

Department of Mathematics and Statistics

**Computational Methods for Solving Optimal Industrial Process  
Control Problems**

**Qinqin Chai**

**This thesis is presented for the Degree of  
Doctor of Philosophy  
of  
Curtin University**

**March 2013**

---

---

# Declaration

---

To the best of my knowledge and belief, this thesis contains no material previously published by any other person except where due acknowledgment has been made.

This thesis contains no material which has been accepted for the award of any other degree or diploma in any university.

.....  
Qinqin Chai  
March 2013

---

---

# Abstract

---

In this thesis, we develop new computational methods for three classes of dynamic optimization problems: (i) A parameter identification problem for a general nonlinear time-delay system; (ii) an optimal control problem involving systems with both input and output delays, and subject to continuous inequality state constraints; and (iii) a max-min optimal control problem arising in gradient elution chromatography.

In the first problem, we consider a parameter identification problem involving a general nonlinear time-delay system, where the unknown time delays and system parameters are to be identified. This problem is posed as a dynamic optimization problem, where its cost function is to measure the discrepancy between predicted output and observed system output. The aim is to find unknown time-delays and system parameters such that the cost function is minimized. We develop a gradient-based computational method for solving this dynamic optimization problem. We show that the gradients of the cost function with respect to these unknown parameters can be obtained via solving a set of auxiliary time-delay differential systems from  $t = 0$  to  $t = T$ . On this basis, the parameter identification problem can be solved as a nonlinear optimization problem and existing optimization techniques can be used. Two numerical examples are solved using the proposed computational method. Simulation results show that the proposed computational method is highly effective. In particular, the convergence is very fast even when the initial guess of the parameter values is far away from the optimal values.

Unlike the first problem, in the second problem, we consider a time delay identification problem, where the input function for the nonlinear time-delay system is piecewise-constant. We assume that the time-delays—one involving the state variables and the other involving the input variables—are unknown and need to be estimated using experimental data. We also formulate the problem of estimating the unknown delays as a nonlinear optimization problem in which the cost function measures the least-squares error between predicted output and measured system output. This estimation problem can be viewed as a switched system optimal control problem with time-delays. We show that the gradient of the cost function with respect to the unknown state delay can be obtained via solving a auxiliary time-delay differential system. Furthermore, the gradient of the cost function with respect to the unknown input delay can be obtained via solving an auxiliary time-delay differential system with jump conditions at the delayed control switching time points. On this basis, we develop a heuristic computational algorithm for

solving this problem using gradient based optimization algorithms. Time-delays in two industrial processes are estimated using the proposed computational method. Simulation results show that the proposed computational method is highly effective.

For the third problem, we consider a general optimal control problem governed by a system with input and output delays, and subject to continuous inequality constraints on the state and control. We focus on developing an effective computational method for solving this constrained time delay optimal control problem. For this, the control parameterization technique is used to approximate the time planning horizon  $[0, T]$  into  $N$  subintervals. Then, the control is approximated by a piecewise constant function with possible discontinuities at the pre-assigned partition points, which are also called the switching time points. The heights of the piecewise constant function are decision variables which are to be chosen such that a given cost function is minimized. For the continuous inequality constraints on the state, we construct approximating smooth functions in integral form. Then, the summation of these approximating smooth functions in integral form, which is called the constraint violation, is appended to the cost function to form a new augmented cost function. In this way, we obtain a sequence of approximate optimization problems subject to only boundedness constraints on the decision variables. Then, the gradient of the augmented cost function is derived. On this basis, we develop an effective computational method for solving the time-delay optimal control problem with continuous inequality constraints on the state and control via solving a sequence of approximate optimization problems, each of which can be solved as a nonlinear optimization problem by using existing gradient-based optimization techniques. This proposed method is then used to solve a practical optimal control problem arising in the study of a real evaporation process. The results obtained are highly satisfactory, showing that the proposed method is highly effective.

The fourth problem that we consider is a max-min optimal control problem arising in the study of gradient elution chromatography, where the manipulative variables in the chromatographic process are to be chosen such that the separation efficiency is maximized. This problem has three non-standard characteristics: (i) The objective function is non-smooth; (ii) each state variable is defined over a different time horizon; and (iii) the order of the final times for the state variable, the so-called retention times, are not fixed. To solve this problem, we first introduce a set of auxiliary decision variables to govern the ordering of the retention times. The integer constraints on these auxiliary decision variables are approximated by continuous boundedness constraints. Then, we approximate the control by a piecewise constant function, and apply a novel time-scaling transformation to map the retention times and control switching times to fixed points in a new time horizon. The retention times and control switching times become decision variables in the new time horizon. In addition, the max-min objective function is approximated by a minimization problem subject to an additional constraint. On this basis, the optimal control problem is

reduced to an approximate nonlinear optimization problem subject to smooth constraints, which is then solved using a recently developed exact penalty function method. Numerical results obtained show that this approach is highly effective.

Finally, some concluding remarks and suggestions for further study are made in the conclusion chapter.

---

---

# List of publications

---

The following papers were published or accepted for publication during the PhD candidature:

- Q. Chai, R. Loxton, K. L. Teo, and C. H. Yang, “A class of optimal state-delay control problems,” *Nonlinear Analysis: Real World Applications*, vol. 14, no. 3, pp. 1536-1550, 2013.
- Q. Chai, R. Loxton, K. L. Teo, and C. H. Yang, “A unified parameter identification method for nonlinear time-delay systems,” *Journal of Industrial and Management Optimization*, vol. 9, no. 2, pp. 471-486, 2013.
- Q. Chai, R. Loxton, K. L. Teo, and C. H. Yang, “Time-delay estimation for nonlinear systems with piecewise-constant input,” *Applied Mathematics and Computation*, vol. 219, pp. 9543-9560, 2013.
- Q. Chai, R. Loxton, K. L. Teo, and C. H. Yang, “A max-min control problem arising in gradient elution chromatography,” *Industrial and Engineering Chemistry Research*, vol. 51, no. 17, pp. 6137-6144, 2012.
- Q. Chai, C. H. Yang, K. L. Teo, and W. H. Gui, “Optimal control of an industrial-scale evaporation process: sodium aluminate solution,” *Control Engineering Practice*, vol. 20, no. 6, pp. 618-628, 2012.
- Y. G. Li, W. H. Gui, K. L. Teo, H. Q. Zhu, and Q. Q. Chai, “Optimal control for zinc solution purification based on interaction CSTR models,” *Journal of Process Control*, vol. 22, no. 10, pp. 1878-1889, 2012.
- R. Loxton, Q. Chai, and K. L. Teo, “A class of max-min optimal control problems with applications to chromatography,” *Proceedings of the 5th International Conference on Optimization and Control with Applications*, Beijing, China, December 4-8, 2012.

---

---

# Acknowledgements

---

The research reported in this thesis was carried out from September 2008 to December 2012. During this period, I was a PhD student in the Department of Mathematics and Statistics, Curtin University and the School of Information Science and Engineering, Central South University. I wish to express my appreciation for all kinds of help I received from my supervisors, families, friends, and colleagues during this period of time.

I would like to express my heartfelt thanks to my supervisor, Prof. Kok Lay Teo, and his wife, Mrs. Lye-Hen Teo. Prof. Teo has guided my research during the past two years with remarkable patience and enthusiasm. During the two years of my stay in Australia, he was not only a great supervisor of my research, but also a gracious mentor of my life.

I would like to thank Dr. Ryan Loxton, my co-supervisor in the Department of Mathematics and Statistics, Curtin University. He is an excellent mathematician, and he is always willing to help whenever I have difficulties in my research. He has shared his experience with me on the skill of writing papers. I really appreciate all the unselfish help I received from him.

I would like to thank Prof. Chunhua Yang, my co-supervisor in the School of Information Science and Engineering, Central South University. I have known Prof. Yang since I became a post-graduate student in the School of Information Science and Engineering, Central South University in 2006. It was her who led me to the road of research and helped me to apply for a scholarship from China Scholarship Council, so that I could come to Australia and start a new stage of my research.

I would like to give thanks to Prof. Weihua Gui, in the School of Information Science and Engineering, Central South University. He is a great engineer. He has shared with me many practical experiences, which have inspired me in my research.

I also thank Prof. Yonghong Wu, Postgraduate Coordinator in the Department of Mathematics and Statistics, Curtin University. He is the Chair of my thesis committee and he has helped me on many occasions.

I would like to give thanks to the people and friends that I have worked with in Australia, especially, Dr. Qun Lin, Dr. Bin Li, Dr. Jingyang Zhou and his wife Yujing Wang, Dr. Changjun Yu, Prof. Honglei Xu, Yufei Sun, Yanyan Yin, Qian Sun, Mingliang Xue, Yanli Zhou and the department's academic visitors that I had the opportunity to meet, Dr. Canghua Jiang, A/Prof. Chuanjiang Li, A/Prof. Tieqiao Tang, Prof. Xuegang Hu, A/Prof. Lingling Xu, and A/Prof. Xiangyu Gao.

In addition, I thank each of my friends and teachers in Changsha, Prof. Yalin Wang, Prof. Yongfang Xie, Prof. Jianjun He, A/Prof. Yonggang Li, Dr. Hongqiu Zhu, Dr. Can Huang, Dr. Huifeng Ren, Dr. Canhui Xu, Dr. Cao Shi, Jianhua Liu, and Jingjing Qin. I really enjoyed the time that I spent with them.

I thank all of the staff in the Department of Mathematics and Statistics for contributing to a friendly working environment. The administrative staff—Joyce Yang, Shuie Liu, Lisa Holling, Jeannie Darmago, and Carey Ryken Rapp—deserve special thanks for providing kind and professional help on numerous occasions.

Finally, on a more personal note, I sincerely thank my husband Xinmei Lin and everyone in my family for their love, understanding and support during my PhD candidature in Australia.



---

---

# Contents

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation and background . . . . .	1
1.2	Nonlinear programming problems . . . . .	2
1.3	Numerical techniques for optimal control problems . . . . .	4
1.4	Parameter identification for time-delay systems . . . . .	10
1.5	Control methods for time-delay systems . . . . .	12
1.5.1	Optimal control for time-delay systems . . . . .	12
1.5.2	Model predictive control for time-delay systems . . . . .	12
1.6	Max-Min optimal control problems . . . . .	15
1.7	Overview of the thesis . . . . .	17
<b>2</b>	<b>Identification of time-delays and parameters for nonlinear systems</b>	<b>21</b>
2.1	Introduction . . . . .	21
2.2	Problem formulation . . . . .	22
2.3	Preliminaries . . . . .	23
2.4	Gradient computation . . . . .	30
2.4.1	State variation with respect to time-delays . . . . .	30
2.4.2	State variation with respect to system parameters . . . . .	34
2.4.3	Gradient computation algorithm . . . . .	37
2.5	Numerical examples . . . . .	38
2.5.1	Example 2.1 . . . . .	38
2.5.2	Example 2.2 . . . . .	40
2.6	Conclusion . . . . .	42
<b>3</b>	<b>Identification of time-delays for nonlinear systems with piecewise-constant input</b>	<b>44</b>
3.1	Introduction . . . . .	44
3.2	Problem formulation . . . . .	45
3.3	State variation . . . . .	48
3.3.1	State variation with respect to the state-delay . . . . .	48
3.3.2	State variation with respect to the input-delay . . . . .	49
3.4	Computation algorithm . . . . .	62

3.5	Numerical examples . . . . .	64
3.5.1	Example 1: Zinc sulphate purification . . . . .	64
3.5.2	Example 2: Sodium aluminate evaporation . . . . .	67
3.6	Conclusions . . . . .	70
<b>4</b>	<b>Time-delay optimal control application problem: an industrial evaporation process</b>	<b>71</b>
4.1	Introduction . . . . .	71
4.2	Problem statement . . . . .	72
4.3	Solution method . . . . .	74
4.3.1	Control parameterization . . . . .	74
4.3.2	Constraints transformation . . . . .	76
4.3.3	Convergence analysis . . . . .	78
4.3.4	Computational method . . . . .	85
4.4	Application: optimal control of an evaporation process . . . . .	91
4.4.1	The evaporation system . . . . .	92
4.4.2	Mathematical model of the evaporation system . . . . .	94
4.4.3	Optimal control problem formulation . . . . .	99
4.4.4	Numerical results . . . . .	102
4.5	Conclusion . . . . .	107
<b>5</b>	<b>A max-min control problem arising in gradient elution chromatography</b>	<b>108</b>
5.1	Introduction . . . . .	108
5.2	Problem statement . . . . .	111
5.3	Problem transformation . . . . .	112
5.4	A computational method . . . . .	118
5.5	Numerical examples . . . . .	125
5.5.1	Example 5.1 . . . . .	125
5.5.2	Example 5.2 . . . . .	127
5.6	Conclusions . . . . .	128
<b>6</b>	<b>Summary and future research directions</b>	<b>130</b>
6.1	Main contributions of the thesis . . . . .	130
6.2	Future research directions . . . . .	132
	<b>Bibliography</b>	<b>133</b>

---

---

# CHAPTER 1

---

## Introduction

### 1.1 Motivation and background

In an optimal control problem, there is a governing dynamic system whose trajectory, called the state, is influenced by an adjustable variable, called the control. Its aim is to find a control such that a performance index is optimized subject to some specified constraints. The performance index, which is also called the cost function, could represent energy consumption, wastage of consumable materials, or the time taken to achieve a given task, just to name a few examples. The specified constraints could arise due to design specifications, safety operation conditions or engineering requirements. Optimal control has many successful real world applications in areas ranging from engineering to economics. Many of these real world problems are too complicated to admit analytical solutions. Thus, it is unavoidable to rely on numerical methods to deal with these problems. There are now many computational methods available in the literature for solving various classes of optimal control problems. Most of these methods are for control problems in which the governing dynamic system does not involve time-delay. However, time-delays arise in many real world applications, such as chemical tank reactors [86], aerospace engineering [143], chromatography [107], and power converters [153]. The effects of time-delays must not be ignored, because it is known [28, 123] that time-delays in a dynamic system could cause instability in the system concerned. However, techniques, theory and methods for optimal control problems without time-delay are often not applicable to optimal control problems with time-delays. Consequently, it has attracted a considerable interest amongst mathematicians and engineers, especially process engineers, to develop effective computational methods for solving optimal control problems involving nonlinear time-delay systems [25, 74, 77, 131].

In this thesis, we will formulate several optimal control problems arising in practical industrial processes. Then, we will develop effective computational methods for solving these real world optimal control problems.

## 1.2 Nonlinear programming problems

Nonlinear programming (NLP) problem is an optimization problem with nonlinear objective function and/or nonlinear constraints. NLP problems arise in many real world applications. A general NLP problem can be stated as follows:

**Problem** ( $\tilde{P}_1$ ).

$$\begin{aligned} & \text{Minimize} && f(\mathbf{x}) \\ & \text{s.t.} && g_i(\mathbf{x}) = 0, \quad i = 1, \dots, m_1, \\ & && g_i(\mathbf{x}) \leq 0, \quad i = m_1 + 1, \dots, m_1 + m_2. \end{aligned}$$

where  $\mathbf{x} \in \mathbb{R}^r$  is the decision vector;  $f(\mathbf{x})$  is the objective function;  $g_i(\mathbf{x})$ ,  $i = 1, \dots, m_1$ , are given equality constraint functions; and  $g_i(\mathbf{x})$ ,  $i = m_1 + 1, \dots, m_1 + m_2$ , are given inequality constraint functions. Assume that the objective function and the constraint functions are twice continuously differentiable.

For NLP problem, a vector  $\mathbf{x}$  is called a feasible solution if it satisfies all the constraint functions of Problem ( $\tilde{P}_1$ ). The set containing all the feasible solutions is called the feasible region. It is denoted by  $\mathcal{F}$ . Furthermore, the  $j$ th inequality constraint is said to be active at the point  $\mathbf{x}$  if  $g_j(\mathbf{x}) = 0$ .

**Definition 1.1.** (*Active set.*) The active set  $\mathcal{A}(\mathbf{x})$  at point  $\mathbf{x}$  is the set of indices defined by

$$\mathcal{A}(\mathbf{x}) = \{j \in \{m_1 + 1, \dots, m_1 + m_2\} | g_j(\mathbf{x}) = 0\}.$$

Let  $\bar{\mathcal{A}}(\mathbf{x}) = \mathcal{A}(\mathbf{x}) \cup \{1, \dots, m_1\}$  denote the set of indices of the active constraints and the equality constraints at  $\mathbf{x}$ .

A feasible point  $\mathbf{x}^* \in \mathcal{F}$  is called a global minimum if  $f(\mathbf{x}^*) \leq f(\mathbf{x})$ ,  $\forall \mathbf{x} \in \mathcal{F}$ . Moreover, a feasible point  $\mathbf{x}^* \in \mathcal{F}$  is called a local minimum if there exists an  $\varepsilon > 0$  such that

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{N}_\varepsilon(\mathbf{x}^*) = \{\mathbf{x}^* \in \mathcal{F} | \|\mathbf{x}^* - \mathbf{x}\| \leq \varepsilon\}.$$

We say that the linearly independent constraint qualification (LICQ) holds at a point  $\hat{\mathbf{x}} \in \mathcal{F}$  if the following set at point  $\hat{\mathbf{x}}$  is linearly independent:

$$\{\nabla g_j(\hat{\mathbf{x}}) | j \in \bar{\mathcal{A}}(\hat{\mathbf{x}})\},$$

where  $\nabla g_j(\hat{\mathbf{x}})$  denotes the partial derivative (i.e., gradient) of  $g_j(\mathbf{x})$  evaluated at  $\mathbf{x} = \hat{\mathbf{x}}$ . In addition, a point at which the LICQ holds is called a regular point.

Consider Problem  $(\tilde{P}_1)$ . The Lagrangian is:

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{i=1}^{m_1+m_2} \lambda_i g_i(\mathbf{x}),$$

where  $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_{m_1+m_2}]^\top$  is the vector of Lagrange multipliers.

Let  $\nabla \mathbf{h}(\mathbf{x})$  be the partial derivatives of function  $h$  with respect to  $\mathbf{x}$  and let  $\nabla^2 \mathbf{h}(\mathbf{x})$  be the second partial derivatives (Hessian matrix) of  $h$  with respect to  $\mathbf{x}$ , i.e.,

$$\nabla \mathbf{f}(\mathbf{x}) = \left[ \frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right]^\top,$$

$$\nabla^2 \mathbf{f}(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \ddots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_n} \end{bmatrix}.$$

Similarly,  $\nabla L(\mathbf{x})$  and  $\nabla^2 L(\mathbf{x}, \boldsymbol{\lambda})$  denote, respectively, the partial derivatives and the second partial derivatives (Hessian matrix) of  $L$  with respect to  $\mathbf{x}$ .

The first-order optimality conditions, called the Karush-Kuhn-Tucker (KKT) conditions, for local optimal solutions are given below.

**Theorem 1.1.** (*KKT conditions*). *Suppose that  $\mathbf{x}^* \in \mathcal{F}$  is a local optimal solution of Problem  $(\tilde{P}_1)$  and that the LICQ holds at  $\mathbf{x}^*$ . Then, there exists a non trivial vector  $\boldsymbol{\lambda}^*$  such that the following conditions are satisfied*

$$\begin{aligned} \nabla L(\mathbf{x}^*, \boldsymbol{\lambda}^*) &= \mathbf{0}, \\ g_i(\mathbf{x}^*) &= 0, \quad i = 1, \dots, m_1, \\ g_i(\mathbf{x}^*) &\leq 0, \quad i = m_1 + 1, \dots, m_1 + m_2, \\ \lambda_i^* &\geq 0, \quad i = m_1 + 1, \dots, m_1 + m_2, \\ \lambda_i^* g_i(\mathbf{x}^*) &= 0, \quad i = 1, \dots, m_1 + m_2. \end{aligned}$$

For NLP problem, the global optimal solution is difficult to find, and hence the focus is on finding a local optimal solution. There are many methods available for solving various NLP problems. Examples include convex programming [71], separable programming [49], fractional programming [22], quadratic programming [23], and sequential quadratic programming [54]. In real word applications, the optimization problems are usually nonlinear, where both the objective function and constraint functions are nonlinear. For such problems, sequential quadratic programming (SQP) method is known to be effective. For more details on theory and computational algorithms, see, for example, [73, 104].

The main idea of SQP method is to solve a sequence of quadratic programming (QP)

subproblems, each of which is a quadratic model with quadratic objective function and linearized constraints. A QP subproblem is solved in each iteration step, giving rise to a search direction for the NLP for the current iterate  $\mathbf{x}(k)$ , where  $k$  denotes the  $k$ th iteration. More specifically, the objective function  $f$  is approximated by its local quadratic approximation

$$f(\mathbf{x}(k) + d(\mathbf{x})) \approx f(\mathbf{x}(k)) + \nabla f(\mathbf{x}(k))d(\mathbf{x}) + \frac{1}{2}d(\mathbf{x})^\top \nabla^2 L(\mathbf{x}(k), \boldsymbol{\lambda}(k))d(\mathbf{x}(k)),$$

where  $d(\mathbf{x}) = \mathbf{x} - \mathbf{x}(k)$ , and the constraint functions  $g_i$ ,  $i = 1, \dots, m_1 + m_2$ , are approximated by their local affine approximations

$$g_i(\mathbf{x}(k) + d(\mathbf{x})) \approx g_i(\mathbf{x}(k)) + \nabla g_i(\mathbf{x}(k))^\top d(\mathbf{x}), \quad i = 1, \dots, m_1 + m_2.$$

Let  $\mathbf{B}(k) = \nabla^2 L(\mathbf{x}(k), \boldsymbol{\lambda}(k))$ . The QP subproblem is:

**Problem ( $\tilde{P}_2$ ).**

$$\begin{aligned} \min \quad & f(\mathbf{x}(k)) + (\nabla f(\mathbf{x}(k), \boldsymbol{\lambda}(k)))^\top d(\mathbf{x}) + \frac{1}{2}d(\mathbf{x})^\top \mathbf{B}(k)d(\mathbf{x}), \\ \text{s.t.} \quad & g_i(\mathbf{x}(k)) + (\nabla g_i(\mathbf{x}(k)))^\top d(\mathbf{x}) = 0, \quad i = 1, \dots, m_1, \\ & g_i(\mathbf{x}(k)) + (\nabla g_i(\mathbf{x}(k)))^\top d(\mathbf{x}) \leq 0, \quad i = m_1 + 1, \dots, m_1 + m_2. \end{aligned}$$

Problem ( $\tilde{P}_2$ ) is solved as a quadratic programming problem with the active set strategy. For further details, see [62, 80].

To summarize, the SQP method is an iterative method for solving nonlinear optimization problems. It has been implemented in many software packages, such as NPSOL, NLPQL, OPSYC, OPTIMA, and SQP in MATLAB. It is important to note that the gradients of the cost function and constraint functions are essential information needed in the optimization process. Thus, this thesis pertains to the development of theory and methods for computing gradients of the cost function and the constraint functions. These gradients are then used in conjunction with the SQP iterative method for solving the optimization problems under considerations.

## 1.3 Numerical techniques for optimal control problems

In engineering, a mathematical model is commonly used to describe the behavior of a dynamic process. This mathematical model is often expressed in terms of a system of

ordinary differential equations as given below:

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t)), \quad t \in [0, T], \quad (1.1)$$

$$\mathbf{x}(0) = \mathbf{x}^0, \quad (1.2)$$

where  $T > 0$  is the given terminal time of the time planning horizon  $[0, T]$ ;  $\mathbf{x} \in \mathbb{R}^n$  is the state vector;  $\mathbf{u} \in \mathbb{R}^r$  is the control vector;  $\mathbf{f} \in \mathbb{R}^n$  is a given function describing the evolution of the states; and  $\mathbf{x}^0$  is the initial state vector at the initial time  $t = 0$ . Note that here, the value of a function  $\omega$  at time  $t$  is denoted by  $\omega(t)$ . The control can change its values from  $t = 0$  to  $t = T$ . For a given control, the state evolves according to the system of ordinary differential equations (1.1) with initial condition (1.2) over the time planning horizon  $[0, T]$ .

In practice, the control strategy for (1.1)-(1.2) cannot be completely arbitrary, because the control strategy is limited by the capacity of the equipment used. For example, the feed flow rate into a 30 m<sup>3</sup> reactor tank must be less than or equal to  $\frac{30}{T}$  m<sup>3</sup>. In other words, the control is subject to physical limitations, often expressed mathematically as the following control restraint set:

$$\mathbf{U} = \{\mathbf{v} = [v_1, \dots, v_r]^\top : a_i \leq v_i \leq b_i, i = 1, \dots, r\},$$

where  $a_i$  and  $b_i$ ,  $i = 1, \dots, r$  are given constants; and the superscript  $\top$  denotes the transpose. A measurable function  $\mathbf{u}$  such that  $\mathbf{u}(t) \in \mathbf{U}$  for almost all  $t \in [0, T]$  is called an admissible control. Let  $\mathcal{U}$  be the set which consists of all such admissible controls. It is called the set of admissible controls.

Note that the state is influenced by the control through system (1.1). For an optimal control problem, it is required to choose an admissible control such that a performance measure, which could be energy consumption, wastage of consumable materials, etc., is optimized. This performance measure is also called an objective function or a cost function. In general, there are two terms in a cost function—a terminal cost and an integral cost. For example, in an evaporation process, the energy usage should be minimized while the solution level in each of the evaporators must be as close as possible to specific given values at the terminal time. Let  $x_j$ ,  $j = 1, \dots, 7$ , be the states describing the levels; and let  $\hat{x}_j$ ,  $j = 1, \dots, 7$ , be the target levels at terminal time. Furthermore, let  $u$  be the control representing the flow rate of the high temperature steam, and let  $W$  be a function of  $u$  and  $\mathbf{x} = [x_1, \dots, x_7]^\top$  representing the water evaporated from the process. A typical cost function of the form is given below:

$$J = \sum_{j=1}^7 (x_j(T) - \hat{x}_j)^2 + \int_0^T \frac{u(t)^2}{W(u(t), \mathbf{x}(t))^2} dt. \quad (1.3)$$

The first term on the right hand side of (1.3) measures the differences between the real levels and the target levels at the terminal time. It is the terminal cost. The second term on the right hand side of (1.3) evaluates the energy usage during the whole period of the time horizon.

The cost function given by (1.3) is a special case of a general cost function given as follows:

$$J(\mathbf{u}) = \Phi(\mathbf{x}(T)) + \int_0^T \mathcal{L}(t, \mathbf{x}(t), \mathbf{u}(t)) dt, \quad (1.4)$$

where  $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}$  is given function measuring the terminal cost, and  $\mathcal{L} : \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^r \rightarrow \mathbb{R}$  measures the cost during the whole time horizon. An admissible control which minimizes the cost function (1.4) is called an optimal control. Note that the cost function  $J$  depends entirely on  $\mathbf{u}$ , as  $\mathbf{x}$  is implicitly determined by  $\mathbf{u}$  through the system of ordinary differential equations (1.1)-(1.2).

We may now state formally a simple optimal control problem in the following.

**Problem ( $\tilde{P}_3$ ).** *Given the dynamic system (1.1)-(1.2), find a control  $\mathbf{u} \in \mathcal{U}$  such that the cost function (1.4) is minimized.*

Problem ( $\tilde{P}_3$ ) is called a Bolza problem. It can, in principle, be solved by using Pontryagin Minimum Principle or Bellman's Principle of Optimality.

Let us look at how the Pontryagin Minimum Principle is used to solve Problem ( $\tilde{P}_3$ ). For this, we introduce the Hamiltonian function for Problem ( $\tilde{P}_3$ ) given below:

$$H(t, \mathbf{x}(t), \mathbf{u}(t), \boldsymbol{\lambda}(t)) = \mathcal{L}(t, \mathbf{x}(t), \mathbf{u}(t)) + \boldsymbol{\lambda}^\top(t) \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t)), \quad (1.5)$$

where  $\boldsymbol{\lambda}$  is called the co-state, which satisfies the following system of differential equations:

$$\dot{\boldsymbol{\lambda}}(t) = - \left[ \frac{\partial H(t, \mathbf{x}(t), \mathbf{u}(t), \boldsymbol{\lambda}(t))}{\partial \mathbf{x}} \right]^\top, \quad t \in [0, T], \quad (1.6)$$

with boundary condition

$$\boldsymbol{\lambda}(T) = \left[ \frac{\partial \Phi(\mathbf{x}(T))}{\partial \mathbf{x}} \right]^\top. \quad (1.7)$$

System (1.6)-(1.7) is called the co-state system. If  $\mathbf{u}^*$  is an optimal control, and  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  are the corresponding state and co-state, respectively, then it can be shown [110] that

$$H(t, \mathbf{x}^*(t), \mathbf{u}^*(t), \boldsymbol{\lambda}^*(t)) = \min_{\mathbf{v} \in \mathcal{U}} H(t, \mathbf{x}^*(t), \mathbf{v}, \boldsymbol{\lambda}^*(t)), \quad (1.8)$$

for all  $t \in [0, T]$ , except possibly on a finite subset of  $[0, T]$ . This condition is known as the Pontryagin Minimum Principle. By solving the Pontryagin Minimum Principle, the optimal control can, in principle, be obtained as a function of time, state, and co-state. If



such a control is obtained, we can substitute it into the state system (1.1)-(1.2) and the co-state system (1.6)-(1.7), yielding a two-point boundary-value (TPBV) problem. The optimal control can be obtained through solving this TPBV problem. This is, however, a very difficult task (even solving it numerically is difficult, let alone analytically).

We now look at Bellman's Principle of Optimality. Let the system (1.1) be evolved starting at time point  $t \in [0, T]$  from a given state  $\mathbf{x}$ , and let the corresponding solution of the system be denoted as  $\mathbf{y}(s|t, \mathbf{x})$ , where  $s \in [t, T]$ . Then, define the value function  $V : [0, T] \times \mathbb{R}^n \rightarrow \mathbb{R}$  as follows:

$$V(t, \mathbf{x}) = \inf_{\mathbf{u} \in \mathcal{U}} \left\{ \Phi(\mathbf{y}(T|t, \mathbf{x})) + \int_t^T \mathcal{L}(s, \mathbf{y}(s|t, \mathbf{x}), \mathbf{u}(s)) ds \right\}, \quad (t, \mathbf{x}) \in [0, T] \times \mathbb{R}^n.$$

By Bellman's Principle of Optimality, it can be shown that the following partial differential equation is satisfied:

$$\frac{\partial V(t, \mathbf{x})}{\partial t} + \inf_{\mathbf{v} \in U} \left\{ \frac{\partial V(t, \mathbf{x})}{\partial \mathbf{x}} \mathbf{f}(t, \mathbf{x}(t), \mathbf{v}) + \mathcal{L}(t, \mathbf{x}(t), \mathbf{v}) \right\} = 0, \quad (t, \mathbf{x}) \in [0, T] \times \mathbb{R}^n, \quad (1.9)$$

with the boundary condition

$$V(T, \mathbf{x}) = \Phi(\mathbf{x}(T)), \quad \mathbf{x} \in \mathbb{R}^n. \quad (1.10)$$

Equation (1.9) is known as the Hamilton-Jacobi-Bellman (HJB) equation. The solution of the HJB equation (1.9) with boundary condition (1.10) can be used to construct an optimal feedback control to Problem  $(\tilde{P}_3)$ . However, the HJB equation can only be solved analytically for problems involving linear dynamics and quadratic cost function. Furthermore, to solve it numerically, the dimension of the problem must be small, because the numerical solution of HJB equation is computationally very demanding.

