

A detailed proof of both criterion can be found in [Nise, 2015]. The introduction of these stability evaluation techniques may seem redundant, given that the closed-loop BIBO stability can still be characterized from Theorem 2.10. However, their graphical nature allows for an easy understanding of how disturbances can affect the stability of state-feedback systems. For instance, a system subject to a sinusoidal disturbance of constant magnitude, but with the same frequency as the natural frequency of the system, will show a response with higher magnitude for the phase crossover frequency than the one visualized in the Bode plot. This phenomenon is widely known in Physics as “resonance”. Therefore, the influence of disturbances can inflict instability to a stable system.

Of course, not all disturbances observed in real operations are strong enough to bring any reasonable stable controller to an unstable condition. However, depending on the choice of the gain  $\mathbf{K}$ , some closed-loop systems can be more prone to these undesired problems than others. This motivates the discussion on stability margins.

**Definition 3.7.** (Stability Margins) Given a closed-loop system with transfer function  $T(s)$ , the *Gain Margin (GM)* is defined as a factor of how much a gain can be increased before the system becomes unstable, and is equated as:

$$GM = \frac{1}{|T(j\omega_{pc})|}, \quad (3.32)$$

where  $\omega_{pc}$  is the phase crossover frequency such that  $\angle T(j\omega_{pc}) = -180^\circ$  (or the point where a Nyquist diagram crosses the real axis for  $-1 < s < 0$ ). In addition, the *Phase Margin (PM)* is defined as how much phase lag can be added to  $T(s)$  the system becomes unstable, and is equated as:

$$PM = \angle T(j\omega_{gc}), \quad (3.33)$$

where  $\omega_{gc}$  is the *gain crossover frequency* such that  $|T(j\omega_{gc})| = 0dB$  (or the angle when a Nyquist diagram crosses the unit circle centered at  $s = 0$ ).

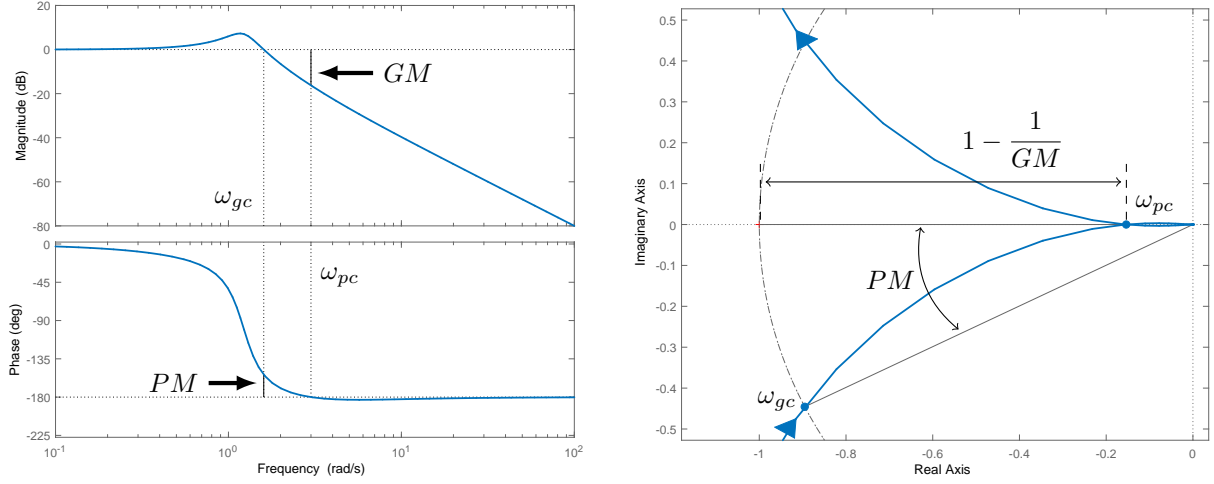


Figure 3.8: Stability margins visualizations given Bode plots (left) and Nyquist diagrams (right) of closed-loop dynamical systems.

A graphical representation of these margins, for both Bode plots and Nyquist diagrams is shown in Fig. 3.8, for the same system as the one from last figure. If a closed-loop system has small Gain Margin, it is clear that its stability is not robust to gain uncertainties, while a small Phase Margins implies that its stability is not robust to time delay uncertainties in the control actions. From this motivation, the use of high-gain controllers have been a working example of how to achieve robust control [Khalil, 2005, Mahmoud and Khalil, 2002], but is usually challenging for the instrumentation in real-world operations. Therefore, the design of state-feedback gains must account for these quantities, and a trade-off between performance and robustness is always necessary for this formulation.





## Chapter 4

# Optimal Control and Estimation

This chapter introduces the vast field of optimal control and optimal estimation of dynamical systems. The developments are focused in optimizing a cost function that produces an optimal sequence of control actions to achieve the desired control objective, under some conditions. The chapter starts by discussing a general formulation of the optimal control problem and then specializes the formulation to a case of a linear system with quadratic cost functions. After that, the dual optimal state estimation problem is discussed and a result that merges the optimal control with the optimal estimation is presented. Finally, the main stability and robustness properties of the controllers derived in this chapter are accessed.

### 4.1 Formulation

The last chapter introduced the notion of controller synthesis as an engineering procedure to be done “by hand” from a designer with some knowledge about the system and the control environment. This approach, despite being very popular and practical, is rather time-consuming and demands the designer to consider what are the best pole locations given some hard specifications. In addition to that, the design techniques previously presented becomes more involving when applied to MIMO systems. Thus, a more general and automatic procedure to control synthesis is desirable.

With this motivation, the concept of *Optimal Control* [Anderson and Moore, 1990, Kirk, 1998, Liberzon, 2012, Bertsekas, 2017] was introduced as an alternative strategy for controlling dynamical systems. Optimal controllers determines the necessary action by minimizing a cost function (or maximizing a reward function). These formulations produces feedback controllers that autonomously requires optimal input signals given a desired objective and possible restrictions, begin easily generalized to different systems.

**Definition 4.1.** (Optimal Control) Given a system in State-Space formulation, with state signal  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^n$ , and a reference signal  $\mathbf{r} : \mathbb{R} \rightarrow \mathbb{R}^n$ , the input signal  $\mathbf{u}(t) \in \mathbb{R}^r$ , for any time  $t$ , is optimal if an optimal control law  $\pi^* : \mathbb{R}^{n \times n \times 1} \rightarrow \mathbb{R}^r$  can be found as:

$$\mathbf{u}(t) = \pi^*(\mathbf{x}, \mathbf{r}, t) = \min_{\mathbf{u}} J(\mathbf{x}, \mathbf{r}, t), \quad (4.1)$$

where  $J : \mathbb{R}^{n \times n \times 1} \rightarrow \mathbb{R}$  is known as a *cost function* of the states and reference signals.

Note that this optimization can be converted to maximizing a function  $V(\cdot)$  by evaluating  $V(\cdot) = -J(\cdot)$ . Thus, this document will only refers to optimization as minimizing some cost function. The problem of finding an optimal control law, or optimal control policy,  $\pi^*(\cdot)$  depends

on the choice of the cost function  $J(\cdot)$ , and several optimization techniques can be used to determine the value of  $\mathbf{u}(t)$  that achieve its minimum. This problem differs from standard optimization problems because it is dependent of time, there can be a decision variable and it is constrained by the dynamics of the model and by the optimal policy action. Solutions to this optimization can be obtained from several methods, however, this work specifically focus on the concept of Dynamic Programming, first introduced by [Bellman, 1954].

To facilitate the discussion and analysis of optimal controllers, consider a subclass of these controllers (that is still very general) defined below by a specific choice of functional for the cost function.

**Definition 4.2.** (Finite-Horizon Optimal Regulators) Consider a controller setup as in Definition 4.1. A *Finite-Horizon Optimal Regulator* is defined as any controller whose optimal policy over a time interval  $t \in [t_0, T]$  minimizes the cost functional:

$$J(\mathbf{x}, \mathbf{u}, t_0) = \int_{t_0}^T l(\mathbf{x}, \mathbf{u}, \tau) d\tau + l_f(\mathbf{x}, T), \quad (4.2)$$

where  $l(\cdot) : \mathbb{R}^{n \times r \times 1} \rightarrow \mathbb{R}$  and  $l_f(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}$  are, respectively, the *trajectory* and *terminal loss functions*. In the case that  $t_0 = 0$ ,  $T$  is also known as the *control horizon*.

This formulation presents a notion of optimizing for a sequence of control actions  $[\mathbf{u}(t_0), \dots, \mathbf{u}(T)]$  that can cause a state trajectory  $[\mathbf{x}(t_0), \dots, \mathbf{x}(T)]$ , which is optimal given the loss functions. As will be shown later, this is a very feasible strategy for controllable linear systems. In an optimization notation, this problem can also be presented as the following *program*:

$$\begin{aligned} & \text{minimize} && \int_{t_0}^T l(\mathbf{x}, \mathbf{u}, \tau) d\tau + l_f(\mathbf{x}, T) \\ & \text{s.t.} && \dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}, \mathbf{u}, t) \\ & && \mathbf{u}(t) = \mathbf{p}\mathbf{i}(\mathbf{x}(t_0), \dots, \mathbf{x}(t)) \end{aligned} \quad (4.3)$$

Now, the discussion turns to how to solve this general problem. In this work, the solution for the optimal control will follow a dynamic programming formulation. The first necessary effort, then, is to define an important partial differential equation known as the Hamilton-Jacobi equation, given below.

**Theorem 4.1.** (Hamilton-Jacobi equation) Consider a finite-horizon cost function in the form of Definition 4.2, for a system described by the state-equation  $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u}, t)$ , restated as:

$$V(\mathbf{x}, \mathbf{u}, t_0) = \int_{t_0}^T l(\mathbf{x}, \mathbf{u}, \tau) d\tau + l_f(\mathbf{x}(T)) \quad (4.4)$$

Consider also that the loss  $l(\cdot)$  and state function  $\mathbf{f}$  are smooth on their parameters. Then, minimizing any functional in the form of  $V(\cdot)$  is equivalent to determining the solution of the Hamilton-Jacobi equation, which is given by the partial differential equation:

$$\frac{\partial V^*}{\partial t} = - \min_{\mathbf{u}(t)} \left[ l(\mathbf{x}, \mathbf{u}, t) + \left[ \frac{\partial V^*}{\partial \mathbf{x}} \right]^T \mathbf{f}(\mathbf{x}, \mathbf{u}, t) \right] \quad (4.5)$$

and the boundary condition:

$$V^*(\mathbf{x}, T) = l_f(\mathbf{x}(T)). \quad (4.6)$$

*Proof.* First of all, consider the cost functional  $V(\cdot)$  from (4.4). Minimizing this functional with respect to the control inputs  $\mathbf{u}(t_0), \dots, \mathbf{u}(T)$  consists of evaluating the optimal cost:

$$V^*(\mathbf{x}, t_0) = \min_{\mathbf{u}(t_0), \dots, \mathbf{u}(T)} \left[ \int_{t_0}^T l(\mathbf{x}, \mathbf{u}, \tau) d\tau + l_f(\mathbf{x}(T)) \right]. \quad (4.7)$$

Now, consider any  $t \in [t_0, T]$  and  $t_r \in [t, T]$ . Since the original control action  $\mathbf{u}(t), \dots, \mathbf{u}(T)$  can be obtained through the concatenation of  $\mathbf{u}(t), \dots, \mathbf{u}(t_r)$  and  $\mathbf{u}(t_r), \dots, \mathbf{u}(T)$ , the optimal cost in this interval can be represented as:

$$\begin{aligned} V^*(\mathbf{x}, t) &= \min_{\mathbf{u}(t), \dots, \mathbf{u}(T)} \left[ \int_t^T l(\mathbf{x}, \mathbf{u}, \tau) d\tau + l_f(\mathbf{x}(T)) \right] \\ &= \min_{\mathbf{u}(t), \dots, \mathbf{u}(t_r)} \left\{ \min_{\mathbf{u}(t_r), \dots, \mathbf{u}(T)} \left[ \int_t^{t_r} l(\mathbf{x}, \mathbf{u}, \tau) d\tau + \int_{t_r}^T l(\mathbf{x}, \mathbf{u}, \tau) d\tau + l_f(\mathbf{x}(T)) \right] \right\} \\ &= \min_{\mathbf{u}(t), \dots, \mathbf{u}(t_r)} \left\{ \int_t^{t_r} l(\mathbf{x}, \mathbf{u}, \tau) d\tau + \min_{\mathbf{u}(t_r), \dots, \mathbf{u}(T)} \left[ \int_{t_r}^T l(\mathbf{x}, \mathbf{u}, \tau) d\tau + l_f(\mathbf{x}(T)) \right] \right\} \\ &= \min_{\mathbf{u}(t), \dots, \mathbf{u}(t_r)} \left\{ \int_t^{t_r} l(\mathbf{x}, \mathbf{u}, \tau) d\tau + V^*(\mathbf{x}, t_r) \right\}. \end{aligned} \quad (4.8)$$

Notice that this is a recursive formula, since the optimal cost  $V^*(\mathbf{x}, t)$  depends on  $V^*(\mathbf{x}, t_r)$ , which itself is an optimal cost for the future time  $t_r \geq t$ . Without loss of generalization, let  $t_r = t + \delta t$ , where  $\delta t$  is a small number. Since  $l(\cdot)$  is a smooth function, the right-hand side of the recursive form above can be expanded by a Taylor series expansion:

$$\begin{aligned} V^*(\mathbf{x}, t) &= \min_{\mathbf{u}(t), \dots, \mathbf{u}(t+\delta t)} \left\{ l(\mathbf{x}, \mathbf{u}, t + \delta t) \delta t \right. \\ &\quad \left. + V^*(\mathbf{x}, t) + \left[ \frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \right]^T \frac{d\mathbf{x}(t)}{dt} \delta t + \frac{\partial V^*(\mathbf{x}, t)}{\partial t} \delta t + \mathcal{O}(\delta t)^2 \right\}, \end{aligned} \quad (4.9)$$

where  $\mathcal{O}(\delta t)^2$  denotes high order terms. Since the terms  $V^*(\mathbf{x}, t)$  and  $(\partial V^*/\partial t)\delta t$  do not depend on  $\mathbf{u}(t)$ , they can be disregarded of the minimization. Rearranging the terms and substituting  $d\mathbf{x}/dt = \mathbf{f}(\mathbf{x}, \mathbf{u}, t)$  results in:

$$\frac{\partial V^*}{\partial t}(\mathbf{x}, t) = - \min_{\mathbf{u}(t), \dots, \mathbf{u}(t+\delta t)} \left\{ l(\mathbf{x}, \mathbf{u}, t + \delta t) + \left[ \frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \right]^T \mathbf{f}(\mathbf{x}, \mathbf{u}, t) + \mathcal{O}(\delta t)^2 \right\}. \quad (4.10)$$

Finally, letting  $\delta t \rightarrow 0$  yields:

$$\frac{\partial V^*}{\partial t} = - \min_{\mathbf{u}(t)} \left\{ l(\mathbf{x}, \mathbf{u}, t) + \left[ \frac{\partial V^*}{\partial \mathbf{x}} \right]^T \mathbf{f}(\mathbf{x}, \mathbf{u}, t) \right\}. \quad (4.11)$$

To prove the theorem it remains to derive the boundary condition. This result, however, is direct from the form of the cost function, since  $V^*(\mathbf{x}, T) = l_f(\mathbf{x}(T))$  can not be changed through any more control action inside the time horizon.  $\square$

The Hamilton-Jacobi equation implies that finite-horizon optimal controllers can be optimized in a recursive manner, starting from the boundary condition at  $t = T$  to the beginning of the horizon at  $t = t_0$ . This results from the fact that the last action  $\mathbf{u}(T)$  depends only on  $l_f(\cdot)$ , so it can be directly calculated for a desired final state  $\mathbf{x}(T)$  and used to calculate the remaining optimal costs backward in time. This recursive property, directly evidenced in (4.8), is a statement of the Bellman's Principle of Optimality [Bellman, 1954] and, for this reason, the Hamilton-Jacobi equation in the context of optimal control theory is known as the Hamilton-Jacobi-Bellman (HJB) equation. An illustration of this principle is shown in Fig. 4.1

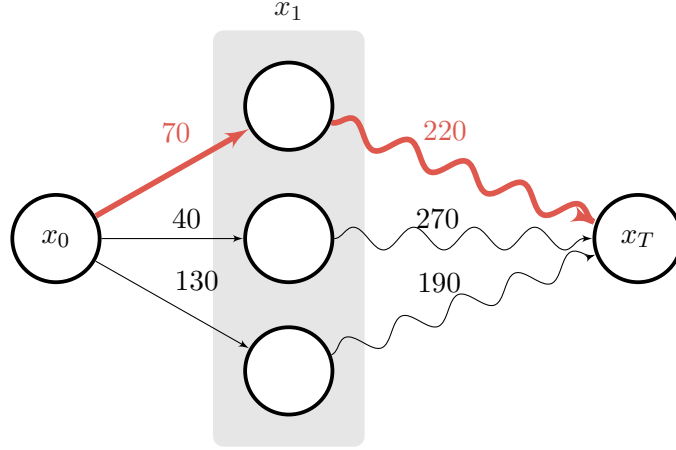


Figure 4.1: Illustration of the Bellman's Principle of Optimality for a system with discrete set of states for a discrete time evolution. Each column represents a time instance and each node represents a possible state. The straight lines represents state transitions given an action with associated costs, whereas the curves represents the trajectory from that state to the terminal state with associated optimal cost. The optimal trajectory between the initial and terminal state is shown in red.

## 4.2 Linear Quadratic (LQ) Controllers

The last section introduced a general condition for solving a finite-horizon optimal control problem. Developing an analytical solution of that condition for any arbitrary loss function  $l(\cdot)$  and state-equation  $\mathbf{f}(\cdot)$  is usually intractable. However, there is a choice of loss function that, under a linear system, allows for a nice closed-form solution to the Hamilton-Jacobi-Bellman equation. This defines the popular class of optimal controllers known as the Linear Quadratic Controllers. Since these controllers are realizations of state-feedback controllers, this section discuss the two possible control operations presented in the last chapter, namely regulation and reference tracking.

### Regulation

**Definition 4.3.** (Linear Quadratic Regulator) Given a linear State-Space system in the form:

$$\begin{cases} \dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) \\ \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t) \end{cases} \quad (4.12)$$

A *Linear Quadratic Regulator* (LQR) for this system is an optimal controller defined by the quadratic cost function:

$$J(\mathbf{x}, \mathbf{u}, t_0) = \int_{t_0}^T (\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}) dt + \mathbf{x}^T(T) \mathbf{Q}_f \mathbf{x}(T), \quad (4.13)$$

where is assumed that  $\mathbf{Q}, \mathbf{Q}_f \succ 0$  and  $\mathbf{R} \succ 0$  are matrices penalizing, respectively, the state-vector magnitude and the control effort.

The LQR optimal controller was a major breakthrough in the development of optimal control theory [Bryson, 1996], and has been a foundation for more complex systems ever since. From

an engineering point of view, this choice of cost function has the advantage that it breaks the controller design to simply choose matrices  $\mathbf{Q}$  and  $\mathbf{R}$  as a trade-off between performance and actuator restrictions. From a mathematical point of view, the cost function has the advantage of being a quadratic function of the state and input signals, which is desirable since quadratic optimization problems have been widely studied in the literature [Boyd and Vandenberghe, 2004]. Now, it is necessary to develop a solution for the Hamilton-Jacobi equation in the light of this formulation. First of all, consider the following theorem.

**Theorem 4.2.** *Consider a continuous cost function  $V : \mathbb{R}^{n \times r \times 1} \rightarrow \mathbb{R}$  given by the LQR cost function defined in Definition 4.3 for a linear system. The optimal cost  $V^*(\mathbf{x}, t)$  has the quadratic form:*

$$V^*(\mathbf{x}, t) = \mathbf{x}^T(t) \mathbf{P}(t) \mathbf{x}(t) \quad (4.14)$$

for any (possibly symmetric) matrix  $\mathbf{P}(t)$  of appropriate dimensions. More precisely, the optimal cost  $V^*(\mathbf{x}, t)$  satisfies the necessary and sufficient conditions for a quadratic function given as, for any  $\lambda \in \mathbb{R}$ :

$$V^*(\lambda \mathbf{x}, t) = \lambda^2 V^*(\mathbf{x}, t) \quad (4.15)$$

$$V^*(\mathbf{x}_1, t) + V^*(\mathbf{x}_2, t) = \frac{1}{2} (V^*(\mathbf{x}_1 + \mathbf{x}_2, t) + V^*(\mathbf{x}_1 - \mathbf{x}_2, t)). \quad (4.16)$$

*Proof.* Since  $V^*(\mathbf{x}, t)$  is the minimum value of  $V(\mathbf{x}, \mathbf{u}^*, t)$ , given an optimal input  $\mathbf{u}^*(t)$  for  $t \in [t, T]$ , then any deviation  $\lambda \in \mathbb{R}$  in the parameters  $(\mathbf{x}, \mathbf{u}(t))$  will result in a greater value of the cost. A direct result from this and the fact that  $V(\mathbf{x}, \mathbf{u}, t)$  is a quadratic function of  $\mathbf{x}(t)$  and  $\mathbf{u}(t)$  is that:

$$V^*(\mathbf{x}, t) \leq V(\lambda \mathbf{x}, \lambda \mathbf{u}^*, t) = \lambda^2 V^*(\mathbf{x}, t) \leq \lambda^2 V(\mathbf{x}, \lambda^{-1} \mathbf{u}^*, t) = V^*(\mathbf{x}, t), \quad (4.17)$$

or, simply:

$$V^*(\mathbf{x}, t) \leq \lambda^2 V^*(\mathbf{x}, t) \leq V^*(\mathbf{x}, t), \quad (4.18)$$

which directly implies  $V^*(\mathbf{x}, t) = \lambda^2 V^*(\mathbf{x}, t)$  and establishes (4.15). Similarly:

$$\begin{aligned} V^*(\mathbf{x}_1, t) + V^*(\mathbf{x}_2, t) &= \frac{1}{4} (V^*(2\mathbf{x}_1, t) + V^*(2\mathbf{x}_2, t)) \\ &\leq \frac{1}{4} (V(2\mathbf{x}_1, \mathbf{u}_{x_1+x_2}^* + \mathbf{u}_{x_1-x_2}^*, t) + V(2\mathbf{x}_2, \mathbf{u}_{x_1+x_2}^* - \mathbf{u}_{x_1-x_2}^*, t)) \\ &= \frac{1}{2} (V(\mathbf{x}_1 + \mathbf{x}_2, \mathbf{u}_{x_1+x_2}^*, t) + V(\mathbf{x}_1 - \mathbf{x}_2, \mathbf{u}_{x_1-x_2}^*, t)) \\ &= \frac{1}{2} (V^*(\mathbf{x}_1 + \mathbf{x}_2, t) + V^*(\mathbf{x}_1 - \mathbf{x}_2, t)) \leq V^*(\mathbf{x}_1, t) + V^*(\mathbf{x}_2, t), \end{aligned} \quad (4.19)$$

which implies  $V^*(\mathbf{x}_1, t) + V^*(\mathbf{x}_2, t) = (1/2) (V^*(\mathbf{x}_1 + \mathbf{x}_2, t) + V^*(\mathbf{x}_1 - \mathbf{x}_2, t))$  and establishes (4.16). Therefore, the optimal cost function has a quadratic form  $V^*(\mathbf{x}, t) = \mathbf{x}^T(t) \mathbf{P}(t) \mathbf{x}(t)$ .  $\square$

This theorem shows that the fact that the cost function  $V(\mathbf{x}, \mathbf{u}, t)$  is quadratic directly implies that the optimal cost  $V^*(\mathbf{x}, t)$  is also quadratic. Therefore, it is possible to define the optimal action  $\mathbf{u}^*(t)$  that produce this optimal cost by evaluating its relationship with the matrix  $\mathbf{P}$ . Since the quadratic cost just evaluated is defined for a finite-horizon optimal controller, it has to obey the condition imposed by the Hamilton-Jacobi equation, and the optimal controller can be solved as shown in the following theorem.

**Theorem 4.3.** (*LQR Control Action*) Given a Linear Quadratic Regulator as defined in Definition 4.3, the optimal action produced by this optimal controller at any time  $t \in [t_0, T]$  is given by:

$$\mathbf{u}^*(t) = -\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}(t)\mathbf{x}(t), \quad (4.20)$$

where  $\mathbf{P}(t)$  is the solution of the matrix Riccati differential equation:

$$-\dot{\mathbf{P}}(t) = \mathbf{A}^T\mathbf{P}(t) + \mathbf{P}(t)\mathbf{A} - \mathbf{P}(t)\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}(t) + \mathbf{Q}, \quad (4.21)$$

with terminal condition  $\mathbf{P}(T) = \mathbf{Q}_f$ .

*Proof.* Consider the Hamilton-Jacobi equation from (4.5) restated below for a quadratic loss function  $l(\mathbf{x}, \mathbf{u}, t)$  and a linear system with state equation  $\mathbf{f}(t)$  as given by Definition 4.3:

$$\begin{aligned} \frac{\partial(\mathbf{x}^T\mathbf{P}\mathbf{x})}{\partial t} &= -\min_{\mathbf{u}(t)} \left\{ \mathbf{x}^T\mathbf{Q}\mathbf{x} + \mathbf{u}^T\mathbf{R}\mathbf{u} + \left[ \frac{\partial(\mathbf{x}^T\mathbf{P}\mathbf{x})}{\partial \mathbf{x}} \right]^T (\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}) \right\}, \\ \mathbf{x}^T\dot{\mathbf{P}}\mathbf{x} &= -\min_{\mathbf{u}(t)} \{ \mathbf{x}^T\mathbf{Q}\mathbf{x} + \mathbf{u}^T\mathbf{R}\mathbf{u} + (2\mathbf{x}^T\mathbf{P})(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}) \} \end{aligned} \quad (4.22)$$

where, without loss of generalization,  $\mathbf{P}(t)$  was assumed to be symmetric. Because of the quadratic nature, it is possible to calculate the minimum of the right-hand side of the equation by taking its derivative with respect to the control action and equating it to zero:

$$\begin{aligned} 0 &= \frac{\partial}{\partial \mathbf{u}(t)} (\mathbf{x}^T\mathbf{Q}\mathbf{x} + \mathbf{u}^T\mathbf{R}\mathbf{u} + (2\mathbf{x}^T\mathbf{P})(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u})) \\ 0 &= 2\mathbf{u}^T(t)\mathbf{R} + (2\mathbf{x}^T(t)\mathbf{P}(t))\mathbf{B} \\ \mathbf{u}(t) &= -\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}(t)\mathbf{x}(t) \end{aligned} \quad (4.23)$$

To solve for  $\mathbf{P}$  first note that the following identity can be found by completing the squares:

$$\begin{aligned} \mathbf{x}^T\mathbf{Q}\mathbf{x} + \mathbf{u}^T\mathbf{R}\mathbf{u} + (2\mathbf{x}^T\mathbf{P})(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}) &= (\mathbf{u} + \mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{x})^T\mathbf{R}(\mathbf{u} + \mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{x}) \\ &\quad + \mathbf{x}^T(\mathbf{A}^T\mathbf{P} + \mathbf{P}\mathbf{A} - \mathbf{P}\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P} + \mathbf{Q})\mathbf{x}. \end{aligned} \quad (4.24)$$

Thus, substituting this identity and substituting the optimal control action results in the matrix Riccati equation:

$$\begin{aligned} \mathbf{x}^T\dot{\mathbf{P}}\mathbf{x} &= \mathbf{x}^T(\mathbf{A}^T\mathbf{P} + \mathbf{P}\mathbf{A} - \mathbf{P}\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P} + \mathbf{Q})\mathbf{x} \\ \dot{\mathbf{P}} &= \mathbf{A}^T\mathbf{P} + \mathbf{P}\mathbf{A} - \mathbf{P}\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P} + \mathbf{Q} \end{aligned} \quad (4.25)$$

Finally, the terminal condition is direct from the fact that  $V^*(\mathbf{x}, T) = \mathbf{x}^T(T)\mathbf{Q}_f\mathbf{x}(T)$ .  $\square$

A closed-form solution for the LQR problem makes this controller a very appealing solution in several applications, since not only it is an optimal controller but the action can be calculated analytically. Several others optimal controllers may offer performance improvements, but usually depends on iterative optimization procedures, which can be computationally expensive to calculate for high-dimensional systems. Furthermore, the time complexity of the LQR control action computation is governed by the solution of the Riccati differential equation which can be done very efficiently [Davison and Maki, 1973, Opanuga et al., 2015].

Now, consider the optimal control action  $\mathbf{u}^*(t)$ . Applying this action to a linear system in State-Space representation results in the following state-equation:

$$\begin{aligned}\dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t) + \mathbf{B}(-\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}(t)\mathbf{x}(t)) \\ &= (\mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}(t))\mathbf{x}(t) \\ &= (\mathbf{A} - \mathbf{B}\mathbf{K}(t))\mathbf{x}(t)\end{aligned}\quad (4.26)$$

This formulation indicates that the Linear Quadratic Regulator solution follows the exact same form of a general state-feedback regulator, represented by the linear but time-variant feedback system  $\mathbf{A}_{cl}(t) = (\mathbf{A} - \mathbf{B}\mathbf{K}(t))$ . Notice, also, that both the computation of  $\mathbf{K}(t)$  and  $\mathbf{P}(t)$  does not explicitly depends on  $\mathbf{x}(t)$ , meaning that, for a specific control horizon, they can be calculated off-line and then provided to the controller actuator for the on-line operation. For this reason, despite being a closed-loop controller with corrective action, the LQR can be considered an open-loop optimizer, since the optimization solution itself does not depend on the state signal, but only on the system dynamical model.

### Reference Tracking

As discussed in the previous chapter, the class of regulation controllers are broad but still restricted to a zero-state reference signal. To expand the range of application for the optimal controller just derived, it is possible to introduce integral action to the feedback.

**Definition 4.4.** (Linear Quadratic Servo) Given a linear State-Space system represented by matrices  $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$ , augmented with state  $\dot{\mathbf{x}}_a(t) = \mathbf{r}(t) - \mathbf{C}\mathbf{x}(t)$ :

$$\begin{cases} \begin{bmatrix} \dot{\mathbf{x}}(t) \\ \dot{\mathbf{x}}_a(t) \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{A} & \mathbf{0} \\ -\mathbf{C} & \mathbf{0} \end{bmatrix}}_{\tilde{\mathbf{A}}} \underbrace{\begin{bmatrix} \mathbf{x}(t) \\ \mathbf{x}_a(t) \end{bmatrix}}_{\tilde{\mathbf{x}}(t)} + \underbrace{\begin{bmatrix} \mathbf{B} \\ \mathbf{0} \end{bmatrix}}_{\tilde{\mathbf{B}}} \mathbf{u}(t) + \begin{bmatrix} \mathbf{0} \\ \mathbf{I} \end{bmatrix} \mathbf{r}(t) \\ \mathbf{y}(t) = \begin{bmatrix} \mathbf{C} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}(t) \\ \mathbf{x}_a(t) \end{bmatrix} \end{cases} \quad (4.27)$$

A *Linear Quadratic Servo* (LQ-Servo) for this system is an optimal controller defined by the quadratic cost function:

$$J(\mathbf{x}, \mathbf{u}, t_0) = \int_{t_0}^T \left( \tilde{\mathbf{x}}^T \tilde{\mathbf{Q}} \tilde{\mathbf{x}} + \mathbf{u}^T \mathbf{R} \mathbf{u} \right) dt + \tilde{\mathbf{x}} \tilde{\mathbf{Q}}_f \tilde{\mathbf{x}}(T), \quad (4.28)$$

where is assumed that  $\tilde{\mathbf{Q}}, \tilde{\mathbf{Q}}_f \succ 0$  and  $\mathbf{R} \succ 0$  are matrices penalizing, respectively, the state-vector magnitude and the control effort.

As before, the integral action turns the tracking problem into a regulation problem. In this case, the optimal controller will try to regulate towards a zero-state for  $\mathbf{x}_a(t)$ , yielding a zero error between the output  $\mathbf{y}(t)$  and reference signal  $\mathbf{r}(t)$ . The new  $p$  diagonal entries of the augmented matrix  $\tilde{\mathbf{Q}}$  can be interpreted as weights penalizing the state deviation from the reference signal. The solution for this controller is the same as the one derived in Theorem 4.3, with  $\mathbf{A} = \tilde{\mathbf{A}}$  and  $\mathbf{B} = \tilde{\mathbf{B}}$ . The resulting control action is equal to:

$$\mathbf{u}(t) = \tilde{\mathbf{K}}(t)\tilde{\mathbf{x}}(t) = [\mathbf{K}(t) \quad \mathbf{K}_a(t)] \begin{bmatrix} \mathbf{x}(t) \\ \mathbf{x}_a(t) \end{bmatrix}. \quad (4.29)$$



### Infinite-Horizon LQ Controllers

To conclude the section, it is also worth mentioning that, for the LQR cost function, it is possible to determine an optimal controller for an infinite-horizon operation, that is, when  $T \rightarrow \infty$ . The fact that the control horizon is now infinite results in a stationary state-feedback gain  $\mathbf{K}$ , as stated below.

**Theorem 4.4.** (*Infinite-Horizon LQR*) Consider a Linear Quadratic Regulator as defined in Definition 4.3, but with  $T = \infty$ . The optimal control action produced by this optimal controller at any time  $t \in [t_0, \infty]$  is given by:

$$\mathbf{u}^*(t) = -\mathbf{K}\mathbf{x}(t) = -\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{x}(t), \quad (4.30)$$

where  $\mathbf{K} = \mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}$  and  $\mathbf{P}$  is the solution of the matrix algebraic Riccati equation:

$$\mathbf{0} = \mathbf{A}^T\mathbf{P} + \mathbf{P}\mathbf{A} - \mathbf{P}\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P} + \mathbf{Q}. \quad (4.31)$$

A detailed proof of this result can be found in [Anderson and Moore, 1990]. The possibility for an infinite-horizon LQR is desirable in the sense that it doesn't restrict the operation to a fixed time interval, however, the resulting optimal controller suffers from a worst performance when compared to finite-time horizon controllers.

## 4.3 Optimal State Estimators

The state-vector  $\mathbf{x}(t)$  is a necessary information to calculate  $\mathbf{u}(t)$  in state-feedback controllers, but may also be inaccessible in practice. In Chapter 3, a deterministic observer, namely the Luenberger observer, was derived to deal with this problem. Despite being possible to utilize the Luenberger observer together with optimal controllers, this method suffers from the same design limitations of the Pole-Placement method, since they are dual problems. Therefore, this section develops an optimal state estimation approach which, from the estimation nature of the problem, relies on a statistical interpretation of the dynamical system and its observations.

### General State Estimators

First of all, consider a system perturbed by disturbances in the state-response and in the measurements, respectively  $\mathbf{w} : \mathbb{R} \rightarrow \mathbb{R}^n$  and  $\mathbf{v} : \mathbb{R} \rightarrow \mathbb{R}^p$ . This configuration is illustrated at Fig. 4.2. If these disturbances are assumed to be zero-mean Gaussian variables, i.e.,  $\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{kf})$  and  $\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_{kf})$  where  $\mathbf{Q}_{kf}$  and  $\mathbf{R}_{kf}$  are covariances matrices of appropriate sizes, then it is possible to motivate a stochastic formulation of the linear State-Space model. For the sake of simplicity, the discussion is focused on the discrete-time case.

**Definition 4.5.** (Stochastic Discrete-Time State-Space) Given an additive process noise variable  $\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{kf})$  with covariance  $\mathbf{Q}_{kf} \in \mathbb{R}^{n \times n}$  and an additive measurement noise variable  $\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_{kf})$  with covariance  $\mathbf{R}_{kf} \in \mathbb{R}^{p \times p}$ , the stochastic discrete-time State-Space model is given by the equations:

$$\begin{cases} \mathbf{x}_{k+1} = \mathbf{A}_d\mathbf{x}_k + \mathbf{B}_d\mathbf{u}_k + \mathbf{w}_k \\ \mathbf{y}_k = \mathbf{C}_d\mathbf{x}_k + \mathbf{v}_k \end{cases}, \quad (4.32)$$

for the initial condition  $\mathbf{x}_0 \sim p(\mathbf{x}_0)$ .

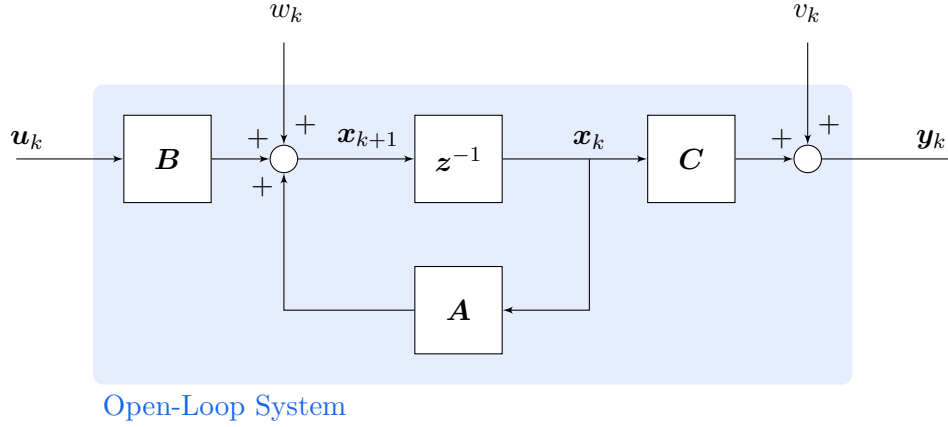


Figure 4.2: Block diagram of a stochastic discrete State-Space system.

With this formulation, the noise variables can actually represent both the effect of external disturbances and the model uncertainty. This is important because dynamical models, even those developed according to a first-principle methodology, are not the true representation of the actual dynamical system or process. The motivation for a discussion of only a discrete State-Space model is to facilitate the interpretation of the statistical properties of the model.

A direct result from Definition 4.5 is that, since  $w_k$  and  $v_k$  are random variables,  $x_k$  and  $y_k$  are also random variables. More precisely, since the noises variables are assumed to have a Gaussian distribution, both variables  $x_k \sim \mathcal{N}(\mu_{x_k}, \Sigma_{x_k})$  and  $y_k \sim \mathcal{N}(\mu_{y_k}, \Sigma_{y_k})$  are also Gaussian, where  $\mu_{x_k}, \mu_{y_k}$  are their respective mean vectors and  $\Sigma_{x_k}, \Sigma_{y_k}$  are the respective covariance matrices. This is because the Gaussian distribution is closed to linear operations. For this reason, each state and output trajectory over a discrete time interval is a collection of time-indexed random variables following the Gaussian distribution, i.e., a stochastic process known as a Gaussian Process [Rasmussen and Williams, 2006].