In view of the difficulties involved in the direct application of the Pontryagin Minimum Principle and Bellman's Principle of Optimality, even for a simple class of optimal control problems, it is inevitable to rely on computational algorithms to solve optimal control problems, especially for real world problems which are often very complicated and subject to various constraints arising from engineering limitations and design specifications. Consequently, numerous families of computational approaches have been developed, such as the direct collocation approach [30, 75], the iterative dynamic programming technique [130], the leap-frog algorithm [36], the switching time computation method [35], the sequential gradient-restoration methods [135], the multiple shooting methods [44, 59], and the control parameterization methods [80].

The direct collocation approach is to approximate a constrained optimal control problem by a finite dimensional nonlinearly constrained optimization problem (NLP), where the entire time horizon is divided into a finite number of subintervals. In each subinter-

val, the controls are approximated by a continuous and piecewise linear function, and the state variables are approximated by a continuously differentiable piecewise cubic function. Clearly, the dimension of the discretised problem is directly related to the number of partition points of the time horizon. A sequence of refinement steps is applied resulting in a sequence of NLPs of increasing size. Each of these NLPs is solved by standard sequential quadratic programming (SQP) methods (see, for example, [29, 48]). This approach (see, for example, [114] and [113]) is applicable to solve optimal control problems with nonlinear dynamical equations subject to nonlinear constraints. Some convergence properties are derived in [138], which are further used to obtain reliable estimates of the co-state variables. Various applications of the method are demonstrated in [56, 102]. Although it can readily solve small scale problems, the discretization of both control and state variables in the direct collocation approach can lead to excessive computation cost for large scale optimal control problems, especially if a reasonably accuracy is required to be met.

The iterative dynamic programming (IDP) technique is derived in [130] based on the Dynamic Programming Principle. It is refined in [128] and [96] to improve the efficiency of the computational procedure. This technique constructs a grid structure for the discretization of both the state and the control. The grid of the state defines accessible points in the state trajectory, while the grid of the control defines admissible control values. The grids are refined iteratively until a satisfactory control policy is obtained. Initial development employed piecewise-constant controls and this was later extended to piecewise linear control policies [129]. Constraints are handled by using a penalty function approach to incorporate them into the objective function. The IDP technique has been successfully applied to a wide range of optimal control problems in chemical engineering, see, for example, [12, 154]. However, as mentioned in [129], there exist many algorithmic parameters, which include the region contraction factor, the number of allowable values for each control variable, the number of grid points, the initial region size and the restoration factor. Proper determination of these parameters is not an easy task.

The leap-frog algorithm is initially developed in [36], where it is used to solve a special type of Two Point Boundary Value Problem arising in geodesics. In [108] and [37], the algorithm is further developed and implemented to handle general nonlinear systems with unbounded and bounded controls. A description of the algorithm is presented in [34], while some theoretical analysis of the algorithm is presented in [37] for a class of optimal control problems with bounded controls in the plane.

The switching time computation (STC) method, proposed in [35], is a computational procedure to find optimal locations of switching points for single-input nonlinear systems. A concatenation of constant-input arcs is used to take the system from a given initial point to the target. In [137] and [39], the STC method is used in the development of a time optimal bang-bang control algorithm. However, this approach is rather restrictive. It is not directly applicable to many types of constraints which appear in practice.

The sequential gradient-restoration algorithms [5, 6] are applicable to optimal control problems involving differentiable constraints, non-differentiable constraints, and terminal constraints. This family of algorithms involves a sequence of two-phase cycles, where each cycle includes a gradient phase and a restoration phase. There is a function, called the augmented function, that consists of the original cost function and the constraints violations. In the gradient phase, the value of the augmented function is decreased; in the restoration phase, the constraint error is decreased, while avoiding excessive change in the value of the cost function. In the complete gradient-restoration cycle, the value of the cost function is decreased, while the constraints are satisfied to a predetermined accuracy. Hence a succession of suboptimal solutions is obtained. It is further enhanced by the dual version [7] of the algorithm.

The multiple shooting approach is proposed in [59]. It divides the time horizon into many subintervals. Then, at each subinterval, the shooting method is used to solve the co-state dynamic system based on an initial guess of the solution of the co-state system. Re-estimating the co-states is continued based on the mismatches until the conditions of the minimum principle are satisfied. This approach is rather sensitive [142] to the initial guess of the co-states at the initial time point.

For the control parameterization method, it is done by partitioning the time horizon of an optimal control problem into several subintervals such that each control can be approximated by a piecewise-constant function (or piecewise linear function or piecewise smooth function) which is consistent with the corresponding partition. The partition points are often referred to as control switching times. The heights of the approximating piecewise-constant function are decision variables, known as control parameters. For the continuous inequality constraints on the state and/or control, the constraint transcription method was first proposed in [82] to approximate the continuous inequality constraints on the state by constraints in integral form, called constraints in canonical form. The constraint transcription method is later extended in [80] to handle continuous inequality constraints on the state and control. Thus, by using the control parameterization method together with the constraint transcription technique, an optimal control problem subject to continuous inequality constraints on the state and/or control is approximated by a sequence of optimal parameter selection problems subject to canonical constraints, where the cost function is to be minimized with respect to the control parameters subject to the constraints being satisfied. Each of the resulting optimal parameter selection problems can be regarded as a mathematical programming problem solvable by gradient-based optimization techniques, and hence many existing optimization software packages can be readily used. The control parameterization technique is used extensively in the literature (see [45, 93, 149]). In [146], a survey on developments of the technique is presented. It is observed that the technique is applicable to a wide range of optimal control problems. In particular, several computational algorithms supported by sound theoretical convergence

analysis are presented in [127] and [16] for dealing with a variety of different classes of optimal control problems.

It is intuitively clear that the switching times of the control should also be taken as decision variables. However, it is numerically sensitive if the gradients of the cost function and the canonical constraints with respect to these switching times are used in the optimization process (see Chapter 5 of [80]). Thus, a time-scaling transform is introduced in [66,81], where it is called the control parameterization enhancing transform (CPET), to map these switching times to fixed knot points in a new time horizon via introducing a new control variable, called the time-scaling control, and an additional differential equation describing the relationship between the original time variable and the new one. Thus, under the time-scaling transform, the optimal control problem subject to continuous inequality constraints on the state and/or control is also approximated by a sequence of optimal parameter selection problems subject to canonical constraints. Each of these problems can be solved as a nonlinear mathematical programming problem.

The optimal control software package MISER [91] is an implementation of algorithms based on the control parameterization technique [80] and the constraint transcription method [89]. It can be used in conjunction with the time-scaling transform [81]. NUDOC-CCS [19] is another optimal control software based on the control parameterization approach. It is used for both the simulation and the optimization of dynamical systems. While it lacks some of the flexibility of the MISER 3.3 package, it has some additional features such as an adaptive grid refinement strategy and efficient posterior sensitivity analysis.

## 1.4 **Parameter identification for time-delay systems**

For a given practical process, it is required to construct a mathematical model to describe the interactions between various factors so as to behave as a whole unit. In practice, there are delay effects on the process of the system. For example, in an imperfectly mixed system [63,64], it takes some time for the changes of mass and energy in particular parts of the vessel to reach the rest of the vessel. Time-delays do arise in many real world situations, including chemical tank reactors [86], aerospace engineering [143], chromatography [107], and power converters and batteries [153]. For such systems, the mathematical model is often expressed as a system of ordinary differential equations involving time-delay arguments.

A typical nonlinear time-delay system is given below:

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{x}(t - \alpha_1), \dots, \mathbf{x}(t - \alpha_m), \mathbf{u}(t), \mathbf{u}(t - \beta_1), \dots, \mathbf{u}(t - \beta_p), \boldsymbol{\zeta}), \quad t \in (0, T], \quad (1.11)$$

$$\mathbf{x}(t) = \boldsymbol{\phi}(t), \quad t \leq 0, \quad (1.12)$$

where  $T > 0$  is the terminal time;  $\boldsymbol{\zeta} \in \mathbb{R}^q$  is a vector of system parameters;  $\mathbf{x}(t) \in \mathbb{R}^n$  is the state of the system;  $\mathbf{x}(t - \alpha_i) \in \mathbb{R}^n$ ,  $i = 1, \dots, m$ , are delayed states, meaning that if  $\mathbf{x}(t)$  is the value of  $\mathbf{x}$  at the time point  $t$ , then  $\mathbf{x}(t - \alpha_i)$  denotes the value of  $\mathbf{x}$  at the time point  $t - \alpha_i$ ;  $\mathbf{u}(t) \in \mathbb{R}^r$  is the control;  $\mathbf{u}(t - \beta_k) \in \mathbb{R}^r$ ,  $k = 1, \dots, p$ , are the delayed controls. Let all the state time-delays and all the control time-delays be referred to collectively as  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_m]^\top$  and  $\boldsymbol{\beta} = [\beta_1, \dots, \beta_p]^\top$ , respectively. Furthermore,  $\mathbf{f} : \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^{nm} \times \mathbb{R}^r \times \mathbb{R}^{pr} \times \mathbb{R}^q \rightarrow \mathbb{R}^n$  is a given function, and  $\boldsymbol{\phi} : \mathbb{R} \rightarrow \mathbb{R}^n$  is a given function representing the state before the time  $t = 0$ . In practice, the state of the system may not be available for observation. Instead, what we can observe is some function of the state variables given by

$$\mathbf{y}(t) = \Psi(\mathbf{x}(t)), \quad (1.13)$$

where  $\mathbf{y}(t) \in \mathbb{R}^m$  with  $m < n$ .

In a practical process, the system parameters and the time-delays are often not known. However, once suitable values for the system parameters and time-delays are chosen, the time-delay system (1.11) can be solved and hence the corresponding output can be computed from the output system (1.13). On the other hand, the output of the real process can be measured at certain time points over the time horizon. The sum of the differences between the computed output and real output at these time points is called the *mismatch*. To identify the system parameters and time-delays, it amounts to minimize the mismatch with respect to the system parameters and time-delays. This is, in fact, an optimal parameter selection problem.

We now return to the time-delay system (1.11)-(1.12). Once the system parameters and the time-delays are identified, then the mathematical model is completely specified. It is then possible to construct optimal control algorithms based on the control parameterization method to synthesize an optimal control for the system (e.g. [79, 94]). The problem of identifying the system parameters and the time-delays based on a given time series data is a key problem in the study of time-delay systems [74]. Such problems are known as *parameter identification problems*.

There are many results available in the literature pertaining to parameter identification problems. An exact least squares algorithm for single time-delay estimation is studied in [116]. Algebraic techniques [85] and the steepest descent algorithm [134] are proposed to determine the input delays. In [95], information theory is used to identify multiple time-

delays from a time series. In [95], it is assumed that each nonlinear term in the dynamical system contains at most one unknown delay. In [57], a genetic algorithm is utilized for the identification of multiple time-delays. However, this method depends critically on the initial guess of the parameters that are to be identified. Lyapunov functions are used to design delay estimators in [141], where it is shown that for an appropriately chosen guess for each of the delays, the approximate delays will converge to the real ones after a finite number of iterations. In [125], a class of systems is considered, where the system dynamics are expressed as the sum of a finite number of nonlinear terms and each nonlinear term involves only one unknown state delay. There is no unknown system parameters involved. The gradient formula of the objective function with respect to the time-delays is derived. This gradient formula is expressed in terms of the solutions of the dynamical system and an auxiliary delay-differential system, both of which are to be solved forward in time over the time horizon. With this gradient formula, existing gradient-based optimization techniques can be incorporated to solve the parameter identification problem for this time-delay system.

## 1.5 Control methods for time-delay systems

### 1.5.1 Optimal control for time-delay systems

Consider the process evolving according to system (1.11)-(1.12) over the time horizon  $[0, T]$ , where system parameters  $\zeta$  and time-delays  $\alpha$  and  $\beta$  are assumed to be given. Let  $\mathbf{U}$  be a given compact and convex subset of  $\mathbb{R}^r$ , and let  $\gamma : [-\beta, 0) \rightarrow \mathbb{R}^r$ . A function  $\mathbf{u} : [-\beta, T] \rightarrow \mathbb{R}^r$  such that  $\mathbf{u}(t) = \gamma(t)$  on  $[-\beta, 0)$  and  $\mathbf{u}(t) \in \mathbf{U}$  for almost all  $t \in [0, T]$  is called an admissible control. Let  $\mathcal{U}$  be the class of all such admissible controls. A simple optimal control problem for time-delay systems may now be stated as follows:

**Problem ( $\tilde{P}_4$ ).** *Consider system (1.11)-(1.12), find a control  $\mathbf{u} \in \mathcal{U}$  such that the cost function (1.4) is minimized.*

### 1.5.2 Model predictive control for time-delay systems

Model predictive control (MPC) is also called the receding horizon control. It was first introduced in 1960s [50], and has since become popular in areas such as chemical processes and paper industries. This is because the MPC algorithm is simple and intuitive. The main structure of MPC is show in Figure 1.1. In Figure 1.1,  $T_s$  is the sample time;  $t_k$  is the present time;  $n_c T_s$  and  $n_p T_s$  are, respectively, the length of the control horizon and the prediction horizon, where  $n_c \leq n_p$  are positive integers. We can see that, in each sample time prered, MPC involves the prediction of the system output over a finite prediction period  $[t_k, t_k + n_p T_s]$  by using process model based on empirical data fitting

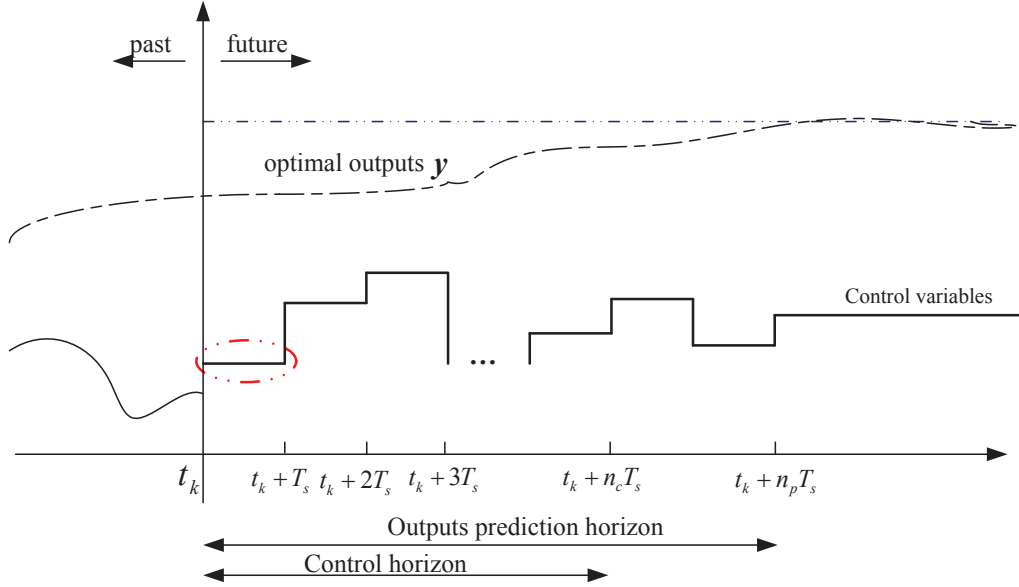


Figure 1.1: Structure of MPC

or dynamic model based on fundamental mass and energy balances. The optimal control is obtained through solving the following open-loop constrained optimal control problem over the prediction period  $[t_k, t_k + n_p T_s]$ .

**Problem ( $\tilde{P}_5$ ).** Consider system (1.11) over the time horizon  $[t_k, t_k + n_p T_s]$  with given system parameters  $\zeta$  and time-delays  $\alpha$  and  $\beta$ , the following optimal control problem is minimized

$$\begin{aligned} \min \quad & J_{\mathbf{u}(t_k), \dots, \mathbf{u}(t_k + n_c T_s)} = \sum_{i=1}^{n_p} \{\mathbf{y}(t_k + iT_s) - \mathbf{y}^*\}^2 + \sum_{i=1}^{n_c} \Delta \mathbf{u}(t_k + iT_s)^\top \mathbf{R} \Delta \mathbf{u}(t_k + iT_s) \\ \text{s.t.} \quad & \dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t)), \quad t \in [t_k, t_k + n_p T_s], \\ & \underline{\alpha} \leq \Delta \mathbf{u}(t) \leq \bar{\alpha}, \\ & \underline{\beta} \leq \mathbf{u}(t) \leq \bar{\beta}, \\ & \mathbf{g}(t, \mathbf{x}, \mathbf{u}) \leq 0, \end{aligned}$$

with initial conditions

$$\mathbf{u}(t) = \boldsymbol{\gamma}(t), \quad \mathbf{x}(t) = \boldsymbol{\phi}(t), \quad t \leq t_k,$$

and the output function

$$\mathbf{y}(t) = \boldsymbol{\Psi}(\mathbf{x}(t)), \quad t \in [t_k, t_k + n_p T_s],$$

where  $\mathbf{x} = [x_1, \dots, x_n]^\top \in \mathbb{R}^n$  is the state vector and  $\boldsymbol{\phi} = [\phi_1, \dots, \phi_n]^\top \in \mathbb{R}^n$  is its

initial condition;  $\mathbf{u} = [u_1, \dots, u_r]^\top \in \mathbb{R}^r$  is the control vector and  $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_r]^\top \in \mathbb{R}^r$  is its initial condition;  $\mathbf{f} = [f_1, \dots, f_n]^\top \in \mathbb{R}^n$  is the dynamic system function;  $\Delta \mathbf{u} = [\Delta u_1, \dots, \Delta u_r]^\top \in \mathbb{R}^r$  is the change rate of the control at time  $t_k + iT_s$ ,  $i = 1, \dots, n_c$ , with  $\Delta \mathbf{u}(t_k + iT_s) = \mathbf{u}(t_k + iT_s) - \mathbf{u}(t_k + (i-1)T_s)$ ;  $\mathbf{R} \in \mathbb{R}^{rr}$  is a given weighting matrix;  $\mathbf{g} = [g_1, \dots, g_{n_g}]^\top \in \mathbb{R}^{n_g}$  denotes the constraint function;  $\underline{\boldsymbol{\beta}} = [\underline{\beta}_1, \dots, \underline{\beta}_r]^\top \in \mathbb{R}^r$  and  $\bar{\boldsymbol{\beta}} = [\bar{\beta}_1, \dots, \bar{\beta}_r]^\top \in \mathbb{R}^r$  are the lower and upper bounds of the control, respectively;  $\underline{\boldsymbol{\alpha}} = [\underline{\alpha}_1, \dots, \underline{\alpha}_r]^\top \in \mathbb{R}^r$  and  $\bar{\boldsymbol{\alpha}} = [\bar{\alpha}_1, \dots, \bar{\alpha}_r]^\top \in \mathbb{R}^r$  are the lower and upper bounds of the change rate of the control, respectively;  $\mathbf{y} = [y_1, \dots, y_m]^\top \in \mathbb{R}^m$  is the system output;  $\boldsymbol{\Psi} = [\Psi_1, \dots, \Psi_m]^\top \in \mathbb{R}^m$  is the given output function;  $\mathbf{y}^* = [y_1^*, \dots, y_m^*]^\top \in \mathbb{R}^m$  is the given reference trajectory. Moreover,  $\mathbf{u}(t) = \mathbf{u}(t_k + n_c T_s)$  if  $t > t_k + n_c T_s$ .

Problem ( $\tilde{P}_5$ ) is solvable by gradient based algorithms. Note that for MPC, only the control adjustments for the next instant  $t_k + T_s$  of the optimal control over  $[t_k, t_k + T_s]$  is used. Then, the time is moved forwards to  $t_k := t_k + T_s$  and a new optimal control problem in the form of Problem ( $\tilde{P}_5$ ) is solved for subsequent sampling periods. This process is repeated.

Amongst the existing literature, the MPC can be divided into two categories: (i) Linear MPC; and (ii) nonlinear MPC. For the linear MPC, it involves solving a sequence of open-loop optimal control problems, each of which contains a linear system dynamics, a quadratic cost function and linear constraints over a future time horizon. It has been widely used in practice (see, for example, [1, 68]). In the linear MPC, the resulting optimal control problems can be solved as quadratic programming problems involving linear dynamical systems. Since the control to be applied over the next instant is to be calculated online, the MPC can only be used for systems with slow dynamics and long sampling time. The linear MPC is extended to nonlinear systems in [15], where the process model is approximately linear over a small operating range. Then, the linearization of the nonlinear process model in a small operating range can be carried out. However, it achieves poor performance when the process model could not be approximated accurately by linear model.

This problem is overcome in the nonlinear model predictive control (NMPC) [136]. The NMPC is characterized by nonlinear system model, nonlinear constraints, and non-quadratic cost function. The NMPC techniques have been successfully implemented in many large scale industrial processes since 1980s (see, for example, [31, 61, 139] for time-delay industrial processes). However, it may not yet be suitable for on-line implementation. To be more specific, let us mention an example, where the NMPC is applied to an evaporation process in [61]. In this example, it takes nearly one hour of computational time to compute the control by using the NMPC on a workstation with dual Pentium III Xeon processors for one hour of simulation of the evaporation process using the control obtained. The computational time grows as the prediction horizon is increased. Since many practical processes have large time-delays [32, 61], a short prediction horizon may



not be able to capture the effect on the changes of the state for such processes. Thus, a long prediction horizon is required, but then the consequent computational burden will be much increased so that it will become impractical for on-line implementation. Other studies also show that for optimal control problems involving nonlinear dynamics with high dimension and state constraints, especially continuous state constraints, NMPC requires a large computation time in each sampling period. Indeed, the computation time needed for NMPC is much more than that of the linear MPC. The high computational burden has limited its application to optimal control problems involving nonlinear dynamics and subject to constraints on the state and control [144].

To overcome the drawback of requiring a large computation time, many improved methods have been proposed. The most important one is the linear matrix inequalities based MPC [106, 158]. As the linear matrix inequalities based optimization can be solved in polynomial time, the computation time for MPC is reduced. However, this approach is only applicable to problems involving linear dynamics and linear constraints, while the cost function is quadratic. A decreasing horizon MPC is introduced in [87]. As the control horizon decreases from one iteration to the next, the computation time will clearly be decreased. In [17], the off-time calculation is carried on a discrete set of constraints and terminal costs. The results are then utilized in the implementation of the MPC. However, this method can only be used for systems with lower dimension. Some on-line optimization methods, such as swarm-staring method and grid method, have been proposed to speed-up the computation time for MPC [156]. However, they are yet to be realized in practice.

## 1.6 Max-Min optimal control problems

A max-min optimal control problem involves finding a control such that the minimum of a set of objectives is maximized. Thus, its cost function is a max-min function which is non-smooth and non-differentiable. This problem has extensive real world applications in engineering optimality design [115], electronic circuit design [21], robust control design [4], economics [18], and the staffing problems in call centers [9]. Let us look at an example arising in the study of a separation process, where several different kinds of products are drained from a single outlet. It is required to ensure the purification of each product. For this, it is required that one product is totally drained at one time. Thus, the duration time between two successive outlet times should be maximized. However, large duration times will increase the total operation time and decrease the productivity. In practice, the total operation time should be as short as possible. Thus, the optimal control problem for this separation process can be formulated as: find a control such that the minimum ratio of the duration time between each pair of two successive outlet times to the total operation time is maximized subject to some constraints due to engineering specifications.

As solving a max-min optimal control problem basically involves solving a sequence of approximate max-min optimization problems, we will mention some approaches proposed for solving max-min optimization problems. One way of solving a max-min optimization problem is to view the max-min objective as multiple objectives. Then, the aim is to find an efficient solution over a set of possible solutions, which is called the Pareto set. Thus, it is required to first find the Pareto solution set. Then, a sorting procedure is applied to evaluate the solutions. The Pareto solutions could, in principle, be found by Pareto-set based optimization methods, such as genetic algorithm [3], evolution algorithm [4, 10] and particle swarm optimization algorithm [155]. However, these methods are heuristic methods.

It is known that max-min problems and min-max problems are equivalent. For closely related minmax problems, smoothing techniques are proposed to convert min-max problems into simple, smooth, unconstrained or constrained optimization problems. Then, standard unconstrained or constrained minimization techniques can be used to solve these converted problems. More specifically, the smoothing technique [53] is used to approximate a max function by a smooth function, while the maximum entropy function is used in [72] to approximate the inner maximum objective function by a continuous smooth function. However, if the accuracy requirement on the approximation is high, then the smooth approximating problems will become ill-conditioned. Hence, when applied to these problems, the unconstrained optimization techniques may experience numerical difficulties, leading to slow convergence—and in some cases no convergence at all. In addition to this, two distinct search directions based algorithms are proposed in [21] to solve min-max problems directly, while interior-point method is used in [43, 52].

For max-min optimal control problems, we wish to mention the following two methods. For the first method reported in [8], the max-min optimal control problem is first converted to min-max optimal control problem, then a smoothing technique is used to transform the min-max optimal control problem into a standard optimal control problem in the form of Bolza with additional inequality constraints on the state/or control variables. This standard optimal control problem can be solved by existing methods and theories. This idea has been widely adopted, see, for example, [16, 88, 97]. For another method reported in [90], the parametrization technique is used to approximate the max-min optimal control problem by a sequence of max-min optimal parameter selection problems. Then, each of these approximate problems is shown to be equivalent to a standard min-max optimization problem. Hence, it is solvable by existing optimization software for solving minmax optimization problems, such as FFSQP 3.7 [76] and CONSOL-OPTCAD [99].

## 1.7 Overview of the thesis

In previous sections, we presented brief surveys on computational methods for solving optimal control problems, time-delay optimal control problems and maxmin optimal control problems. Furthermore, a brief introduction to MPC for time-delay optimal control problems is also given.

The purpose of this thesis is to present new computational methods for several classes of practical optimal control problems. They are briefly mentioned below.

In Chapter 2, we consider a general nonlinear time-delay system described by (1.11)-(1.12), where not only the unknown system parameters, but also the unknown time-delays need to be identified. Furthermore, the input is smooth function. Let  $\alpha \in \mathcal{T}$ ,  $\beta \in \mathcal{B}$ , and  $\zeta \in \mathcal{Z}$  denote, respectively, the state time-delays, input time-delays and system parameters, where  $\mathcal{T}$ ,  $\mathcal{B}$  and  $\mathcal{Z}$  are the set of candidate parameter vectors for the corresponding variables. Our parameter identification problem for time-delay system can be described as follows:

**Problem ( $\tilde{P}_6$ ).** Choose  $\alpha \in \mathcal{T}$ ,  $\beta \in \mathcal{B}$ , and  $\zeta \in \mathcal{Z}$  such that the cost function

$$J(\alpha, \beta, \zeta) = \sum_{l=1}^M |\mathbf{y}(t_l | \tau, \beta, \zeta) - \hat{\mathbf{y}}^l|^2. \quad (1.14)$$

is minimized, where  $|\cdot|$  denotes the usual Euclidean norm. Here,  $\hat{\mathbf{y}}^l$  denotes the value of the output function (1.13) at the sample time  $t_l$ ,  $l = 1, \dots, M$ . The cost function measures the discrepancy between predicted and observed system output.

We develop a unified gradient-based computational approach that involves solving Problem ( $\tilde{P}_6$ ). Since the delays and parameters influence the cost function *implicitly* through the dynamic system, in this computational method, the gradients of the cost function with respect to the delays and system parameters are derived. They are obtained by solving a set of auxiliary delay-differential systems from  $t = 0$  to  $t = T$ . Then, the delays and parameters are determined simultaneously through solving a dynamic optimization problem using existing optimization techniques. Two nonlinear parameter identification problems involving time-delay systems are solved by using the method proposed. From numerical simulations, it is clearly indicated that this algorithm is effective.

In Chapter 3, we consider a more difficult time-delay identification problem, where the input function of the nonlinear time-delay system is piecewise-constant. The main difficulties with this problem are: i) Since the input function is discontinuous, the dynamics are clearly discontinuous with respect to the input delay. Thus, the results obtained in Chapter 2 cannot be used to determine the state variation with respect to input time-delay; and ii) Problem ( $\tilde{P}_6$ ) in this case can be restated as a switched system optimal control problem. Unfortunately, the well-known *time-scaling transform* technique for solving optimal

control problems involving switched systems (see [118, 124, 132, 151]) is not applicable to time-delay systems such as system (1.11)-(1.12) defined above (see the discussion in [33]). Thus, we propose a new computational approach, which is based on a novel derivation of the cost function's gradient. We then apply this approach to estimate the time-delays in two industrial chemical processes—a zinc sulphate purification process and a sodium aluminate evaporation process. Numerical simulations demonstrate the effectiveness of this algorithm.

In Chapter 4, we consider a general class of optimal control problems for time-delay systems subject to continuous inequality constraints on states and controls. In other words, these constraints must be satisfied for each  $t \in [0, T]$ . This problem can be stated as follows.

**Problem ( $\tilde{P}_7$ ).** *Given system (1.11)-(1.12) with given  $\alpha$ ,  $\beta$ ,  $\zeta$ , find a control  $\mathbf{u} \in \mathcal{U}$  such that the cost function (1.4) is minimized subject to the following constraints:*

$$g_i(t, \mathbf{x}(t), \mathbf{u}(t)) \leq 0, \quad i = 1, \dots, N_c, \quad \forall t \in [0, T].$$

An efficient gradient-based computational method is devised for solving this optimal control problem. In this method, the control parameterization technique is used to approximate the control by a piecewise-constant function (it could also be approximated by a piecewise linear or piecewise smooth function) with possible discontinuities at the  $N - 1$  pre-assigned partition points. The heights of the piecewise-constant function are regarded as decision variables, which are referred to collectively as the control parameter vector. Then, the constrained time-delay optimal control problem is approximated by a sequence of optimal parameter selection problems involving time-delay dynamical system and subject to continuous inequality constraints on the state and boundedness constraints on the control parameter vector. For the continuous inequality constraints, they are transformed, by using the constraint transcription method, into equivalent equality constraints in integral form. However, the integrands of these equality constraints are nonlinear and nonsmooth. Thus, a local smoothing technique is used to approximate these nonsmooth integrands by smooth functions. Then, the equality constraints in integral form are approximated by inequality constraints in integral form, where their integrands are approximating smooth functions. There are two parameters involved in the inequality constraints in integral form—one controls the accuracy of the approximation and the other controls the feasibility of the original constraints satisfaction. Now, by using penalty function ideas, the summation of these inequality constraints in integral form is appended to the cost function to form an augmented cost function. In this way, the constrained time-delay optimal control problem is approximated by a sequence of optimal parameter selection problems involving time-delay dynamical system and subject to only boundedness constraints on the control parameter vector. The gradient of the

augmented cost function with respect to the control parameter vector is derived. On this basis, an effective gradient-based optimization method is developed for solving each of these optimal parameter selection problems with simple boundedness constraints on the control parameter vector. The optimal control parameter vector obtained can then be used to construct a piecewise-constant control for the original constrained time-delay optimal control problem. Supporting convergence results are established. In particular, it is shown that when the penalty factor is sufficiently large, then the optimal piecewise-constant control obtained will satisfy the continuous inequality constraints of the original problem. Furthermore, when the number of partition points is increased, the cost corresponding to the optimal piecewise-constant control will converge to the true optimal cost. The computational method proposed is applied to an optimal control problem for an industrial-scale evaporation process with long time-delays, where the mass units of live steam consumption used for evaporating one unit of water is minimized subject to the requirements for the product concentrations and the levels in the evaporators being satisfied.

In Chapter 5, we consider the following set of differential equations:

$$\begin{aligned}\dot{x}_1(t) &= f_1(t, \mathbf{x}(t), \mathbf{u}(t)), & t \in [0, t_1], \\ &\vdots \\ \dot{x}_n(t) &= f_n(t, \mathbf{x}(t), \mathbf{u}(t)), & t \in [0, t_n].\end{aligned}$$

where  $t_i > 0$ ,  $i = 1, \dots, n$ , are the unknown final times for the  $i$ th state dynamic function, respectively; and if  $i \neq j$ , then  $t_i \neq t_j$ .

This problem arises in the study of gradient elution chromatography, which is used to separate different kinds of components in a solution. A typical chromatography system consists of a column containing an adsorbent (called the stationary phase) and a liquid that flows through the column (called the mobile phase). The mixture to be separated is injected into the mobile phase and flows through the column. Since different components are attracted to the adsorbent at different grades, they exit the gradient elution column at different times. The main concern of this process is to maximize separation efficiency so as to ensure the purification of each component. Thus, our goal is to find a control such that the minimum duration between successive final times is maximized. This max-min problem has three non-standard characteristics: (i) The objective function is non-smooth; (ii) each state variable is defined over a different time horizon; and (iii) the ordering of the final times is unknown. In this max-min optimal control problem, there are multiple characteristics times, max-min objective function, and binary decision variables. We will propose an efficient gradient-based computational method to solve this complicated optimal control problem. To demonstrate the effectiveness of the computational method proposed, two numerical examples are solved. One of these two examples is an optimal

control problem involving a real chromatography process. The results obtained are highly promising.

Finally, in Chapter 6, we summarize the main contributions of the thesis and discuss some interesting directions for future research.

---

---

# CHAPTER 2

---

## Identification of time-delays and parameters for nonlinear systems

### 2.1 Introduction

In this chapter, we consider a general nonlinear delay-differential system with unknown time-delays and unknown system parameters. We formulate the problem of identifying these unknown quantities as a nonlinear optimization problem in which the cost function measures the least-square error between predicted output and observed system output. This type of parameter identification problem is previously considered in [125] for a simple class of systems in which each nonlinear component contains at most one unknown delay with no unknown system parameters. However, in many real-world systems, such as the purification process of zinc sulphate solution [93], the nonlinear terms contain *both* delays and parameters that need to be identified. Our goal in this chapter is to develop an efficient method to identify the unknown delays and the unknown parameters in a complicated time-delay dynamical system. We will introduce a set of auxiliary delay-differential systems. Then, we will show that the gradient of the least-square cost function can be expressed in terms of the solutions of these auxiliary systems. A numerical integration is used to solve the auxiliary systems, and thereby we obtain the gradient of the cost function, which is the main information needed to solve the parameter identification problem as a nonlinear optimization problem by using numerical optimization techniques. Based on this idea, a computational algorithm is developed for identifying the unknown time-delays and system parameters in a general nonlinear system. We demonstrate the effectiveness of the proposed algorithm on two nonlinear parameter identification problems, one of which is the parameter identification problem for the zinc sulphate purification process.

## 2.2 Problem formulation

Consider the following nonlinear time-delay system:

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \boldsymbol{\zeta}), \quad t \in [0, T], \quad (2.1)$$

$$\mathbf{x}(t) = \boldsymbol{\phi}(t), \quad t \leq 0, \quad (2.2)$$

where  $T > 0$  is a given *terminal time*;  $\mathbf{x}(t) = [x_1(t), \dots, x_n(t)]^\top \in \mathbb{R}^n$  is the *state*;  $\tilde{\mathbf{x}}(t) = [\mathbf{x}(t - \tau_1)^\top, \dots, \mathbf{x}(t - \tau_m)^\top]^\top \in \mathbb{R}^{nm}$  is the *delayed state*; and  $\boldsymbol{\zeta} = [\zeta_1, \dots, \zeta_r]^\top \in \mathbb{R}^r$  is a vector of unknown *system parameters*. Furthermore,  $\mathbf{f} : \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^{nm} \times \mathbb{R}^r \rightarrow \mathbb{R}^n$  and  $\boldsymbol{\phi} : \mathbb{R} \rightarrow \mathbb{R}^n$  are given functions.

The time-delays in (2.1)-(2.2) are unknown quantities that need to be determined. We assume that the  $i$ th time-delay is in the interval  $[a_i, b_i]$ , where  $a_i$  and  $b_i$  are given constants such that  $0 \leq a_i < b_i$ . Hence, the unknown time-delays satisfy the following bound constraints:

$$a_i \leq \tau_i \leq b_i, \quad i = 1, \dots, m. \quad (2.3)$$

Any vector  $\boldsymbol{\tau} = [\tau_1, \dots, \tau_m]^\top \in \mathbb{R}^m$  that satisfies (2.3) is called a *candidate time-delay vector*. Let  $\mathcal{T}$  denote the set of all such candidate time-delay vectors.

In addition to the time-delays, the system parameters in (2.1)-(2.2) are also unknown quantities that need to be determined. We suppose that

$$c_j \leq \zeta_j \leq d_j, \quad j = 1, \dots, r, \quad (2.4)$$

where  $c_j$  and  $d_j$  are given real numbers such that  $0 \leq c_j < d_j$ . Note that there is no loss of generality in assuming that  $c_j \geq 0$ ; if  $c_j < 0$ , then we may replace  $\zeta_j$  with  $\zeta_j + c_j$ . Any vector  $\boldsymbol{\zeta} = [\zeta_1, \dots, \zeta_r]^\top \in \mathbb{R}^r$  that satisfies (2.4) is called a *candidate parameter vector*. Let  $\mathcal{Z}$  denote the set of all such candidate parameter vectors.

The output of system (2.1)-(2.2) is given by

$$\mathbf{y}(t) = \mathbf{g}(\mathbf{x}(t), \boldsymbol{\zeta}), \quad t \in [0, T], \quad (2.5)$$

where  $\mathbf{g} : \mathbb{R}^n \times \mathbb{R}^r \rightarrow \mathbb{R}^p$  is a given function.

We assume that the following conditions are satisfied.

**(2.A.1).** The given functions  $\mathbf{f}$  and  $\mathbf{g}$  are continuously differentiable, and  $\boldsymbol{\phi}$  is twice continuously differentiable.

**(2.A.2).** There exists a real number  $L_1 > 0$  such that

$$|\mathbf{f}(t, \mathbf{x}, \tilde{\mathbf{x}}, \boldsymbol{\zeta})| \leq L_1(1 + |\mathbf{x}| + |\tilde{\mathbf{x}}| + |\boldsymbol{\zeta}|), \quad (t, \mathbf{x}, \tilde{\mathbf{x}}, \boldsymbol{\zeta}) \in \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^{nm} \times \mathbb{R}^r,$$



where  $|\cdot|$  denotes the Euclidean norm.

On the basis of assumptions (2.A.1) and (2.A.2), the dynamic system (2.1)-(2.2) admits a unique solution corresponding to each pair  $(\boldsymbol{\tau}, \boldsymbol{\zeta}) \in \mathcal{T} \times \mathcal{Z}$  [111]. We denote this solution by  $\boldsymbol{x}(\cdot|\boldsymbol{\tau}, \boldsymbol{\zeta})$ . Substituting  $\boldsymbol{x}(\cdot|\boldsymbol{\tau}, \boldsymbol{\zeta})$  into (2.5) gives  $\boldsymbol{y}(\cdot|\boldsymbol{\tau}, \boldsymbol{\zeta})$ , the predicted system output corresponding to  $(\boldsymbol{\tau}, \boldsymbol{\zeta}) \in \mathcal{T} \times \mathcal{Z}$ . More formally,

$$\boldsymbol{y}(t|\boldsymbol{\tau}, \boldsymbol{\zeta}) = \boldsymbol{g}(\boldsymbol{x}(t|\boldsymbol{\tau}, \boldsymbol{\zeta}), \boldsymbol{\zeta}), \quad t \leq T. \quad (2.6)$$

Suppose that the output from system (2.1)-(2.2) has been measured experimentally at times  $t = t_l$ ,  $l = 1, \dots, q$ , where each  $t_l \in [0, T]$ . Let  $\hat{\boldsymbol{y}}^l \in \mathbb{R}^p$  denote the measured output at time  $t = t_l$ . Then the problem of identifying the unknown time-delays and system parameters can be formulated mathematically as follows.

**Problem (P).** *Choose  $\boldsymbol{\tau} \in \mathcal{T}$  and  $\boldsymbol{\zeta} \in \mathcal{Z}$  such that the following cost function:*

$$J(\boldsymbol{\tau}, \boldsymbol{\zeta}) = \sum_{l=1}^q |\boldsymbol{y}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta}) - \hat{\boldsymbol{y}}^l|^2 \quad (2.7)$$

*is minimized, where  $|\cdot|$  denotes the usual Euclidean norm.*

Problem (P) is a nonlinear dynamic optimization problem whose decision variables are the delays and system parameters in system (2.1)-(2.2). We need to select optimal values for these delays and parameters so that the predicted system output best fits the experimental data. There are very few optimization techniques available in the literature for time-delay systems. In the existing literature, the delays are often assumed to be fixed and known (see, for example, [79, 119, 150]). Problem (P) is unique in that the delays are not fixed, but are decision variables to be chosen optimally. The cost function in Problem (P) is also highly non-standard, as it depends on the system's state at a set of discrete time points, not just at the terminal time. Such cost functions have been considered in [120, 121] for non-delay systems, and in [93] for systems with fixed delays. However, the computational techniques developed in these references are not applicable to Problem (P) because the time-delays in system (2.1)-(2.2) are decision variables to be identified.

## 2.3 Preliminaries

Throughout this subsection, let  $k \in \{1, \dots, m\}$  and  $(\boldsymbol{\tau}, \boldsymbol{\zeta}) \in \mathcal{T} \times \mathcal{Z}$  be arbitrary but fixed. For simplicity, we write  $\boldsymbol{x}(t)$  instead of  $\boldsymbol{x}(t|\boldsymbol{\tau}, \boldsymbol{\zeta})$ , and  $\boldsymbol{x}^\epsilon(t)$  instead of  $\boldsymbol{x}(t|\boldsymbol{\tau} + \epsilon \boldsymbol{e}^k, \boldsymbol{\zeta})$ , where  $\boldsymbol{e}^k$  denotes the  $k$ th unit basis vector in  $\mathbb{R}^m$ .

Define

$$I = [a_k - \tau_k, b_k - \tau_k].$$

Note that  $I \neq \emptyset$  and  $0 \in I$ . Clearly,

$$\epsilon \in I \iff \tau + \epsilon e^k \in \mathcal{T}.$$

For each  $\epsilon \in I$ , define

$$\varphi^\epsilon(t) = \mathbf{x}^\epsilon(t) - \mathbf{x}(t), \quad t \leq T,$$

and

$$\boldsymbol{\theta}^{\epsilon,i}(t) = \mathbf{x}^\epsilon(t - \tau_i - \epsilon \delta_{ki}) - \mathbf{x}(t - \tau_i), \quad t \leq T, \quad i = 1, \dots, m,$$

where  $\delta_{ki}$  denotes the Kronecker delta function. Furthermore, let

$$\boldsymbol{\theta}^\epsilon(t) = [(\boldsymbol{\theta}^{\epsilon,1}(t))^\top, \dots, (\boldsymbol{\theta}^{\epsilon,m}(t))^\top]^\top \in \mathbb{R}^{nm}, \quad t \leq T.$$

Clearly,

$$\boldsymbol{\theta}^{\epsilon,i}(t) = \varphi^\epsilon(t - \tau_i), \quad t \leq T, \quad i \neq k, \quad (2.8)$$

$$\varphi^\epsilon(t) = \mathbf{0}, \quad t \leq 0. \quad (2.9)$$

In the sequel, we will use the notation  $\frac{\partial}{\partial \tilde{\mathbf{x}}^i}$  to denote partial differentiation with respect to the  $i$ th delayed state in  $\tilde{\mathbf{x}}(t)$  (i.e. partial differentiation with respect to  $\mathbf{x}(t - \tau_i)$ ).

Now, define

$$\boldsymbol{\chi}^\epsilon(t) = \begin{cases} \dot{\phi}(t), & \text{if } t \leq 0, \\ \mathbf{f}(t, \mathbf{x}^\epsilon(t), \mathbf{x}^\epsilon(t - \tau_1 - \epsilon \delta_{k1}), \dots, \mathbf{x}^\epsilon(t - \tau_m - \epsilon \delta_{km}), \boldsymbol{\zeta}), & \text{if } t \in (0, T]. \end{cases} \quad (2.10)$$

We immediately see that for almost all  $t \in (-\infty, T]$ ,

$$\dot{\boldsymbol{x}}^\epsilon(t) = \boldsymbol{\chi}^\epsilon(t). \quad (2.11)$$

Let  $\bar{b} > 0$  be a fixed constant such that

$$\bar{b} = \max_{i=1, \dots, m} \{b_i\}.$$

We have the following lemma.

**Lemma 2.1.** *There exists a positive real number  $L_2 > 0$  such that for each  $\epsilon \in I$ ,*

$$|\mathbf{x}^\epsilon(t)|, |\boldsymbol{\chi}^\epsilon(t)| \leq L_2, \quad t \in [-\bar{b}, T], \quad (2.12)$$

*Proof.* Recall that  $\mathbf{x}(s) = \boldsymbol{\phi}(s)$  is given for all  $s \leq 0$ . By (2.A.1),  $\boldsymbol{\phi}(s)$  is twice differentiable. Clearly, there exist positive numbers  $\alpha_1$  and  $\alpha_2$  such that

$$|\boldsymbol{\phi}(s)| \leq \alpha_1, \quad s \in [-\bar{b}, 0], \quad (2.13)$$

$$|\dot{\boldsymbol{\phi}}(s)| \leq \alpha_2, \quad s \in [-\bar{b}, 0]. \quad (2.14)$$

For brevity, we denote

$$\tilde{\mathbf{x}}^\epsilon(t) = [\mathbf{x}^\epsilon(t - \tau_1 - \epsilon\delta_{k1})^\top, \dots, \mathbf{x}^\epsilon(t - \tau_m - \epsilon\delta_{km})^\top]^\top.$$

For each  $s \in [0, T]$ , we have

$$\mathbf{x}^\epsilon(t) = \mathbf{x}(0) + \int_0^t \mathbf{f}(s, \mathbf{x}^\epsilon(s), \tilde{\mathbf{x}}^\epsilon(s), \boldsymbol{\zeta}) ds, \quad t \in [0, T]. \quad (2.15)$$

Since  $\boldsymbol{\zeta}$  is bounded in  $\mathcal{Z}$ , applying (2.A.2) to (2.15) it follows that

$$\begin{aligned} |\mathbf{x}^\epsilon(t)| &\leq \alpha_1 + \int_0^t L_1(1 + |\mathbf{x}^\epsilon(s)| + |\tilde{\mathbf{x}}^\epsilon(s)| + |\boldsymbol{\zeta}|) ds \\ &\leq \alpha_1 + L_1 T + \int_0^t L_1(|\mathbf{x}^\epsilon(s)| + |\tilde{\mathbf{x}}^\epsilon(s)| + |\boldsymbol{\zeta}|) ds \\ &\leq \alpha_1 + L_1 T + \int_0^t L_1(|\mathbf{x}^\epsilon(s)| + |\boldsymbol{\zeta}|) ds + \sum_{i=1}^m \int_{-\tau_i}^{t-\tau_i} L_1 |\mathbf{x}^\epsilon(s)| ds \\ &\leq L_0 + \int_0^t L_1 |\mathbf{x}^\epsilon(s)| ds + \sum_{i=1}^m \int_{-\bar{b}}^t L_1 |\mathbf{x}^\epsilon(s)| ds, \quad t \in [0, T], \end{aligned}$$

where  $L_0 = \alpha_1 + L_1 T + r L_1 \bar{d} T$  and  $\bar{d} = \max_{j=1, \dots, r} \{\zeta_j\}$ . Therefore, by using (2.13),

$$|\mathbf{x}^\epsilon(t)| \leq L_0 + m \alpha_1 \bar{b} + (m+1) L_1 \int_0^t |\mathbf{x}^\epsilon(s)| ds, \quad t \in [0, T].$$

Then, by using (2.10) and Gronwall-Bellman's lemma [111], we obtain

$$|\mathbf{x}^\epsilon(t)| \leq \rho_1 \exp((m+1)L_1 T), \quad t \in [0, T], \quad (2.16)$$

where  $\rho_1 = L_0 + m \alpha_1 \bar{b}$ .

In addition, for each  $s \in [0, T]$ , consider (2.10)-(2.11), then by using (2.A.2),

$$|\boldsymbol{\chi}^\epsilon(t)| \leq L_1(1 + |\mathbf{x}^\epsilon(t)| + |\tilde{\mathbf{x}}^\epsilon(t)| + |\boldsymbol{\zeta}|), \quad t \in [0, T],$$

Clearly,

$$|\boldsymbol{\chi}^\epsilon(t)| \leq L_1 + rL_1\bar{d} + L_1|\boldsymbol{x}^\epsilon(t)| + L_1 \sum_{i=1}^m |\boldsymbol{x}^\epsilon(s - \tau_i - \epsilon\delta_{ki})|. \quad (2.17)$$

Thus, by applying (2.16) to (2.17), we obtain

$$|\boldsymbol{\chi}^\epsilon(t)| \leq L_2, \quad (2.18)$$

where  $L_2 = L_1 + rL_1\bar{d} + (m+1)\rho_1L_1 \exp((m+1)L_1T)$ . Combining (2.13), (2.14), (2.16) and (2.18), the conclusion of the lemma follows readily.  $\square$

Define

$$\Xi = \{\boldsymbol{\omega} \in \mathbb{R}^n : |\boldsymbol{\omega}| \leq L_2\}. \quad (2.19)$$

Then, it follows from Lemma 2.1 that  $\boldsymbol{x}^\epsilon(t) \in \Xi$  for all  $t \in [-\bar{b}, T]$  and  $\epsilon \in I$ . We have the following lemma.

**Lemma 2.2.** *There exists a positive real number  $L_3 > 0$  such that for all  $\epsilon \in I$ ,*

$$|\boldsymbol{\varphi}^\epsilon(t)|, |\boldsymbol{\chi}^\epsilon(t) - \boldsymbol{\chi}^0(t)|, \max_{i=1, \dots, m} |\boldsymbol{\theta}^{\epsilon, i}(t)| \leq L_3|\epsilon|, \quad t \in [0, T]. \quad (2.20)$$

*Proof.* For each  $t \in [0, T]$ , we have

$$|\boldsymbol{\varphi}^\epsilon(t)| = |\boldsymbol{x}^\epsilon(t) - \boldsymbol{x}(t)| \leq \int_0^t |\boldsymbol{\chi}^\epsilon(s) - \boldsymbol{\chi}^0(s)| ds, \quad t \in [0, T]. \quad (2.21)$$

By using (2.A.1), the function  $\boldsymbol{f}$  is Lipschitz continuous on  $\Xi \times \Xi$ . Hence, there exists a real number  $\rho > 0$  such that

$$|\boldsymbol{\chi}^\epsilon(s) - \boldsymbol{\chi}^0(s)| \leq \rho|\boldsymbol{\varphi}^\epsilon(s)| + \rho \sum_{i=1}^m |\boldsymbol{\theta}^{\epsilon, i}(s)|, \quad s \in [0, T]. \quad (2.22)$$

Let  $\epsilon \in I$  be arbitrary but fixed. For each  $s \in [0, T]$ ,

$$\begin{aligned} |\boldsymbol{\theta}^{\epsilon, k}(s)| &= |\boldsymbol{x}^\epsilon(s - \tau_k - \epsilon) - \boldsymbol{x}(s - \tau_k)| \\ &\leq |\boldsymbol{x}^\epsilon(s - \tau_k - \epsilon) - \boldsymbol{x}^\epsilon(s - \tau_k)| + |\boldsymbol{x}^\epsilon(s - \tau_k) - \boldsymbol{x}(s - \tau_k)|. \end{aligned}$$

Hence, by (2.11),

$$|\boldsymbol{\theta}^{\epsilon, k}(s)| \leq \int_{\alpha(s)}^{\beta(s)} |\boldsymbol{\chi}^\epsilon(\eta)| d\eta + |\boldsymbol{\varphi}^\epsilon(s - \tau_k)|, \quad s \in [0, T], \quad (2.23)$$

where

$$\alpha(s) = \min\{s - \tau_k, s - \tau_k - \epsilon\}, \quad \beta(s) = \max\{s - \tau_k, s - \tau_k - \epsilon\}.$$

Clearly,

$$|\beta(s) - \alpha(s)| = \epsilon, \quad s \in [0, T], \quad (2.24)$$

and

$$[\alpha(s), \beta(s)] \subset [-\bar{b}, T], \quad s \in [0, T]. \quad (2.25)$$

Substituting (2.24)-(2.25) into (2.23), and using Lemma 2.1, it yields

$$|\theta^{\epsilon, k}(s)| \leq L_2|\epsilon| + |\varphi^\epsilon(s - \tau_k)|, \quad s \in [0, T]. \quad (2.26)$$

Substituting (2.26) into (2.22) gives

$$|\chi^\epsilon(s) - \chi^0(s)| \leq \rho|\varphi^\epsilon(s)| + m\rho L_2|\epsilon| + \rho \sum_{i=1}^m |\varphi^\epsilon(s - \tau_i)|, \quad s \in [0, T]. \quad (2.27)$$

Substituting (2.27) into (2.21), it follows that

$$\begin{aligned} |\varphi^\epsilon(t)| &\leq m\rho L_2|\epsilon|T + \int_0^t \rho|\varphi^\epsilon(s)|ds + \rho \sum_{i=1}^m \int_0^t |\varphi^\epsilon(s - \tau_i)|ds \\ &\leq m\rho L_2|\epsilon|T + \int_0^t \rho|\varphi^\epsilon(s)|ds + \rho \sum_{i=1}^m \int_{-\tau_i}^{t-\tau_i} |\varphi^\epsilon(s)|ds \\ &\leq m\rho L_2|\epsilon|T + \rho \int_0^t |\varphi^\epsilon(s)|ds + \rho m \int_{-\bar{b}}^t |\varphi^\epsilon(s)|ds, \quad t \in [0, T]. \end{aligned}$$

Recall that  $\varphi^\epsilon(s) = 0$  for all  $s \leq 0$ . Therefore,

$$|\varphi^\epsilon(t)| \leq m\rho L_2|\epsilon|T + (m+1)\rho \int_0^t |\varphi^\epsilon(s)|ds, \quad t \in [0, T].$$

Thus, by Gronwall-Bellman's lemma,

$$|\varphi^\epsilon(t)| \leq \rho_2|\epsilon|, \quad t \in [0, T],$$

where  $\rho_2 = m\rho L_2T \exp\{(m+1)\rho T\}$ . Since  $\varphi^\epsilon(s) = 0$  for all  $s \leq 0$ , it follows that

$$|\varphi^\epsilon(t)| \leq \rho_2|\epsilon|, \quad t \leq T. \quad (2.28)$$

Substituting (2.28) into (2.26) yields

$$|\boldsymbol{\theta}^{\epsilon,k}(s)| \leq (L_2 + \rho_2)|\epsilon|, \quad s \in [0, T].$$

Substituting (2.28) into (2.27), we obtain

$$|\boldsymbol{\chi}^\epsilon(s) - \boldsymbol{\chi}^0(s)| \leq m\rho L_2|\epsilon| + \rho\rho_2|\epsilon| + \rho \sum_{i=1}^m \rho_2|\epsilon| = (m\rho L_2 + (m+1)\rho\rho_2)|\epsilon|, \quad s \in [0, T].$$

Choose  $L_3 = \max\{m\rho L_2 + (m+1)\rho\rho_2, L_2 + \rho_2, \rho_2\}$ . The proof is completed.  $\square$

To proceed further, we need the following lemma.

**Lemma 2.3.** *For almost all  $t \in [0, T]$ , it holds that*

$$\lim_{\epsilon \rightarrow 0} \frac{\boldsymbol{\theta}^{\epsilon,k}(t) - \boldsymbol{\varphi}^\epsilon(t - \tau_k)}{\epsilon} = -\boldsymbol{\chi}^0(t - \tau_k). \quad (2.29)$$

*Proof.* Let  $t \in [0, T] \setminus \tau_k$  be arbitrary but fixed. Then, for each  $\epsilon \in I \setminus 0$ ,

$$\boldsymbol{\theta}^{\epsilon,k}(t) - \boldsymbol{\varphi}^\epsilon(t - \tau_k) = \boldsymbol{x}^\epsilon(t - \tau_k - \epsilon) - \boldsymbol{x}^\epsilon(t - \tau_k).$$

Hence, by (2.11),

$$\frac{\boldsymbol{\theta}^{\epsilon,k}(t) - \boldsymbol{\varphi}^\epsilon(t - \tau_k)}{\epsilon} = \frac{1}{\epsilon} \int_{t-\tau_k}^{t-\tau_k-\epsilon} \boldsymbol{\chi}^\epsilon(s) ds.$$

We can write this equation as follows:

$$\begin{aligned} \frac{\boldsymbol{\theta}^{\epsilon,k}(t) - \boldsymbol{\varphi}^\epsilon(t - \tau_k)}{\epsilon} &= -\boldsymbol{\chi}^0(t - \tau_k) + \boldsymbol{\chi}^0(t - \tau_k) + \frac{1}{\epsilon} \int_{t-\tau_k}^{t-\tau_k-\epsilon} \boldsymbol{\chi}^\epsilon(s) ds \\ &= -\boldsymbol{\chi}^0(t - \tau_k) + \rho(\epsilon), \end{aligned} \quad (2.30)$$

where

$$\rho(\epsilon) = \frac{1}{\epsilon} \int_{t-\tau_k}^{t-\tau_k-\epsilon} \{\boldsymbol{\chi}^\epsilon(s) - \boldsymbol{\chi}^0(t - \tau_k)\} ds.$$

Then, by using triangle inequality,

$$|\rho(\epsilon)| \leq \frac{1}{|\epsilon|} \int_{\alpha_1}^{\beta_1} |\boldsymbol{\chi}^\epsilon(s) - \boldsymbol{\chi}^0(s)| ds + \frac{1}{|\epsilon|} \int_{\alpha_1}^{\beta_1} |\boldsymbol{\chi}^0(s) - \boldsymbol{\chi}^0(t - \tau_k)| ds, \quad (2.31)$$

where

$$\alpha_1 = \min\{t - \tau_k, t - \tau_k - \epsilon\}, \quad \beta_1 = \max\{t - \tau_k, t - \tau_k - \epsilon\}.$$

Clearly,  $\beta_1 - \alpha_1 = |\epsilon|$ . Thus, by Lemma 2.2 and (2.31), it follows that, for each  $\epsilon \in I$ ,

$$|\rho(\epsilon)| \leq L_3|\epsilon| + \frac{1}{|\epsilon|} \int_{\alpha_1}^{\beta_1} |\boldsymbol{\chi}^0(s) - \boldsymbol{\chi}^0(t - \tau_k)| ds. \quad (2.32)$$

Since  $t \neq \tau_k$ , we need to consider the following two cases. *Case 1:*  $t < \tau_k$  and *Case 2:*  $t > \tau_k$ .

*Case 1:*  $t < \tau_k$ . Clearly,

$$\epsilon \in I, \quad |\epsilon| < \tau_k - t \quad \Rightarrow \quad [\alpha_1, \beta_1] \subset [-\bar{b}, 0]. \quad (2.33)$$

Now, since  $\boldsymbol{f}$  is continuously differentiable (recall (2.A.1)), and  $\boldsymbol{\zeta}$  and  $\boldsymbol{x}$  are bounded on  $[-\bar{b}, T]$  (recall (2.4) and (2.12)), we can show that  $\boldsymbol{\chi}^0$  is Lipschitz continuous on  $[-\bar{b}, 0]$ . Hence, there exists a real number  $\eta_1 > 0$  such that

$$|\boldsymbol{\chi}^0(s) - \boldsymbol{\chi}^0(t - \tau_k)| \leq \eta_1 |s - t + \tau_k|, \quad s \in [-\bar{b}, 0] \quad (2.34)$$

It follows from (2.33) and (2.34) that

$$|\boldsymbol{\chi}^0(s) - \boldsymbol{\chi}^0(t - \tau_k)| \leq \eta_1 |s - t + \tau_k| \leq \eta_1 (\beta_1 - \alpha_1) = \eta_1 |\epsilon|, \quad s \in [\alpha_1, \beta_1],$$

when  $\epsilon \in I$  is sufficiently small. Substituting this inequality into (2.32) gives

$$|\rho(\epsilon)| \leq (L_3 + \eta_1) |\epsilon|.$$

This shows that

$$\lim_{\epsilon \rightarrow 0} \rho(\epsilon) = 0, \quad \epsilon \in I \setminus 0, \quad s \in [-\bar{b}, 0] \quad (2.35)$$

*Case 2:*  $t > \tau_k$

Suppose  $t > \tau_k$ . Clearly,

$$\epsilon \in I, \quad |\epsilon| < t - \tau_k \quad \Rightarrow \quad [\alpha_1, \beta_1] \subset (0, T]. \quad (2.36)$$

Similarly, since  $\boldsymbol{f}$  is continuously differentiable (recall (2.A.1)), and  $\boldsymbol{\zeta}$  and  $\boldsymbol{x}$  are bounded on  $[-\bar{b}, T]$  (recall (2.4) and (2.12)), we can show that  $\boldsymbol{\chi}^0$  is Lipschitz continuous on  $(0, T]$ . Hence, there exists a real number  $\eta_2 > 0$  such that

$$|\boldsymbol{\chi}^0(s) - \boldsymbol{\chi}^0(t - \tau_k)| \leq \eta_2 |s - t + \tau_k|, \quad s \in (0, T] \quad (2.37)$$

It follows from (2.36) and (2.37) that

$$|\boldsymbol{\chi}^0(s) - \boldsymbol{\chi}^0(t - \tau_k)| \leq \eta_2 |s - t + \tau_k| \leq \eta_2(\beta_1 - \alpha_1) = \eta_2 |\epsilon|, \quad s \in [\alpha_1, \beta_1],$$

when  $\epsilon \in I$  is sufficiently small.