**Definition 4.6.** (Gaussian Process) Given a Gaussian process defined as a collection of random variables:

$$\mathbf{X} = \{x_k \sim \mathcal{N}(\mu_{x_k}, \Sigma_{x_k}); k \in [0, 1, \dots, K]\}. \quad (4.33)$$

In the case that the variables represents states of a causal system, the joint probability of the random variables can be modeled as:

$$p(\mathbf{X}) = p(x_0, x_1, \dots, x_K) = \prod_{k=0}^K p(x_k | x_{k-1}, \dots, x_1, x_0). \quad (4.34)$$

The defined joint probability  $p(x_0, \dots, x_K)$  is interpreted as the probability of occurring a specific state trajectory. The conditional probability  $p(x_k | x_{k-1}, \dots, x_0)$  states that the probability associated with a random state variable  $x_k$  can be conditioned on the state trajectory  $(x_{k-1}, \dots, x_0)$  from the previous time instants. These two probabilities can be related through the well-known *Chain Rule of Probability*, which, in the context of stochastic processes, is a realization of the *Markov Property* [Stratonovich, 1968]. A nice property of Gaussian distributions, in this context, is that the joint and condition distribution of Gaussian random variables are still Gaussian [Ross, 2010].

### Kalman Filter

The dynamic response  $\mathbf{x}_{k+1}$  of a system for a time instant  $k + 1$ , given its trajectory until the actual state  $\mathbf{x}_k$ , is given by the conditional probability  $p(\mathbf{x}_{k+1}|\mathbf{x}_k, \dots, \mathbf{x}_1, \mathbf{x}_0)$ . One can assume, however, that the an actual state stores enough information from the state trajectory until that point, that is possible to determine the next state response directly through its probability. Therefore, the conditional probability could be restated as:

$$p(\mathbf{x}_{k+1}|\mathbf{x}_k, \dots, \mathbf{x}_1, \mathbf{x}_0) = p(\mathbf{x}_{k+1}|\mathbf{x}_k). \quad (4.35)$$

A process that exhibits this property is known as a First-Order Markov process. With this formulation, it is always possible to marginalize the probability distribution of a random variable  $\mathbf{x}_{k+1}$  given the initial state  $\mathbf{x}_0$ , which are related through the equation:

$$p(\mathbf{x}_{k+1}) = p(\mathbf{x}_{k+1}|\mathbf{x}_k)p(\mathbf{x}_k|\mathbf{x}_{k-1}) \cdots p(\mathbf{x}_1|\mathbf{x}_0)p(\mathbf{x}_0) = p(\mathbf{x}_0) \prod_{i=0}^k p(\mathbf{x}_{i+1}|\mathbf{x}_i). \quad (4.36)$$

Combining this assumption with Definition 4.6, in order to describe a Gauss-Markov Process, and using the model from Definition 4.5 it is possible to describe a stochastic dynamical system as the following probabilistic State-Space model.

**Definition 4.7.** (Controlled Hidden Markov Model) Given a stochastic State-Space model as in Definition 4.5, and assuming the Markov property, a *Controlled Hidden Markov Model* for this dynamical system at a time instant  $k$  is defined through the two distributions:

$$\begin{cases} p(\mathbf{x}_k|\mathbf{x}_{k-1}, \mathbf{u}_{k-1}) & \text{(Transition Distribution)} \\ p(\mathbf{y}_k|\mathbf{x}_k) & \text{(Emission Distribution)} \end{cases}, \quad (4.37)$$

for the initial condition given by  $p(\mathbf{x}_0)$ . Furthermore, since the noises  $\mathbf{w}_k$  and  $\mathbf{v}_k$  are zero-mean Gaussian noises, these distributions are modeled as the Normal conditional distributions:

$$\begin{cases} p(\mathbf{x}_k|\mathbf{x}_{k-1}) = \mathcal{N}(\mathbf{x}_k|\mathbf{A}\mathbf{x}_{k-1} + \mathbf{B}\mathbf{u}_{k-1}, \mathbf{Q}_{kf}) \\ p(\mathbf{y}_k|\mathbf{x}_k) = \mathcal{N}(\mathbf{y}_k|\mathbf{C}\mathbf{x}_k, \mathbf{R}_{kf}) \end{cases}. \quad (4.38)$$

A Hidden-Markov Model is a statistical interpretation of dynamical systems that is well-known in the literature [Barber, 2012]. An illustration of the dependence between these variables as a graphical model is depicted in Fig. 4.3. Notice that, in this formulation, the state-vector  $\mathbf{x}_k$  actual value is never known (since it is a latent variable), but can be observed through the emission variable  $\mathbf{y}_k$ , under the assumption that system model is perfectly known. This raises the problem of state estimation as a common inference problem known as the *filtering*, which consists on determining statistically the actual value of  $\mathbf{x}_k$  given the observations  $\mathbf{y}_0, \dots, \mathbf{y}_k$  and control actions  $\mathbf{u}_0, \dots, \mathbf{u}_k$ . Thus, the state can be obtained through the distribution:

$$p(\mathbf{x}_k|\mathbf{y}_k, \mathbf{y}_{k-1}, \dots, \mathbf{y}_0, \mathbf{u}_{k-1}, \mathbf{u}_{k-2}, \dots, \mathbf{u}_0), \quad (4.39)$$

from which a representative value, such as the expected value  $\mathbb{E}[p(\mathbf{x}_k|\mathbf{y}_k, \dots, \mathbf{y}_0, \mathbf{u}_{k-1}, \dots, \mathbf{u}_0)]$  can be selected.

There are several approaches to solve the filtering and other related problems in statistical state estimation [Särkkä, 2013], but the most popular method has been the Kalman filter, presented in the seminal paper [Kalman, 1960a]. The discrete formulation for a linear Kalman Filter is given below.

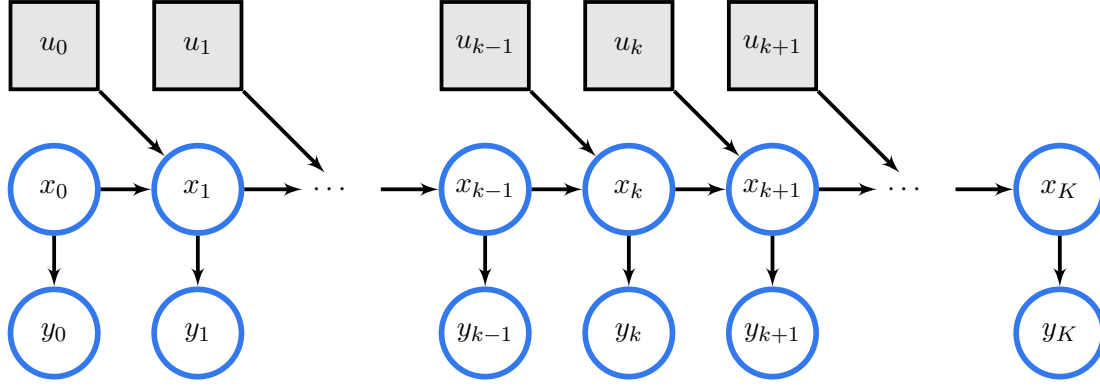


Figure 4.3: A Controlled Hidden Markov Model for a discrete-time system, where the shaded boxes indicates that the inputs are not random variables.

**Theorem 4.5.** (*Kalman Filter*) Given a Controlled Hidden Markov Model as in Definition 4.7 and initial state distribution  $\mathbf{x}_0 \sim \mathcal{N}(\bar{\mathbf{x}}_0, \mathbf{P}_0)$ , the filtering distribution can be solved in closed-form as:

$$p(\mathbf{x}_k | \mathbf{y}_k, \dots, \mathbf{y}_0, \mathbf{u}_{k-1}, \dots, \mathbf{u}_0) = \mathcal{N}(\bar{\mathbf{x}}_k, \mathbf{P}_k), \quad (4.40)$$

where  $\bar{\mathbf{x}}_k$ , which denotes the mean of  $\mathbf{x}_k$ , and covariance  $\mathbf{P}_k$  are recursively computed through the prediction and update steps below.

- *Prediction Step:*

$$\begin{aligned} \hat{\mathbf{x}}_k &= \mathbf{A}\bar{\mathbf{x}}_{k-1} + \mathbf{B}\mathbf{u}_{k-1}, \\ \hat{\mathbf{P}}_k &= \mathbf{A}\mathbf{P}_{k-1}\mathbf{A}^T + \mathbf{Q}_{kf}. \end{aligned} \quad (4.41)$$

- *Update Step:*

$$\begin{aligned} \mathbf{K}_k &= \hat{\mathbf{P}}_k \mathbf{C}^T (\mathbf{R}_{kf} + \mathbf{C}\hat{\mathbf{P}}_k\mathbf{C}^T)^{-1}, \\ \bar{\mathbf{x}}_k &= \hat{\mathbf{x}}_k + \mathbf{K}_k(\mathbf{y}_k - \mathbf{C}\hat{\mathbf{x}}_k), \\ \mathbf{P}_k &= \hat{\mathbf{P}}_k - \mathbf{K}_k(\mathbf{R}_{kf} + \mathbf{C}\hat{\mathbf{P}}_k\mathbf{C}^T)\mathbf{K}_k^T. \end{aligned} \quad (4.42)$$

*Proof.* Consider the transition and emission distributions given as:

$$\begin{cases} p(\mathbf{x}_k | \mathbf{x}_{k-1}) = \mathcal{N}(\mathbf{x}_k | \mathbf{A}\bar{\mathbf{x}}_{k-1} + \mathbf{B}\mathbf{u}_{k-1}, \mathbf{Q}_{kf}) \\ p(\mathbf{y}_k | \mathbf{x}_k) = \mathcal{N}(\mathbf{y}_k | \mathbf{C}\bar{\mathbf{x}}_k, \mathbf{R}_{kf}) \end{cases}. \quad (4.43)$$

Consider, now, the distribution of  $(\mathbf{x}_{k-1}, \mathbf{x}_k)$  conditioned on measurements  $(\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{k-1})$ . Since the joint and condition distributions of a Gaussian are themselves Gaussian, this density can be expressed as:

$$\begin{aligned} p(\mathbf{x}_k, \mathbf{x}_{k-1} | \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{k-1}) &= p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{k-1}) \\ &= \mathcal{N}(\mathbf{x}_k | \mathbf{A}\bar{\mathbf{x}}_{k-1} + \mathbf{B}\mathbf{u}_{k-1}, \mathbf{Q}_{kf}) \mathcal{N}(\mathbf{x}_{k-1} | \bar{\mathbf{x}}_{k-1}, \mathbf{P}_{k-1}) \\ &= \mathcal{N} \left( \begin{bmatrix} \mathbf{x}_{k-1} \\ \mathbf{x}_k \end{bmatrix} \middle| \begin{bmatrix} \bar{\mathbf{x}}_{k-1} \\ \mathbf{A}\bar{\mathbf{x}}_{k-1} + \mathbf{B}\mathbf{u}_{k-1} \end{bmatrix}, \begin{bmatrix} \mathbf{P}_{k-1} & \mathbf{P}_{k-1}\mathbf{A}^T \\ \mathbf{A}\mathbf{P}_{k-1} & \mathbf{A}\mathbf{P}_{k-1}\mathbf{A}^T + \mathbf{Q}_{kf} \end{bmatrix} \right). \end{aligned} \quad (4.44)$$

Marginalizing the joint distribution, it is possible to obtain the distribution of  $\mathbf{x}_k$  conditioned

in the measurement history as:

$$p(\mathbf{x}_k | \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{k-1}) = \mathcal{N}(\underbrace{\mathbf{A}\bar{\mathbf{x}}_{k-1} + \mathbf{B}\mathbf{u}_{k-1}}_{\hat{\mathbf{x}}_k}, \underbrace{\mathbf{A}\mathbf{P}_{k-1}\mathbf{A}^T + \mathbf{Q}_{kf}}_{\hat{\mathbf{P}}_k}), \quad (4.45)$$

which completes the *prediction step* of the Kalman filter. For the update step, consider the joint distribution between  $(\mathbf{x}_k, \mathbf{y}_k)$ , the actual state and actual measurement, which applying Lemma A.1 results in:

$$\begin{aligned} p(\mathbf{x}_k, \mathbf{y}_k | \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{k-1}) &= p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{k-1}) \\ &= \mathcal{N}(\mathbf{y}_k | \mathbf{C}\bar{\mathbf{x}}_k, \mathbf{R}_{kf}) \mathcal{N}(\mathbf{x}_k | \hat{\mathbf{x}}_k, \hat{\mathbf{P}}_k) \\ &= \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_k \\ \mathbf{y}_k \end{bmatrix} \middle| \begin{bmatrix} \hat{\mathbf{x}}_k \\ \mathbf{C}\hat{\mathbf{x}}_k \end{bmatrix}, \begin{bmatrix} \hat{\mathbf{P}}_k & \hat{\mathbf{P}}_k \mathbf{C}^T \\ \mathbf{C}\hat{\mathbf{P}}_k & \mathbf{C}\hat{\mathbf{P}}_k \mathbf{C}^T + \mathbf{R}_{kf} \end{bmatrix}\right). \end{aligned} \quad (4.46)$$

Therefore, using Lemma A.2, the filtering distribution can be obtained as:

$$p(\mathbf{x}_k | \mathbf{y}_k, \mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{k-1}) = \mathcal{N}(\mathbf{x}_k | \bar{\mathbf{x}}_k, \mathbf{P}_k), \quad (4.47)$$

where  $\hat{\mathbf{x}}_k$  and  $\mathbf{P}_k$  are the computations of the update step, given as:

$$\begin{cases} \bar{\mathbf{x}}_k = \hat{\mathbf{x}}_k + \hat{\mathbf{P}}_k \mathbf{C}^T (\mathbf{R}_{kf} + \mathbf{C}\hat{\mathbf{P}}_k \mathbf{C}^T)^{-1} (\mathbf{y}_k - \mathbf{C}\hat{\mathbf{x}}_k) \\ \mathbf{P}_k = \hat{\mathbf{P}}_k - \hat{\mathbf{P}}_k \mathbf{C}^T (\mathbf{R}_{kf} + \mathbf{C}\hat{\mathbf{P}}_k \mathbf{C}^T)^{-1} \mathbf{C}\hat{\mathbf{P}}_k \end{cases}. \quad (4.48)$$

Finally, making  $\mathbf{K}_k = \hat{\mathbf{P}}_k \mathbf{C}^T (\mathbf{R}_{kf} + \mathbf{C}\hat{\mathbf{P}}_k \mathbf{C}^T)^{-1}$  concludes the proof.  $\square$

Notice that the Kalman filter solution implies that an “open-loop” estimation of the mean of the states,  $\hat{\mathbf{x}}_k$ , can be corrected by the equation:

$$\bar{\mathbf{x}}_k = \hat{\mathbf{x}}_k + \mathbf{K}_k (\mathbf{y}_k - \mathbf{C}\hat{\mathbf{x}}_k), \quad (4.49)$$

which, by expanding  $\hat{\mathbf{x}}_k$ , results in a similar state-equation expression for a system with a deterministic observer as defined earlier, but for a time-varying gain  $\mathbf{K}_k$ :

$$\begin{aligned} \bar{\mathbf{x}}_k &= \mathbf{A}\bar{\mathbf{x}}_k + \mathbf{B}\mathbf{u}_k + \mathbf{K}_k (\mathbf{y}_k - \mathbf{C}\hat{\mathbf{x}}_k) \\ &= (\mathbf{A} - \mathbf{K}_k \mathbf{C}) \bar{\mathbf{x}}_k + \mathbf{B}\mathbf{u}_k + \mathbf{K}_k \mathbf{y}_k. \end{aligned} \quad (4.50)$$

Notice that, despite the similarity with Luenberger observers, the Kalman gain has the advantage that the gain matrix  $\mathbf{K}_k$  is not resulting from a designer procedure but rather from solving an inference problem whose only additional information needed are the covariances  $\mathbf{Q}_{kf}$  and  $\mathbf{R}_{kf}$ , which can be assumed or estimated.

### Controller-Estimator Duality

In addition to that, it can be shown that the recursive steps that solves the filtering problem are also responsible for iteratively minimizing the trace of the covariance matrix  $\mathbf{P}_k$ , which is given as:

$$\mathbf{P}_k = \mathbb{E}[(\mathbf{x}_k - \bar{\mathbf{x}}_k)(\mathbf{x}_k - \bar{\mathbf{x}}_k)^T]. \quad (4.51)$$

Since this method minimizes some cost function, it is indeed an optimal estimator. In fact, the Kalman filter can be formulated through an optimal control formulation, which motivates the development of a continuous-time version of this estimator, as shown below.

**Theorem 4.6.** (*Kalman-Bucy Filter*) Consider a continuous-time State-Space linear system subject to additive process noise variable  $\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{kf})$  and measurement noise variable  $\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_{kf})$ , where the covariances  $\mathbf{Q}_{kf} \in \mathbb{R}^{n \times n}$  and  $\mathbf{R}_{kf} \in \mathbb{R}^{p \times p}$  represents the power spectral density of the noises. In this case, for an estimated state  $\bar{\mathbf{x}}(t)$  at time  $t$ , the error covariance:

$$\mathbf{J}(\mathbf{x}, \bar{\mathbf{x}}, t) = \mathbb{E} \{ [\mathbf{x}(t) - \bar{\mathbf{x}}(t)][\mathbf{x}(t) - \bar{\mathbf{x}}(t)]^T \} \quad (4.52)$$

is minimized by  $\bar{\mathbf{x}}(t)$  obtained through the system:

$$\dot{\bar{\mathbf{x}}}(t) = \mathbf{A}\bar{\mathbf{x}}(t) + \mathbf{B}\mathbf{u}(t) + \mathbf{K}_e(t)(\mathbf{y}(t) - \mathbf{C}\bar{\mathbf{x}}(t)), \quad (4.53)$$

where  $\mathbf{K}_e(t) = \mathbf{P}_e(t)\mathbf{C}\mathbf{R}^{-1}$ , being  $\mathbf{P}_e(t)$  the solution of the Riccati differential matrix equation:

$$\dot{\mathbf{P}}_e(t) = \mathbf{A}\mathbf{P}_e(t) + \mathbf{P}_e(t)\mathbf{A}^T - \mathbf{P}_e(t)\mathbf{C}^T\mathbf{R}_{kf}^{-1}\mathbf{C}\mathbf{P}_e(t) + \mathbf{Q}_{kf}, \quad (4.54)$$

with initial condition  $\mathbf{P}_e(t_0) = \mathbb{E} \{ [\mathbf{x}(t_0) - \bar{\mathbf{x}}(t_0)][\mathbf{x}(t_0) - \bar{\mathbf{x}}(t_0)]^T \}$  for  $t_0 > -\infty$ .

A detailed proof of this theorem can be found in [Crassidis and Junkins, 2011]. The formulation of an optimal state estimator further emphasizes the duality between controllers and observers. Notice that the gain  $\mathbf{K}_e(t)$  of a Kalman-Bucy filter depends on a matrix  $\mathbf{P}_e(t)$  that solves the exact same Riccati differential equation as the one in Theorem 4.3, but forward in time for matrices  $\mathbf{A}^T$  and  $\mathbf{C}^T$  instead of  $\mathbf{A}$  and  $\mathbf{B}$ . Notice, however, that these filters are not limited to a terminal time  $T$ , since the boundary condition only requires information on the initial time  $t_0 > -\infty$ , which implies that they could also be used in forward infinite-horizon operations. Another direct result of the duality is that is possible to derive an estimator considering  $t_0 \rightarrow -\infty$ , which, as in Theorem 4.4, results in a time-invariant filter with gain  $\mathbf{K}_e$ .

## 4.4 Linear Quadratic Gaussian (LQG)

This chapter presented formulation for controllers and state estimators that are optimal, in the sense that they minimize some cost function based on information about the model and measurements. In the previous chapter, the Theorem 3.5 stated that a controller and an estimator can be designed separately, and a resulting control by feedback of estimated states is always feasible. In this section, a similar configuration is shown for optimal control and estimation operations, as depicted in Fig. 4.4.

The control configuration that connects an optimal control action of a Linear Quadratic Regulator together with the estimation states from a Kalman filter (or Kalman-Bucy filter) is known as a *Linear Quadratic Gaussian* (LQG) controller. Since the controller is independent of the estimator, it is also possible to include integral action to the configuration, as shown in the block diagram, presenting this architecture as a general configuration that can deal with a broad range of applications. A formal definition is given below.

**Definition 4.8.** (Linear Quadratic Gaussian) Consider a stochastic system in State-Space representation just as in Definition 4.5, rewritten below:

$$\begin{cases} \dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) + \mathbf{w}(t) \\ \mathbf{y}(t) &= \mathbf{C}\mathbf{x}(t) + \mathbf{v}(t) \end{cases}, \quad (4.55)$$

whose estimated state-vector  $\hat{\mathbf{x}}(t)$  is determined by a Kalman-Bucy filter and whose optimal input signal  $\mathbf{u}(t)$  is calculated through a finite-horizon LQR. The Linear Quadratic Gaussian (LQG) control for the horizon  $t \in [t_0, T]$ , with  $-\infty < t_0 \leq T < \infty$ , is defined as:

$$\dot{\hat{\mathbf{x}}}(t) = [\mathbf{A} - \mathbf{K}_e(t)\mathbf{C} - \mathbf{B}\mathbf{K}(t)]\hat{\mathbf{x}}(t) + \mathbf{K}_e(t)\mathbf{y}(t), \quad (4.56)$$

where  $\mathbf{K}(t) = \mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}(t)$  and  $\mathbf{K}_e(t) = \mathbf{P}_e(t)\mathbf{C}\mathbf{R}^{-1}$  are, respectively, the LQR and Kalman-Bucy gains for matrices  $\mathbf{P}(t)$  and  $\mathbf{P}_e(t)$  that solve the Riccati differential equations:

$$\begin{cases} -\dot{\mathbf{P}}(t) = \mathbf{A}^T\mathbf{P}(t) + \mathbf{P}(t)\mathbf{A} - \mathbf{P}(t)\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}(t) + \mathbf{Q} \\ \dot{\mathbf{P}}_e(t) = \mathbf{A}\mathbf{P}_e(t) + \mathbf{P}_e(t)\mathbf{A}^T - \mathbf{P}_e(t)\mathbf{C}^T\mathbf{R}_k^{-1}\mathbf{C}\mathbf{P}_e(t) + \mathbf{Q}_{kf} \end{cases} \quad (4.57)$$

for boundary conditions  $\mathbf{P}(T) = \mathbf{Q}_f$  and  $\mathbf{P}_e(t_0) = \mathbb{E}\{[\mathbf{x}(t_0) - \bar{\mathbf{x}}(t_0)][\mathbf{x}(t_0) - \bar{\mathbf{x}}(t_0)]^T\}$ , respectively.

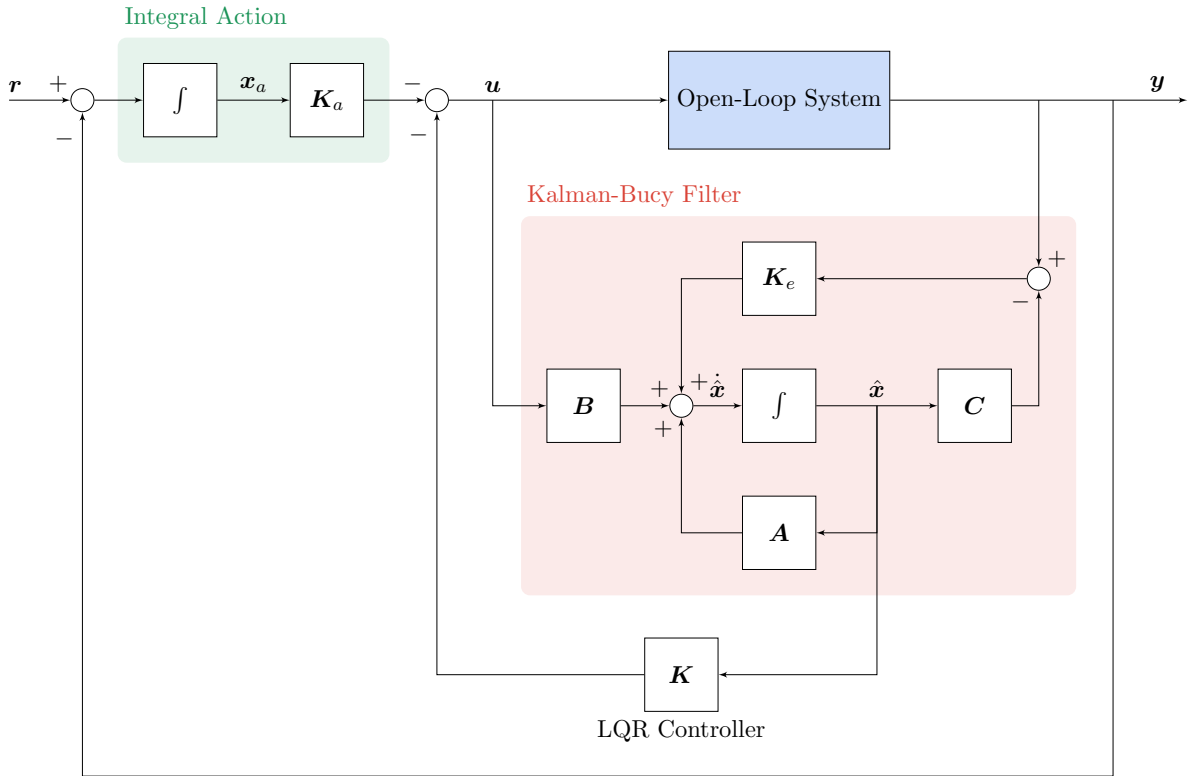


Figure 4.4: Block diagram of a Linear Quadratic Gaussian control configuration.

Notice that another advantage of this controller-estimator configuration is that the solution for both Riccati differential equations does not depend on the states at any time  $\mathbf{x}(t)$ , with the exception of the boundary condition  $\mathbf{P}_e(t_0)$  which can in any case be obtained before the system operation. Therefore, the values for matrices  $\mathbf{P}(t)$  and  $\mathbf{P}_e(t)$ , and consequently  $\mathbf{K}(t)$  and  $\mathbf{K}_e(t)$ , for  $t \in [t_0, T]$ , can be evaluated beforehand and then applied in the on-line operation.

The LQG controller can also be used for tracking non-constant references, given that the system can be augmented while maintaining controllability. The augmented system is used to calculate the controller gains  $\mathbf{K}(t)$  using Definition 4.4, while the estimator gains  $\mathbf{K}_e(t)$  are optimized through the original system. Given that it can optimally solve both the regulator and tracking problems, while still accounting for uncertainty in the system, the Linear Quadratic

Gaussian poses as a high performing multi-purpose controller. The resulting augmented system can be represented by the closed-loop state equation:

$$\begin{cases} \begin{bmatrix} \dot{\hat{\mathbf{x}}}(t) \\ \dot{\hat{\mathbf{x}}}_a(t) \end{bmatrix} = \begin{bmatrix} \mathbf{A} - \mathbf{BK}(t) - \mathbf{K}_e(t)\mathbf{C} & -\mathbf{BK}_a(t) \\ -\mathbf{C} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{x}}(t) \\ \hat{\mathbf{x}}_a(t) \end{bmatrix} + \begin{bmatrix} \mathbf{K}_e(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{y}(t) \\ \mathbf{r}(t) \end{bmatrix} \\ \mathbf{y}(t) = \begin{bmatrix} \mathbf{C} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}(t) \\ \mathbf{x}_a(t) \end{bmatrix} \end{cases} \quad (4.58)$$

**Example 4.1.** For the sake of illustration, consider the same system as in Example 3.1. Consider, now, that it is perturbed by noises  $\mathbf{w}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{kf})$  and  $\mathbf{v}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_{kf})$ , with covariances:

$$\mathbf{Q}_{kf} = \begin{bmatrix} 0.001 & 0.001 \\ 0.001 & 0.001 \end{bmatrix} \quad \mathbf{R}_{kf} = 0.001. \quad (4.59)$$

A LQG controllers for reference tracking is obtained by solving the controller Ricatti differential equation for weights  $\mathbf{Q} = \text{diag}(10, 10, 1)$  and  $\mathbf{R} = 0.1$ , considering the terminal weight  $\mathbf{P}(T) = \mathbf{Q}$ . The estimator is solved forward in time for initial condition  $\mathbf{P}_e(t_0) = \mathbf{Q}_{kf}$ .

The resulting simulations from Example 4.1 are shown in Fig. 4.5. Notice that, since the integral action is now calculated over the observations, which are noisy, the controller produced a noisy input signal. The reference tracking was still achieved, but with very oscillatory behavior. A solution to this problem could be using the integral action as a feedback over an output signal obtained from the estimated states, which are probably smoother. Another solution would be to use a *smoothing filter* directly to the input signal to avoid this oscillatory action.

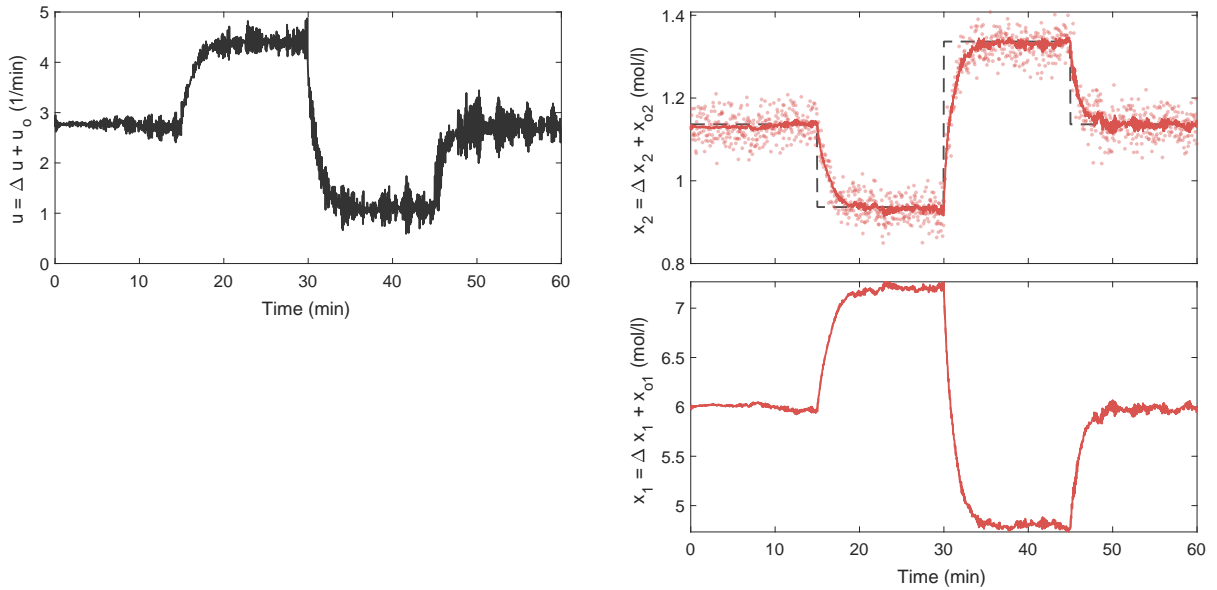


Figure 4.5: Simulation of the LQG controllers showing the input signal (left) and correspondent actual state responses (right). The dots denotes represents the measurements from the real system.



## 4.5 Stability and Robustness Analysis

In the last chapter, the poles of the feedback controlled system were arbitrarily assigned using the pole-placement method. In this chapter, however, the controllers result from optimization procedures, so it is necessary to discuss whether or not these controllers are stable and robust given the criteria already presented.

### Stability

First of all, the stability of the closed-loop system via Linear Quadratic Regulator is discussed. One might conclude that, since this controller is obtained from an optimization process, the resulting controller can not be unstable. Consider, however, the popular counter-intuition: consider a scalar system  $\dot{x}(t) = x(t) + u(t)$ , with  $R = 1$  and  $Q = 0$ , such that its associated LQR cost function is:

$$J(x, u, t_0) = \int_{t_0}^{\infty} u^2(t) dt. \quad (4.60)$$

In this case, the optimal control action is always  $u(t) = 0$ , for  $t \in [t_0, \infty]$ . However, this control action results in the system  $\dot{x}(t) = x(t)$  which is clearly BIBO unstable. Therefore, there are conditions under which the optimization process may result in a stable or unstable closed-loop system, as shown by the following theorem.

**Theorem 4.7.** (*LQR Asymptotic Stability*) Consider the infinite-horizon Linear Quadratic Regulator as defined in Theorem 4.4. If the system is observable, then the closed-loop matrix  $\mathbf{A}_{cl} = (\mathbf{A} - \mathbf{B}\mathbf{K}) = (\mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P})$  is BIBO stable.

*Proof.* Consider, without loss of generality, that  $\mathbf{Q} = \mathbf{C}^T\mathbf{C}$ , which is still positive definite. In this case  $\mathbf{C} = \mathbf{Q}^{1/2}$ , given that  $\mathbf{Q}^{1/2}\mathbf{Q}^{1/2} = \mathbf{C}$ , which is a nonsingular matrix and the observability assumption holds. Furthermore, consider the LQR cost function given an optimal input trajectory  $\mathbf{u}^* : \mathbb{R} \rightarrow \mathbb{R}^r$ :

$$J(\mathbf{x}, \mathbf{u}^*, t_0) = \int_{t_0}^{\infty} (\mathbf{x}^T(\mathbf{C}^T\mathbf{C})\mathbf{x} + (\mathbf{u}^*)^T\mathbf{R}\mathbf{u}^*) dt = \int_{t_0}^{\infty} (\mathbf{y}^T\mathbf{y} + (\mathbf{u}^*)^T\mathbf{R}\mathbf{u}^*) dt. \quad (4.61)$$

The optimal output trajectory resulting from this control action is given by  $\mathbf{y}^*(t) = \mathbf{C}\mathbf{x}^*(t)$ , where  $\mathbf{x}^*(t)$  is the corresponding optimal state trajectory. Since  $\mathbf{u}^*(t)$  was obtained by solving the algebraic Riccati equation on Theorem 4.4, a well-known result is that the optimal cost function is bounded:

$$\mathbf{V}^*(\mathbf{x}, t_0) \leq \mathbf{x}^T(t_0)\mathbf{P}\mathbf{x}(t_0), \quad (4.62)$$

and, therefore, (4.61) must also be bounded. This is only possible if  $\mathbf{y}^*(t) \rightarrow 0$  and  $\mathbf{u}^*(t) \rightarrow 0$  as  $t \rightarrow \infty$ , which, since the system is observable, directly implies  $\mathbf{x}^*(t) \rightarrow 0$  as  $t \rightarrow \infty$ , concluding that the system is BIBO stable.  $\square$

Notice that, in contrast to the Pole-Placement method where the stability must be explicitly defined through the assignment of the eigenvalues, the LQR controller guarantees a stable closed-loop system directly, as long as the system is observable, even when the open-loop system is unstable. Since the BIBO stability (Definition 2.7) discusses the magnitude of a system response as  $t \rightarrow \infty$ , the results are derived for infinite-horizon LQR. Despite of this, it is also possible to interpret some notion of stable response for finite-horizon LQR, since the cost function is always bounded by the optimal cost as conditioned by the Hamilton-Jacobi equation.

### Robustness

Now, the discussion turns to the robustness properties of these optimal controllers. It can be shown that the LQR has a gain margin  $GM = \infty$  and a phase margin  $PM \geq 60^\circ$  [Anderson and Moore, 1990]. This assumption can be verified from a visual inspection of the Nyquist plot, as shown in Fig. 4.6. The reason for these values comes from the fact that the Nyquist diagram of an optimal regulator never crosses the unit circle centered at  $s = -1$  (implied by the Return Difference Equality). Therefore, the diagram never crosses the point  $s = -1$  no matter how much the gain  $K$  is changed, which establishes the infinite gain margin. In addition, the angle that the closest permissible point of a unit circle centered in  $s = 0$  makes with the point  $s = -1$  is exactly  $60^\circ$ , which establishes the lower bound on the gain phase.

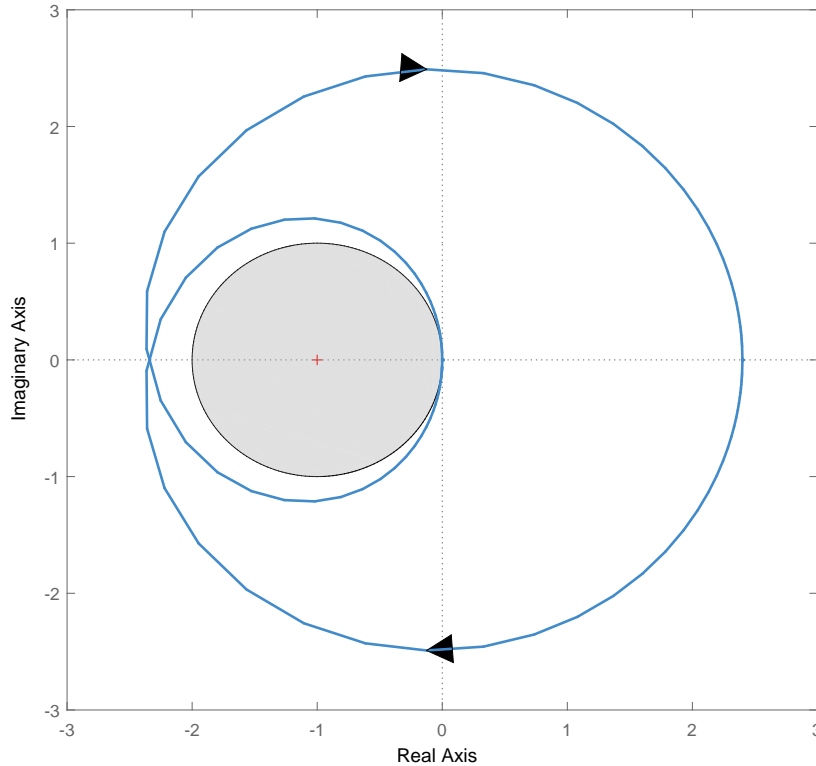


Figure 4.6: Nyquist diagram of a single state system with an open-loop unstable pole but which is closed-loop stable. The gain calculated by a LQR stabilizing controller never crosses the unit circle centered at  $s = -1 + j0$ .

Of course, there is no physical support of saying that a controller device exhibits an infinite gain margin, since there is always uncertainty about the model, the instrumentation and the environment. However, the fact that the mathematical model displays infinite gain margin allows for the assumption that the physical controller will have an extremely large margin. In fact, the uncertainty just mentioned is the main motivation into using optimal state estimators to reconstruct the state-vector from data, and it was shown by [Doyle, 1978] that an optimal controller-estimator configuration, i.e., the Linear Quadratic Gaussian, has no guaranteed margins.



# Chapter 5

## Methodology

This chapter details the methodology used in order to attest the discussion presented until now. The experiments were performed in a simulated environment, using the [MATLAB, 2018] software as the framework to compute such simulations. The results obtained are presented and discussed in the next chapter.

### 5.1 Non-Isothermal Continuous Stirred Tank

The system used for the experiments was the non-isothermal Continuous Stirred Tank Reactor (CSTR) presented by [Klatt and Engell, 1998]. This system is a realization of a class of processes that characterize a wide range of industrial applications, hence being considered a classical benchmark for reactor systems. In particular, the system proposed for the experiment is a multiple-input multiple-output (MIMO) system which is highly nonlinear, with non-minimum phase behavior and unmeasurable states, thus makes it very challenging to control. A schematic of this system is shown in Fig. 5.1.