Substituting this inequality into (2.32) gives

$$|\rho(\epsilon)| \leq (L_3 + \eta_2) |\epsilon|, \quad s \in (0, T].$$

This shows that

$$\lim_{\epsilon \rightarrow 0} \rho(\epsilon) = 0, \quad \epsilon \in I \setminus 0, \quad s \in (0, T] \quad (2.38)$$

Applying (2.35) and (2.38) to (2.30) completes the proof.  $\square$

## 2.4 Gradient computation

Problem (P) involves choosing a finite number of decision variables to minimize the cost function (2.7). Thus, in principle, Problem (P) can be viewed as a nonlinear programming problem. Standard algorithms for solving nonlinear programming problems—for example, sequential quadratic programming or interior-point methods [79]—typically require the gradient of the cost function, which is difficult to determine in Problem (P) because the delays and parameters influence (2.7) *implicitly* through the dynamic system (2.1)-(2.2). The aim of this section is to develop an efficient computational method for computing the gradient of the cost function in Problem (P). This method, which is inspired by earlier works in [117, 125, 126], can be integrated with a standard nonlinear programming algorithm to solve Problem (P).

### 2.4.1 State variation with respect to time-delays

The solution of system (2.1)-(2.2) is normally viewed as a function of time, with  $\boldsymbol{\tau}$  and  $\boldsymbol{\zeta}$  being fixed vectors. By fixing  $t \in (-\infty, T]$  while allowing  $\boldsymbol{\tau}$  and  $\boldsymbol{\zeta}$  to vary, we obtain a new function  $\boldsymbol{x}(t|\cdot, \cdot) : \mathcal{T} \times \mathcal{Z} \rightarrow \mathbb{R}^n$  whose value at  $(\boldsymbol{\tau}, \boldsymbol{\zeta}) \in \mathcal{T} \times \mathcal{Z}$  is  $\boldsymbol{x}(t|\boldsymbol{\tau}, \boldsymbol{\zeta})$ . In the following theorem, we show that  $\boldsymbol{x}(t|\cdot, \cdot)$  is differentiable with respect to the time-delays. This result is central to the development of a computational procedure for solving Problem (P).

**Theorem 2.1.** *Let  $t \in (0, T]$  be a fixed time point. Then,  $\boldsymbol{x}(t|\cdot, \cdot)$  is differentiable with respect to  $\tau_k$  on  $\mathcal{T} \times \mathcal{Z}$ . In fact, for each  $(\boldsymbol{\tau}, \boldsymbol{\zeta}) \in \mathcal{T} \times \mathcal{Z}$ ,*

$$\frac{\partial \boldsymbol{x}(t|\boldsymbol{\tau}, \boldsymbol{\zeta})}{\partial \tau_k} = \boldsymbol{\Lambda}^k(t|\boldsymbol{\tau}, \boldsymbol{\zeta}), \quad k = 1, \dots, m, \quad (2.39)$$



where  $\Lambda^k(\cdot|\tau, \zeta)$  satisfies the auxiliary time-delay system

$$\begin{aligned} \dot{\Lambda}^k(t) &= \frac{\partial \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \zeta)}{\partial \mathbf{x}} \Lambda^k(t) + \sum_{i=1}^m \frac{\partial \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \zeta)}{\partial \tilde{\mathbf{x}}^i} \Lambda^k(t - \tau_i) \\ &\quad - \frac{\partial \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \zeta)}{\partial \tilde{\mathbf{x}}^k} \chi(t - \tau_k) \end{aligned} \quad (2.40)$$

with initial condition

$$\Lambda^k(t) = \mathbf{0}, \quad t \leq 0. \quad (2.41)$$

*Proof.* Let  $k \in \{1, \dots, m\}$  and  $(\tau, \zeta) \in \mathcal{T} \times \mathcal{Z}$  be arbitrary but fixed. As in Section 2.3, we write  $\mathbf{x}^\epsilon(t)$  instead of  $\mathbf{x}(t|\tau + \epsilon \mathbf{e}^k, \zeta)$ , and  $\mathbf{x}(t)$  instead of  $\mathbf{x}(t|\tau, \zeta)$ .

For each  $\epsilon \in I \setminus \{0\}$ , define

$$\rho(\epsilon) = \int_0^T |\epsilon^{-1} \boldsymbol{\theta}^{\epsilon, k}(s) - \epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s - \tau_k) + \boldsymbol{\chi}(s - \tau_k)| ds. \quad (2.42)$$

It follows from (2.9), (2.12), and (2.20) that for each  $\epsilon \in I \setminus \{0\}$ ,

$$|\epsilon^{-1} \boldsymbol{\theta}^{\epsilon, k}(s) - \epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s - \tau_k) + \boldsymbol{\chi}(s - \tau_k)| \leq L_2 + 2L_3, \quad s \in [0, T].$$

Hence, the integrand in (2.42) is uniformly bounded with respect to  $\epsilon \in I \setminus \{0\}$ . Furthermore, it follows from (2.29) that  $\epsilon^{-1} \boldsymbol{\theta}^{\epsilon, k}(s) - \epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s - \tau_k) + \boldsymbol{\chi}(s - \tau_k)$  converges to zero almost everywhere on  $[0, T]$  as  $\epsilon \rightarrow 0$ . Thus, from the Lebesgue dominated convergence theorem,

$$\lim_{\epsilon \rightarrow 0} \rho(\epsilon) = \lim_{\epsilon \rightarrow 0} \int_0^T |\epsilon^{-1} \boldsymbol{\theta}^{\epsilon, k}(s) - \epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s - \tau_k) + \boldsymbol{\chi}(s - \tau_k)| ds = 0.$$

Now, keeping  $\epsilon \in I \setminus \{0\}$  fixed for the time being, we define

$$\bar{\mathbf{f}}(s, \alpha) = \mathbf{f}(s, \mathbf{x}(s) + \alpha \boldsymbol{\varphi}^\epsilon(s), \tilde{\mathbf{x}}(s) + \alpha \boldsymbol{\theta}^\epsilon(s), \zeta), \quad (s, \alpha) \in [0, T] \times [0, 1].$$

Then, by the chain rule,

$$\frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \alpha} = \frac{\partial \mathbf{f}(s, \alpha)}{\partial \mathbf{x}} \boldsymbol{\varphi}^\epsilon(s) + \sum_{i=1}^m \frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \tilde{\mathbf{x}}^i} \boldsymbol{\theta}^{\epsilon, i}(s), \quad (2.43)$$

where

$$\frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}(s, \mathbf{x}(s) + \alpha \boldsymbol{\varphi}^\epsilon(s), \tilde{\mathbf{x}}(s) + \alpha \boldsymbol{\theta}^\epsilon(s), \zeta)}{\partial \mathbf{x}}, \quad (2.44)$$

$$\frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \tilde{\mathbf{x}}^i} = \frac{\partial \mathbf{f}(s, \mathbf{x}(s) + \alpha \boldsymbol{\varphi}^\epsilon(s), \tilde{\mathbf{x}}(s) + \alpha \boldsymbol{\theta}^\epsilon(s), \zeta)}{\partial \tilde{\mathbf{x}}^i}. \quad (2.45)$$

We can rewrite (2.43) as follows:

$$\begin{aligned} \frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \alpha} &= \Delta_1(s, \alpha) + \Delta_2(s, \alpha) + \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \mathbf{x}} \boldsymbol{\varphi}^\epsilon(s) + \sum_{i=1}^m \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^i} \boldsymbol{\varphi}^\epsilon(s - \tau_i) \\ &\quad + \sum_{i=1}^m \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^i} \{\boldsymbol{\theta}^{\epsilon, i}(s) - \boldsymbol{\varphi}^\epsilon(s - \tau_i)\}, \end{aligned}$$

where

$$\Delta_1(s, \alpha) = \left\{ \frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \mathbf{x}} - \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \mathbf{x}} \right\} \boldsymbol{\varphi}^\epsilon(s), \quad (2.46)$$

$$\Delta_2(s, \alpha) = \sum_{i=1}^m \left\{ \frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \tilde{\mathbf{x}}^i} - \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^i} \right\} \boldsymbol{\theta}^{\epsilon, i}(s). \quad (2.47)$$

Applying (2.8) gives

$$\begin{aligned} \frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \alpha} &= \Delta_1(s, \alpha) + \Delta_2(s, \alpha) + \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \mathbf{x}} \boldsymbol{\varphi}^\epsilon(s) \\ &\quad + \sum_{i=1}^m \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^i} \boldsymbol{\varphi}^\epsilon(s - \tau_i) + \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^k} \{\boldsymbol{\theta}^{\epsilon, k}(s) - \boldsymbol{\varphi}^\epsilon(s - \tau_k)\}. \end{aligned} \quad (2.48)$$

Now,

$$\boldsymbol{\varphi}^\epsilon(t) = \mathbf{x}^\epsilon(t) - \mathbf{x}(t) = \int_0^t \{\bar{\mathbf{f}}(s, 1) - \bar{\mathbf{f}}(s, 0)\} ds.$$

Thus, by the fundamental theorem of calculus,

$$\boldsymbol{\varphi}^\epsilon(t) = \int_0^t \{\bar{\mathbf{f}}(s, 1) - \bar{\mathbf{f}}(s, 0)\} ds = \int_0^t \int_0^1 \frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \alpha} d\alpha ds. \quad (2.49)$$

Substituting (2.48) into (2.49) yields

$$\begin{aligned} \boldsymbol{\varphi}^\epsilon(t) &= \int_0^t \int_0^1 \{\Delta_1(s, \alpha) + \Delta_2(s, \alpha)\} d\alpha ds + \int_0^t \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \mathbf{x}} \boldsymbol{\varphi}^\epsilon(s) ds \\ &\quad + \int_0^t \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^k} \{\boldsymbol{\theta}^{\epsilon, k}(s) - \boldsymbol{\varphi}^\epsilon(s - \tau_k)\} ds + \sum_{i=1}^m \int_0^t \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^i} \boldsymbol{\varphi}^\epsilon(s - \tau_i) ds. \end{aligned} \quad (2.50)$$

Now, by using (2.44) and (2.45), we can write the auxiliary system (2.40) as follows:

$$\dot{\Lambda}^k(s) = \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \mathbf{x}} \Lambda^k(s) + \sum_{i=1}^m \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^i} \Lambda^k(s - \tau_i) - \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^k} \boldsymbol{\chi}(s - \tau_k).$$

Hence,

$$\begin{aligned} \mathbf{\Lambda}^k(t) &= \int_0^t \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \mathbf{x}} \mathbf{\Lambda}^k(s) ds + \sum_{i=1}^m \int_0^t \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^i} \mathbf{\Lambda}^k(s - \tau_i) ds \\ &\quad - \int_0^t \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^k} \boldsymbol{\chi}(s - \tau_k) ds. \end{aligned} \quad (2.51)$$

Now, since  $\mathbf{f}$  is continuously differentiable (recall (2.A.1)), and  $\mathbf{x}$  is bounded on  $[-\bar{b}, T]$  (recall (2.12)), there exists constants  $M_1 > 0$  and  $M_2 > 0$  such that

$$\left| \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \mathbf{x}} \right| \leq M_1, \quad \left| \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \tilde{\mathbf{x}}^i} \right| \leq M_2, \quad s \in [0, T],$$

where  $|\cdot|$  denotes the usual Euclidean norm on  $\mathbb{R}^{n \times n}$ . Thus, by multiplying (2.50) by  $\epsilon^{-1}$ , subtracting (2.51), and then finally taking the norm of both sides, we obtain

$$\begin{aligned} |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) - \mathbf{\Lambda}^k(t)| &\leq M_2 \rho(\epsilon) + \int_0^t M_1 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \mathbf{\Lambda}^k(s)| ds \\ &\quad + \sum_{i=1}^m \int_0^t M_2 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s - \tau_i) - \mathbf{\Lambda}^k(s - \tau_i)| ds \\ &\quad + |\epsilon|^{-1} \int_0^t \int_0^1 \{|\Delta_1(s, \alpha)| + |\Delta_2(s, \alpha)|\} d\alpha ds, \end{aligned} \quad (2.52)$$

where  $\rho(\epsilon)$  is as defined in (2.42). The second integral term on the right-hand side of (2.52) can be simplified as follows:

$$\begin{aligned} \sum_{i=1}^m \int_0^t M_2 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s - \tau_i) - \mathbf{\Lambda}^k(s - \tau_i)| ds &= \sum_{i=1}^m \int_{-\tau_i}^{t-\tau_i} M_2 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \mathbf{\Lambda}^k(s)| ds \\ &\leq \sum_{i=1}^m \int_0^t M_2 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \mathbf{\Lambda}^k(s)| ds \\ &= \int_0^t m M_2 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \mathbf{\Lambda}^k(s)| ds. \end{aligned}$$

Hence, (2.52) becomes

$$\begin{aligned} |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) - \mathbf{\Lambda}^k(t)| &\leq M_2 \rho(\epsilon) + \int_0^t \bar{M} |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \mathbf{\Lambda}^k(s)| ds \\ &\quad + |\epsilon|^{-1} \int_0^t \int_0^1 \{|\Delta_1(s, \alpha)| + |\Delta_2(s, \alpha)|\} d\alpha ds, \end{aligned} \quad (2.53)$$

where  $\bar{M} = M_1 + m M_2$ . Since  $\mathbf{f}$  is continuously differentiable and  $\mathbf{x}^\epsilon$  is uniformly bounded with respect to  $\epsilon$ , both  $\frac{\partial \bar{\mathbf{f}}}{\partial \mathbf{x}}$  and  $\frac{\partial \bar{\mathbf{f}}}{\partial \tilde{\mathbf{x}}^i}$  are uniformly continuous on  $[0, T] \times [0, 1]$ . Furthermore, by (2.20),  $\mathbf{x}(s) + \alpha \boldsymbol{\varphi}^\epsilon(s) \rightarrow \mathbf{x}(s)$  and  $\tilde{\mathbf{x}}(s) + \alpha \boldsymbol{\theta}^\epsilon(s) \rightarrow \tilde{\mathbf{x}}(s)$  uniformly on  $[0, T] \times [0, 1]$  as

$\epsilon \rightarrow 0$ . Thus, for each  $\delta > 0$ , there exists an  $\epsilon' > 0$  such that for all  $\epsilon$  satisfying  $|\epsilon| < \epsilon'$ ,

$$\begin{aligned} \left| \frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \mathbf{x}} - \frac{\partial \bar{\mathbf{f}}(s, 0)}{\partial \mathbf{x}} \right| &< \delta, \quad (s, \alpha) \in [0, T] \times [0, 1], \\ \left| \frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \tilde{\mathbf{x}}^i} - \frac{\partial \bar{\mathbf{f}}(s, \alpha)}{\partial \tilde{\mathbf{x}}^i} \right| &< \delta, \quad (s, \alpha) \in [0, T] \times [0, 1]. \end{aligned}$$

By taking the norm of (2.46) and (2.47), and then using these inequalities together with (2.20), we obtain

$$|\Delta_1(s, \alpha)| \leq \delta L_3 |\epsilon|, \quad |\Delta_2(s, \alpha)| \leq \delta m L_3 |\epsilon|,$$

where  $|\epsilon| < \epsilon'$ . Substituting these inequalities into (2.53) yields,

$$|\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) - \boldsymbol{\Lambda}^k(t)| \leq M_2 \rho(\epsilon) + (L_3 T + mL_3 T) \delta + \int_0^t \bar{M} |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \boldsymbol{\Lambda}^k(s)| ds.$$

Now, recall that  $\rho(\epsilon) \rightarrow 0$  as  $\epsilon \rightarrow 0$ . Hence, there exists an  $\epsilon'' > 0$  such that  $\rho(\epsilon) < \delta$  whenever  $|\epsilon| < \epsilon''$ . Thus, for all  $\epsilon$  such that  $|\epsilon| < \min\{\epsilon', \epsilon''\}$ ,

$$|\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) - \boldsymbol{\Lambda}^k(t)| \leq M_2 \delta + (L_3 T + mL_3 T) \delta + \int_0^t \bar{M} |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \boldsymbol{\Lambda}^k(s)| ds.$$

Applying the Gronwall-Bellman Lemma [111] gives

$$|\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) - \boldsymbol{\Lambda}^k(t)| \leq \delta (M_2 + L_3 T + mL_3 T) \exp\{\bar{M} T\},$$

where  $|\epsilon| < \min\{\epsilon', \epsilon''\}$ . Since  $\delta$  is arbitrary, this shows that  $\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) \rightarrow \boldsymbol{\Lambda}^k(t)$  as  $\epsilon \rightarrow 0$ . It follows that

$$\frac{\partial \mathbf{x}(t | \boldsymbol{\tau}, \boldsymbol{\zeta})}{\partial \tau_k} = \lim_{\epsilon \rightarrow 0} \frac{\mathbf{x}^\epsilon(t) - \mathbf{x}(t)}{\epsilon} = \lim_{\epsilon \rightarrow 0} \epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) = \boldsymbol{\Lambda}^k(t),$$

as required. □

## 2.4.2 State variation with respect to system parameters

In Theorem 2.1, we derive formulae for the gradient of the state with respect to the time-delays. We now turn our attention to the gradient of the state with respect to the system parameters.

Let  $w$  be a new state variable with dynamics

$$\dot{w}(t) = 1, \quad t \in [0, T], \tag{2.54}$$

$$w(t) = t, \quad t \leq 0. \tag{2.55}$$

Clearly,  $w(t) = t$  for all  $t \in (-\infty, T]$ . Thus, we can express the system parameters in (2.1)-(2.2) in terms of the new state  $w$  as follows:

$$\zeta_j = t - w(t - \zeta_j), \quad j = 1, \dots, r. \quad (2.56)$$

Substituting (2.56) into the original system (2.1)-(2.2) gives

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), t - w(t - \zeta_1), \dots, t - w(t - \zeta_r)), \quad t \in [0, T], \quad (2.57)$$

$$\mathbf{x}(t) = \boldsymbol{\phi}(t), \quad t \leq 0. \quad (2.58)$$

The system parameters  $\zeta_j$ ,  $j = 1, \dots, r$ , are now time-delays in the enlarged system consisting of (2.54)-(2.55) and (2.57)-(2.58). Thus, to determine the state variation with respect to the system parameters in system (2.1)-(2.2), we just need to apply Theorem 2.1 to the enlarged system consisting of (2.54)-(2.55) and (2.57)-(2.58). It is important to note that each system parameter is bounded below by zero (see the problem formulation in Section 2.2). Thus, the enlarged system considered here is a valid time-delay system with all time-delays being non-negative.

Let  $\mathbf{z}(t) \in \mathbb{R}^{n+1}$  and  $\tilde{\mathbf{z}}(t) \in \mathbb{R}^{(n+1)(m+r)}$  denote, respectively, the state and delayed state vectors for the enlarged system, where

$$\mathbf{z}(t) = [x_1(t), \dots, x_n(t), w(t)]^\top$$

and

$$\tilde{\mathbf{z}}(t) = [\mathbf{z}(t - \tau_1)^\top, \dots, \mathbf{z}(t - \tau_m)^\top, \mathbf{z}(t - \zeta_1)^\top, \dots, \mathbf{z}(t - \zeta_r)^\top]^\top.$$

The enlarged system consisting of (2.54)-(2.55) and (2.57)-(2.58) can be written as follows:

$$\dot{\mathbf{z}}(t) = \hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t)), \quad t \in [0, T], \quad (2.59)$$

$$\mathbf{z}(t) = \hat{\boldsymbol{\phi}}(t), \quad t \leq 0, \quad (2.60)$$

where

$$\hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t)) = \begin{bmatrix} \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), t - w(t - \zeta_1), \dots, t - w(t - \zeta_r)) \\ 1 \end{bmatrix}$$

and

$$\hat{\boldsymbol{\phi}}(t) = \begin{bmatrix} \boldsymbol{\phi}(t) \\ t \end{bmatrix}.$$

Define

$$\hat{\boldsymbol{\chi}}(t) = \begin{cases} \dot{\hat{\boldsymbol{\phi}}}(t), & \text{if } t \leq 0, \\ \hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t)), & \text{if } t \in (0, T]. \end{cases}$$

Let  $j \in \{1, \dots, r\}$  and  $(\tau, \zeta) \in \mathcal{T} \times \mathcal{Z}$ . Then the auxiliary system for (2.59)-(2.60) corresponding to the system parameter  $\zeta_j$  is

$$\begin{aligned} \begin{bmatrix} \dot{\Gamma}^j(t) \\ \dot{\gamma}^j(t) \end{bmatrix} &= \frac{\partial \hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t))}{\partial \mathbf{z}} \begin{bmatrix} \Gamma^j(t) \\ \gamma^j(t) \end{bmatrix} + \sum_{i=1}^m \frac{\partial \hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t))}{\partial \tilde{\mathbf{z}}^i} \begin{bmatrix} \Gamma^j(t - \tau_i) \\ \gamma^j(t - \tau_i) \end{bmatrix} \\ &+ \sum_{i=1}^r \frac{\partial \hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t))}{\partial \tilde{\mathbf{z}}^{m+i}} \begin{bmatrix} \Gamma^j(t - \zeta_i) \\ \gamma^j(t - \zeta_i) \end{bmatrix} - \frac{\partial \hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t))}{\partial \tilde{\mathbf{z}}^{m+j}} \hat{\mathbf{x}}(t - \zeta_j) \end{aligned} \quad (2.61)$$

with the initial conditions

$$\begin{bmatrix} \Gamma^j(t) \\ \gamma^j(t) \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \quad t \leq 0. \quad (2.62)$$

Here,  $\Gamma^j(t) : (-\infty, T] \rightarrow \mathbb{R}^n$  is the variation of the original state  $\mathbf{x}$  with respect to  $\zeta_j$  and  $\gamma^j(t) : (-\infty, T] \rightarrow \mathbb{R}^n$  is the variation of the new state  $w$  with respect to  $\zeta_j$ . Note that

$$\begin{aligned} \frac{\partial \hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t))}{\partial \mathbf{z}} &= \begin{bmatrix} \frac{\partial \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \zeta)}{\partial \mathbf{x}} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix}, \\ \frac{\partial \hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t))}{\partial \tilde{\mathbf{z}}^i} &= \begin{bmatrix} \frac{\partial \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \zeta)}{\partial \tilde{\mathbf{x}}^i} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix}, \quad i = 1, \dots, m, \\ \frac{\partial \hat{\mathbf{f}}(t, \mathbf{z}(t), \tilde{\mathbf{z}}(t))}{\partial \tilde{\mathbf{z}}^i} &= \begin{bmatrix} \mathbf{0} & -\frac{\partial \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \zeta)}{\partial \zeta_{i+m}} \\ \mathbf{0} & 0 \end{bmatrix}, \quad i = 1, \dots, r. \end{aligned}$$

Furthermore, it is clear that  $\gamma^j(t) = 0$  for all  $t \in (-\infty, T]$ . Thus, the auxiliary system (2.61)-(2.62) becomes

$$\begin{aligned} \dot{\Gamma}^j(t) &= \frac{\partial \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \zeta)}{\partial \mathbf{x}} \Gamma^j(t) + \sum_{i=1}^m \frac{\partial \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \zeta)}{\partial \tilde{\mathbf{x}}^i} \Gamma^j(t - \tau_i) \\ &+ \frac{\partial \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \zeta)}{\partial \zeta_j} \end{aligned} \quad (2.63)$$

with initial conditions

$$\Gamma^j(t) = \mathbf{0}, \quad t \leq 0. \quad (2.64)$$

Applying Theorem 2.1 to the enlarged system consisting of (2.54)-(2.55) and (2.57)-(2.58) yields the following result.

**Theorem 2.2.** *Let  $t \in (0, T]$  be a fixed time point. Then  $\mathbf{x}(t|\cdot, \cdot)$  is differentiable with respect to  $\zeta_j$  on  $\mathcal{T} \times \mathcal{Z}$ . In fact, for each  $(\tau, \zeta) \in \mathcal{T} \times \mathcal{Z}$ ,*

$$\frac{\partial \mathbf{x}(t|\tau, \zeta)}{\partial \zeta_j} = \Gamma^j(t|\tau, \zeta), \quad j = 1, \dots, r, \quad (2.65)$$

where  $\Gamma^j(\cdot|\boldsymbol{\tau}, \boldsymbol{\zeta})$  satisfies the auxiliary time-delay system (2.63)-(2.64).

### 2.4.3 Gradient computation algorithm

We are now ready to derive formulae for the gradients of the cost function in Problem (P). By using Theorems 2.1 and 2.2 and the chain rule of differentiation, we obtain

$$\frac{\partial J(\boldsymbol{\tau}, \boldsymbol{\zeta})}{\partial \tau_k} = 2 \sum_{l=1}^q (\mathbf{y}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta}) - \hat{\mathbf{y}}^l)^\top \frac{\partial \mathbf{g}(\mathbf{x}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta}), \boldsymbol{\zeta})}{\partial \mathbf{x}} \boldsymbol{\Lambda}^k(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta}), \quad (2.66)$$

$$\begin{aligned} \frac{\partial J(\boldsymbol{\tau}, \boldsymbol{\zeta})}{\partial \zeta_j} &= 2 \sum_{l=1}^q (\mathbf{y}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta}) - \hat{\mathbf{y}}^l)^\top \frac{\partial \mathbf{g}(\mathbf{x}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta}), \boldsymbol{\zeta})}{\partial \mathbf{x}} \boldsymbol{\Gamma}^j(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta}) \\ &\quad + 2 \sum_{l=1}^q (\mathbf{y}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta}) - \hat{\mathbf{y}}^l)^\top \frac{\partial \mathbf{g}(\mathbf{x}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta}), \boldsymbol{\zeta})}{\partial \zeta_j}. \end{aligned} \quad (2.67)$$

We now present the following algorithm for computing the cost function (2.7) and its gradient at a given pair  $(\boldsymbol{\tau}, \boldsymbol{\zeta}) \in \mathcal{T} \times \mathcal{Z}$ .

#### Algorithm 2.1.

- Step 1.* Obtain  $\mathbf{x}(\cdot|\boldsymbol{\tau}, \boldsymbol{\zeta})$ ,  $\boldsymbol{\Lambda}^k(\cdot|\boldsymbol{\tau}, \boldsymbol{\zeta})$ ,  $k = 1, \dots, m$ , and  $\boldsymbol{\Gamma}^j(\cdot|\boldsymbol{\tau}, \boldsymbol{\zeta})$ ,  $j = 1, \dots, r$ , by solving the enlarged time-delay system consisting of the original system (2.1)-(2.2) and the auxiliary systems (2.40)-(2.41) and (2.63)-(2.64).
- Step 2.* Use the state values  $\mathbf{x}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta})$ ,  $l = 1, \dots, q$ , to compute  $\mathbf{y}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta})$  through equation (2.6).
- Step 3.* Use  $\mathbf{y}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta})$ ,  $l = 1, \dots, q$ , to compute  $J(\boldsymbol{\tau}, \boldsymbol{\zeta})$  through equation (2.7).
- Step 4.* Use  $\mathbf{x}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta})$ ,  $\mathbf{y}(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta})$ ,  $\boldsymbol{\Lambda}^k(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta})$ , and  $\boldsymbol{\Gamma}^j(t_l|\boldsymbol{\tau}, \boldsymbol{\zeta})$ ,  $l = 1, \dots, q$ , to compute  $\frac{\partial J(\boldsymbol{\tau}, \boldsymbol{\zeta})}{\partial \tau_k}$ ,  $k = 1, \dots, m$ , and  $\frac{\partial J(\boldsymbol{\tau}, \boldsymbol{\zeta})}{\partial \zeta_j}$ ,  $j = 1, \dots, r$ , through equations (2.66) and (2.67).

This gradient computation algorithm can be integrated with a standard gradient-based optimization method (e.g. sequential quadratic programming) to solve Problem (P) as a nonlinear programming problem.

In some applications, the governing dynamic system includes input-delays as well as state-delays. For example, consider the following system:

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \mathbf{u}(t), \tilde{\mathbf{u}}(t), \boldsymbol{\zeta}), \quad t \in [0, T], \quad (2.68)$$

$$\mathbf{x}(t) = \boldsymbol{\phi}(t), \quad t \leq 0, \quad (2.69)$$

where  $\mathbf{u}(t) = [u_1(t), \dots, u_v(t)]^\top \in \mathbb{R}^v$  is the *control input* of system (2.68)-(2.69);  $\tilde{\mathbf{u}}(t) = [\mathbf{u}(t-\lambda_1)^\top, \dots, \mathbf{u}(t-\lambda_d)^\top]^\top \in \mathbb{R}^{vd}$  is the *delayed control*; and  $\lambda_i, i = 1, \dots, d$  are unknown *control delays*. The other symbols are as defined in Section 2.2.

In (2.68)-(2.69),  $\tilde{\mathbf{u}}$  is assumed to be a known input function. Thus, we can write (2.68)-(2.69) in the form of (2.1)-(2.2) as follows:

$$\dot{\mathbf{x}}(t) = \bar{\mathbf{f}}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \boldsymbol{\zeta}, \boldsymbol{\lambda}), \quad t \in [0, T], \quad (2.70)$$

where  $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_d]^\top$  is a parameter vector containing the control delays. If the input function  $\mathbf{u}$  is continuously differentiable, then  $\bar{\mathbf{f}}$  is also continuously differentiable, and thus the approach outlined above for solving Problem (P) is applicable. Hence, our identification method can also be applied to systems with input delay (assuming that the input is smooth).