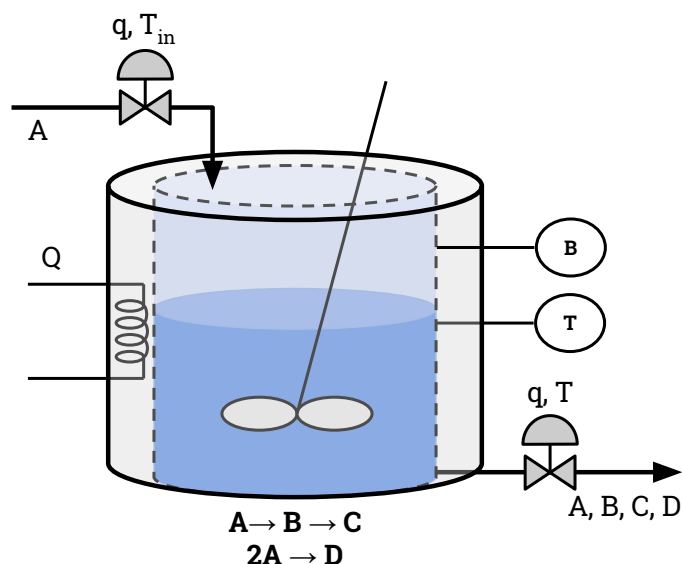
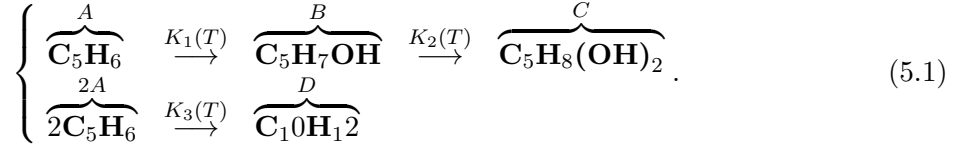


Figure 5.1: The non-isothermal continuous stirred tank reactor proposed.

This system is comprised by a tank containing a dilute solution of cyclopentadiene ( $C_5H_6$ ), which suffers reactions together with water molecules and together with the side-products of these reactions. The reaction scheme follows the same Van de Vusse scheme explored earlier in

this report:



Since the reaction is non-isothermal, the kinetics rates  $K_1(T)$ ,  $K_2(T)$  and  $K_3(T)$  are function of the temperature inside the tank, and they are assumed to follow the Arrhenius equation:

$$K_i(T) = K_{i0} e^{\frac{-E_i/R}{T+273.15}} \quad (5.2)$$

In order to control the temperature, ensuring some safety constraints to the process, a coolant system is coupled to the tank reactor to perform heat exchange by conduction. This scenario also assumes that only the concentration of chemical  $\mathbf{C}_5\mathbf{H}_7\mathbf{OH}$  and the temperature inside the tank can be measured. Additionally, the liquid inflow into the tank carries only the chemical  $\mathbf{C}_5\mathbf{H}_6$  with a concentration  $\rho_{in}^{(A)}$  and temperature  $T_{in}$ .

A nonlinear dynamical model, which was obtained from the first principles approach, is proposed to represent the evolution of the important states:

$$\left\{ \begin{array}{l} \frac{d\rho_A}{dt} = q(\rho_{in}^{(A)} - \rho_A) - (K_1(T)\rho_A + K_3(T)\rho_A^2) \\ \frac{d\rho_B}{dt} = -q\rho_B + K_1(T)\rho_A - K_2(T)\rho_B \\ \frac{dT}{dt} = q(T_{in} - T) + \frac{k_W A_r}{\rho C_p V_r} (T_C - T) \\ \quad - \frac{1}{\rho C_p} (K_1(T)\rho_A \Delta H_{AB} + K_2(T)\rho_B \Delta H_{BC} + K_3(T)\rho_A^2 \Delta H_{AC}) \\ \frac{dT_C}{dt} = \frac{1}{m_K C_{pK}} Q + \frac{k_W A_r}{m_K C_{pK}} (T - T_C) \end{array} \right. \quad (5.3)$$

The meaning and value of each system variable is shown in Table 5.1. Additionally, the real system has some constraints in respect to the manipulated variables, namely:

$$5hr^{-1} \leq F \leq 35hr^{-1}; \quad -8500 \frac{kJ}{h} \leq Q \leq 0 \frac{kJ}{h}. \quad (5.4)$$

## 5.2 Simulation

The process and all the control structures were simulated using the MATLAB software. In the first part of the experiments, the dynamics of the process were analyzed in respect to the mathematical structure of the models and through the simulation results.

In respect to the control, the second part of the experiments focused on the design of controllers using the optimal control techniques. The controllers were developed and analyzed in the same order that they were presented in this document:

1. Linear Quadratic Regulator, assuming a nonzero initial state and perfect measurement of the states;
2. Linear Quadratic Regulator with Integral Action, for a non-constant reference signals and assuming perfect measurement of the states;
3. Linear Quadratic Gaussian, assuming a nonzero initial state and reconstruction of the state-vector through observations from a stochastic linear system;

System Variable	Symbol	Value	Unit
Concentration $A$ in the reactor solution	$\rho_A$	—	mol/l
Concentration $B$ in the reactor solution	$\rho_B$	—	mol/l
Temperature inside the reactor	$T$	—	$^{\circ}C$
Temperature of the coolant jacket	$T_C$	—	$^{\circ}C$
Inflow concentration of $A$	$\rho_{in}^{(A)}$	5.1	mol/l
Inflow liquid temperature	$T_{in}$	130	$^{\circ}C$
Nominal kinetic rate for $A \rightarrow B$	$K_{10}$	$(1.287) \times 10^{12}$	$h^{-1}$
Nominal kinetic rate for $B \rightarrow D$	$K_{20}$	$(1.287) \times 10^{12}$	$h^{-1}$
Nominal kinetic rate for $2A \rightarrow D$	$K_{30}$	$(9.043) \times 10^9$	$l \text{ mol}^{-1}h$
Activation energy for $A \rightarrow B$	$E_1/R$	9758.3	-
Activation energy for $B \rightarrow C$	$E_2/R$	9758.3	-
Activation energy for $2A \rightarrow D$	$E_3/R$	8560.0	-
Enthalpy for $A \rightarrow B$	$\Delta H_{AB}$	4.2	$kJ/mol$
Enthalpy for $B \rightarrow C$	$\Delta H_{BC}$	-11.0	$kJ/mol$
Enthalpy for $2A \rightarrow D$	$\Delta H_{AD}$	-41.85	$kJ/mol$
Total density of the reactor solution	$\varrho$	$(0.9342) \times 10^{-4}$	$kg/l$
Heat capacity of the solution	$C_p$	3.01	$kJ \text{ kg}^{-1}K$
Heat capacity of the coolant	$C_{pK}$	2.0	$kJ \text{ kg}^{-1}K^{-1}$
Surface area of the tank reactor	$A_R$	0.215	$m^2$
Volume of the tank reactor	$V_R$	10.01	$m^3$
Coolant mass	$m_K$	5.0	$kg$
Heat transfer coefficient	$k_W$	4032.0	$kJ \text{ h}^{-1}m^2K$

Table 5.1: Model parameters and variables nomenclature.

4. Linear Quadratic Gaussian with Integral Action, for a non-constant reference and assuming a reconstruction of the state-vector through observations from a stochastic linear system;

The simulation results are mainly analyzed through graphical visualizations of the data obtained. The experiments were focused on validating the modeling and control procedures for the system presented, instead of comparing and ranking the controllers based on the most performing one. For the sake of completeness, however, some metrics are used to evaluate each controller performance and relate them to the visualizations. The metrics used were the *Integral of Squared Error* (ISE), *Integral of Absolute Error* (IAE), *Integral of Time-Weighted Squared Error* (ITSE) and *Integral of Time-Weighted Absolute Error* (ITAE), which are given as:

$$\begin{aligned}
 \text{IAE} &= \int_{t_0}^T |\mathbf{r}(\tau) - \mathbf{y}(\tau)| d\tau; & \text{ITAE} &= \int_{t_0}^T \tau |\mathbf{r}(\tau) - \mathbf{y}(\tau)| d\tau; \\
 \text{ISE} &= \int_{t_0}^T [\mathbf{r}(\tau) - \mathbf{y}(\tau)]^2 d\tau; & \text{ITSE} &= \int_{t_0}^T \tau [\mathbf{r}(\tau) - \mathbf{y}(\tau)]^2 d\tau.
 \end{aligned} \tag{5.5}$$



## Chapter 6

# Results and Discussion

This chapter presents the results of the experiments described in the last chapter. The sections starts by presenting the simulation of the dynamical model and analyzing its properties. After that, several optimal controllers are formulated and simulated in closed-loop conditions, both for regulation and tracking. The control simulations are grouped accordingly to the control objective that they impose to the system.

### 6.1 Dynamical Analysis

#### Nonlinear model simulation

The behavior of the process is accessed directly by simulating the nonlinear model presented at (5.3). The main goal is to observe how the system responds to variations in each of its manipulated variables. In the first simulation, shown in Fig. 6.1, the process is started at the steady-state condition, and then a pulse is applied to the liquid inflow to the tank.

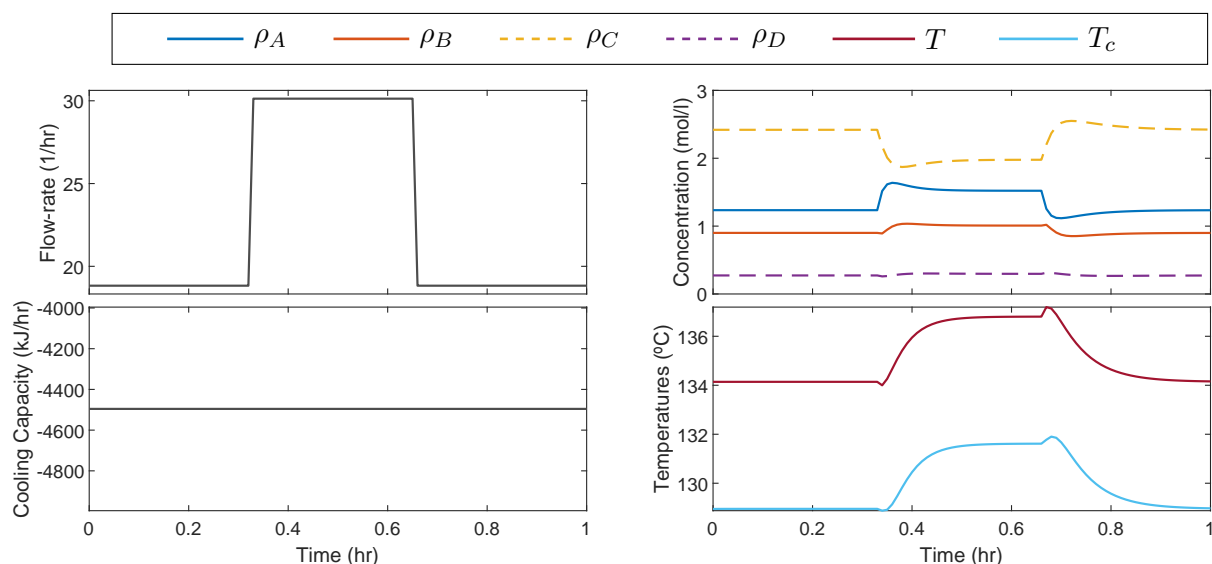


Figure 6.1: Non-linear simulation showing the manipulated variables (left) and correspondent response of the system variables (right) for a pulse change in the input flow-rate.

The simulation shows that both compounds  $A$  and  $B$  grows together when there is an increase in the input flow-rate, in which case the first exhibits a greater response than the other. Of course, this conclusion is direct from the fact that the input carries concentrations of this very same chemical. In respect to the temperatures, notice that, since the process is non-isothermal,



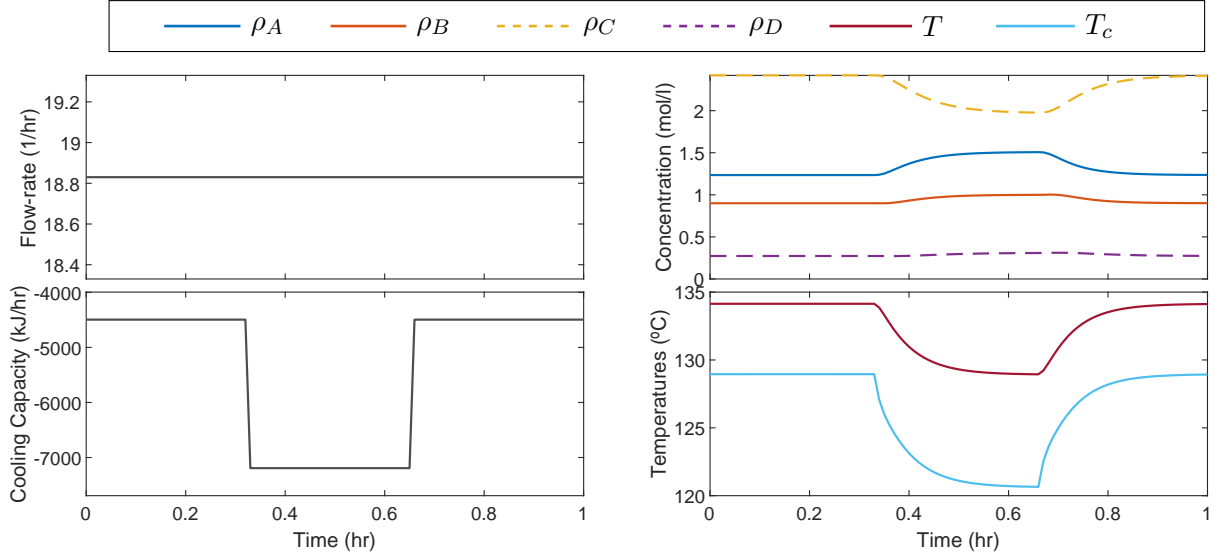


Figure 6.2: Non-linear simulation showing the manipulated variables (left) and correspondent response of the system variables (right) for a pulse change in the cooling capacity.

the temperatures raised together with the increase in the reagents concentrations. Given that the temperature of the liquid inflow is actually lower than the tank temperature at the steady-state, it is possible to relate the increase in the temperature directly to the occurrence of the chemical reactions. This fact characterizes the process as an exothermic process, as it can be also noted from the first-principles model.

The second simulation was performed similarly, but applying a negative pulse to the jacket cooling capacity. The result is shown in Fig. 6.2. As expected, the change in the cooling capacity affects the temperatures inside the reactor tank and inside the cooling jacket. The visualization also shows that the change in the temperature slightly changed the concentration of the substances inside the reactor. Therefore, it is possible to consider an indirect effect of this manipulated variable to these system variables.

Now, the goal was to visualize if it is possible to increase the concentration of the chemicals through the reactions, while maintaining the tank temperature close to the steady-state value. For this reason, the last simulation considers a pulse to increase the input flow-rate while another pulse increases the magnitude of the cooling capacity. The result of this simulation is shown at Fig. 6.3. Note that the tank temperature remained restricted to values closer to the steady-state point, whereas the concentrations were increased by the flow-rate input. Thus, the possibility of increasing the production of a chemical without breaking safety constraints holds.

### Linearized model realization and simulation

The discussion throughout the document focused on control design methods for linear systems. For this reason, the nonlinear model can not directly be used to develop the controllers, and a linearization (Theorem 2.2) of this model is required. First of all, consider the following choice for the state, input and output variables:

$$\mathbf{x}(t) = \begin{bmatrix} \rho_A(t) \\ \rho_B(t) \\ T(t) \\ T_c(t) \end{bmatrix}; \quad \mathbf{u}(t) = \begin{bmatrix} q(t) \\ Q(t) \end{bmatrix}; \quad \mathbf{y}(t) = \begin{bmatrix} \rho_B(t) \\ T(t) \end{bmatrix}. \quad (6.1)$$

By evaluating the derivative of each differential state equation in relation to each state

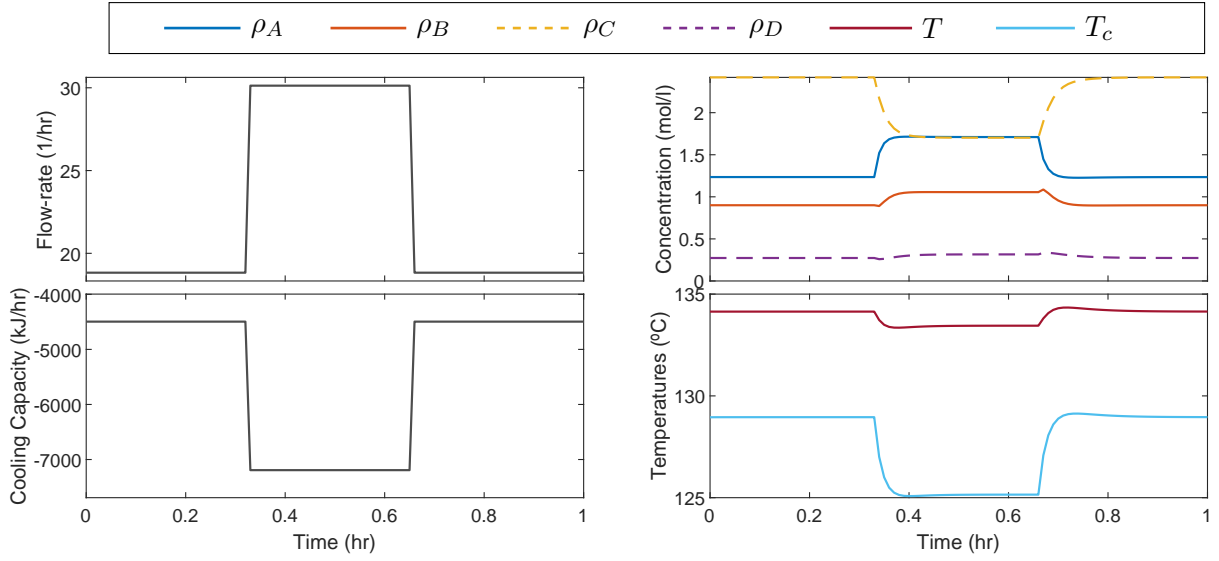


Figure 6.3: Non-linear simulation showing the manipulated variables (left) and correspondent response of the system variables (right) for a change in the input flow-rate and cooling capacity.

variable, the obtained Jacobian is

$$\mathbf{A} = \begin{bmatrix} -u_{1o} - K_1(x_{3o}) - 2K_3(x_{3o})x_{1o} & 0 & \frac{K_1(x_{3o})}{dx_{3o}}x_{1e} + \frac{K_3(x_{3o})}{dx_{3o}}x_{1o}^2 & 0 \\ K_1(x_{3o}) & -u_{1o} - K_2(x_{3o}) & \frac{K_1(x_{3o})}{dx_{3o}}x_{1e} \frac{K_2(x_{3o})}{dx_{3o}}x_{2e} & 0 \\ \frac{\partial H(x_{1o}, x_{2o}, x_{3o})}{\partial x_{1o}} & \frac{\partial H(x_{1o}, x_{2o}, x_{3o})}{\partial x_{2o}} & -u_{1e} - \alpha + \frac{\partial H(x_{1o}, x_{2o}, x_{3o})}{\partial x_{3o}} & \alpha \\ 0 & 0 & \beta & -\beta \end{bmatrix}, \quad (6.2)$$

where the parameters  $\alpha$  and  $\beta$  are given by

$$\alpha = \frac{k_W A_r}{\rho C_p V_r}; \quad \beta = \frac{k_W A_r}{m_K C_{pK}}, \quad (6.3)$$

and the derivatives are calculated as

$$\begin{aligned} \frac{dK_i}{dx_{3o}} &= \frac{E_1/R}{(x_{3o} + 273.15)^2} K_{i0} e^{\frac{-E_i/R}{x_{3o} + 273.15}}; \\ \frac{\partial H}{\partial x_{1o}} &= \frac{-1}{m_K C_{pK}} (K_1(x_{3o}) \Delta H_{AB} + 2K_3(x_{3o}) x_{1o} \Delta H_{AD}); \\ \frac{\partial H}{\partial x_{2o}} &= \frac{-1}{m_K C_{pK}} (K_2(x_{3o}) \Delta H_{BC}); \\ \frac{\partial H}{\partial x_{3o}} &= \frac{-1}{m_K C_{pK}} \left( \frac{dK_1}{dx_{3o}} x_{1o} \Delta H_{AB} + \frac{dK_2}{dx_{3o}} x_{2o} \Delta H_{BC} + \frac{dK_3}{dx_{3o}} x_{1o}^2 \Delta H_{AC} \right). \end{aligned} \quad (6.4)$$

Similarly, the derivative of each differential equation in relation to each input variable is

evaluated, which results in the simpler Jacobian

$$\mathbf{B} = \begin{bmatrix} \rho_{in}^{(A)} - x_{1o} & 0 \\ -x_{2o} & 0 \\ T_{in} - x_{3o} & 0 \\ 0 & \frac{1}{m_K C_{pK}} \end{bmatrix}. \quad (6.5)$$

Regarding the matrices  $\mathbf{C}$  and  $\mathbf{D}$ , recall that only the variables  $\rho_B$  and  $T$  are measured in the process. The measurement, however, is assumed to be perfect, so that these matrices can be defined as

$$\mathbf{C} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}; \quad \mathbf{D} = \mathbf{0}. \quad (6.6)$$

The linear model presented in (6.2) is an approximation of the nonlinear model for any point in the state space, given that it is a steady-state point. To obtain feasible linearization points, the differential equation system is evaluated for all points where  $\dot{\rho}_A = \dot{\rho}_B = \dot{T} = \dot{T}_c = 0$ . To solve this system, the input variables are selected to be the interval defined by the physical variable restrictions. The result is shown in Fig. 6.4 for the two measured variables. In this work, it is assumed the steady-state point given by  $\mathbf{x}_o = [1.235, 0.90, 134.14, 128.95]^T$  and  $\mathbf{x}_o = [18.83, -4495.7]^T$ , indicated in the figure by a white cross.

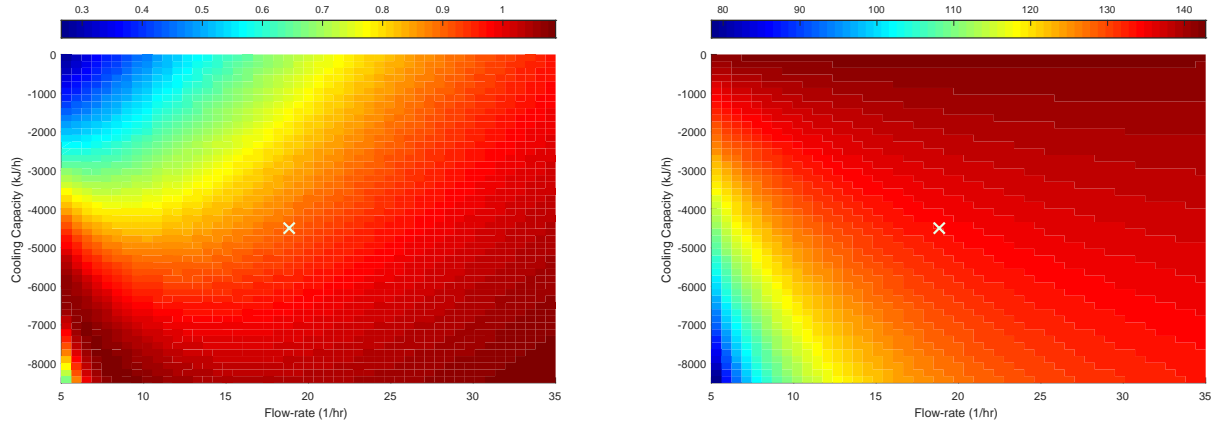


Figure 6.4: Visualization of the steady-state values for the linearized model.

Using the steady-state values just mentioned, a realization for the linearized model is obtained:

$$\left\{ \begin{array}{l} \mathbf{A} = \begin{bmatrix} -86.0962 & 0 & -4.2010 & 0 \\ 50.6146 & -69.4446 & 0.9958 & 0 \\ 172.2263 & 197.9985 & -36.7597 & 30.7977 \\ 0 & 0 & 86.6880 & -86.6880 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} 3.865 & 0 \\ -0.900 & 0 \\ -4.140 & 0 \\ 0 & 0.100 \end{bmatrix} \\ \mathbf{C} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \quad \mathbf{D} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \end{array} \right. \quad (6.7)$$

In order to demonstrate that the linear model is a valid approximation of the nonlinear model, both systems were simulated in parallel for the same input signal. The simulation, as shown in Fig. 6.5, starts at the steady-state values and then apply a sinusoidal change in the input to drag the system away from this initial point. As expected, the results shows that the linearized model is an accurate representation of the system, except when it is far from the steady-state values.

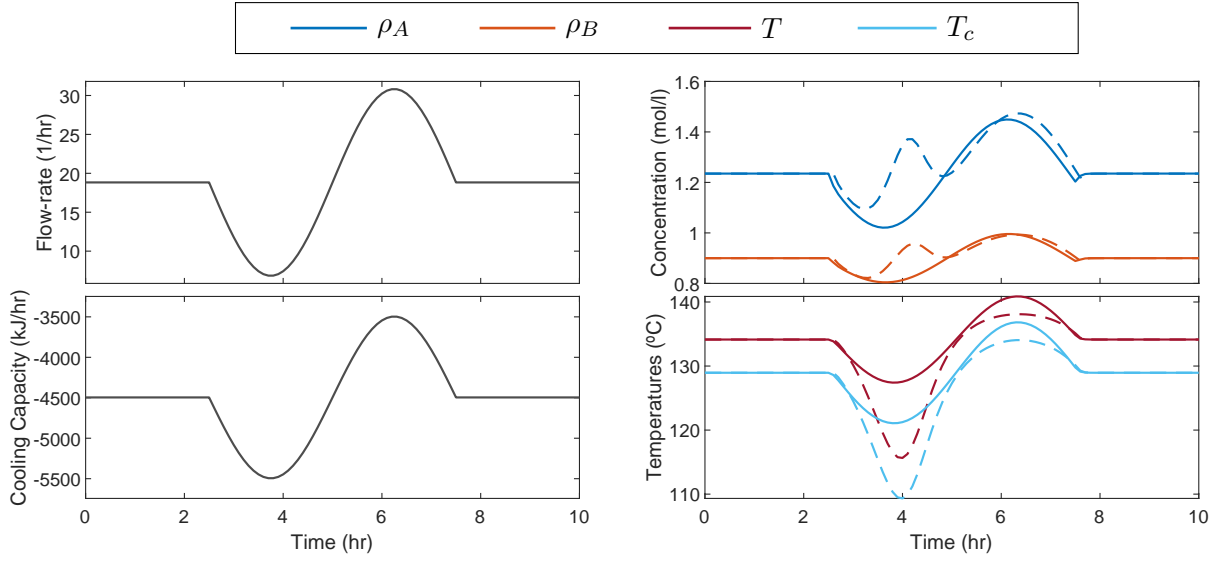


Figure 6.5: Simulation of both the nonlinear and linearized State-Space models. The nonlinear response is represented by the dashed lines.

### Linearized model properties

Considering a State-Space realization, it is possible to access the system properties. First of all, the BIBO stability of the system was determined through the poles location in the complex plane. In the case of (6.7), the eigenvalues of  $\mathbf{A}$  were calculated as  $\boldsymbol{\lambda} = [-16.79, -54.84, -86.33, -121.01]$ , which are all in right-hand side of the plane, i.e.,  $\text{Re}[\lambda_i], \forall i \in [1, 4]$ . It is possible to conclude that, in the light of this model, the system is stable.

The next property to be accessed consists in the controllability of the system. Using the definition, the controllability matrix of this model was calculated as:

$$\mathbf{C} = \begin{bmatrix} 3.8650 & 0 & -315.3696 & 0 & 24465 & -12.9 & 2271 \\ -0.9000 & 0 & 254.0027 & 0 & -32964 & 3.1 & -1250 \\ -4.1400 & 0 & 639.6409 & 3.0798 & -38589 & -380.2 & 43720 \\ 0 & 0.1000 & -358.8883 & -8.6688 & 86561 & 1018.5 & -121250 \end{bmatrix}. \quad (6.8)$$

This matrix has a row rank equal to 4, which is exactly the number of state variables. Therefore, it is concluded that this system is indeed controllable, meaning that it is possible to drive the system to any position in the space using a specific input signal. It is worth mentioning that this matrix exhibits some entries with high values, which can become numerically unstable in the case of augmenting the system.

Moreover, the last property of the system to be accessed concerns the observability of the system. Similarly to the controllability case, the observability matrix was calculated as:

$$\mathbf{O} = \begin{bmatrix} 0 & 1.0000 & 0 & 0 \\ 0 & 0 & 1.0000 & 0 \\ 50.6146 & -69.4446 & 0.9958 & 0 \\ 172.2263 & 197.9985 & -36.7597 & 30.7977 \\ -7701.1 & 5019.7 & -318.4 & 30.7 \\ -11137 & -21028 & 3495 & -3802 \\ 862270 & -411630 & 51710 & -1246 \\ 496400 & 2152200 & -432200 & 437200 \end{bmatrix}. \quad (6.9)$$

This matrix has a column rank equal to 4, which, again, is exactly the number of state variables. Thus, it is possible to conclude that the system is observable, meaning that the information of any state can be reconstructed given only the outputs.

## 6.2 Optimal Control

### Regulation

The results for the controllers arrangements are now presented. The first controller simulated was the Linear Quadratic Regulator (LQR), as shown in Fig. 6.6. The simulation started with initial condition  $x_1(0) = x_2(0) = 0$  (or, equivalently,  $\Delta x_1(0) = -x_{1o}$  and  $\Delta x_2(0) = -x_{2o}$ ) while maintaining the states  $x_3$  and  $x_4$  at the steady-state value. The goal of this simulation was to represent the scenario where the solution in the tank does not actually contain any concentration of the reagents or products. This scenario, thus, requires that the controller “starts” the process by providing the flow of reagent necessary to achieve the desirable steady-state concentration of the controlled product. The weights used, regarding the LQR cost function, were:

$$\mathbf{Q} = \begin{bmatrix} 0.25 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0.01 \end{bmatrix} ; \quad \mathbf{R} = \begin{bmatrix} 0.25 & 0 \\ 0 & 1.1e-5 \end{bmatrix}. \quad (6.10)$$

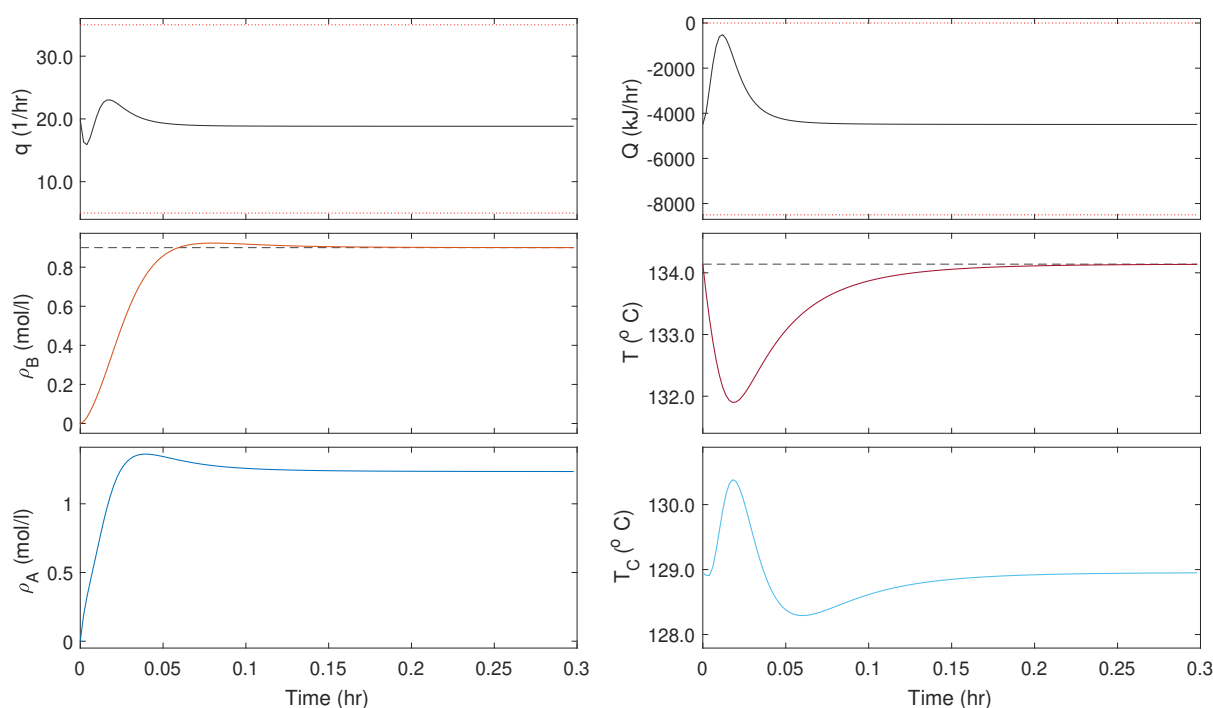


Figure 6.6: Simulation of the  $LQR_1$  controller for “starting” the tank reactor to the steady-state point. The reference is indicated by the black dashed lines, and the manipulated variables restrictions are indicated by the red dotted lines.

Notice that the resulting reactions will change the temperature inside the tank, requiring also that the controller manipulates the cooling capacity in order to control the reactor temperature. As expected, the controller momentarily increases the input flow-rate in order to accelerate the production of compound  $B$ , returning to the steady-state value to achieve the desired reference. Notice that the input flow-rate does not grow further in the first minutes in order to avoid overshoots, which occurs in the results presented but with small magnitude. In respect to the reactor temperature, the controller was able to respond to the sudden change by decreasing the absolute value of the cooling capacity, causing the coolant to increase its temperature and, consequently, returns the reactor temperature to the reference. Therefore, the simulation shows that the proposed optimal controller is able to drive the process to a desired constant reference in an optimal manner, not causing overshoots and/or oscillations, in a small amount of time.

Further testing the capabilities of this controller, the same experiment was considered, but while subjecting the system to disturbances in some system variables. Specifically, disturbances were applied to the inflow concentration and temperature,  $\rho_{in}^{(A)}$  and  $T_{in}$ , respectively, to represent the scenario where the fluid entering the tank is perturbed by previous processes or by instrumentation limitations. To simulate these cases, pulse signals of values  $w_1 = \pm 0.5$  and  $w_2 = \pm 2$  were applied to the nominal value of the inflow concentration and temperature, respectively. Each pulse starts at a time value that is an odd multiple of  $\lfloor T/4 \rfloor$  and lasts for  $0.02T$  hours, where  $T$  is the total size of the control horizon. The results are shown in Fig. 6.7.

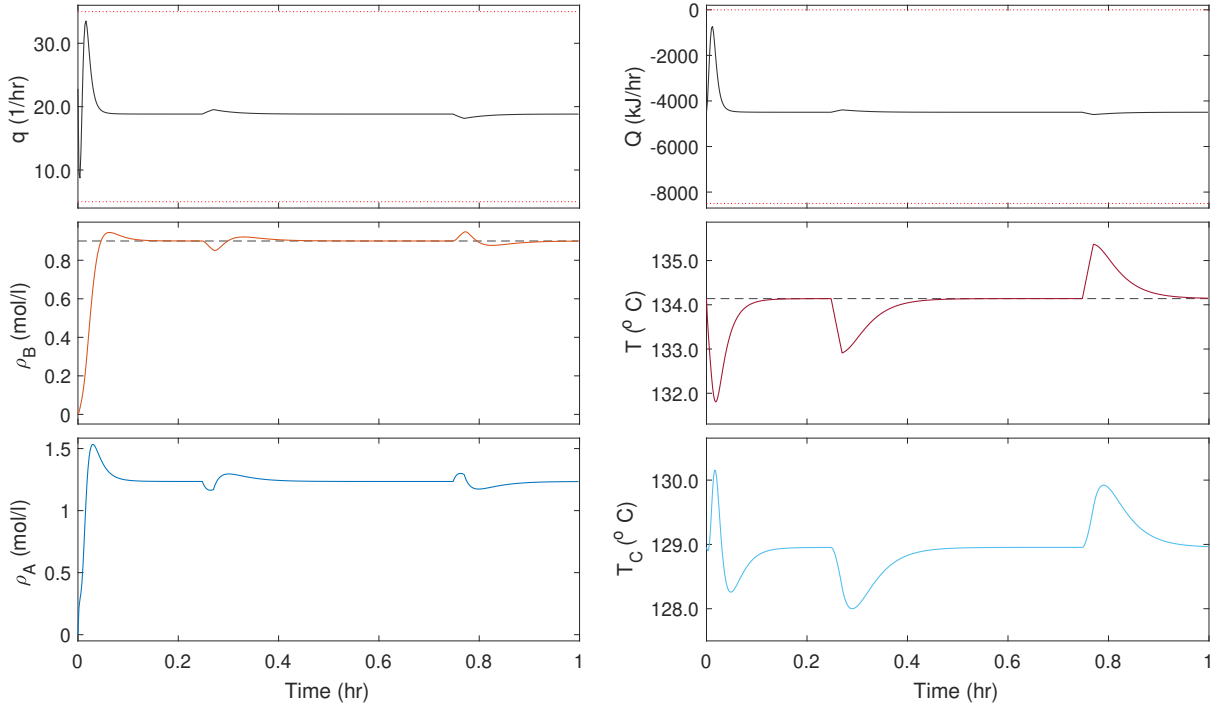


Figure 6.7: Simulation of the  $LQR_2$  controller for “starting” the tank reactor to the steady-state point, but subjected to two disturbances pulses in the system variables.

First of all, it was observed that the control horizon had to be increased in order to allow the feedback controller to perceive the disturbance and corrects the system response accordingly. As a direct result of the larger horizon, the controller performed more aggressively at the first minutes. Despite of this, the input signals did not violated the defined constraints for the manipulated variables. It was observed that the regulator is capable of react to the disturbances and provide a corrective action to return to the steady-state point.

Until now, the experiments were done assuming that the optimal controller had direct access to the real system states. In a more realistic scenario, the states can only be reconstructed through noisy measurements. A Linear Quadratic Gaussian (LQG) controller is proposed to deal with this limitation, considering the same control objectives of the previous experiments and using the same weighting matrices as in (6.10). In the experiments it was assumed a process noise  $\mathbf{w} : \mathbb{R} \rightarrow \mathbb{R}^4$  and measurement noise  $\mathbf{z} : \mathbb{R} \rightarrow \mathbb{R}^2$  described by the Gaussian distributions:

$$\mathbf{w}(t) \sim \mathcal{N} \left( \begin{bmatrix} x_1(t) \\ x_2(t) \\ x_3(t) \\ x_4(t) \end{bmatrix} \middle| \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0.1 & 0 & 0 & 0 \\ 0 & 0.1 & 0 & 0 \\ 0 & 0 & 0.1 & 0 \\ 0 & 0 & 0 & 0.1 \end{bmatrix} \right) \quad \mathbf{z}(t) \sim \mathcal{N} \left( \begin{bmatrix} u_1(t) \\ u_2(t) \end{bmatrix} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0.0001 & 0 \\ 0 & 0.1 \end{bmatrix} \right) \quad (6.11)$$

The covariance matrix of the noise  $\mathbf{w}(t)$  tries to represent the uncertainty of the non-linear model behavior in respect to the linearized model, while the covariance matrix of  $\mathbf{z}(t)$  aims

to represent realistic limitation of measurement devices in respect to these variables. It is assumed that the Kalman-Bucy estimator can estimate these covariances fairly accurately, before optimizing each estimator gain  $\mathbf{K}_e(t)$ . The simulation results, for the scenario of starting the process, is shown in Fig. 6.8.