## 2.5 Numerical examples

### 2.5.1 Example 2.1

We now apply the solution method developed in Section 2.4 to the industrial purification process described in [93,94]. The purpose of this process is to remove harmful cobalt and cadmium ions from a zinc sulphate electrolyte by adding zinc powder to induce deposition. This is a key step in the production of zinc.

The changes in concentrations of cobalt and cadmium ions in the electrolyte are described by the following differential equations:

$$V\dot{x}_1(t) = Qx_1^0 - Qx_1(t-\tau) - \alpha u(t)x_1(t-\tau) + cx_2(t-\tau), \quad (2.71)$$

$$V\dot{x}_2(t) = Qx_2^0 - Qx_2(t-\tau) - \beta v(t)x_2(t-\tau) + dx_1(t-\tau), \quad (2.72)$$

and

$$x_1(t) = 3.3 \times 10^{-4}, \quad x_2(t) = 4.0 \times 10^{-3} \quad t \leq 0, \quad (2.73)$$

where  $x_1$  is the concentration of cobalt ions;  $x_2$  is the concentration of cadmium ions; and  $u$  and  $v$  are control variables representing the zinc powder reaction surface areas for two metallic impurities ions which depend on the amount of zinc powder added to the reaction tank.

Furthermore,  $V$  is the volume of the reaction tank ( $V = 400$ );  $Q$  is the flux of solution ( $Q = 200$ );  $\alpha, \beta, c, d$ , are system parameters; and  $x_1^0$  and  $x_2^0$  are the concentrations of cobalt and cadmium ions at the inlet of the reaction tank, respectively ( $x_1^0 = 6 \times 10^{-4}$ ,  $x_2^0 = 9 \times 10^{-3}$ ). Reference [94] considers the parameter identification problem for system



(2.71)-(2.73) with a given time-delay of  $\tau = 2$ . Here, we consider the problem of identifying the time-delay. We assume that  $\beta$ ,  $c$ , and  $d$  are equal to the optimal values reported in [94]:

$$\beta = 2.823 \times 10^{-4}, \quad c = 16.67, \quad d = 7.107 \times 10^2. \quad (2.74)$$

These values were obtained using data from a real zinc production factory in China. We assume that the terminal time is  $T = 8$ . We set the input variables  $u$  and  $v$  as equal to the optimal control functions obtained in [94].

$$u(t) = \sum_{k=1}^8 \sigma^k \psi_{[\bar{t}_{k-1}, \bar{t}_k)}(t), \quad t \in [0, 8], \quad (2.75)$$

$$v(t) = \sum_{k=1}^8 \bar{\sigma}^k \psi_{[\bar{t}_{k-1}, \bar{t}_k)}(t), \quad t \in [0, 8], \quad (2.76)$$

where the values of  $\bar{t}_k$ ,  $\sigma^k$ , and  $\bar{\sigma}^k$ ,  $k = 1, \dots, 8$ , are listed in Table 2.1, and

$$\psi_{[\bar{t}_{k-1}, \bar{t}_k)}(t) = \begin{cases} 1, & \text{if } t \in [\bar{t}_{k-1}, \bar{t}_k), \\ 0, & \text{otherwise.} \end{cases}$$

The output of the system is the concentration of cadmium ions,  $y(t) = x_2(t)$ .

Given system (2.71)-(2.73), with data (2.74)-(2.76), our goal is to identify the system parameter  $\alpha$  and the delay  $\tau$ .

Table 2.1: Control values for Example 2.1.

$k$	1	2	3	4	5	6	7	8
$\bar{t}_k$	1	2	3	4	5	6	7	8
$\sigma_k \times 10^{-5}$	1.08	1.57	1.24	1.56	1.59	1.43	1.25	1.25
$\bar{\sigma}_k \times 10^{-5}$	5.20	4.70	4.97	4.60	4.53	4.64	4.74	4.62

We simulate system (2.71)-(2.73) with  $\tau = \hat{\tau} = 2$  and  $\alpha = \hat{\alpha} = 7.828 \times 10^{-4}$  to generate the observed data in Problem (P). The sample times are  $t_l = l/2$ ,  $l = 1, \dots, 16$ , and

$$\hat{y}^l = x_2(t_l | \hat{\tau}, \hat{\alpha}).$$

Our identification problem is: choose  $\tau$  and  $\alpha$  to minimize

$$J(\tau, \alpha) = \sum_{l=1}^{16} |y(t_l | \tau, \alpha) - \hat{y}^l|^2 = \sum_{l=1}^{16} |x_2(t_l | \tau, \alpha) - x_2(t_l | \hat{\tau}, \hat{\alpha})|^2$$

subject to the dynamic system (2.71)-(2.73).

Note that this problem cannot be solved using the identification method in [125], as the third term on the right-hand side of (2.71) is a nonlinear term containing both an unknown parameter and an unknown delay. The identification method in [125] is only applicable when each nonlinear term contains a single delay and no unknown parameters. We instead solve this problem using a Matlab program that integrates the SQP optimization method with the gradient computation algorithm described in Section 2.4.3.

Computational results for different initial guesses are shown in Table 2.2. The convergence of the output trajectory for the initial guess  $\tau = 3$  and  $\alpha = 0$  is displayed in Figure 2.1. This figure shows the output trajectory at intermediate iterations of the algorithm, as well as the final (converged) trajectory. In Table 2.2 and Figure 2.1,  $\tau^i$  and  $\alpha^i$  are the values of  $\tau$  and  $\alpha$  at the  $i$ th iteration during the optimization process ( $i = 0$  signifies the initial guess). We can see from Table 2.2 and Figure 2.1 that the optimal trajectory converges to the observed data well, regardless of the initial guess. Thus, the algorithm easily recovers the true values of the delay and parameter for this problem.

Table 2.2: Convergence of the cost values in Example 2.1.

No.	Initial guess		Cost value at $i$ th iteration			
	$\tau^0$	$\alpha^0$	$i = 0$	$i = 5$	$i = 10$	$i = 70$
1	0.5	0.5	$9.111 \times 10^{33}$	$5.392 \times 10^{-6}$	$5.157 \times 10^{-9}$	$7.751 \times 10^{-15}$
2	1.0	1.0	$4.558 \times 10^{20}$	$5.106 \times 10^{-6}$	$7.709 \times 10^{-10}$	$1.088 \times 10^{-13}$
3	1.5	0.5	$3.346 \times 10^{10}$	$1.722 \times 10^{-6}$	$1.496 \times 10^{-6}$	$1.700 \times 10^{-13}$
4	3.0	0.0	$7.094 \times 10^{-5}$	$2.536 \times 10^{-5}$	$2.209 \times 10^{-5}$	$3.341 \times 10^{-14}$
5	3.0	1.0	$8.533 \times 10^3$	$2.589 \times 10^{-5}$	$2.180 \times 10^{-5}$	$2.050 \times 10^{-14}$

For comparison, we also solve this problem using the genetic algorithm (GA) in [57]. The parameters of GA are: the size of population is 20, the crossover probability is 0.8, the selection rate is 0.9, the mutation probability is 0.01, the number of bits for each individual is 14, and the maximum number of iterations is 1000. It takes about 40 minutes for GA to solve this problem, which is more than 20 times longer than the computation time taken by our method. Moreover, the cost value obtained by GA is  $1.3787 \times 10^{-9}$  with corresponding parameter estimates  $\tau = 2.0026$  and  $\alpha = 7.9351 \times 10^{-4}$ . Clearly, the results obtained by our new method are better than those from GA. This is not surprising, as our method exploits the gradient of the cost function to achieve fast convergence.

## 2.5.2 Example 2.2

We now demonstrate the applicability of our approach to systems with multiple delays. Consider the dynamic system given below:

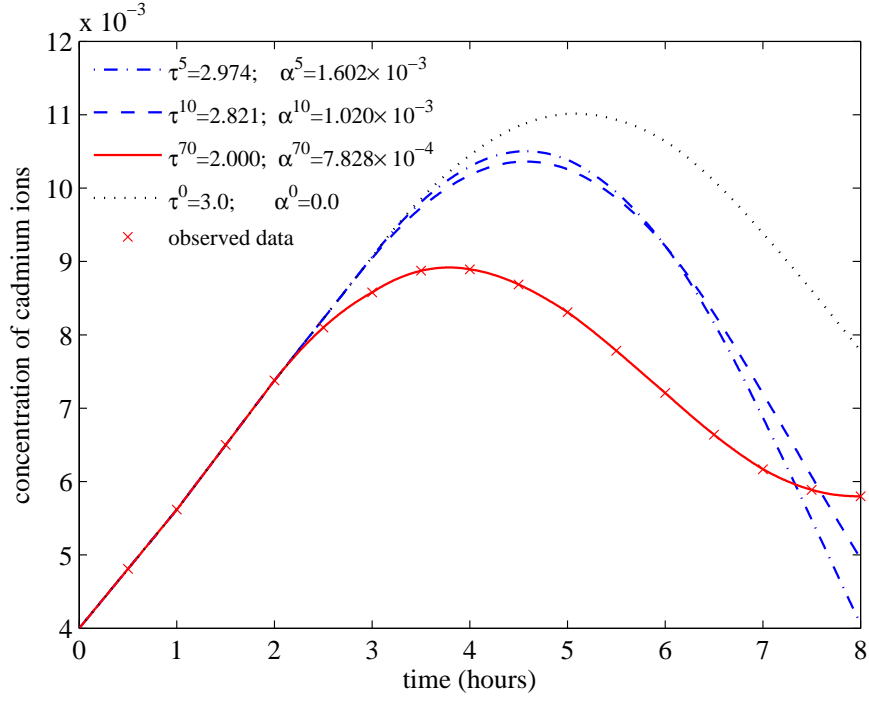


Figure 2.1: Convergence of the output trajectory in Example 2.1 for initial guess No.4.

$$\begin{aligned} \dot{x}_1(t) = & -2x_1(t) + 0.1(1 - x_1(t - \tau_1)) \exp \left\{ \frac{20x_2(t)}{20 + x_2(t)} \right\} \\ & + 0.1x_1(t - \tau_1)x_2(t - \tau_2) + u(t - \tau_3), \end{aligned} \quad (2.77)$$

$$\begin{aligned} \dot{x}_2(t) = & -2.5x_2(t) + 0.8(1 - x_1(t - \tau_1)) \exp \left\{ \frac{20x_2(t)}{20 + x_2(t)} \right\} \\ & + 0.1x_2(t - \tau_1)x_2(t - \tau_2) + u(t - \tau_3), \end{aligned} \quad (2.78)$$

with initial condition

$$x_1(t) = 1, \quad x_2(t) = 1, \quad t \leq 0. \quad (2.79)$$

Here,  $\tau_1$  and  $\tau_2$  are unknown state-delays, and  $\tau_3$  is an unknown input delay. Assume that the terminal time of this system is  $T = 10$ . The input function is given by

$$u(t) = 0.1 \sin(t), \quad t \leq 10.$$

Furthermore, the output is

$$y(t) = x_2(t), \quad t \leq 10.$$

We use the output trajectory of (2.77)-(2.79) with  $[\hat{\tau}_1, \hat{\tau}_2, \hat{\tau}_3] = [2.4, 1.8, 1.1]^\top$  to generate the observed data for Problem (P). We set

$$\hat{y}^l = x_2(t_l | \hat{\tau}_1, \hat{\tau}_2, \hat{\tau}_3), \quad l = 1, \dots, 20,$$

where  $t_l = l/2$ ,  $l = 1, \dots, 20$ . Thus, our identification problem is: choose  $\tau_1$ ,  $\tau_2$ , and  $\tau_3$  to minimize

$$J(\boldsymbol{\tau}) = \sum_{l=1}^{20} |y(t_l | \tau_1, \tau_2, \tau_3) - \hat{y}^l|^2 = \sum_{l=1}^{20} |x_2(t_l | \tau_1, \tau_2, \tau_3) - x_2(t_l | \hat{\tau}_1, \hat{\tau}_2, \hat{\tau}_3)|^2$$

subject to the dynamics (2.77)-(2.79).

We solved this problem using the same Matlab program that was used to solve Example 2.1. The convergence process of the program is shown in Table 2.3 for four sets of initial guesses. The convergence of the output trajectory for the initial guess  $\boldsymbol{\tau}^0 = [3.0, 3.0, 3.0]^\top$  is shown in Figure 2.2. In Table 2.3 and Figure 2.2,  $\boldsymbol{\tau}^i = [\tau_1^i, \tau_2^i, \tau_3^i]^\top$  is the values of  $\boldsymbol{\tau}$  at the  $i$ th iteration, while  $i = 0$  signifies the initial guess. We also solve this problem using GA with the same parameters as in Example 2.1. The optimal cost obtained by GA is  $1.3 \times 10^{-4}$ . Moreover, the computation time is much longer than our new method. As with Example 2.1, we see that the optimization results converge from all initial guesses to the optimal solution.

Table 2.3: Convergence of the cost values in Example 2.2.

No.	Initial guess			Cost value at $i$ th iteration			
	$\tau_1^0$	$\tau_2^0$	$\tau_3^0$	$i = 0$	$i = 5$	$i = 10$	$i = 30$
1	0.5	0.5	0.5	0.4922	0.0188	$5.667 \times 10^{-3}$	$6.661 \times 10^{-15}$
2	1.5	1.5	1.5	0.1386	0.0035	$3.357 \times 10^{-6}$	$6.618 \times 10^{-15}$
3	2.5	2.5	2.5	0.0747	0.0083	$4.405 \times 10^{-4}$	$1.534 \times 10^{-14}$
4	3.0	3.0	3.0	0.1710	0.0298	$2.780 \times 10^{-3}$	$6.656 \times 10^{-15}$

## 2.6 Conclusion

In this chapter, we have developed a gradient-based computational method for determining unknown time-delays and unknown parameters in a general nonlinear system. This method is unified in the sense that the delays and parameters are determined simultaneously by solving a dynamic optimization problem. The gradient of the cost function in this problem is obtained by solving a set of auxiliary delay-differential systems from  $t = 0$  to  $t = T$ . The numerical simulations in Section 2.5 demonstrate that this approach

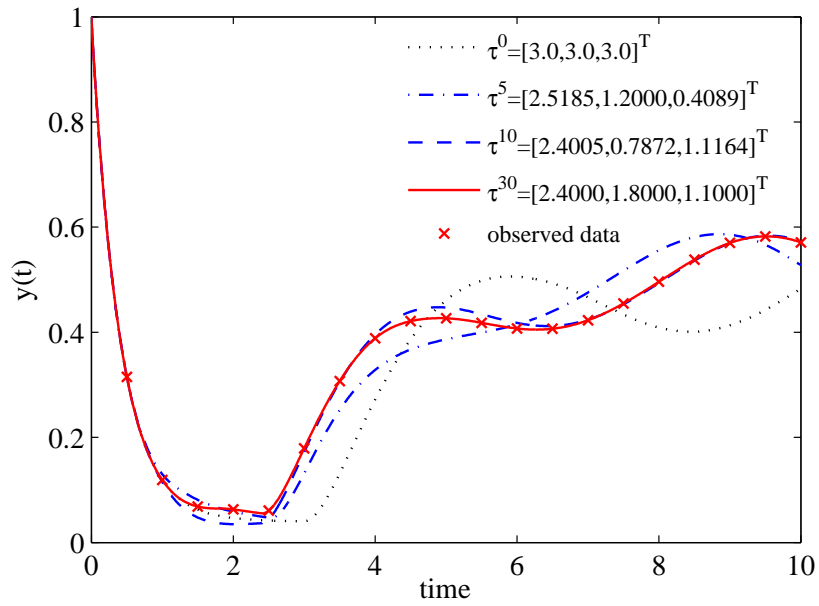


Figure 2.2: Convergence of the output trajectory in Example 2.2 for initial guess No.4.

is highly effective. In particular, it converges quickly even when the initial estimates for the delays and parameters are far away from the optimal values.

---

---

# CHAPTER 3

---

## Identification of time-delays for nonlinear systems with piecewise-constant input

### 3.1 Introduction

The optimization-based approach developed in Chapter 2 is designed for nonlinear delayed systems with smooth inputs. For systems with input-delays, if the input function is smooth, then the system dynamics will be continuously differentiable with respect to the input-delays, and thus the approach proposed in Chapter 2 can be easily modified to estimate the input-delays in this case. Unfortunately, the input function is often non-smooth in practical applications. Examples include biodiesel production [58], evaporation process [119], chromatography process [133], quadruple-tank process [55], batch reactor [105], and distillation column [60]. Since flow rate transmission, sensors, transfer delays of sensor-to-controller that are involved in control loops unavoidably introduce input-delays. As such, the estimation method in Chapter 2 and [125] is not applicable in such situations.

Time-delays identification for nonlinear delay systems with input-delays has been an interest research topic. However, the vast majority of delay estimation methods for delayed systems with piecewise inputs are only applicable to simple systems with linear dynamics and a single delay, see for example, step input and system parameters identification method [103], annihilation and integration based identification method [103]. In this chapter, we consider the time-delay estimation problem for nonlinear systems in which the input function is piecewise-constant. Such estimation problems arise, for example, in evaporation and purification processes [94, 119]. We assume that the governing system contains one state-delay and one input-delay, both of which are unknown and need to be estimated using experimental data. As with Chapter 2, we formulate the delay estimation problem as a dynamic optimization problem in which the cost function measures the least-squares error between predicted output and observed system output. The main difficulties in solving this problem are: i) The delays are decision variables to be optimized, rather than fixed values. Thus, conventional optimization techniques are

not directly applicable; and ii) since the input function is discontinuous, thus the dynamics are clearly discontinuous with respect to input delays. Hence, the results obtained in Chapter 2 can not be used to determine this state variation with respect to input time-delays. In this chapter, we focus on the derivation of a computational procedure for determining the gradient of the cost function for this problem. This procedure, which involves integrating an auxiliary impulsive system with instantaneous jumps forward in time, is far more complex than the procedure given in Chapter 2, which does not involve any jumps. Moreover, because of the discontinuous nature of the input function, the cost function's gradient does not exist at certain points. We propose a heuristic strategy for dealing this complication. Subsequently, this heuristic strategy can be combined with our gradient computation procedure to solve the estimation problem using standard nonlinear programming algorithms. We then apply this approach to estimate the time-delays in two large-scale industrial engineering systems. The purpose of this chapter is to develop a new method for estimating the time-delays.

## 3.2 Problem formulation

Consider the following nonlinear time-delay system:

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t), \mathbf{u}(t - \beta)), \quad t \in [0, T], \quad (3.1)$$

$$\mathbf{x}(t) = \boldsymbol{\phi}(t), \quad t \leq 0, \quad (3.2)$$

where  $T > 0$  is a given *terminal time*;  $\mathbf{x}(t) = [x_1(t), \dots, x_n(t)]^\top \in \mathbb{R}^n$  is the *state vector*;  $\mathbf{u}(t) = [u_1(t), \dots, u_r(t)]^\top \in \mathbb{R}^r$  is the *input vector*;  $\alpha$  and  $\beta$  are unknown time-delays that need to be determined; and  $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R}^r \times \mathbb{R}^r \rightarrow \mathbb{R}^n$  and  $\boldsymbol{\phi} : \mathbb{R} \rightarrow \mathbb{R}^n$  are given functions. Many dynamic processes in chemical engineering—for example, the distillation process described in [60]—can be modeled by equations (3.1) and (3.2). We assume that  $\mathbf{f}$ ,  $\mathbf{g}$ , and  $\boldsymbol{\phi}$  are continuously differentiable. We also assume that there exists a positive real number  $L_1 > 0$  such that for all  $\mathbf{x}', \mathbf{x}'' \in \mathbb{R}^n$  and  $\mathbf{u}', \mathbf{u}'' \in \mathbb{R}^r$ ,

$$|\mathbf{f}(\mathbf{x}', \mathbf{x}'', \mathbf{u}', \mathbf{u}'')| \leq L_1(1 + |\mathbf{x}'| + |\mathbf{x}''| + |\mathbf{u}'| + |\mathbf{u}''|), \quad (3.3)$$

where  $|\cdot|$  denotes the Euclidean norm. This assumption is standard in the control systems literature [93, 111, 118, 122, 124].

The output  $\mathbf{y}(t)$  of system (3.1)-(3.2) is defined by

$$\mathbf{y}(t) = \mathbf{g}(\mathbf{x}(t)), \quad t \in [0, T], \quad (3.4)$$

where  $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^q$  is a given continuously differentiable function.

We refer to  $\alpha$  as the *state-delay* and  $\beta$  as the *input-delay*. The exact values of these

delays are unknown; the only information we are given is that  $\alpha$  lies within the interval  $[\alpha_{\min}, \alpha_{\max}]$  and  $\beta$  lies within the interval  $[\beta_{\min}, \beta_{\max}]$ , where  $\alpha_{\min} \geq 0$  and  $\beta_{\min} > 0$ . Thus, we have the following bound constraints:

$$\alpha_{\min} \leq \alpha \leq \alpha_{\max}, \quad (3.5)$$

$$\beta_{\min} \leq \beta \leq \beta_{\max}. \quad (3.6)$$

We assume that the input signal  $\mathbf{u}$  is a given piecewise-constant function (this is the case in many engineering systems). Hence,  $\mathbf{u}$  can be expressed as follows:

$$\mathbf{u}(t) = \boldsymbol{\sigma}^i, \quad t \in [t_{i-1}, t_i), \quad i = 1, \dots, p, \quad (3.7)$$

where  $\boldsymbol{\sigma}^i \in \mathbb{R}^r$ ,  $i = 1, \dots, p$ , are given vectors and  $t_i$ ,  $i = 0, \dots, p$ , are given time points such that  $-\beta_{\max} = t_0 < t_1 < \dots < t_p = T$ . Equation (3.7) can be rewritten as

$$\mathbf{u}(t) = \sum_{i=1}^p \boldsymbol{\sigma}^i \chi_{[t_{i-1}, t_i)}(t), \quad t \in [-\beta_{\max}, T], \quad (3.8)$$

where the characteristic function  $\chi_{[t_{i-1}, t_i)} : \mathbb{R} \rightarrow \mathbb{R}$  is defined by

$$\chi_{[t_{i-1}, t_i)}(t) = \begin{cases} 1, & \text{if } t \in [t_{i-1}, t_i), \\ 0, & \text{otherwise.} \end{cases}$$

For each pair  $(\alpha, \beta) \in [\alpha_{\min}, \alpha_{\max}] \times [\beta_{\min}, \beta_{\max}]$ , let  $\mathbf{x}(\cdot | \alpha, \beta)$  denote the corresponding solution of system (3.1)-(3.2). Substituting  $\mathbf{x}(\cdot | \alpha, \beta)$  into equation (3.4) gives  $\mathbf{y}(\cdot | \alpha, \beta)$ , the predicted system output corresponding to  $(\alpha, \beta)$ . Mathematically,

$$\mathbf{y}(t | \alpha, \beta) = \mathbf{g}(\mathbf{x}(t | \alpha, \beta)), \quad t \leq T. \quad (3.9)$$

Suppose that the output from system (3.1)-(3.2) has been measured experimentally at a set of sample times  $\{\tau_j\}_{j=1}^m \subset [0, T]$ . Let  $\hat{\mathbf{y}}^j \in \mathbb{R}^q$  denote the measured output at the  $j$ th sample time. Our goal is to use the experimental data  $\{(\tau_j, \hat{\mathbf{y}}^j)\}_{j=1}^m$  to identify the time-delays  $\alpha$  and  $\beta$ . We formulate this problem mathematically as follows.

**Problem (P).** *Choose the state-delay  $\alpha$  and the input-delay  $\beta$  to minimize the cost function*

$$J(\alpha, \beta) = \sum_{j=1}^m |\mathbf{y}(\tau_j | \alpha, \beta) - \hat{\mathbf{y}}^j|^2 \quad (3.10)$$

*subject to the dynamic system (3.1)-(3.2) and the bound constraints (3.5)-(3.6).*

Problem (P) is a dynamic optimization problem governed by the time-delay system (3.1)-(3.2). The most interesting aspect of Problem (P) is that the time-delays in (3.1)-



(3.2) are actually decision variables to be chosen optimally. This is highly unusual; in most optimization problems involving time-delay systems, the delays are fixed and known, and the control input function is the decision variable to be chosen optimally [74, 93, 119]. In Problem (P), the input function is known, and the delays are the variables that need to be optimized.

We now conclude this section by showing that Problem (P) can be transformed into a *switched system optimal control problem*.

First, from (3.8),

$$\mathbf{u}(t - \beta) = \sum_{i=1}^p \boldsymbol{\sigma}^i \chi_{[t_{i-1}, t_i)}(t - \beta) = \sum_{i=1}^p \boldsymbol{\sigma}^i \chi_{[t_{i-1} + \beta, t_i + \beta)}(t) = \sum_{i=1}^p \boldsymbol{\sigma}^i \chi_{[v_{i-1}, v_i)}(t), \quad (3.11)$$

where  $v_i$ ,  $i = 0, \dots, p$  are new decision variables defined by

$$v_i = t_i + \beta, \quad i = 0, \dots, p. \quad (3.12)$$

It follows from (3.12) that

$$v_i - t_i = v_{i-1} - t_{i-1}, \quad i = 1, \dots, p. \quad (3.13)$$

Substituting (3.11) into (3.1) gives

$$\dot{\mathbf{x}}(t) = \mathbf{f}^i(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t)), \quad t \in [v_{i-1}, v_i) \cap [0, T], \quad i = 1, \dots, p, \quad (3.14)$$

where

$$\mathbf{f}^i(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t)) = \mathbf{f}(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t), \boldsymbol{\sigma}^i).$$

System (3.14) is a *switched system* in which the dynamics change instantaneously at the *switching times*  $v_i$ ,  $i = 1, \dots, p$ .

Problem (P) can now be restated as follows: Choose the state-delay  $\alpha$  and the switching times  $v_i$ ,  $i = 1, \dots, p$  to minimize (3.10) subject to the switched system (3.14), the initial condition (3.2), and the constraints (3.5)-(3.6) and (3.13). This is an example of a switched system optimal control problem. Such problems have been the subject of active research over the last decade (see, for example, [24, 26, 69, 100] and the references cited therein). In particular, the well-known *time-scaling transformation* is a powerful tool for solving switched system optimal control problems (see [118, 124, 132, 151]). Unfortunately, the time-scaling transformation is not applicable to time-delay systems such as system (3.14) defined above. Thus, a new method is needed to solve Problem (P).

### 3.3 State variation

Our goal is to solve Problem (P) using nonlinear optimization techniques. To do this, we need the partial derivatives of  $J$  with respect to the decision variables  $\alpha$  and  $\beta$ . However, since  $J$  is not an explicit function of  $\alpha$  and  $\beta$ , these partial derivatives cannot be determined using standard differentiation rules. To derive formulae for the partial derivatives of  $J$ , we first need to consider the *state variation* with respect to  $\alpha$  and  $\beta$ .

#### 3.3.1 State variation with respect to the state-delay

Define

$$\boldsymbol{\psi}(t) = \begin{cases} \dot{\boldsymbol{\phi}}(t), & \text{if } t \leq 0, \\ \mathbf{f}(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t), \mathbf{u}(t - \beta)), & \text{if } t \in (0, T]. \end{cases}$$

Furthermore, let  $\frac{\partial}{\partial \tilde{\mathbf{x}}}$  denote differentiation with respect to the delayed state argument. We will use this notation frequently throughout this chapter.

The solution  $\mathbf{x}(\cdot|\alpha, \beta)$  of system (3.1)-(3.2) is normally viewed as a function of time, with  $\alpha$  and  $\beta$  being fixed values. By instead of fixing  $t \in (-\infty, T]$ , while allowing  $\alpha$  and  $\beta$  to vary, we obtain function  $\mathbf{x}(t|\cdot, \cdot) : [\alpha_{\min}, \alpha_{\max}] \times [\beta_{\min}, \beta_{\max}] \rightarrow \mathbb{R}^n$  whose value at  $(\alpha, \beta)$  is  $\mathbf{x}(t|\alpha, \beta)$ . The partial derivative of  $\mathbf{x}(t|\cdot, \cdot)$  with respect to  $\alpha$  is called the *state variation* with respect to  $\alpha$ . The following result, which can be proved in a similar manner to the proof of Theorem 2.1 in Chapter 2, gives a method for determining this state variation.

**Theorem 3.1.** *Let  $t \in (0, T]$  be a fixed time point. Then  $\mathbf{x}(t|\cdot, \cdot)$  is differentiable with respect to the state-delay  $\alpha$ . In fact, for each  $(\alpha, \beta) \in [\alpha_{\min}, \alpha_{\max}] \times [\beta_{\min}, \beta_{\max}]$ ,*

$$\frac{\partial \mathbf{x}(t|\alpha, \beta)}{\partial \alpha} = \boldsymbol{\Lambda}(t|\alpha, \beta), \quad (3.15)$$

where  $\boldsymbol{\Lambda}(\cdot|\alpha, \beta)$  satisfies the auxiliary time-delay system

$$\begin{aligned} \dot{\boldsymbol{\Lambda}}(t) &= \frac{\partial \mathbf{f}(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t), \mathbf{u}(t - \beta))}{\partial \mathbf{x}} \boldsymbol{\Lambda}(t) \\ &\quad + \frac{\partial \mathbf{f}(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t), \mathbf{u}(t - \beta))}{\partial \tilde{\mathbf{x}}} \boldsymbol{\Lambda}(t - \alpha) \\ &\quad - \frac{\partial \mathbf{f}(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t), \mathbf{u}(t - \beta))}{\partial \tilde{\mathbf{x}}} \boldsymbol{\psi}(t - \alpha) \end{aligned} \quad (3.16)$$

with initial condition

$$\boldsymbol{\Lambda}(t) = \mathbf{0}, \quad t \leq 0. \quad (3.17)$$

According to Theorem 3.1, the state variation with respect to  $\alpha$  can be computed by solving the auxiliary time-delay system (3.16)-(3.17). This result is a simple extension of

the main result in Chapter 2, which pertains to systems with multiple state-delays but no input delays. To solve Problem (P), we also need the state variation with respect to  $\beta$ . Unfortunately, the results in Chapter 2, which are based on the assumption that the system dynamics are continuous with respect to the time-delays, cannot be used to determine this state variation. Indeed, since the input function  $\mathbf{u}$  is discontinuous, the dynamics (3.1) are clearly discontinuous with respect to  $\beta$ . In the next subsection, we describe a new method for computing the state variation with respect to  $\beta$ .