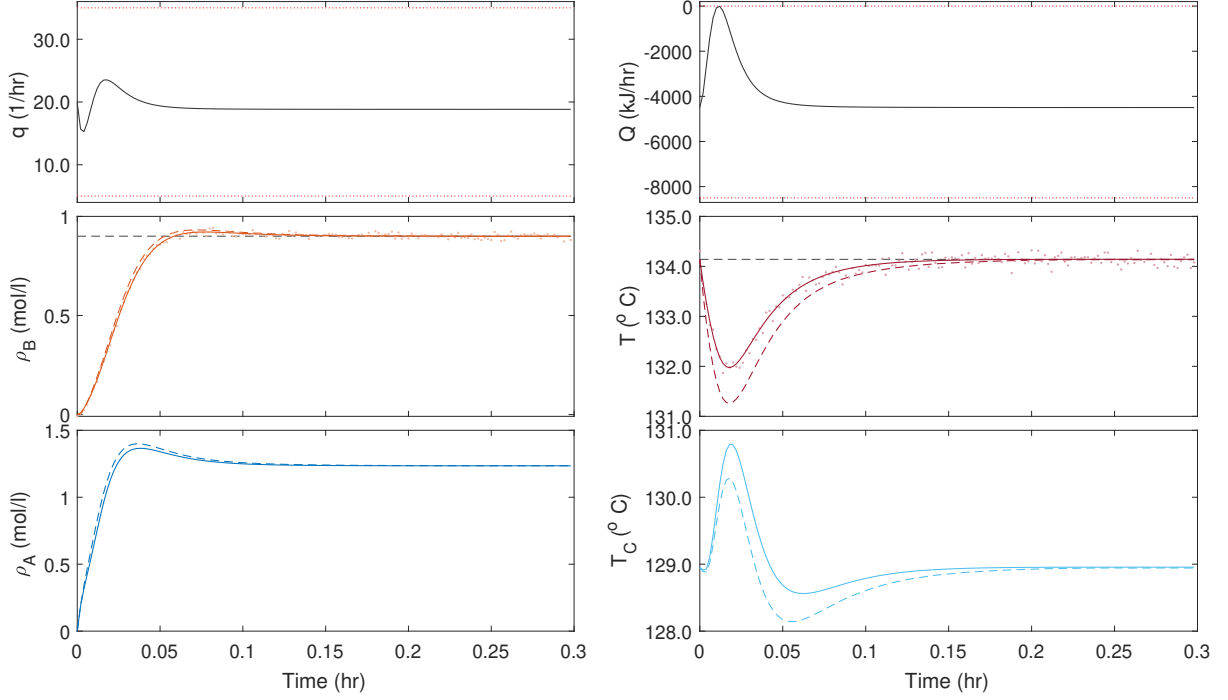


Figure 6.8: Simulation of the  $LQG_1$  controller for “starting” the tank reactor to the steady-state point. The dots represents the measurements of the output variables, and the colored dashed lines represents the state estimations.

The estimation of the covariance matrices gave the estimator enough information to accurately reconstruct the state vector from the observations. As it can be seen in the visualizations, this control-estimation operation is very similar to the one shown in Fig. 6.6, when the states were directly known. Thus, the LQG controller proved to be optimal for this more realistic situation. Furthermore, the same controller was also evaluated for the case where disturbances in the system variables are present. In this experiment, the same disturbance signal used in the LQR simulation is provided to the system. The results are shown in Fig. 6.9. The system responses are very similar. However, the corrective actions for the disturbances were smoother and are barely visible in the scale of the visualization. Finally, a summary of the results in this subsection are displayed at Table 6.1.

Controller	ISE	IAE	ITSE	ITAE
$LQR_1$	7.2678	12.2719	0.0015	0.2592
	82.9199	60.6988	0.1045	2.8347
$LQG_1$	7.2765	12.0436	0.0015	0.2299
	67.6334	48.4748	0.0596	1.7535
$LQR_2$	7.2713	14.7600	0.0270	1.9614
	137.0858	134.5661	24.5280	51.1356
$LQG_2$	7.3481	15.1579	0.0270	2.0193
	155.0745	160.7632	32.2785	65.5359

Table 6.1: Comparison between the LQR and LQG controllers obtained. The pair  $(LQR_2, LQG_2)$  comprises the simulations in which the system was subjected to disturbances.

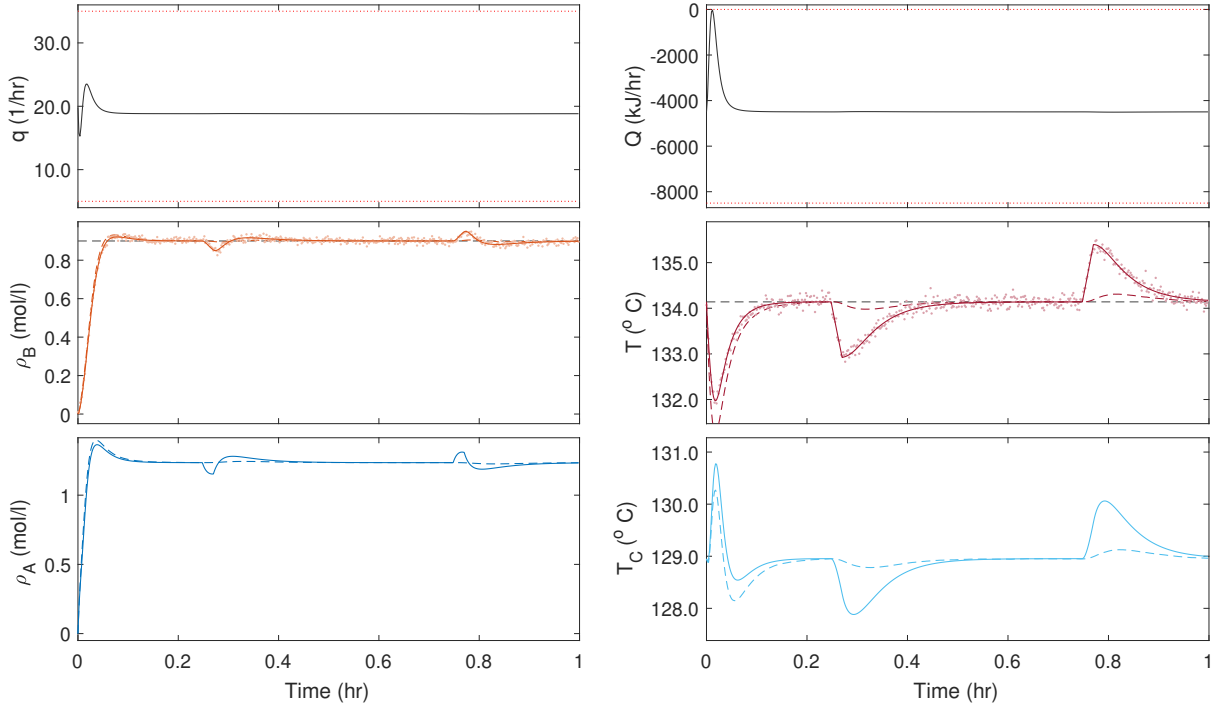


Figure 6.9: Simulation of the  $LQG_2$  controller for “starting” the tank reactor to the steady-state point, but subjected to two disturbances pulses in the system variables.

### Tracking

Now, the experiments for controllers that track non-constant references is presented. As in the previous subsection, the first controller simulated consists in the Linear Quadratic Regulator with Integral Action (LQRI), also known as Linear Quadratic Servo (LQ-Servo). The simulations starts with initial states equal to the steady-state values. The weights used were the same from (6.10), but including the weights for the augmented states:

$$Q = \begin{bmatrix} 0.25 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0 & 0 & 10 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.5 \end{bmatrix} ; \quad R = \begin{bmatrix} 0.16 & 0 \\ 0 & 1e-6 \end{bmatrix}. \quad (6.12)$$

The proposed scenario for this controller is to increase and then decrease the amount of compound  $B$  that is produced, while maintaining the reactor temperature at a constant value. For the concentration control, the reference signal is comprised by a sequence of step functions which varies from approximately 12% of the operation value, which is a realistic set-point change in such processes. The results of the simulation are shown in Fig. 6.10.

After the introduction of the integral action, two more weights were needed to be selected for the controller optimization. The MIMO case implies a trade-off between reaching the reference for the concentration and for the temperature, in the case that they do not represent a valid steady-state. Therefore, the choices of weights to the integral action has to take in account the fact that an aggressive set-point tracking may result in oscillatory, or even unstable, systems. To avoid these situations, the weights were selected to cause a smooth tracking, and the control horizon was expanded to allow the controller to operate for more hours.

Despite the mentioned difficulties, the controller was able to drive the system as close as possible to the reference, while still obeying the manipulated variables constraints. Notice



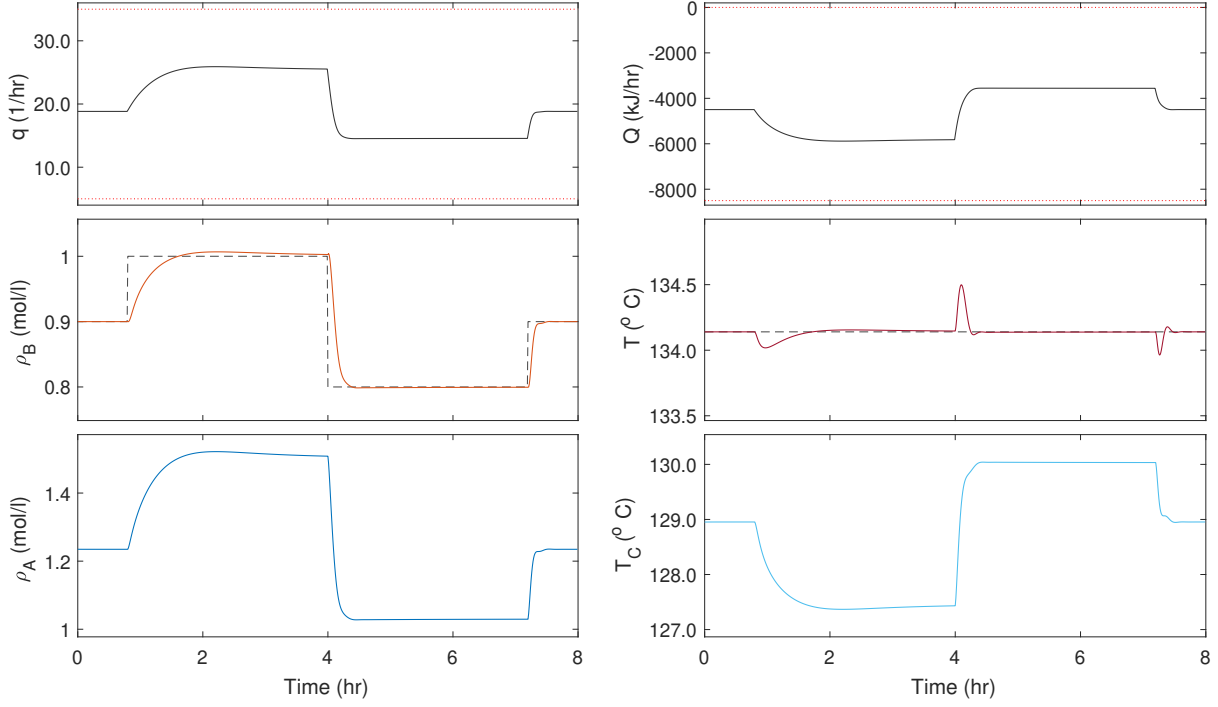


Figure 6.10: Simulation of the  $LQRI_1$  for reference tracking.

that, in the moments where the reference changes, the change in the concentration resulted in a momentarily change in the temperature, which is promptly corrected by the feedback action of the controller. This change, however, is very small and the temperature remained within an interval of  $0.5^\circ\text{C}$  around the reference signal, which is an excellent result. Another nice response observed in the visualization comprises in the change in the coolant temperature, which follows closely the shape of the change in the cooling capacity. In the last simulations, this variable evolved in a similar manner to the reactor temperature, but, in this case, the first was directly manipulated to allow for the latter to stay at the steady-state value.

Following the same case of the previous subsection, the same controller was used to simulate a closed-loop operation in the case when disturbances to the real system are occurring. In this case, a pulse signal is applied to the same variables as in the last case, but for each time instant multiple of  $\lceil T/4 \rceil$  with a duration of  $0.02T$ , which is a longer step in this case since  $T$  is greater. The results are presented at Fig. 6.11. The first interesting phenomena to notice is that the system responded more aggressively to the disturbances that occurred between the time interval  $t \in [0.8, 4]$  for which  $r(t) = 1$ , when compared to the disturbances occurring in the interval  $t \in [4, 7.2]$  for which  $r(t) = 0.8$ . This is expected from the fact that the system has non-linear terms associated with the concentration  $\rho_A$ , which is greater at the first interval. The controller, however, was able to correct the disturbances and drive the system back to the reference.

In the direction for a more realistic scenario, the same reference tracking just proposed was used together with a Linear Quadratic Gaussian, augmented with integral action, to perform the feedback over the estimated states. In this case, a different weighting for the augmented states were modified to still result in a similar operation to the one in Fig. 6.6 and 6.7.

$$Q = \begin{bmatrix} 0.25 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0 & 0 & 100 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0.1 \end{bmatrix} ; \quad R = \begin{bmatrix} 1 & 0 \\ 0 & 1.1e-5 \end{bmatrix}. \quad (6.13)$$

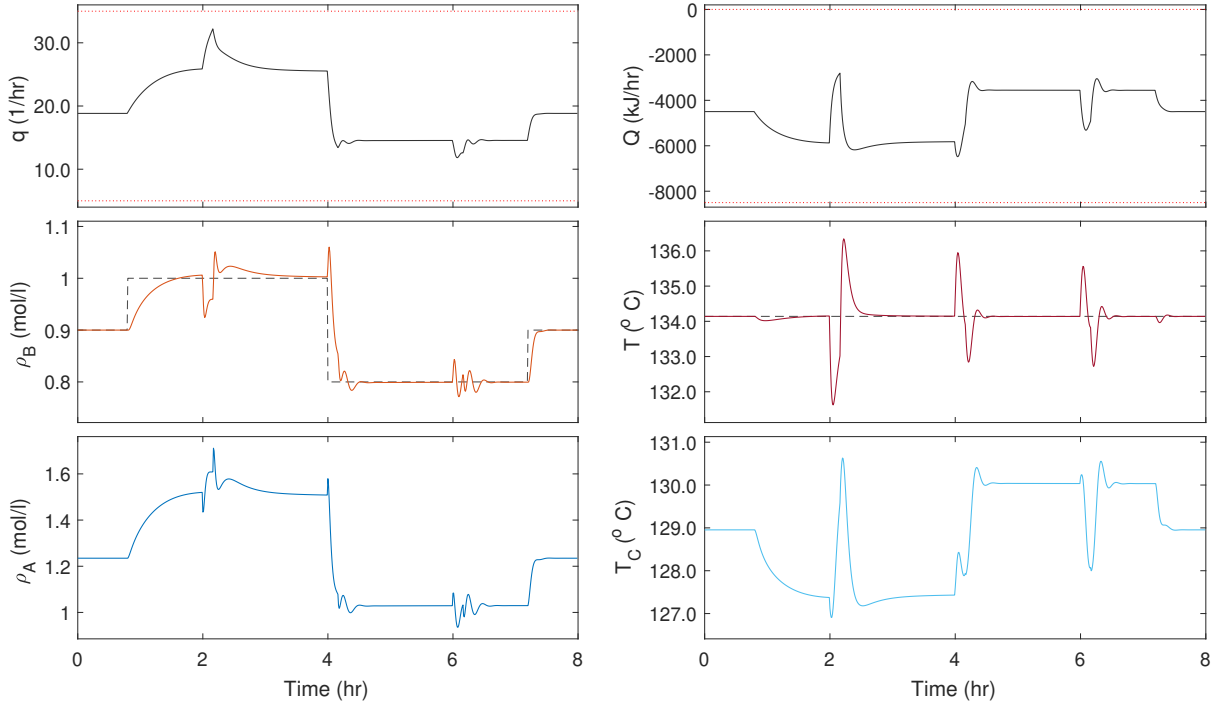


Figure 6.11: Simulation of the  $LQRI_2$  controller for reference tracking, but subjected to three pulse of disturbances in the system variables.

In this case, the process and measurement noises used,  $\mathbf{w} : \mathbb{R} \rightarrow \mathbb{R}^4$  and  $\mathbf{z} : \mathbb{R} \rightarrow \mathbb{R}^2$ , respectively, are described by the Gaussian distributions:

$$\mathbf{w}(t) \sim \mathcal{N} \left( \begin{bmatrix} x_1(t) \\ x_2(t) \\ x_3(t) \\ x_4(t) \end{bmatrix} \middle| \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0.4 & 0 & 0 & 0 \\ 0 & 0.4 & 0 & 0 \\ 0 & 0 & 0.5 & 0 \\ 0 & 0 & 0 & 0.5 \end{bmatrix} \right) \quad \mathbf{z}(t) \sim \mathcal{N} \left( \begin{bmatrix} u_1(t) \\ u_2(t) \end{bmatrix} \middle| \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0.0001 & 0 \\ 0 & 0.01 \end{bmatrix} \right). \quad (6.14)$$

The greater values for the covariance  $\mathbf{Q}_k$  of the covariance noise reflects the fact that the estimator becomes more uncertain about the system when it follows a reference signal far from the steady-state value. The resulting operation is displayed at Fig. 6.12.

The first difference noted was the fact that the weights for the augmented states for the LQG controller actually needed a value of greater magnitude than those used for the LQR case to be able to yield the same results. There are two important phenomena to emphasize in this result. First of all, the input signals requested from the controller are not smooth, resulting in an oscillatory behavior in the actual response of the system (and in the Kalman-Bucy filter state estimations). This is resulting from the fact, as noted in Example 4.1, that the integral action is calculated using the measurements of the system output, which are noisy. In addition to this, the Kalman-Bucy filter was not able to effectively reconstruct the real states in the cases for the references that are far from the linearization steady-state point. It is visible that the estimates have an offset when the system reaches the steady-state after tracking those references. In contrast, the estimate is far accurate when the real response is close to the linearization steady-state point, as evidenced in the temperature responses. This phenomena is resulting from the fact that the Kalman-Bucy filter is modeled from a linearized realization of the non-linear system, exhibiting the same limitations as were observed in Fig. 2.3.