### 3.3.2 State variation with respect to the input-delay

#### A Preliminaries

Before deriving the state variation with respect to  $\beta$ , we first need to derive several preliminary results. Let  $(\alpha, \beta) \in [\alpha_{\min}, \alpha_{\max}] \times [\beta_{\min}, \beta_{\max}]$  be a fixed pair. Define

$$\mathcal{S} = [\beta_{\min} - \beta, \beta_{\max} - \beta].$$

Note that  $\mathcal{S}$  is a non-empty closed interval of positive measure. Clearly,

$$\epsilon \in \mathcal{S} \iff \beta + \epsilon \in [\beta_{\min}, \beta_{\max}].$$

Now, for each  $\epsilon \in \mathcal{S}$ , define

$$\begin{aligned} \varphi^\epsilon(t) &= \mathbf{x}(t|\alpha, \beta + \epsilon) - \mathbf{x}(t|\alpha, \beta), \quad t \leq T, \\ \mathbf{x}^\epsilon &= \mathbf{x}(t|\alpha, \beta + \epsilon). \end{aligned} \tag{3.18}$$

By (3.2),

$$\varphi^\epsilon(t) = \mathbf{0}, \quad t \leq 0. \tag{3.19}$$

Since the system dynamics satisfy the linear growth condition (3.3), it can be shown (see [80]) that there exists a positive real number  $L_2 > 0$  such that

$$|\mathbf{x}(t|\alpha, \beta + \epsilon)| \leq L_2, \quad t \in [-\alpha_{\max}, T], \quad \epsilon \in \mathcal{S}. \tag{3.20}$$

Our first preliminary result is stated and proved below.

**Lemma 3.1.** *There exists a positive real number  $L_3 > 0$  such that for all  $\epsilon \in \mathcal{S}$  of sufficiently small magnitude,*

$$|\varphi^\epsilon(t)| \leq L_3|\epsilon|, \quad t \in (-\infty, T]. \tag{3.21}$$

*Proof.* Let  $\epsilon \in \mathcal{S}$  be such that

$$|\epsilon| < \frac{1}{2} \min \{t_i - t_{i-1}\}_{i=1}^p.$$

For each  $i = 1, \dots, p$ , define  $I_i = (t_{i-1} + |\epsilon|, t_i - |\epsilon|)$ . Furthermore, for each  $i = 0, \dots, p$ , define

$$J_i = \begin{cases} [t_0, t_0 + |\epsilon|], & i = 0, \\ [t_i - |\epsilon|, t_i + |\epsilon|], & i = 1, \dots, p-1, \\ [t_p - |\epsilon|, t_p], & i = p. \end{cases}$$

Note that  $\{I_i\}_{i=1}^p$  and  $\{J_i\}_{i=0}^p$  form a partition of  $[-\beta_{\max}, T]$ . Also,  $|J_i| \leq 2|\epsilon|$ ,  $i = 0, \dots, p$ , and

$$\mathbf{u}(s) = \mathbf{u}(s - \epsilon) = \boldsymbol{\sigma}^i, \quad s \in I_i, \quad i = 1, \dots, p. \quad (3.22)$$

Now, if  $t \leq 0$ , then  $\boldsymbol{\varphi}^\epsilon(t) = \mathbf{0}$  and the proof is complete. Thus, assume that  $t > 0$ . Then

$$\begin{aligned} |\boldsymbol{\varphi}^\epsilon(t)| &\leq |\mathbf{x}^\epsilon(t) - \mathbf{x}(t)| \\ &\leq \int_0^t \left| \mathbf{f}(\mathbf{x}^\epsilon(s), \mathbf{x}^\epsilon(s - \alpha), \mathbf{u}(s), \mathbf{u}(s - \beta - \epsilon)) - \mathbf{f}(\mathbf{x}(s), \mathbf{x}(s - \alpha), \mathbf{u}(s), \mathbf{u}(s - \beta)) \right| ds, \end{aligned}$$

where  $\mathbf{x}^\epsilon(s) = \mathbf{x}(s|\alpha, \beta + \epsilon)$  and  $\mathbf{x}(s) = \mathbf{x}(s|\alpha, \beta)$ .

Thus, since  $\mathbf{x}^\epsilon$  is uniformly bounded with respect to  $\epsilon \in \mathcal{S}$  (recall (3.20)) and  $\mathbf{f}$  is continuously differentiable, there exists a constant  $M_1 > 0$  such that

$$|\boldsymbol{\varphi}^\epsilon(t)| \leq M_1 \int_0^t |\boldsymbol{\varphi}^\epsilon(s)| ds + M_1 \int_0^t |\boldsymbol{\varphi}^\epsilon(s - \alpha)| ds + M_1 \int_0^t |\mathbf{u}(s - \beta - \epsilon) - \mathbf{u}(s - \beta)| ds.$$

By shifting the time variable in the second and third integrals and then using (3.19), we obtain

$$\begin{aligned} |\boldsymbol{\varphi}^\epsilon(t)| &\leq M_1 \int_0^t |\boldsymbol{\varphi}^\epsilon(s)| ds + M_1 \int_{-\alpha}^{t-\alpha} |\boldsymbol{\varphi}^\epsilon(s)| ds + M_1 \int_{-\beta}^{t-\beta} |\mathbf{u}(s - \epsilon) - \mathbf{u}(s)| ds \\ &\leq 2M_1 \int_0^t |\boldsymbol{\varphi}^\epsilon(s)| ds + M_1 \int_{-\beta}^{t-\beta} |\mathbf{u}(s - \epsilon) - \mathbf{u}(s)| ds \\ &= 2M_1 \int_0^t |\boldsymbol{\varphi}^\epsilon(s)| ds + M_1 \sum_{i=1}^p \int_{I_i \cap (-\beta, t-\beta)} |\mathbf{u}(s - \epsilon) - \mathbf{u}(s)| ds \\ &\quad + M_1 \sum_{i=0}^p \int_{J_i \cap (-\beta, t-\beta)} |\mathbf{u}(s - \epsilon) - \mathbf{u}(s)| ds. \end{aligned}$$

Hence, by (3.22),

$$|\varphi^\epsilon(t)| \leq 2M_1 \int_0^t |\varphi^\epsilon(s)| ds + M_1 M_2 \sum_{i=0}^p |J_i|,$$

where  $M_2 = \max_{j \neq k} |\sigma^j - \sigma^k|$ . Since  $|J_i| \leq 2|\epsilon|$ , we have

$$|\varphi^\epsilon(t)| \leq 2M_1 \int_0^t |\varphi^\epsilon(s)| ds + 2(p+1)M_1 M_2 |\epsilon|.$$

Finally, applying the Gronwall-Bellman Lemma [111] yields

$$|\varphi^\epsilon(t)| \leq 2(p+1)M_1 M_2 \exp\{2M_1 T\} |\epsilon|.$$

This completes the proof.  $\square$

For each  $\epsilon \in \mathcal{S}$ , define

$$\bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma}) = \mathbf{f}(\mathbf{x}(s) + \eta\varphi^\epsilon(s), \mathbf{x}(s - \alpha) + \eta\varphi^\epsilon(s - \alpha), \mathbf{u}(s), \boldsymbol{\sigma}),$$

where, as in the proof of Lemma 3.1, let  $\mathbf{x}(t) = \mathbf{x}(t|\alpha, \beta)$ . Then by the chain rule,

$$\frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \eta} = \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \mathbf{x}} \varphi^\epsilon(s) + \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \tilde{\mathbf{x}}} \varphi^\epsilon(s - \alpha), \quad (3.23)$$

where

$$\frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}(\mathbf{x}(s) + \eta\varphi^\epsilon(s), \mathbf{x}(s - \alpha) + \eta\varphi^\epsilon(s - \alpha), \mathbf{u}(s), \boldsymbol{\sigma})}{\partial \mathbf{x}}, \quad (3.24)$$

$$\frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \tilde{\mathbf{x}}} = \frac{\partial \mathbf{f}(\mathbf{x}(s) + \eta\varphi^\epsilon(s), \mathbf{x}(s - \alpha) + \eta\varphi^\epsilon(s - \alpha), \mathbf{u}(s), \boldsymbol{\sigma})}{\partial \tilde{\mathbf{x}}}. \quad (3.25)$$

We can rewrite (3.23) as follows:

$$\begin{aligned} \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \eta} &= \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})}{\partial \mathbf{x}} \varphi^\epsilon(s) + \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})}{\partial \tilde{\mathbf{x}}} \varphi^\epsilon(s - \alpha) \\ &\quad + \Delta_1(s, \eta, \boldsymbol{\sigma}) + \Delta_2(s, \eta, \boldsymbol{\sigma}), \end{aligned} \quad (3.26)$$

where

$$\Delta_1(s, \eta, \boldsymbol{\sigma}) = \left\{ \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \mathbf{x}} - \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})}{\partial \mathbf{x}} \right\} \varphi^\epsilon(s), \quad (3.27)$$

$$\Delta_2(s, \eta, \boldsymbol{\sigma}) = \left\{ \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \tilde{\mathbf{x}}} - \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})}{\partial \tilde{\mathbf{x}}} \right\} \varphi^\epsilon(s - \alpha). \quad (3.28)$$

Since  $\mathbf{f}$  is continuously differentiable and  $\mathbf{x}$  and  $\mathbf{u}$  are bounded, the following result is easily established.

**Lemma 3.2.** For each  $\boldsymbol{\sigma} \in \mathbb{R}^r$ , there exists a corresponding  $L_4 > 0$  such that

$$\left| \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})}{\partial \mathbf{x}} \right| \leq L_4, \quad \left| \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})}{\partial \tilde{\mathbf{x}}} \right| \leq L_4, \quad s \in [0, T], \quad (3.29)$$

where  $|\cdot|$  denotes the Euclidean norm on  $\mathbb{R}^{n \times n}$ .

We now show that the functions  $\Delta_1$  (defined by (3.27)) and  $\Delta_2$  (defined by (3.28)) are of order  $\epsilon$ .

**Lemma 3.3.** Let  $\delta > 0$  and  $\boldsymbol{\sigma} \in \mathbb{R}^r$  be arbitrary. Then for any  $\epsilon \in \mathcal{S}$  of sufficiently small magnitude,

$$|\Delta_1(s, \eta, \boldsymbol{\sigma})| \leq L_3 \delta |\epsilon|, \quad |\Delta_2(s, \eta, \boldsymbol{\sigma})| \leq L_3 \delta |\epsilon|,$$

where  $L_3 > 0$  is as defined in Lemma 3.1.

*Proof.* By (3.21),  $\mathbf{x}(s) + \eta \boldsymbol{\varphi}^\epsilon(s) \rightarrow \mathbf{x}(s)$  and  $\mathbf{x}(s - \alpha) + \eta \boldsymbol{\varphi}^\epsilon(s - \alpha) \rightarrow \mathbf{x}(s - \alpha)$  uniformly on  $[0, T]$  as  $\epsilon \rightarrow 0$ . Hence, since  $\mathbf{f}$  is continuously differentiable and  $\mathbf{x}^\epsilon$  is uniformly bounded with respect to  $\epsilon$ , there exists an  $\epsilon' > 0$  such that for any  $\epsilon \in \mathcal{S}$  satisfying  $|\epsilon| < \epsilon'$ ,

$$\left| \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \mathbf{x}} - \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})}{\partial \mathbf{x}} \right| < \delta, \\ \left| \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \tilde{\mathbf{x}}} - \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})}{\partial \tilde{\mathbf{x}}} \right| < \delta.$$

By taking the norm of (3.27)-(3.28), and then using the above inequalities together with (3.21), we obtain the desired result.  $\square$

Let  $a$  and  $b$  be given constants such that  $a, b \in [0, T]$ . Define

$$\rho_\epsilon(a, b, \boldsymbol{\sigma}) = \int_a^b \{ \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}) - \bar{\mathbf{f}}^\epsilon(a, 0, \boldsymbol{\sigma}) \} ds. \quad (3.30)$$

Our final preliminary result is stated and proved below.

**Lemma 3.4.** For each  $\boldsymbol{\sigma} \in \mathbb{R}^r$ , there exists a corresponding  $L_5 > 0$  such that for all  $\epsilon \in \mathcal{S}$  of sufficiently small magnitude,

$$|\rho_\epsilon(a, b, \boldsymbol{\sigma})| \leq L_5 |b - a| \cdot |\epsilon| + L_5 (b - a)^2 + L_5 \int_{\min\{a, b\}}^{\max\{a, b\}} |\mathbf{u}(s) - \mathbf{u}(a)| ds.$$

*Proof.* From (3.30),

$$\rho_\epsilon(a, b, \boldsymbol{\sigma}) = \int_a^b \{ \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}) \} ds + \int_a^b \{ \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}) - \bar{\mathbf{f}}^\epsilon(a, 0, \boldsymbol{\sigma}) \} ds.$$

Thus,

$$\begin{aligned} |\rho_\epsilon(a, b, \boldsymbol{\sigma})| &\leq \int_{\min\{a, b\}}^{\max\{a, b\}} |\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})| ds \\ &\quad + \int_{\min\{a, b\}}^{\max\{a, b\}} |\bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}) - \bar{\mathbf{f}}^\epsilon(a, 0, \boldsymbol{\sigma})| ds. \end{aligned} \quad (3.31)$$

Consider the first integrand on the right-hand side of (3.31). Using (3.26) and (3.29) yields

$$\begin{aligned} |\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})| &\leq \int_0^1 \left| \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma})}{\partial \eta} \right| d\eta \\ &\leq \int_0^1 \{|\Delta_1(s, \eta, \boldsymbol{\sigma})| + |\Delta_2(s, \eta, \boldsymbol{\sigma})|\} d\eta \\ &\quad + L_4 |\boldsymbol{\varphi}^\epsilon(s)| + L_4 |\boldsymbol{\varphi}^\epsilon(s - \alpha)|. \end{aligned}$$

By Lemma 3.1 and Lemma 3.3 with  $\delta = 1$ , we see that for any  $\epsilon \in \mathcal{S}$  of sufficiently small magnitude,

$$|\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma})| \leq 2L_3 |\epsilon| + 2L_3 L_4 |\epsilon|. \quad (3.32)$$

Now, consider the second integrand on the right-hand side of (3.31). Since  $\mathbf{f}$  is continuously differentiable and  $\mathbf{x}^\epsilon$  is uniformly bounded with respect to  $\epsilon$  (recall (3.20)), there exists a constant  $M_3 > 0$  such that

$$|\bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}) - \bar{\mathbf{f}}^\epsilon(a, 0, \boldsymbol{\sigma})| \leq M_3 |\mathbf{x}(s) - \mathbf{x}(a)| + M_3 |\mathbf{x}(s - \alpha) - \mathbf{x}(a - \alpha)| + M_3 |\mathbf{u}(s) - \mathbf{u}(a)|.$$

Note that  $\dot{\mathbf{x}}(s) = \boldsymbol{\psi}(s)$  for almost all  $s \in (-\infty, T]$ , where  $\boldsymbol{\psi}$  is as defined in Subsection 3.3.1. Thus,

$$\begin{aligned} |\bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}) - \bar{\mathbf{f}}^\epsilon(a, 0, \boldsymbol{\sigma})| &\leq M_3 \int_{\min\{a, s\}}^{\max\{a, s\}} |\boldsymbol{\psi}(\eta)| d\eta + M_3 \int_{\min\{a, s\} - \alpha}^{\max\{a, s\} - \alpha} |\boldsymbol{\psi}(\eta)| d\eta \\ &\quad + M_3 |\mathbf{u}(s) - \mathbf{u}(a)| \\ &\leq M_3 M_4 |s - a| + M_3 M_4 |s - a| + M_3 |\mathbf{u}(s) - \mathbf{u}(a)|, \end{aligned} \quad (3.33)$$

where  $M_4 = \max_{\eta \in [-\alpha_{\max}, T]} |\boldsymbol{\psi}(\eta)|$ . Substituting (3.32) and (3.33) into (3.31) gives

$$|\rho_\epsilon(a, b, \boldsymbol{\sigma})| \leq (2L_3 + 2L_3 L_4) |b - a| \cdot |\epsilon| + 2M_3 M_4 (b - a)^2 + M_3 \int_{\min\{a, b\}}^{\max\{a, b\}} |\mathbf{u}(s) - \mathbf{u}(a)| ds.$$

Taking  $L_5 = \max\{2L_3 + 2L_3 L_4, 2M_3 M_4, M_3\}$  completes the proof.  $\square$

### B Main result

Equipped with Lemmas 3.1-3.4, we are now ready to derive the state variation with respect to the input-delay  $\beta$ . First, define

$$\mathcal{I} = \{t_i + \beta, i = 0, \dots, p\}.$$

Consider the following auxiliary system:

$$\begin{aligned} \dot{\Gamma}(t) &= \frac{\partial \mathbf{f}(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t), \mathbf{u}(t - \beta))}{\partial \mathbf{x}} \Gamma(t) \\ &+ \frac{\partial \mathbf{f}(\mathbf{x}(t), \mathbf{x}(t - \alpha), \mathbf{u}(t), \mathbf{u}(t - \beta))}{\partial \tilde{\mathbf{x}}} \Gamma(t - \alpha), \end{aligned} \quad (3.34)$$

where, for each  $t \in \mathcal{I} \cap (0, T]$ ,

$$\lim_{t \rightarrow (t_i + \beta)^+} \Gamma(t) = \lim_{t \rightarrow (t_i + \beta)^-} \Gamma(t) + \bar{\mathbf{f}}^\epsilon(t_i + \beta, 0, \boldsymbol{\sigma}^i) - \bar{\mathbf{f}}^\epsilon(t_i + \beta, 0, \boldsymbol{\sigma}^{i+1}), \quad (3.35)$$

and

$$\Gamma(t) = \mathbf{0}, \quad t \leq 0. \quad (3.36)$$

Let  $\Gamma(\cdot | \alpha, \beta)$  denote the unique right continuous solution of (3.34)-(3.36). We have the following important result.

**Theorem 3.2.** *Let  $(\alpha, \beta) \in [\alpha_{\min}, \alpha_{\max}] \times [\beta_{\min}, \beta_{\max}]$  be a fixed pair such that*

$$t_i + \beta \notin \{0\} \cup \{t_j, j = 0, \dots, p\}, \quad i = 0, \dots, p.$$

*Furthermore, consider a fixed time point  $t \in (t_{i-1} + \beta, t_i + \beta) \cap (0, T]$ , where  $i \in \{1, \dots, p\}$ .*

*Then*

$$\lim_{\epsilon \rightarrow 0^+} \epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) = \Gamma(t | \alpha, \beta), \quad (3.37)$$

*where  $\boldsymbol{\varphi}^\epsilon$  is as defined in (3.18).*

*Proof.* Let

$$a_i = \max\{t_{i-1} + \beta, 0\}.$$

Then

$$\mathbf{x}(t) = \mathbf{x}(a_i) + \int_{a_i}^t \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i) ds. \quad (3.38)$$

Let  $\epsilon \in \mathcal{S}$  be sufficiently small so that  $0 < \epsilon < \min\{t_j - t_{j-1}\}_{j=1}^p$  and  $t > t_{i-1} + \beta + \epsilon$ .



Define

$$a_i^\epsilon = \max\{t_{i-1} + \beta + \epsilon, 0\}.$$

Then

$$\mathbf{x}^\epsilon(t) = \mathbf{x}^\epsilon(a_i^\epsilon) + \int_{a_i^\epsilon}^t \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) ds. \quad (3.39)$$

We can write (3.39) as follows:

$$\begin{aligned} \mathbf{x}^\epsilon(t) &= \mathbf{x}^\epsilon(a_i) + \int_{a_i}^{a_i^\epsilon} \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^{i-1}) ds + \int_{a_i^\epsilon}^t \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) ds \\ &= \mathbf{x}^\epsilon(a_i) + \int_{a_i}^{a_i^\epsilon} \{\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^{i-1}) - \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i)\} ds + \int_{a_i^\epsilon}^t \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) ds, \end{aligned} \quad (3.40)$$

where  $\boldsymbol{\sigma}^{i-1}$  is arbitrary if  $i = 1$  (in this case, we must have  $a_i^\epsilon = a_i = 0$  when  $\epsilon$  is sufficiently small, because  $\beta < \beta_{\max}$ ). From (3.38) and (3.40), we have

$$\begin{aligned} \boldsymbol{\varphi}^\epsilon(t) &= \mathbf{x}^\epsilon(t) - \mathbf{x}(t) \\ &= \boldsymbol{\varphi}^\epsilon(a_i) + \int_{a_i}^t \{\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)\} ds + \int_{a_i}^{a_i^\epsilon} \{\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^{i-1}) - \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i)\} ds. \end{aligned}$$

Thus,

$$\begin{aligned} \boldsymbol{\varphi}^\epsilon(t) &= \boldsymbol{\varphi}^\epsilon(a_i) + \int_{a_i}^t \{\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)\} ds - \rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^i) + \rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^{i-1}) \\ &\quad + (a_i^\epsilon - a_i) \{\bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^{i-1}) - \bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^i)\}, \end{aligned}$$

where  $\rho_\epsilon$  is as defined in (3.30). By the fundamental theorem of calculus,

$$\begin{aligned} \boldsymbol{\varphi}^\epsilon(t) &= \boldsymbol{\varphi}^\epsilon(a_i) + \int_{a_i}^t \int_0^1 \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma}^i)}{\partial \eta} d\eta ds - \rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^i) + \rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^{i-1}) \\ &\quad + (a_i^\epsilon - a_i) \{\bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^{i-1}) - \bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^i)\}. \end{aligned}$$

Using (3.26),

$$\begin{aligned} \boldsymbol{\varphi}^\epsilon(t) &= \boldsymbol{\varphi}^\epsilon(a_i) + \int_{a_i}^t \int_0^1 \{\Delta_1(s, \eta, \boldsymbol{\sigma}^i) + \Delta_2(s, \eta, \boldsymbol{\sigma}^i)\} d\eta ds \\ &\quad + \int_{a_i}^t \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)}{\partial \mathbf{x}} \boldsymbol{\varphi}^\epsilon(s) ds + \int_{a_i}^t \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)}{\partial \tilde{\mathbf{x}}} \boldsymbol{\varphi}^\epsilon(s - \alpha) ds - \rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^i) \\ &\quad + \rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^{i-1}) + (a_i^\epsilon - a_i) \{\bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^{i-1}) - \bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^i)\}. \end{aligned} \quad (3.41)$$

We can express the solution of the auxiliary system as follows:

$$\Gamma(t) = \Gamma(a_i^+) + \int_{a_i}^t \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)}{\partial \mathbf{x}} \Gamma(s) ds + \int_{a_i}^t \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)}{\partial \tilde{\mathbf{x}}} \Gamma(s - \alpha) ds. \quad (3.42)$$

Thus, from Lemma 3.2 and equations (3.41) and (3.42),

$$\begin{aligned} |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) - \Gamma(t)| &\leq |\gamma_i(\epsilon)| + \epsilon^{-1} \int_{a_i}^t \int_0^1 \{|\Delta_1(s, \eta, \boldsymbol{\sigma}^i)| + |\Delta_2(s, \eta, \boldsymbol{\sigma}^i)|\} d\eta ds \\ &\quad + \int_{a_i}^t L_4 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \Gamma(s)| ds + \int_{a_i}^t L_4 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s - \alpha) - \Gamma(s - \alpha)| ds \\ &\quad + \epsilon^{-1} |\rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^i)| + \epsilon^{-1} |\rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^{i-1})|, \end{aligned}$$

where  $L_4$  is the constant defined in Lemma 3.2 and

$$\gamma_i(\epsilon) = \epsilon^{-1} \boldsymbol{\varphi}^\epsilon(a_i) - \Gamma(a_i^+) + \epsilon^{-1} (a_i^\epsilon - a_i) \{ \bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^{i-1}) - \bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^i) \}. \quad (3.43)$$

Recall that  $a_i^\epsilon - a_i \leq \epsilon$  and  $t_{i-1} + \beta \neq t_j$  for all  $j$ . Thus, we may assume that  $\epsilon$  is sufficiently small so that  $\mathbf{u}(s) = \mathbf{u}(a_i)$  for all  $s \in [a_i, a_i^\epsilon]$ . It then follows from Lemma 3.4 that

$$|\rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^{i-1})| \leq 2L_5' \epsilon^2, \quad |\rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^i)| \leq 2L_5'' \epsilon^2,$$

where  $L_5'$  and  $L_5''$  are the constants in Lemma 3.4 corresponding to  $\boldsymbol{\sigma}^{i-1}$  and  $\boldsymbol{\sigma}^i$ , respectively. By the above inequalities and Lemma 3.3, assuming that  $\epsilon$  is sufficiently small,

$$\begin{aligned} |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) - \Gamma(t)| &\leq 2TL_3\delta + 2L_5'\epsilon + 2L_5''\epsilon + |\gamma_i(\epsilon)| + \int_{a_i}^t L_4 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \Gamma(s)| ds \\ &\quad + \int_{a_i}^t L_4 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s - \alpha) - \Gamma(s - \alpha)| ds, \end{aligned} \quad (3.44)$$

where  $\delta > 0$  is arbitrary and  $L_3$  is the constant defined in Lemma 3.1. Performing a change of variable in the second integral on the right-hand side of (3.44) yields

$$\begin{aligned} |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(t) - \Gamma(t)| &\leq 2TL_3\delta + 4L_5\epsilon + |\gamma_i(\epsilon)| + \int_{a_i}^t L_4 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \Gamma(s)| ds \\ &\quad + \int_{a_i - \alpha}^{t - \alpha} L_4 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \Gamma(s)| ds \\ &\leq 2TL_3\delta + 4L_5\epsilon + |\gamma_i(\epsilon)| + \mu_i(\epsilon) + \int_{a_i}^t 2L_4 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \Gamma(s)| ds, \end{aligned}$$

where  $L_5 = \max\{L_5', L_5''\}$  and

$$\mu_i(\epsilon) = \int_{a_i - \alpha}^{a_i} L_4 |\epsilon^{-1} \boldsymbol{\varphi}^\epsilon(s) - \Gamma(s)| ds.$$

Assuming that  $\delta$  is sufficiently small so that  $a_i + \delta \leq t$ ,

$$\begin{aligned} |\epsilon^{-1}\varphi^\epsilon(t) - \mathbf{\Gamma}(t)| &\leq 2TL_3\delta + 4L_5\epsilon + \int_{a_i}^{a_i+\delta} 2L_4|\epsilon^{-1}\varphi^\epsilon(s) - \mathbf{\Gamma}(s)|ds \\ &\quad + |\gamma_i(\epsilon)| + \mu_i(\epsilon) + \int_{a_i+\delta}^t 2L_4|\epsilon^{-1}\varphi^\epsilon(s) - \mathbf{\Gamma}(s)|ds. \end{aligned} \quad (3.45)$$

Now, since  $\mathbf{\Gamma}$  is a piecewise continuous function, there exists a constant  $M_1 > 0$  such that

$$|\mathbf{\Gamma}(s)| \leq M_1, \quad s \in (-\infty, T].$$

Therefore, it follows from Lemma 3.1 that for all sufficiently small  $\epsilon > 0$ ,

$$|\epsilon^{-1}\varphi^\epsilon(s) - \mathbf{\Gamma}(s)| \leq L_3 + M_1, \quad s \in (-\infty, T]. \quad (3.46)$$

Substituting (3.46) into (3.45) gives

$$\begin{aligned} |\epsilon^{-1}\varphi^\epsilon(t) - \mathbf{\Gamma}(t)| &\leq 2TL_3\delta + 4L_5\epsilon + |\gamma_i(\epsilon)| + \mu_i(\epsilon) + 2L_4(L_3 + M_1)\delta \\ &\quad + \int_{a_i+\delta}^t 2L_4|\epsilon^{-1}\varphi^\epsilon(s) - \mathbf{\Gamma}(s)|ds. \end{aligned} \quad (3.47)$$

Note that this inequality holds for all  $t \in [a_i + \delta, t_i + \beta)$  and  $t = (t_i + \beta)^-$ , uniformly with respect to  $\epsilon \leq \delta$ . Thus, by the Gronwall-Bellman Lemma [111],

$$|\epsilon^{-1}\varphi^\epsilon(t) - \mathbf{\Gamma}(t)| \leq (2TL_3\delta + 4L_5\epsilon + |\gamma_i(\epsilon)| + \mu_i(\epsilon) + 2L_4(L_3 + M_1)\delta) \exp\{2L_4T\}. \quad (3.48)$$

This inequality holds for all  $\epsilon$  of sufficiently small magnitude.

Now, suppose that  $t \in (t_{i-1} + \beta, t_i + \beta) \cap (0, T]$  for  $i = \min\{j : t_j + \beta > 0\}$ . Then  $a_i = 0$ , and thus by (3.2) and (3.36),

$$\mu_i(\epsilon) = \int_{-\alpha}^0 L_4|\epsilon^{-1}\varphi^\epsilon(s) - \mathbf{\Gamma}(s)|ds = 0.$$

Since by assumption  $t_{i-1} + \beta < 0$ ,  $a_i^\epsilon = a_i = 0$  for all sufficiently small  $\epsilon$ . Thus,

$$\gamma_i(\epsilon) = \epsilon^{-1}\varphi^\epsilon(0) - \mathbf{\Gamma}(0^+) = \mathbf{0}.$$

Substituting  $\mu_i(\epsilon) = 0$  and  $\gamma_i(\epsilon) = \mathbf{0}$  into (3.48) gives

$$|\epsilon^{-1}\varphi^\epsilon(t) - \mathbf{\Gamma}(t)| \leq (2TL_3\delta + 4L_5\epsilon + 2L_4(L_3 + M_1)\delta) \exp\{2L_4T\}. \quad (3.49)$$

Since  $\delta > 0$  was chosen arbitrarily and  $\epsilon$  can be made arbitrarily small, this shows that (3.37) holds for  $i = \min\{j : t_j + \beta > 0\}$ . Moreover, the derivation leading to (3.49) shows

that (3.37) also holds for  $t = (t_i + \beta)^-$ . It is also clear that (3.37) holds for all  $t \in (-\infty, 0]$ .

Now, suppose that (3.37) holds for all  $t \in (-\infty, t_k + \beta) \setminus \{t_j + \beta\}_{j=0}^k$  and  $t = (t_k + \beta)^-$ , where

$$\min\{j : t_j + \beta > 0\} \leq k \leq p - 1. \quad (3.50)$$

We will show that (3.37) holds for all  $t \in (t_k + \beta, t_{k+1} + \beta)$  and  $t = (t_{k+1} + \beta)^-$ . The result will then follow by induction.

Let  $t \in (t_k + \beta, t_{k+1} + \beta)$ , where  $k$  satisfies (3.50). By our inductive hypothesis, for almost all  $s \in (-\infty, t_k + \beta)$ ,

$$\lim_{\epsilon \rightarrow 0^+} \epsilon^{-1} \varphi^\epsilon(s) = \Gamma(s). \quad (3.51)$$

In view of (3.46) and (3.51), applying Lebesgue's dominated convergence theorem gives

$$\lim_{\epsilon \rightarrow 0^+} \mu_{k+1}(\epsilon) = \lim_{\epsilon \rightarrow 0^+} \int_{t_k + \beta - \alpha}^{t_k + \beta} L_4 |\epsilon^{-1} \varphi^\epsilon(s) - \Gamma(s)| ds = 0. \quad (3.52)$$

Furthermore,

$$\begin{aligned} \gamma_{k+1}(\epsilon) &= \epsilon^{-1} \varphi^\epsilon(t_k + \beta) - \Gamma((t_k + \beta)^+) + \bar{\mathbf{f}}^\epsilon(t_k + \beta, 0, \boldsymbol{\sigma}^k) - \bar{\mathbf{f}}^\epsilon(t_k + \beta, 0, \boldsymbol{\sigma}^{k+1}) \\ &= \epsilon^{-1} \varphi^\epsilon(t_k + \beta) - \Gamma((t_k + \beta)^-). \end{aligned}$$

Thus, by our inductive hypothesis,

$$\lim_{\epsilon \rightarrow 0^+} \gamma_{k+1}(\epsilon) = \mathbf{0}. \quad (3.53)$$

By combining equations (3.52) and (3.53) with (3.48) for  $i = k+1$ , we see that (3.37) holds for  $t \in (t_k + \beta, t_{k+1} + \beta)$ . Similar arguments show that (3.37) also holds for  $t = (t_{k+1} + \beta)^-$ . The proof then follows by induction.  $\square$

Theorem 3.2 shows that  $\epsilon^{-1} \varphi^\epsilon \rightarrow \Gamma(\cdot | \alpha, \beta)$  as  $\epsilon \rightarrow 0^+$ . We now derive the analogous result for  $\epsilon \rightarrow 0^-$ .

**Theorem 3.3.** *Let  $(\alpha, \beta) \in [\alpha_{\min}, \alpha_{\max}] \times (\beta_{\min}, \beta_{\max}]$  be a fixed pair such that*

$$t_i + \beta \notin \{0\} \cup \{t_j, j = 0, \dots, p\}, \quad i = 0, \dots, p.$$

*Furthermore, consider a fixed time point  $t \in (t_{i-1} + \beta, t_i + \beta) \cap (0, T]$ , where  $i \in \{1, \dots, p\}$ .*

*Then*

$$\lim_{\epsilon \rightarrow 0^-} \epsilon^{-1} \varphi^\epsilon(t) = \Gamma(t | \alpha, \beta). \quad (3.54)$$

*Proof.* Let  $a_i$  and  $a_i^\epsilon$  be as defined in the proof of Theorem 3.2. Furthermore, let  $\epsilon \in \mathcal{S}$

be such that  $\min\{t_{j-1} - t_j\}_{j=1}^p < \epsilon < 0$  and  $t < t_i + \beta + \epsilon$ . Then

$$\mathbf{x}(t) = \mathbf{x}(a_i^\epsilon) + \int_{a_i^\epsilon}^{a_i} \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^{i-1}) ds + \int_{a_i}^t \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i) ds,$$

where  $\boldsymbol{\sigma}^{i-1}$  is arbitrary if  $i = 1$  (in this case, we must have  $a_i = a_i^\epsilon = 0$ ). Moreover,

$$\mathbf{x}^\epsilon(t) = \mathbf{x}^\epsilon(a_i^\epsilon) + \int_{a_i^\epsilon}^{a_i} \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) ds + \int_{a_i}^t \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) ds.$$

Thus,

$$\begin{aligned} \boldsymbol{\varphi}^\epsilon(t) &= \mathbf{x}^\epsilon(t) - \mathbf{x}(t) \\ &= \boldsymbol{\varphi}^\epsilon(a_i^\epsilon) + \int_{a_i^\epsilon}^{a_i} \{\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^{i-1})\} ds + \int_{a_i}^t \{\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)\} ds. \end{aligned}$$

This equation can be rewritten as follows:

$$\begin{aligned} \boldsymbol{\varphi}^\epsilon(t) &= \boldsymbol{\varphi}^\epsilon(a_i^\epsilon) - \rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^i) + (a_i - a_i^\epsilon) \{\bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^i) - \bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^{i-1})\} \\ &\quad + \int_{a_i^\epsilon}^{a_i} \{\bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^{i-1}) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^{i-1})\} ds + \int_{a_i}^t \{\bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^i) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)\} ds, \end{aligned}$$

where  $\rho_\epsilon$  is as defined in (3.30). Using the fundamental theorem of calculus and (3.26),

$$\begin{aligned} \boldsymbol{\varphi}^\epsilon(t) &= \boldsymbol{\varphi}^\epsilon(a_i^\epsilon) - \rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^i) + \int_{a_i^\epsilon}^{a_i} \{\bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^{i-1}) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^{i-1})\} ds \\ &\quad + (a_i - a_i^\epsilon) \{\bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^i) - \bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^{i-1})\} + \int_{a_i}^t \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)}{\partial \mathbf{x}} \boldsymbol{\varphi}^\epsilon(s) ds \\ &\quad + \int_{a_i}^t \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^i)}{\partial \tilde{\mathbf{x}}} \boldsymbol{\varphi}^\epsilon(s - \alpha) ds + \int_{a_i}^t \int_0^1 \{\Delta_1(s, \eta, \boldsymbol{\sigma}^i) + \Delta_2(s, \eta, \boldsymbol{\sigma}^i)\} d\eta ds. \end{aligned}$$

Note that  $t_{i-1} + \beta \neq t_j$  for all  $j$ . Thus, we may assume that  $\epsilon$  is sufficiently small so that  $\mathbf{u}(s) = \mathbf{u}(a_i)$  for all  $s \in [a_i^\epsilon, a_i]$ . It then follows from Lemma 3.4 that

$$|\rho_\epsilon(a_i, a_i^\epsilon, \boldsymbol{\sigma}^{i-1})| \leq 2L_5\epsilon^2.$$

Furthermore, assuming that  $\epsilon$  is sufficiently small, by using a similar arguments to those in the proof of Lemma 3.4, one can show that there exists a constant  $M_5 > 0$  such that

$$\int_{a_i^\epsilon}^{a_i} |\bar{\mathbf{f}}^\epsilon(a_i, 0, \boldsymbol{\sigma}^{i-1}) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^{i-1})| ds \leq M_5\epsilon^2.$$

Hence, as in the proof of Theorem 3.2,

$$\begin{aligned} |\epsilon^{-1}\varphi^\epsilon(t) - \Gamma(t)| &\leq M_5|\epsilon| + 2L_5|\epsilon| + 2L_3T\delta + \int_{a_i}^t L_4|\epsilon^{-1}\varphi^\epsilon(s) - \Gamma(s)|ds \\ &\quad + |\gamma_i(\epsilon)| + \int_{a_i}^t L_4|\epsilon^{-1}\varphi^\epsilon(s - \alpha) - \Gamma(s - \alpha)|ds, \end{aligned} \quad (3.55)$$

where  $\delta > 0$  is arbitrary and

$$\gamma_i(\epsilon) = \epsilon^{-1}\varphi^\epsilon(a_i^\epsilon) - \Gamma(a_i^+) + \epsilon^{-1}(a_i - a_i^\epsilon) \{ \bar{f}^\epsilon(a_i, 0, \sigma^i) - \bar{f}^\epsilon(a_i, 0, \sigma^{i-1}) \}.$$

Simplifying (3.55) gives

$$|\epsilon^{-1}\varphi^\epsilon(t) - \Gamma(t)| \leq M_5|\epsilon| + 2L_5|\epsilon| + 2L_3T\delta + |\gamma_i(\epsilon)| + \mu_i(\epsilon) + \int_{a_i}^t 2L_4|\epsilon^{-1}\varphi^\epsilon(s) - \Gamma(s)|ds,$$

where

$$\mu_i(\epsilon) = \int_{a_i - \alpha}^{a_i} L_4|\epsilon^{-1}\varphi^\epsilon(s) - \Gamma(s)|ds.$$

Finally, by applying Gronwall's Lemma [111] yields

$$|\epsilon^{-1}\varphi^\epsilon(t) - \Gamma(t)| \leq (M_5|\epsilon| + 2L_5|\epsilon| + 2L_3T\delta + |\gamma_i(\epsilon)| + \mu_i(\epsilon)) \exp\{2L_4T\}. \quad (3.56)$$

In particular, this inequality also holds for all  $t = a_i^+$ , assuming that  $\epsilon$  is of sufficiently small magnitude.

Now, suppose  $t \in (t_{i-1} + \beta, t_i + \beta) \cap (0, T]$  for  $i = \min\{j : t_j + \beta > 0\}$ . Then  $a_i = 0$ , and thus by (3.2) and (3.36),

$$\mu_i(\epsilon) = \int_{-\alpha}^0 L_4|\epsilon^{-1}\varphi^\epsilon(s) - \Gamma(s)|ds = 0.$$

Also, since  $a_i^\epsilon = a_i = 0$ ,  $\gamma_i(\epsilon) = \epsilon^{-1}\varphi^\epsilon(0) - \Gamma(0^+) = \mathbf{0}$ . Substituting  $\mu_i(\epsilon) = 0$  and  $\gamma_i(\epsilon) = \mathbf{0}$  into (3.56) gives

$$|\epsilon^{-1}\varphi^\epsilon(t) - \Gamma(t)| \leq (M_5|\epsilon| + 2L_5|\epsilon| + 2L_3T\delta) \exp\{2L_4T\}.$$

Since  $\delta > 0$  was chosen arbitrarily and  $\epsilon$  can be made arbitrarily small, it follows that (3.54) holds for  $i = \min\{j : t_j + \beta > 0\}$ . It is clear that (3.54) also holds for all  $t \in (-\infty, 0]$ , and for  $t = a_i^+$ .

Now, suppose that (3.54) holds for all  $t \in (-\infty, t_k + \beta) \setminus \{t_j + \beta\}_{j=0}^k$  and  $t = a_k^+$ , where

$$\min\{j : t_j + \beta > 0\} \leq k \leq p - 1. \quad (3.57)$$

We will show that (3.54) holds for all  $t \in (t_k + \beta, t_{k+1} + \beta)$  and  $t = a_{k+1}^+$ . The result will then follow by induction.

Let  $t \in (t_k + \beta, t_{k+1} + \beta)$ , where  $k$  satisfies (3.57) above. By our inductive hypothesis, for almost all  $s \in (-\infty, t_k + \beta)$ ,

$$\lim_{\epsilon \rightarrow 0^-} \epsilon^{-1} \varphi^\epsilon(s) = \Gamma(s). \quad (3.58)$$

Thus, as in the proof of Theorem 3.2, we can apply Lebesgue's dominated convergence theorem to obtain

$$\lim_{\epsilon \rightarrow 0^+} \mu_{k+1}(\epsilon) = \int_{t_k + \beta - \alpha}^{t_k + \beta} L_4 | \epsilon^{-1} \varphi^\epsilon(s) - \Gamma(s) | ds = 0. \quad (3.59)$$

We have

$$\begin{aligned} \gamma_{k+1}(\epsilon) &= \epsilon^{-1} \varphi^\epsilon(t_k + \beta + \epsilon) - \Gamma((t_k + \beta)^+) + \bar{\mathbf{f}}^\epsilon(t_k + \beta, 0, \boldsymbol{\sigma}^k) - \bar{\mathbf{f}}^\epsilon(t_k + \beta, 0, \boldsymbol{\sigma}^{k+1}) \\ &= \epsilon^{-1} \varphi^\epsilon(t_k + \beta + \epsilon) - \Gamma((t_k + \beta)^-). \end{aligned}$$

Hence,

$$\begin{aligned} \gamma_{k+1}(\epsilon) &= \epsilon^{-1} \varphi^\epsilon(a_k) - \Gamma(a_k^+) + \epsilon^{-1} \int_{a_k}^{t_k + \beta + \epsilon} \{ \bar{\mathbf{f}}^\epsilon(s, 1, \boldsymbol{\sigma}^k) - \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^k) \} ds \\ &\quad - \int_{a_k}^{t_k + \beta} \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^k)}{\partial \mathbf{x}} \Gamma(s) ds - \int_{a_k}^{t_k + \beta} \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^k)}{\partial \tilde{\mathbf{x}}} \Gamma(s - \alpha) ds \\ &= \epsilon^{-1} \varphi^\epsilon(a_k) - \Gamma(a_k^+) + \epsilon^{-1} \int_{a_k}^{t_k + \beta + \epsilon} \int_0^1 \frac{\partial \bar{\mathbf{f}}^\epsilon(s, \eta, \boldsymbol{\sigma}^k)}{\partial \eta} d\eta ds \\ &\quad - \int_{a_k}^{t_k + \beta} \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^k)}{\partial \mathbf{x}} \Gamma(s) ds - \int_{a_k}^{t_k + \beta} \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^k)}{\partial \tilde{\mathbf{x}}} \Gamma(s - \alpha) ds. \end{aligned}$$

Using (3.26), we obtain

$$\begin{aligned} \gamma_{k+1}(\epsilon) &= \epsilon^{-1} \varphi^\epsilon(a_k) - \Gamma(a_k^+) + \epsilon^{-1} \int_{a_k}^{t_k + \beta + \epsilon} \int_0^1 \{ \Delta_1(s, \eta, \boldsymbol{\sigma}^k) + \Delta_2(s, \eta, \boldsymbol{\sigma}^k) \} d\eta ds \\ &\quad + \int_{a_k}^{t_k + \beta + \epsilon} \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^k)}{\partial \mathbf{x}} \{ \epsilon^{-1} \varphi^\epsilon(s) - \Gamma(s) \} ds \\ &\quad + \int_{a_k}^{t_k + \beta + \epsilon} \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^k)}{\partial \tilde{\mathbf{x}}} \{ \epsilon^{-1} \varphi^\epsilon(s - \alpha) - \Gamma(s - \alpha) \} ds \\ &\quad - \int_{t_k + \beta + \epsilon}^{t_k + \beta} \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^k)}{\partial \mathbf{x}} \Gamma(s) ds - \int_{t_k + \beta + \epsilon}^{t_k + \beta} \frac{\partial \bar{\mathbf{f}}^\epsilon(s, 0, \boldsymbol{\sigma}^k)}{\partial \tilde{\mathbf{x}}} \Gamma(s - \alpha) ds. \end{aligned}$$

Thus, using Lemma 3.2 and Lemma 3.3,

$$\begin{aligned}
|\gamma_{k+1}(\epsilon)| &= |\epsilon^{-1}\varphi^\epsilon(a_k) - \Gamma(a_k^+)| + 2L_3T\delta + \int_{a_k}^{t_k+\beta+\epsilon} L_4|\epsilon^{-1}\varphi^\epsilon(s) - \Gamma(s)|ds \\
&\quad + \int_{a_k}^{t_k+\beta+\epsilon} L_4|\epsilon^{-1}\varphi^\epsilon(s-\alpha) - \Gamma(s-\alpha)|ds + L_4M_1|\epsilon| + L_4M_1|\epsilon| \\
&\leq |\epsilon^{-1}\varphi^\epsilon(a_k) - \Gamma(a_k^+)| + 2L_3T\delta + 2L_4M_1|\epsilon| \\
&\quad + \int_{a_k}^{t_k+\beta} L_4|\epsilon^{-1}\varphi^\epsilon(s) - \Gamma(s)|ds + \int_{a_k}^{t_k+\beta} L_4|\epsilon^{-1}\varphi^\epsilon(s-\alpha) - \Gamma(s-\alpha)|ds,
\end{aligned}$$

where  $M_1$  is as defined in the proof of Theorem 3.2. Using the Lebesgue's dominated convergence theorem and the induction hypothesis, the two integrals converge to zero as  $\epsilon \rightarrow 0-$ . It follows also from the induction hypothesis that the first term converges to zero as  $\epsilon \rightarrow 0-$ . Thus,

$$\lim_{\epsilon \rightarrow 0-} |\gamma_{k+1}(\epsilon)| = 0. \quad (3.60)$$

Combining equations (3.59) and (3.60) with (3.56), for  $i = k + 1$ , we see that (3.54) holds for all  $t \in (t_k + \beta, t_{k+1} + \beta)$  and  $t = (t_k + \beta)^+$ . The proof then follows by induction.  $\square$

Together, Theorems 3.2 and 3.3 show that the state variation with respect to  $\beta$  is given by  $\Gamma(\cdot|\alpha, \beta)$ . This is stated formally in the following theorem.

**Theorem 3.4.** *Let  $(\alpha, \beta) \in [\alpha_{\min}, \alpha_{\max}] \times [\beta_{\min}, \beta_{\max}]$  be a fixed pair such that*

$$t_i + \beta \notin \{0\} \cup \{t_j, j = 0, \dots, p\}, \quad i = 0, \dots, p.$$

*Furthermore, consider a fixed time point  $t \in (t_{i-1} + \beta, t_i + \beta) \cap (0, T]$ , where  $i \in \{1, \dots, p\}$ .*

*Then*

$$\frac{\partial \mathbf{x}(t|\alpha, \beta)}{\partial \beta} = \Gamma(t|\alpha, \beta). \quad (3.61)$$

### 3.4 Computation algorithm

In this section, based on the results in Section 3.3, we develop a computational algorithm for solving Problem (P). Our approach is to view Problem (P) as a nonlinear programming problem in which  $\alpha$  and  $\beta$  are decision variables to be chosen optimally. On this basis, Problem (P) can, in principle, be solved using standard nonlinear programming algorithms such as the SQP method, which relies on the partial derivatives of the cost function to compute search directions leading to profitable areas of the search space. Thus, to solve Problem (P) as a nonlinear programming problem, we need to derive the partial derivatives of  $J$  with respect to both  $\alpha$  and  $\beta$ .



By using the state variation formulae in Theorems 3.1 and 3.4, we can differentiate  $J$  using the chain rule. However, the state variation with respect to  $\beta$  does not exist for all values of  $\beta$  (recall that Theorem 3.4 is only valid when  $t_i + \beta \notin \{0\} \cup \{t_j, j = 0, \dots, p\}$  for each  $i = 0, \dots, p$ ). Thus, at each stage of the optimization process, we need to check the condition  $t_i + \beta \notin \{0\} \cup \{t_j, j = 0, \dots, p\}$ , and if this condition is not satisfied, then we perturb  $\beta$  by a small amount  $\epsilon$ . More precisely, we first check to see whether  $(\alpha, \beta)$  is an optimal pair of delay estimates for Problem (P) (i.e. if the value of  $J$ , which measures the discrepancy between predicted output and observed system output, is below a desired tolerance). If it is, then we stop. Otherwise, we calculate the following modified value of  $\beta$ :

$$\bar{\beta} = \begin{cases} \beta, & \text{if } t_i + \beta \notin \{0\} \cup \{t_j, j = 0, \dots, p\} \text{ for each } i = 0, \dots, p, \\ \beta + \epsilon, & \text{otherwise,} \end{cases} \quad (3.62)$$

where  $\epsilon$  is a small number chosen to ensure that  $t_i + \beta + \epsilon \notin \{0\} \cup \{t_j, j = 0, \dots, p\}$  for each  $i = 0, \dots, p$ , and  $\beta + \epsilon \in [\beta_{\min}, \beta_{\max}]$ .

Note that the state variation formula in Theorem 3.4 is not applicable at the time points  $t = t_i + \bar{\beta}$ ,  $i = 0, \dots, p$ . Thus, if  $\tau_j \in \{t_i + \bar{\beta}\}_{i=0}^p$  for some  $j$ , where  $\tau_j$  is the  $j$ th sample time, then we will not be able to compute the state variation of  $\mathbf{x}(\tau_j | \alpha, \bar{\beta})$  with respect to the input-delay. In this case, we need to consider a modified cost function in which the experimental data are slightly perturbed. The perturbation procedure is designed to ensure that none of the new sample times coincide with points in  $\{t_i + \bar{\beta}\}_{i=0}^p$ . Details are given below.

After arriving at a new delay pair  $(\alpha, \bar{\beta})$  at some point during the optimization process, we define the  $j$ th perturbed sample time as follows:

$$\bar{\tau}_j = \begin{cases} \tau_j, & \text{if } \tau_j \notin \{t_i + \bar{\beta}\}_{i=0}^p, \\ \tau_j + \delta_j, & \text{if } \tau_j \in \{t_i + \bar{\beta}\}_{i=0}^p, \end{cases} \quad (3.63)$$

where  $\delta_j$  is a small number chosen such that  $\tau_j + \delta_j \in [0, T] \setminus \{t_i + \bar{\beta}\}_{i=0}^p$ . The corresponding output vector is defined as follows:

$$\bar{\mathbf{y}}^j = \begin{cases} \hat{\mathbf{y}}^j, & \text{if } \tau_j \notin \{t_i + \bar{\beta}\}_{i=0}^p, \\ \hat{\mathbf{y}}^j + \boldsymbol{\gamma}^j, & \text{if } \tau_j \in \{t_i + \bar{\beta}\}_{i=0}^p, \end{cases} \quad (3.64)$$

where  $\boldsymbol{\gamma}^j$  is computed using the original experimental data together with an appropriate interpolation technique. Our new objective function is

$$\bar{J}(\alpha, \bar{\beta}) = \sum_{j=1}^m |\mathbf{y}(\bar{\tau}_j | \alpha, \bar{\beta}) - \bar{\mathbf{y}}^j|^2 \approx \sum_{j=1}^m |\mathbf{y}(\tau_j | \alpha, \bar{\beta}) - \hat{\mathbf{y}}^j|^2 = J(\alpha, \bar{\beta}). \quad (3.65)$$

Using Theorem 3.1, the partial derivative of  $\bar{J}$  with respect to the state delay is given by

$$\begin{aligned}\frac{\partial \bar{J}(\alpha, \bar{\beta})}{\partial \alpha} &= 2 \sum_{j=1}^m (\mathbf{g}(\mathbf{x}(\bar{\tau}_j | \alpha, \bar{\beta})) - \bar{\mathbf{y}}^j)^\top \frac{\partial \mathbf{g}(\mathbf{x}(\bar{\tau}_j | \alpha, \bar{\beta}))}{\partial \mathbf{x}} \frac{\partial \mathbf{x}(\bar{\tau}_j | \alpha, \bar{\beta})}{\partial \alpha} \\ &= 2 \sum_{j=1}^m (\mathbf{g}(\mathbf{x}(\bar{\tau}_j | \alpha, \bar{\beta})) - \bar{\mathbf{y}}^j)^\top \frac{\partial \mathbf{g}(\mathbf{x}(\bar{\tau}_j | \alpha, \bar{\beta}))}{\partial \mathbf{x}} \mathbf{\Lambda}(\bar{\tau}_j | \alpha, \bar{\beta}).\end{aligned}\quad (3.66)$$

The partial derivative of  $\bar{J}$  with respect to the input delay can be determined in a similar manner to the derivation of  $\frac{\partial \bar{J}}{\partial \alpha}$  given above:

$$\frac{\partial \bar{J}(\alpha, \bar{\beta})}{\partial \beta} = 2 \sum_{j=1}^m (\mathbf{g}(\mathbf{x}(\bar{\tau}_j | \alpha, \bar{\beta})) - \bar{\mathbf{y}}^j)^\top \frac{\partial \mathbf{g}(\mathbf{x}(\bar{\tau}_j | \alpha, \bar{\beta}))}{\partial \mathbf{x}} \mathbf{\Gamma}(\bar{\tau}_j | \alpha, \bar{\beta}).\quad (3.67)$$

Since it is unlikely that many of the sample times will lie in the set  $\{t_i + \bar{\beta}\}_{i=0}^p$ , there should be no noticeable difference between minimizing  $\bar{J}$  and minimizing  $J$ . Indeed, our numerical results in the next section indicate that this is precisely the case. Also, the input function can be chosen judiciously during experimentation to minimize the chance of one of the sample times lying in the set  $\{t_i + \beta\}_{i=0}^p$ . Our heuristic optimization strategy for descending from a point  $(\alpha, \beta)$  is described below.

*Step 1.* Compute the modified input-delay  $\bar{\beta}$  according to (3.62).

*Step 2.* Compute the new experimental data  $\{(\bar{\tau}_j, \bar{\mathbf{y}}^j)\}_{j=1}^m$  using (3.63) and (3.64).

*Step 3.* Obtain  $\mathbf{x}(\cdot | \alpha, \bar{\beta})$ ,  $\mathbf{\Lambda}(\cdot | \alpha, \bar{\beta})$ , and  $\mathbf{\Gamma}(\cdot | \alpha, \bar{\beta})$  by solving the enlarged time-delay system consisting of the original system (3.1)-(3.2) and the auxiliary systems (3.16)-(3.17) and (3.34)-(3.36).

*Step 4.* Use  $\mathbf{x}(\bar{\tau}_j | \alpha, \bar{\beta})$ ,  $j = 1, \dots, m$  to compute  $\mathbf{y}(\bar{\tau}_j | \alpha, \bar{\beta})$  through equation (3.9).

*Step 5.* Use  $\mathbf{y}(\bar{\tau}_j | \alpha, \bar{\beta})$ ,  $j = 1, \dots, m$  to compute  $\bar{J}(\alpha, \bar{\beta})$  through equation (3.65).

*Step 6.* Use  $\mathbf{x}(\bar{\tau}_j | \alpha, \bar{\beta})$ ,  $\mathbf{y}(\bar{\tau}_j | \alpha, \bar{\beta})$ ,  $\mathbf{\Lambda}(\bar{\tau}_j | \alpha, \bar{\beta})$  and  $\mathbf{\Gamma}(\bar{\tau}_j | \alpha, \bar{\beta})$ ,  $j = 1, \dots, m$  to compute  $\frac{\partial \bar{J}(\alpha, \bar{\beta})}{\partial \alpha}$  and  $\frac{\partial \bar{J}(\alpha, \bar{\beta})}{\partial \beta}$  through equations (3.66) and (3.67).

This procedure can be combined with a standard nonlinear programming software to solve Problem (P) and determine optimal estimates for the time-delays.

## 3.5 Numerical examples

### 3.5.1 Example 1: Zinc sulphate purification

For our first example, we consider the industrial zinc sulphate purification process described in [94]. In this process, zinc powder is added to a zinc sulphate electrolyte to

induce deposition of harmful cobalt and cadmium ions. This is a key step in the production of zinc.

The rates of change of cobalt and cadmium ion concentrations in the electrolyte are described by the following differential equations:

$$V\dot{x}_1(t) = Qx_1^0 - Qx_1(t - \alpha) - c_1u(t - \beta)x_1(t - \alpha) + c_2x_2(t - \alpha), \quad (3.68)$$

$$V\dot{x}_2(t) = Qx_2^0 - Qx_2(t - \alpha) - c_3v(t)x_2(t - \alpha) + c_4x_1(t - \alpha), \quad (3.69)$$

and

$$x_1(t) = 3.3 \times 10^{-4}, \quad x_2(t) = 4.0 \times 10^{-3}, \quad t \leq 0, \quad (3.70)$$

where  $x_1$  is the concentration of cobalt ions;  $x_2$  is the concentration of cadmium ions; and  $u$  and  $v$  are control variables that correspond to the zinc powder reaction surface areas (proportional to the amount of zinc powder added to the reaction tank).

Furthermore,  $V$  is the volume of the reaction tank ( $V = 400 \text{ m}^3$ );  $Q$  is the flux of solution ( $Q = 200 \text{ m}^3/\text{h}$ );  $c_1, c_2, c_3, c_4$ , are model parameters; and  $x_1^0$  and  $x_2^0$  are the concentrations of cobalt and cadmium ions at the inlet of the reaction tank, respectively ( $x_1^0 = 6 \times 10^{-4} \text{ g/L}$ ,  $x_2^0 = 9 \times 10^{-3} \text{ g/L}$ ). Reference [94] considers the parameter identification problem for system (3.68)-(3.70) with a given state-delay of  $\alpha = 2$  and no input delay (i.e.  $\beta = 0$ ). Here, we assume that there is a non-negligible delay in the addition of zinc powder to the tank. We also assume that the model parameters are equal to the optimal values reported in [94]:

$$c_1 = 7.828 \times 10^{-4}, \quad c_2 = 16.67, \quad c_3 = 2.823 \times 10^{-4}, \quad c_4 = 7.107 \times 10^2. \quad (3.71)$$

These values were obtained using data from a real zinc production factory in China. We assume that the terminal time is  $T = 8$ , and we set the input variables  $u$  and  $v$  as equal to the optimal control functions obtained in [94]:

$$u(t) = \sigma^i, \quad t \in [t_{i-1}, t_i), \quad i = 1, \dots, 8, \quad (3.72)$$

$$v(t) = \bar{\sigma}^i, \quad t \in [t_{i-1}, t_i), \quad i = 1, \dots, 8, \quad (3.73)$$

where the values of  $t_i$ ,  $\sigma^i$ , and  $\bar{\sigma}^i$ ,  $i = 1, \dots, 8$ , are given in Table 2.1. The output of the system is the concentration of cadmium ions:

$$y(t) = x_2(t). \quad (3.74)$$

Given system (3.68)-(3.70) with data (3.71) and piecewise-constant inputs (3.72)-(3.73), our goal is to identify the delays  $\alpha$  and  $\beta$ . We simulate system (3.68)-(3.70) with  $[\hat{\alpha}, \hat{\beta}]^\top =$

$[2, 0.25]$  to generate the observed data in Problem (P). The observed data  $\hat{y}^j = x_2(\tau_j|\hat{\alpha}, \hat{\beta})$  is sampled at  $\tau_j = j/5$ ,  $j = 1, \dots, 40$ . Thus, our identification problem is: choose  $\alpha$  and  $\beta$  to minimize

$$J(\alpha, \beta) = \sum_{j=1}^{40} |y(\tau_j|\alpha, \beta) - \hat{y}^j|^2 = \sum_{j=1}^{40} |x_2(\tau_j|\alpha, \beta) - x_2(\tau_j|\hat{\alpha}, \hat{\beta})|^2 \quad (3.75)$$

subject to the dynamic system (3.68)-(3.70).

To solve this problem, we wrote a Matlab program that integrates the SQP optimization method with the gradient computation algorithm described in Section 3.4.