The two problems just mentioned are somehow connected. A possible solution for the oscillatory input could be to use the state estimates and the linearized State-Space model to produce an output estimate; and then use it to calculate the integral action. Since the matrix  $\mathbf{Q}_k$ ,

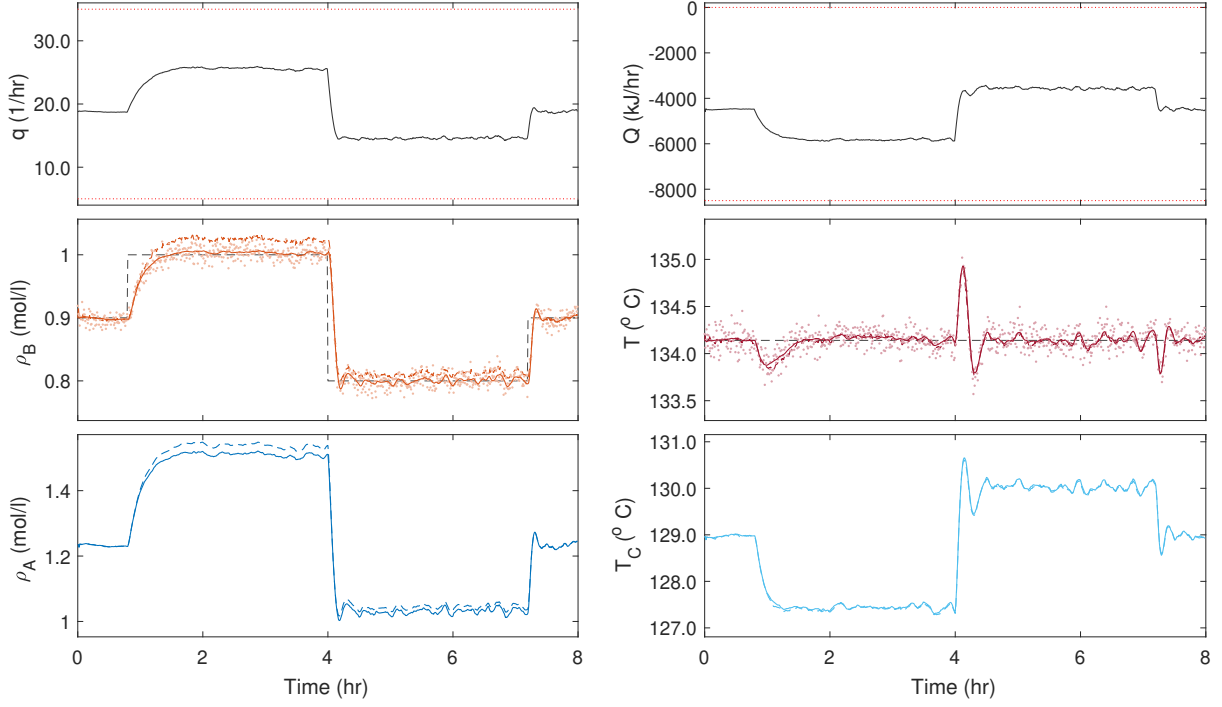


Figure 6.12: Simulation of the  $LQGI_1$  controller for reference tracking.

estimated from the noise covariances, can be underestimated so that the Kalman-Bucy model is considered less uncertain and the estimations are less oscillatory. However, the limitation of the linear Kalman-Bucy filter will result in a poor reference tracking, in respect to the non-linear system response, if the estimated states are used in the integral action. A feasible solution, in this case, would be to explore more performing state estimators for non-linear systems.

For the sake of completeness, the same tracking objective was simulated for the case of disturbances actuating on the real system, just as in the previous cases, as shown in Fig. 6.13. The resulting control shows a more aggressive response in respect to the previous case, but it was still able to correct the disturbance and track the system to the reference. This fact proves the feasibility of using controller-estimator configurations for reference tracking control. Finally, a summary of the results in this subsection are displayed at Table 6.2.

Controller	ISE	IAE	ITSE	ITAE
$LQRI_1$	[0.5929]	[7.4300]	[9.0789]	[23.1565]
	[2.0197]	[15.1220]	[33.7962]	[49.4736]
$LQGI_1$	[0.5060]	[7.0710]	[8.1233]	[25.0444]
	[11.5765]	[47.5308]	[214.0017]	[189.6191]
$LQRI_2$	[0.7860]	[10.1643]	[11.6943]	[32.7671]
	[167.4390]	[128.4992]	[2059.6345]	[450.4048]
$LQGI_2$	[0.8902]	[12.6788]	[17.1311]	[50.5231]
	[687.4735]	[322.6335]	[11317.1011]	[1274.8359]

Table 6.2: Comparison between the  $LQRI$  and  $LQGI$  controllers. The pair  $(LQRI_2, LQGI_2)$  comprises the simulations in which the system was subjected to disturbances.

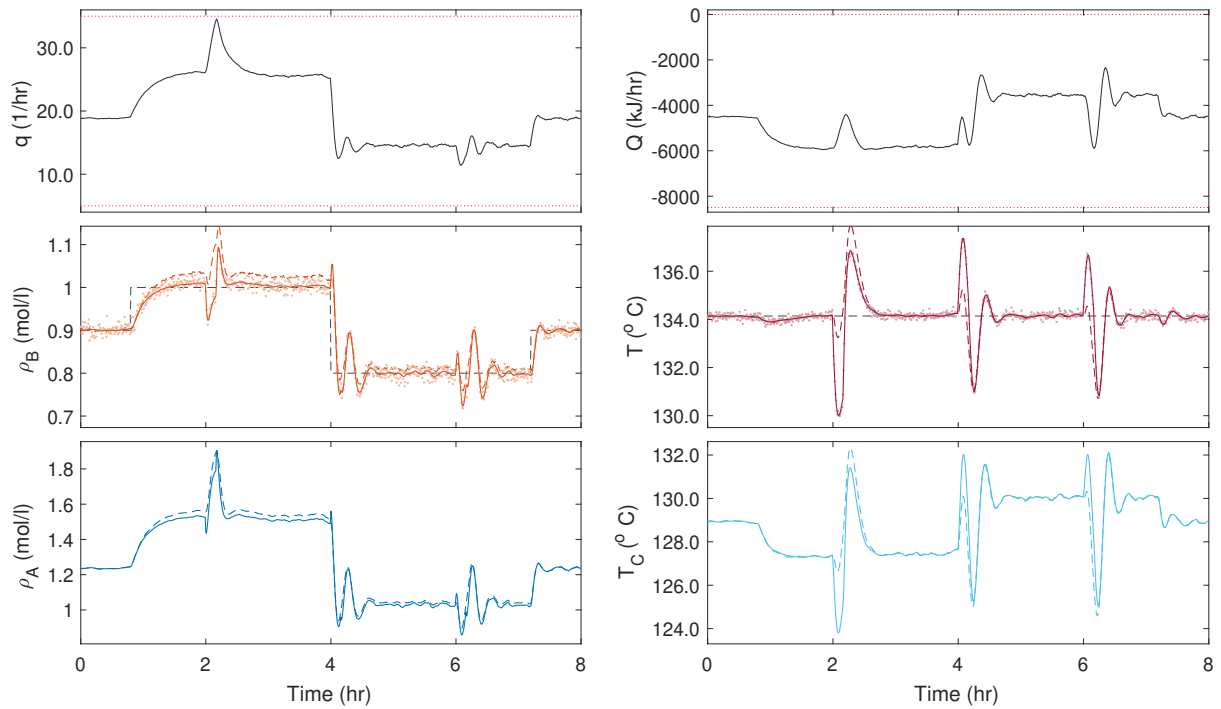


Figure 6.13: Simulation of the  $LQGI_2$  controller for reference tracking, but subjected to three pulse of disturbances in the system variables.



## Chapter 7

# Conclusion

This work has presented a self-contained mathematical framework to model and control chemical reactor systems, focusing on optimal control techniques. In this document, several results were presented and discussed in order to provide a complete understanding of the concepts and methods explored. The discussed methods were applied to a given model of a challenging real-world chemical reactor, which consists of a non-isothermal process following the Van de Vusse reaction scheme. The experiments have shown that optimal control is a feasible solution when the controller has direct access to the internal states of the system, but performed worse when the states are not accessible but estimated. The control simulations, however, still produced practical results that are realistic in respect to real plant of such processes.

The use of optimal controllers greatly reduces the complexity of the controller synthesis procedure, while ensuring an optimal operation. The methods presented can still be expanded to more complex formulations, which would allow for better performance in more adverse cases. Specifically, more advanced state estimation techniques can be considered to mitigate the problem of reference tracking, determining the integral action through an estimate of the controlled variables.

This work also emphasizes the possibility of formulating more advanced techniques, within the scope of optimal control, to achieve better results. As future work, there is the necessity to decouple the control and estimator sample time from the continuous-time in which the real system operates, providing a more realistic configuration of feedback control loops. Additionally, non-linear state estimation techniques should be explored to solve the linearized Kalman filter limitation, and constrained optimal controllers should be explored to mathematically ensure that the hard constraints of the system are being respected.



# Bibliography

- [Anderson and Moore, 1990] Anderson, D. O. B. and Moore, J. B. (1990). *Optimal Control*. Prentice-Hall.
- [Atiyah, 2018] Atiyah, M. (2018). *Introduction to commutative algebra*. CRC Press.
- [Atkins and de Paula, 2011] Atkins, P. and de Paula, J. (2011). *Physical Chemistry for the Life Sciences*. W. H. Freeman, 2nd edition.
- [Barber, 2012] Barber, D. (2012). *Bayesian reasoning and machine learning*. Cambridge University Press.
- [Bellman, 1954] Bellman, R. (1954). The theory of dynamic programming. *Bull. Amer. Math. Soc.*, 60(6):503–515.
- [Bergman et al., 2017] Bergman, T. L., Lavine, A. S., Incropera, F. P., and DeWitt, D. P. (2017). *Fundamentals of Heat and Mass Transfer*. Wiley, 8th edition.
- [Bertsekas, 2017] Bertsekas, D. P. (2017). *Dynamic Programming and Optimal Control*, volume I. Athena Scientific, 4 edition.
- [Bode, 1945] Bode, H. W. (1945). Network analysis and feedback amplifier design.
- [Boyd and Vandenberghe, 2004] Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press.
- [Bryson, 1996] Bryson, A. E. (1996). Optimal control-1950 to 1985. *IEEE Control Systems Magazine*, 16(3):26–33.
- [Cairano et al., 2014] Cairano, S. D., Bernardini, D., Bemporad, A., and Kolmanovsky, I. V. (2014). Stochastic mpc with learning for driver-predictive vehicle control and its application to hev energy management. *IEEE Transactions on Control Systems Technology*, 22(3):1018–1031.
- [Chen, 1998] Chen, C.-T. (1998). *Linear System Theory and Design*. Oxford University Press.
- [Crassidis and Junkins, 2011] Crassidis, J. L. and Junkins, J. L. (2011). *Optimal Estimation of Dynamic Systems*. Chapman & Hall/CRC, 2nd edition.
- [Davison and Maki, 1973] Davison, E. and Maki, M. (1973). The numerical solution of the matrix riccati differential equation. *IEEE Transactions on Automatic Control*, 18(1):71–73.
- [Doyle, 1978] Doyle, J. C. (1978). Guaranteed margins for lqg regulators. *IEEE Transactions on automatic Control*, 23(4):756–757.
- [Eren et al., 2017] Eren, U., Prach, A., Koçer, B. B., Raković, S. V., Kayacan, E., and Açıkmeşe, B. (2017). Model predictive control in aerospace systems: Current state and opportunities. *Journal of Guidance, Control, and Dynamics*, 40(7):1541–1566.



- [Falcone et al., 2007] Falcone, P., Borrelli, F., Asgari, J., Tseng, H. E., and Hrovat, D. (2007). Predictive active steering control for autonomous vehicle systems. *IEEE Transactions on Control Systems Technology*, 15(3):566–580.
- [Franklin et al., 2018] Franklin, G. F., Powell, J. D., Emami-Naeini, A., and Powell, J. D. (2018). *Feedback control of dynamic systems*. Pearson, 8th edition.
- [Goerzen et al., 2009] Goerzen, C., Kong, Z., and Mettler, B. (2009). A survey of motion planning algorithms from the perspective of autonomous uav guidance. *Journal of Intelligent and Robotic Systems*, 57(1):65.
- [González et al., 2016] González, D., Pérez, J., Milanés, V., and Nashashibi, F. (2016). A review of motion planning techniques for automated vehicles. *IEEE Transactions on Intelligent Transportation Systems*, 17(4):1135–1145.
- [Gupta et al., 2012] Gupta, V. K., Ali, I., Saleh, T. A., Nayak, A., and Agarwal, S. (2012). Chemical treatment technologies for waste-water recycling—an overview. *RSC Adv.*, 2:6380–6388.
- [Heller et al., 1978] Heller, H. C., Crawshaw, L. I., and Hammel, H. T. (1978). The thermostat of vertebrate animals. *Scientific American*, 239(2):102–115.
- [Holenda et al., 2008] Holenda, B., Domokos, E., Rédey, A., and Fazakas, J. (2008). Dissolved oxygen control of the activated sludge wastewater treatment process using model predictive control. *Computers & Chemical Engineering*, 32(6):1270 – 1278.
- [Horn and Jackson, 1972] Horn, F. and Jackson, R. (1972). General mass action kinetics. *Archive for rational mechanics and analysis*, 47(2):81–116.
- [Horn and Johnson, 2012] Horn, R. A. and Johnson, C. R. (2012). *Matrix analysis*. Cambridge University Press.
- [Kalman, 1960a] Kalman, R. E. (1960a). A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1):35–45.
- [Kalman, 1960b] Kalman, R. E. (1960b). On the general theory of control systems. In *Proceedings First International Conference on Automatic Control, Moscow, USSR*.
- [Khalil, 2005] Khalil, H. K. (2005). A note on the robustness of high-gain-observer-based controllers to unmodeled actuator and sensor dynamics. *Automatica*, 41(10):1821 – 1824.
- [Kirk, 1998] Kirk, D. E. (1998). *Optimal Control Theory*. Prentice-Hall.
- [Klatt and Engell, 1998] Klatt, K.-U. and Engell, S. (1998). Gain-scheduling trajectory control of a continuous stirred tank reactor. *Computers & Chemical Engineering*, 22(4-5):491–502.
- [Liberzon, 2012] Liberzon, D. (2012). *Calculus of Variations and Optimal Control Theory*. Princeton University Press.
- [Lucia et al., 2013] Lucia, S., Finkler, T., and Engell, S. (2013). Multi-stage nonlinear model predictive control applied to a semi-batch polymerization reactor under uncertainty. *Journal of Process Control*, 23(9):1306 – 1319.
- [Mahmoud and Khalil, 2002] Mahmoud, M. S. and Khalil, H. K. (2002). Robustness of high-gain observer-based nonlinear controllers to unmodeled actuators and sensors. *Automatica*, 38(2):361 – 369.

- [MATLAB, 2018] MATLAB (2018). *version 9.5.0 (R2018b)*. The MathWorks Inc., Natick, Massachusetts.
- [Mnih et al., 2015] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518:529.
- [Moler and Loan, 2003] Moler, C. and Loan, C. V. (2003). Nineteen dubious ways to compute the exponential of a matrix, twenty-five years later. *SIAM Review*, 45(1):3–49.
- [Moore, 1975] Moore, B. C. (1975). On the flexibility offered by state feedback in multivariable systems beyond closed loop eigenvalue assignment. In *1975 IEEE Conference on Decision and Control including the 14th Symposium on Adaptive Processes*, pages 207–214.
- [Mulas et al., 2015] Mulas, M., Tronci, S., Corona, F., Haimi, H., Lindell, P., Heinonen, M., Vahala, R., and Baratti, R. (2015). Predictive control of an activated sludge process: An application to the viikimäki wastewater treatment plant. *Journal of Process Control*, 35:89 – 100.
- [Nise, 2015] Nise, N. S. (2015). *Control Systems Engineering*. Wiley, 7th edition.
- [Nyquist, 1932] Nyquist, H. (1932). Regeneration theory. *Bell system technical journal*, 11(1):126–147.
- [Opanuga et al., 2015] Opanuga, A. A., Edeki, S., Okagbue, H. I., and Akinlabi, G. (2015). A novel approach for solving quadratic riccati differential equations. *International Journal of Applied Engineering Research*, 10(11):29121–29126.
- [Qin and Badgwell, 2003] Qin, S. and Badgwell, T. A. (2003). A survey of industrial model predictive control technology. *Control Engineering Practice*, 11(7):733 – 764.
- [Rasmussen and Williams, 2006] Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning*. The MIT Press.
- [Ross, 2010] Ross, S. M. (2010). *A first course in probability*. Pearson.
- [Särkkä, 2013] Särkkä, S. (2013). *Bayesian Filtering and Smoothing*. Cambridge University Press, New York, NY, USA.
- [Strang, 2016] Strang, G. (2016). *Introduction to linear algebra*, volume 5. Wellesley-Cambridge Press.
- [Stratonovich, 1968] Stratonovich, R. L. (1968). *Conditional Markov Processes and Their Application to the Theory of Optimal Control*. American Elsevier Publisher.
- [Sutton and Barto, 2018] Sutton, R. and Barto, A. (2018). *Reinforcement Learning*. MIT Press, 2nd edition.
- [Syrmos et al., 1997] Syrmos, V., Abdallah, C., Dorato, P., and Grigoriadis, K. (1997). Static output feedback—a survey. *Automatica*, 33(2):125 – 137.
- [Tang et al., 2016] Tang, X., Rupp, B., Yang, Y., Edwards, T. D., Grover, M. A., and Bevan, M. A. (2016). Optimal feedback controlled assembly of perfect crystals. *ACS Nano*, 10(7):6791–6798. PMID: 27387146.
- [Todorov and Jordan, 2002] Todorov, E. and Jordan, M. I. (2002). Optimal feedback control as a theory of motor coordination. *Nature Neuroscience*, 5(11):1226–1235.

- [Van de Vusse, 1964] Van de Vusse, J. (1964). Plug-flow type reactor versus tank reactor. *Chemical Engineering Science*, 19(12):994–996.
- [Vidyasagar, 2002] Vidyasagar, M. (2002). *Nonlinear systems analysis*, volume 42. SIAM.
- [Wang and Su, 2015] Wang, L. and Su, J. (2015). Robust disturbance rejection control for attitude tracking of an aircraft. *IEEE Transactions on Control Systems Technology*, 23(6):2361–2368.

## Appendix A

# Properties of Gaussian Distributions

**Definition A.1.** (Gaussian Distribution) A random variable  $\mathbf{x} \in \mathbb{R}^n$  has a Gaussian distribution with mean  $\boldsymbol{\mu} \in \mathbb{R}^n$  and covariance  $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ , denoted as  $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , if its probability density is described as:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2}|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}. \quad (\text{A.1})$$

**Lemma A.1.** (Joint distribution of Gaussians) Consider the random variables  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^m$ , which have the Gaussian distributions

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{A}\mathbf{x} + \mathbf{u}, \mathbf{L}), \end{aligned} \quad (\text{A.2})$$

then the distribution of  $(\mathbf{x}, \mathbf{y})$  and the marginal distribution of  $\mathbf{y}$  are given as:

$$\begin{aligned} p \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \right) &= \mathcal{N} \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \middle| \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{u} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma} & \boldsymbol{\Sigma}\mathbf{A}^T \\ \mathbf{A}\boldsymbol{\Sigma} & \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T + \mathbf{L} \end{bmatrix} \right), \\ p(\mathbf{y}) &= \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{u}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T + \mathbf{L}). \end{aligned} \quad (\text{A.3})$$

**Lemma A.2.** (Conditional distribution of Gaussians) Consider the random variables  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^m$ , which have the joint Gaussian distribution

$$p \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \right) = \mathcal{N} \left( \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \middle| \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{xy}^T & \boldsymbol{\Sigma}_{yy} \end{bmatrix} \right), \quad (\text{A.4})$$

then the marginal and conditional distributions of  $\mathbf{x}$  and  $\mathbf{y}$  are given as:

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx}), \\ p(\mathbf{y}) &= \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_{yy}), \\ p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_{yy}^{-1}(\mathbf{y} - \boldsymbol{\mu}_y), \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_{yy}^{-1}\boldsymbol{\Sigma}_{xy}^T), \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}_y + \boldsymbol{\Sigma}_{xy}^T\boldsymbol{\Sigma}_{xx}^{-1}(\mathbf{x} - \boldsymbol{\mu}_x), \boldsymbol{\Sigma}_{yy} - \boldsymbol{\Sigma}_{xy}^T\boldsymbol{\Sigma}_{xx}^{-1}\boldsymbol{\Sigma}_{xy}). \end{aligned} \quad (\text{A.5})$$