Computational results for different initial guesses are shown in Table 3.1. The output trajectory for the initial guess  $(\alpha, \beta) = (3, 3)$  is displayed in Figure 3.1. In Table 3.1 and

Table 3.1: Numerical convergence of the cost values in Example 3.1.

No.	Initial guess		Cost value at the $k$ th iteration				$N_1$	$N_2$
	$\alpha^0$	$\beta^0$	$k = 0$	$k = 5$	$k = 10$	$k = 20$		
1	0.5	0.5	$5.865 \times 10^{-5}$	$9.001 \times 10^{-8}$	$8.396 \times 10^{-25}$	$1.151 \times 10^{-27}$	0	0
2	1.0	1.0	$4.171 \times 10^{-5}$	$6.265 \times 10^{-8}$	$2.218 \times 10^{-21}$	$2.287 \times 10^{-34}$	3	0
3	3.0	2.0	$9.169 \times 10^{-5}$	$2.007 \times 10^{-5}$	$6.624 \times 10^{-7}$	$8.209 \times 10^{-27}$	3	0
4	3.0	3.0	$7.828 \times 10^{-5}$	$2.122 \times 10^{-6}$	$2.318 \times 10^{-8}$	$1.141 \times 10^{-26}$	6	0

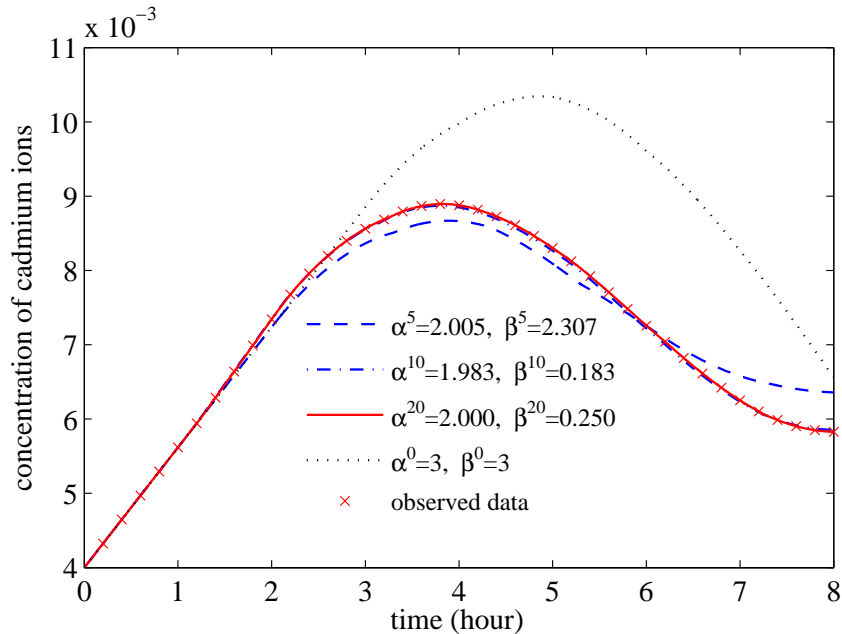


Figure 3.1: Numerical convergence of the output trajectory in Example 3.1 for initial guess No.4.

Figure 3.1,  $\alpha^k$  and  $\beta^k$  are the values of  $\alpha$  and  $\beta$  at the  $k$ th iteration during the optimization process, where  $\alpha^0$  and  $\beta^0$  (i.e.  $k = 0$ ) denote the initial guesses for the values of  $\alpha$  and  $\beta$ . Furthermore,  $N_1$  denotes the number of optimization iterations in which the condition  $t_i + \beta \in \{0\} \cup \{t_j, j = 0, \dots, 8\}$  occurs, and  $N_2$  denotes the number of optimization iterations in which one of the sample times lies in the set  $\mathcal{I} = \{t_j + \beta, j = 0, \dots, 8\}$ . We can see from Table 3.1 and Figure 3.1 that the optimal trajectory converges to the observed data for any initial guess. Note that, at each iteration,  $\tau_j \in \mathcal{I}$  occurs at most  $m$  times (the worst case scenario is when every sample time is in  $\mathcal{I}$ ). However, as expected,  $N_1$  and  $N_2$  are small, and thus the conditions for Theorem 3.4 are satisfied most of the time.

### 3.5.2 Example 2: Sodium aluminate evaporation

We now consider another industrial chemical process—specifically, the evaporation process described in [119]. This process, which takes place in a series of evaporators, is used to process a mother liquor consisting of sodium carbonate, sodium hydroxide, and alumina. The purpose of this process is to improve the concentration of the mother liquor to reach a specific concentration requirement, so that the sodium hydroxide and alumina components can be re-used. Since it takes time for the solution to flow from one reaction vessel to another, changes in the input variables do not cause changes in the evaporation vessel instantaneously—there are delays in the process. For simplicity, we just consider the case where there are two evaporators. The variables that are of interest are the sodium hydroxide concentration, temperature, and level of solution in each evaporation vessel. The dynamics in these two evaporation vessels can be described by the following differential equations:

$$\frac{dh_1(t)}{dt} = \frac{F_{01}(t - \beta)\rho_0(t - \alpha) + F_2(t - \beta)\rho_2(t - \alpha) - F_1(t)\rho_1(t) + V_0}{A_1\rho_1(t)}, \quad (3.76)$$

$$\frac{dh_2(t)}{dt} = \frac{F_0(t - \beta)\rho_0(t - \alpha) - F_2(t)\rho_2(t)}{A_1\rho_2(t)}, \quad (3.77)$$

$$\frac{dC_1(t)}{dt} = \frac{F_{01}(t - \beta)C_0(t - \alpha) + F_2(t - \beta)C_2(t - \alpha) - F_1(t)C_1(t)}{A_1h_1(t)} - \frac{dh_1(t)}{dt} \frac{C_1(t)}{h_1(t)}, \quad (3.78)$$

$$\frac{dC_2(t)}{dt} = \frac{F_0(t - \beta)C_0(t - \alpha) - F_2(t)C_2(t)}{A_2h_2(t)} - \frac{dh_2(t)}{dt} \frac{C_2(t)}{h_2(t)}, \quad (3.79)$$

$$\frac{dT_i(t)}{dt} = \frac{\Delta Q_i(t)}{A_i h_i(t) c_i^p(t)} - \frac{T_i(t) dc_i^p}{c_i^p(t) dt} - \frac{T_i(t) \rho_i(t) dh_i(t)}{h_i(t) dt}, \quad i = 1, 2, \quad (3.80)$$

where  $i$ ,  $i = 1, 2$ , refers to the  $i$ th evaporator;  $h_i$ ,  $T_i$ , and  $C_i$  are the state variables representing the level, temperature, and concentration of the solution in the  $i$ th evaporation vessel, respectively;  $A_i$  is the cross-sectional area of the  $i$ th evaporation vessel;  $F_i$  is the flow rate of the solution output from the  $i$ th evaporator;  $V_0$  is the amount of vapor from

other heat sources mixed with the solution;  $F_0$  is the flow rate of the feed;  $C_0$  is the concentration of the feed;  $\Delta Q_i$  is the heat change in the  $i$ th evaporation vessel (depends on the live steam flow rate);  $c_i^p$  and  $\rho_i$  are the specific heat capacity and the density of the solution in the  $i$ th evaporation vessel, respectively, and depend on the concentration and temperature. Note that  $c^p$ ,  $\rho$ , and  $\Delta Q$  are calculated by using the formulae given in [119].

The state vector for this system is  $\mathbf{x}(t) = [h_1(t), h_2(t), C_1(t), C_2(t), T_1(t), T_2(t)]^\top$ . The initial condition is

$$\mathbf{x}(t) = [1.91, 110.6, 73.0, 1.91, 97.2, 54.5]^\top, \quad t \leq 0. \quad (3.81)$$

Here, the inputs are  $\mathbf{u} = [u_1, u_2, u_3, u_4, u_5]^\top = [F_1, F_2, F_0, F_{01}, V]^\top$ , where  $V$  denotes the live steam flow rate. Also,  $\alpha$  is an unknown state-delay, and  $\beta$  is an unknown input-delay. Assume that the terminal time of this system is  $T = 240$  minutes. The input functions are

$$u_l(t) = \sigma_l^i, \quad t \in [t_{i-1}, t_i], \quad l = 1, 2, 3, 5, \quad i = 1, \dots, 24, \quad (3.82)$$

and

$$u_4(t) = 0.165, \quad t \in [0, 240], \quad (3.83)$$

where  $\sigma_l^i$ ,  $i = 1, \dots, 24$ , are given control heights shown in Figure 3.2. The output is  $\mathbf{y}(t) = [C_1(t), C_2(t)]^\top$ . We use the output trajectory of (3.76)-(3.81) with  $[\hat{\alpha}, \hat{\beta}] = [15, 6]^\top$  to generate the observed data for Problem (P). We set

$$\hat{\mathbf{y}}^j = [x_3(\tau_j | \hat{\alpha}, \hat{\beta}), x_4(\tau_j | \hat{\alpha}, \hat{\beta})]^\top, \quad j = 1, \dots, 24,$$

where  $\tau_{j+1} - \tau_j = 10$ . Thus, our identification problem is: choose  $\alpha$  and  $\beta$  to minimize

$$\begin{aligned} J(\alpha, \beta) &= \sum_{j=1}^{24} |\mathbf{y}(\tau_j | \alpha, \beta) - \hat{\mathbf{y}}^j|^2 \\ &= \sum_{j=1}^{24} |x_3(\tau_j | \alpha, \beta) - x_3(\tau_j | \hat{\alpha}, \hat{\beta})|^2 + \sum_{j=1}^{24} |x_4(\tau_j | \alpha, \beta) - x_4(\tau_j | \hat{\alpha}, \hat{\beta})|^2 \end{aligned} \quad (3.84)$$

subject to the dynamics (3.76)-(3.81).

As with Example 1, we solve this problem using a Matlab program that integrates the SQP optimization method with the gradient computation algorithm described in Section 3.4. The convergence process of the program is shown in Table 3.2 for four sets of initial guesses. The convergence process corresponding to the initial guess of

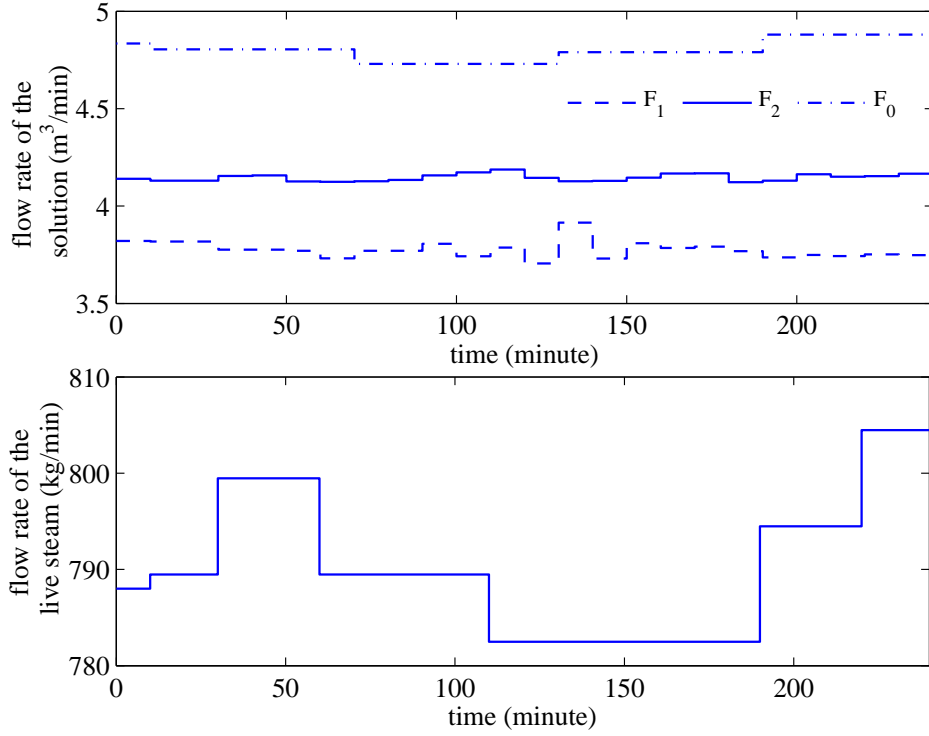


Figure 3.2: Control input variables for Example 3.2

$(\alpha, \beta) = (36, 36)$  is shown in Figure 3.3. In Table 3.2 and Figure 3.3,  $\alpha^k$  and  $\beta^k$  are the values of  $\alpha$  and  $\beta$  at the  $k$ th iteration, where  $k = 0$  indicates the initial guess of  $(\alpha, \beta)$ . Furthermore,  $N_1$  denotes the number of optimization iterations in which the condition  $t_i + \beta \in \{0\} \cup \{t_j, j = 0, \dots, 24\}$  occurs, and  $N_2$  denotes the number of optimization iterations in which one of the sample times lies in the set  $\{t_j + \beta, j = 0, \dots, 24\}$ . We observe that excellent convergence results are achieved for all the initial guesses.

Table 3.2: Numerical convergence of the cost values in Example 3.2.

No.	Initial guess		Cost value at the $k$ th iteration				$N_1$	$N_2$
	$\alpha^0$	$\beta^0$	$k = 0$	$k = 5$	$k = 10$	$k = 20$		
1	12	12	0.077	$4.462 \times 10^{-3}$	$2.782 \times 10^{-8}$	$2.446 \times 10^{-8}$	1	0
2	24	24	0.559	$4.496 \times 10^{-4}$	$3.259 \times 10^{-7}$	$4.684 \times 10^{-8}$	2	0
3	30	30	0.939	$5.300 \times 10^{-3}$	$3.259 \times 10^{-7}$	$4.524 \times 10^{-8}$	1	0
4	36	36	1.170	$8.751 \times 10^{-1}$	$1.728 \times 10^{-4}$	$8.455 \times 10^{-8}$	0	0
5	48	48	1.628	$8.862 \times 10^{-4}$	$6.432 \times 10^{-7}$	$4.802 \times 10^{-8}$	0	0

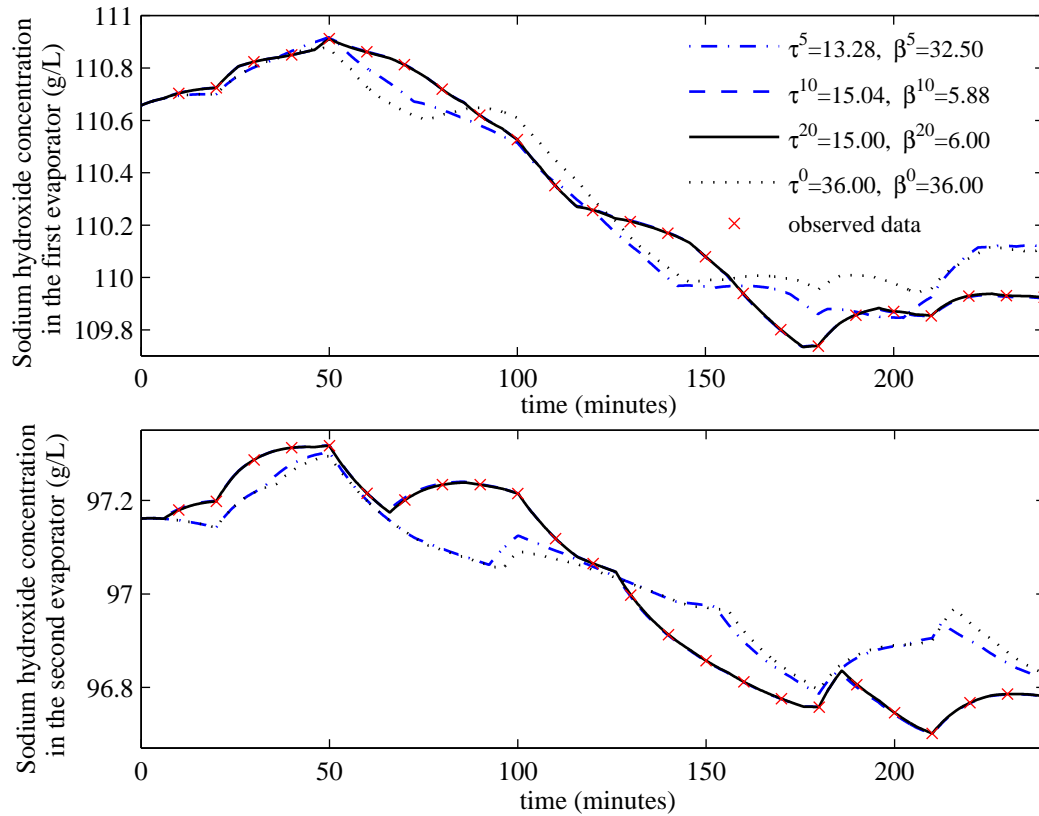


Figure 3.3: Numerical convergence of the output trajectory in Example 3.2 for initial guess No.4.

## 3.6 Conclusions

In this chapter, we have developed a gradient-based computational method for solving a time-delay identification problem for nonlinear systems in which the input function is piecewise-constant. We assume that there are two unknown time-delays in the system—a state-delay and an input-delay. The problem of determining optimal delay estimates is formulated as a dynamic optimization problem. The gradient of the cost function in this problem is obtained by solving two auxiliary delay-differential systems from  $t = 0$  to  $t = T$ . The auxiliary system corresponding to the input-delay is more complicated and involves jump conditions at the delayed control switching time points. The industrial examples demonstrate that our approach is highly effective. In particular, it converges to optimal delay estimates even when the initial estimates are far away from the optimal values.



---

---

# CHAPTER 4

---

## Time-delay optimal control application problem: an industrial evaporation process

### 4.1 Introduction

Time-delay dynamics are encountered in many real-world systems ranging from engineering to economics, such as those reported in [93, 143, 153]. As time-delays will influence the interaction between various system components, control theory and methods developed for dynamical systems without time-delays are not applicable. Thus, new theory and methods have been active research areas over the years and some fundamental and interesting results are now available in the literature (see, for example, [51, 74]). For optimal control problems involving time-delay systems, they have also been extensively studied in the literature such as in [70, 84, 86, 140]. However, most of the results obtained (see, for example, [47, 83, 101]) are for linear dynamical systems. On the computational issues, computational methods are proposed in [78] for optimal control problems with single time-delay. Note that there is no constraint on the state variables in the problem formulation considered in [78], while constraints on the state variables are included in the problem formulation in [119].

In this chapter, we consider a class of optimal control problems involving dynamical systems with multiple time-delays and subject to constraints on the state and/or control variables. Some of these constraints, which are expressed in the form of inequalities, are to be satisfied for all time point over the time planning horizon  $[0, t_f]$ . These constraints are called the continuous inequality constraints. The objective is to find a control such that a cost function is minimized subject to the given constraints. The focus of this chapter is to develop an effective computational method for solving this difficult constrained time-delay optimal control problem. To begin, the control parameterization technique is used to subdivide the time planning horizon  $[0, t_f]$  into  $N$  subintervals. Then, the control is approximated by a piecewise-constant function with possible discontinuities at these partition points. The time-delay optimal control problem is thus approximated by a sequence of time-delay optimal parameter selection problems subject to constraints on the

state and control variables. Amongst these various constraints, the continuous inequality constraints, often involving only state variables, are very difficult to handle directly. Thus, the constraint transcription technique introduced in [80] is used to convert each of them into an equivalent equality constraint in integral form. However, the integrand of each of these equality constraints is non-smooth. Thus, a local smoothing method is used to approximate the non-smooth functions by smooth functions. Consequently, each of these continuous inequality constraints is approximated by a sequence of inequality constraints in integral form, where the integrands are smooth approximating functions. These inequality constraints are known as the inequality constraints in canonical form, as they appear in the same form as the cost function. Then, by using the idea of the penalty function, these inequality constraints are appended to the cost function, forming an augmented cost function. Thus, the constrained time-delay optimal control problem is approximated by a sequence of time-delay optimal parameter selection problems subject to simple bounds on the control parameter vector. Each of them is to be solved as a nonlinear programming problem by using a gradient-based optimization technique, such as the sequential quadratic programming approximation scheme with active set strategy (see, for example, [80,147]). For this, the gradient formula of the augmented cost function with respect to the control parameter vector is derived. On this basis, an effective computational algorithm is developed for the time-delay constrained optimal control problem through solving a sequence of optimal parameter selection problems subject to simple bounds, each of which is regarded as a nonlinear programming problem.

The rest of the chapter is organized as follows. The problem formulation is given in Section 4.2. For the solution method, which consists of constraint transformation, problem approximation, convergence analysis, and computation method, is presented in Section 4.3. For real world application, we consider an optimal control problem of practical alumina evaporation process in Section 4.4, where the objective is to find a control such that the specific requirements of the industrial sodium aluminate solution are met and the solution level for each effect is maintained to within its operation limits with the least energy consumption. In Section 4.5, some concluding remarks are made.

## 4.2 Problem statement

Consider a process that evolves over the time horizon  $[0, t_f]$  as described below:

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \tilde{\mathbf{x}}(t), \mathbf{u}(t), \tilde{\mathbf{u}}(t)), \quad t \in [0, t_f], \quad (4.1)$$

where  $t_f > 0$  is the *terminal time*;  $\mathbf{x}(t) = [x_1(t), \dots, x_n(t)]^\top \in \mathbb{R}^n$  is called the *state*;  $\tilde{\mathbf{x}}(t) = [(\mathbf{x}(t - \alpha_1))^\top, \dots, (\mathbf{x}(t - \alpha_m))^\top]^\top \in \mathbb{R}^{nm}$  is called the *delayed state*.  $\mathbf{u}(t) = [u_1(t), \dots, u_r(t)]^\top \in \mathbb{R}^r$  is the *control*; and  $\tilde{\mathbf{u}}(t) = [(\mathbf{u}(t - \beta_1))^\top, \dots, (\mathbf{u}(t - \beta_p))^\top]^\top \in \mathbb{R}^{pr}$

is the *delayed control*. Furthermore,  $\mathbf{f} : \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^{nm} \times \mathbb{R}^r \times \mathbb{R}^{pr} \rightarrow \mathbb{R}^n$  is a given function.  $\alpha_i$ ,  $i = 1, \dots, m$ , are given state-delays satisfying  $0 < \alpha_1 < \dots < \alpha_m < t_f$ ;  $\beta_j$ ,  $j = 1, \dots, p$ , are given control-delays satisfying  $0 < \beta_1 < \dots < \beta_p < t_f$ . For brevity, let these time-delays, which are sorted in ascending order, be referred to as  $\tau_k$ ,  $k = 1, \dots, M$ . Note that  $\tau_k$ ,  $k = 1, \dots, M$ , are not necessarily equal.

The initial functions for the time-delayed differential equations (4.1) are:

$$\mathbf{x}(t) = \boldsymbol{\phi}(t), \quad t \in [-\tau_M, 0), \quad (4.2)$$

$$\mathbf{x}(0) = \mathbf{x}^0, \quad (4.3)$$

where  $\boldsymbol{\phi}(t) = [\phi_1(t), \dots, \phi_n(t)]^\top$  is a given continuously differentiable function from  $[-\tau_M, 0)$  into  $\mathbb{R}^n$ ; and  $\mathbf{x}^0 \in \mathbb{R}^n$  is a given vector. In addition, the initial condition for the control is:

$$\mathbf{u}(t) = \boldsymbol{\psi}(t), \quad t \in [-\tau_M, 0), \quad (4.4)$$

where  $\boldsymbol{\psi}(t) = [\psi_1(t), \dots, \psi_r(t)]^\top$  is a given piecewise continuous function from  $[-\tau_M, 0)$  into  $\mathbb{R}^r$ .

Define

$$\mathbf{U} = \{\mathbf{u} = [u_1, \dots, u_r]^\top \in \mathbb{R}^r : a_l \leq u_l \leq b_l, \quad l = 1, \dots, r\}, \quad (4.5)$$

where  $a_l$  and  $b_l$ ,  $l = 1, \dots, r$ , are given constants. Clearly,  $\mathbf{U}$  is a compact and convex subset of  $\mathbb{R}^r$ . Any measurable function  $\mathbf{u} = [u_1, \dots, u_r]^\top : [-\tau_M, t_f] \rightarrow \mathbb{R}^r$  such that  $\mathbf{u}(t) = \boldsymbol{\psi}(t)$ ,  $\forall t \in [-\tau_M, 0)$  and  $\mathbf{u}(t) \in \mathbf{U}$  for almost all  $t \in [0, t_f]$ , is called an admissible control. Let  $\mathcal{U}$  be the set which consists of all such admissible controls.

We assume that the following conditions are satisfied.

**(4.A.1).** The function  $\mathbf{f}$  is continuously differentiable with respect to  $\mathbf{x}$  and  $\mathbf{u}$  for each  $t \in [0, t_f]$ , while it is piecewise differentiable with respect to  $t$  for each  $(\mathbf{x}, \mathbf{u}) \in \mathbb{R}^n \times \mathbb{R}^r$ .

**(4.A.2).** The function  $\boldsymbol{\phi}$  is twice continuously differentiable.

**(4.A.3).** There exists a real number  $L_1 > 0$  such that

$$|\mathbf{f}(t, \mathbf{x}, \tilde{\mathbf{x}}, \mathbf{u}, \tilde{\mathbf{u}})| \leq L_1(1 + |\mathbf{x}| + |\tilde{\mathbf{x}}| + |\mathbf{u}| + |\tilde{\mathbf{u}}|),$$

$$(t, \mathbf{x}, \tilde{\mathbf{x}}, \mathbf{u}, \tilde{\mathbf{u}}) \in \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^{nm} \times \mathbb{R}^r \times \mathbb{R}^{pr},$$

where  $|\cdot|$  denotes the usual Euclidean norm.

On the basis of assumptions (4.A.1)-(4.A.3), the dynamic system (4.1)-(4.4) admits a unique solution corresponding to each control  $\mathbf{u} \in \mathcal{U}$  [111]. Let  $\mathbf{x}(\cdot|\mathbf{u})$  denote the corresponding solution of system (4.1) with initial conditions (4.2)-(4.4).

Assume that the following continuous state inequality constraints are to be satisfied:

$$h_i(t, \mathbf{x}(t|\mathbf{u})) \leq 0, \quad i = 1, \dots, N_c, \quad \forall t \in [0, t_f]. \quad (4.6)$$

Now, a general class of time-delay optimal control problems with continuous state inequality constraints may be described as follows:

**Problem (Q).** *Given system (4.1) with initial conditions (4.2)-(4.4), find a control  $\mathbf{u} \in \mathcal{U}$  such that the cost function*

$$J_0(\mathbf{u}) = \Phi_0(\mathbf{x}(t_f|\mathbf{u})) + \int_0^{t_f} \mathcal{L}_0(t, \mathbf{x}(t|\mathbf{u}), \tilde{\mathbf{x}}(t|\mathbf{u}), \mathbf{u}(t), \tilde{\mathbf{u}}(t)) dt, \quad (4.7)$$

*is minimized subject to the continuous state inequality constraints (4.6), where  $\Phi_0$  is the terminal cost.*

We assume that the following conditions are satisfied.

**(4.A.4).** For each  $i = 1, \dots, N_c$ , the function  $h_i$  is continuously differentiable with respect to  $(t, \mathbf{x}) \in [0, t_f] \times \mathbb{R}^n$ .

**(4.A.5).** The function  $\Phi_0$  satisfies the assumption (4.A.2), while the function  $\mathcal{L}_0$  satisfies the assumption (4.A.1).

## 4.3 Solution method

### 4.3.1 Control parameterization

To solve Problem (Q), we apply the control parametrization scheme to approximate the control  $\mathbf{u} \in \mathcal{U}$  by a piecewise-constant function with possible discontinuities at the partition points called the switching times. The heights of the piecewise-constant function are decision variables. More specifically, for each  $l = 1, \dots, r$ ,

$$u_i^N(t) = \sum_{q=1}^N \sigma_l^{N,q} \chi_{I_q}(t), \quad t \in [0, t_f], \quad (4.8)$$

where

$$\chi_{I_q}(t) = \begin{cases} 1, & \text{if } t \in I_q, \\ 0, & \text{otherwise,} \end{cases} \quad (4.9)$$

and

$$I_q = \begin{cases} [t_{q-1}, t_q), & \text{if } q \in \{1, 2, \dots, N-1\}, \\ [t_{q-1}, t_q], & \text{if } q = N. \end{cases} \quad (4.10)$$

Here,  $t_q$ ,  $q = 1, \dots, N$ , are pre-assigned such that  $t_{q-1} < t_q$  with  $t_0 = 0$  and  $t_N = t_f$ . For each  $l = 1, \dots, r$ ,  $\sigma_l^q$ ,  $q = 1, \dots, N$ , are the heights of the piecewise-constant control component  $u_l^N$ .

Let

$$\boldsymbol{\sigma}^{N,q} = [\sigma_1^{N,q}, \dots, \sigma_r^{N,q}]^\top, \quad q = 1, \dots, N,$$

and let

$$\boldsymbol{\sigma}^N = [(\boldsymbol{\sigma}^{N,1})^\top, \dots, (\boldsymbol{\sigma}^{N,N})^\top]^\top.$$

A function  $\mathbf{u}^N = [u_1^N, \dots, u_r^N]^\top$  with  $u_l^N$ ,  $l = 1, \dots, r$ , given by (4.8) and  $\boldsymbol{\sigma}^{N,q} \in \mathbf{U}$ , is called an admissible piecewise-constant control. Let  $\mathcal{U}_N$  be the set of all such admissible piecewise-constant controls, and let  $\Xi_N$  be the set containing all the corresponding  $\boldsymbol{\sigma}^N$ , i.e.,

$$\Xi_N = \{\boldsymbol{\sigma}^N = [(\boldsymbol{\sigma}^{N,1})^\top, \dots, (\boldsymbol{\sigma}^{N,N})^\top]^\top \in \mathbb{R}^{N \times r} : \boldsymbol{\sigma}^{N,q} \in \mathbf{U}, \quad q = 1, \dots, N\}. \quad (4.11)$$

Clearly, each  $\mathbf{u} \in \mathcal{U}_N$  corresponds to a unique  $\boldsymbol{\sigma}^N \in \Xi_N$  and vice versa. For a  $\boldsymbol{\sigma}^N \in \Xi_N$ , let  $\mathbf{u}^N(\cdot|\boldsymbol{\sigma}^N)$  denote the corresponding piecewise-constant control in  $\mathcal{U}_N$ . For each  $\mathbf{u} = \mathbf{u}^N \in \mathcal{U}_N$ , let  $\mathbf{x}(\cdot|\mathbf{u}^N)$  be the corresponding solution of system (4.1) with initial conditions (4.3) and (4.4). Furthermore, let  $\tilde{\mathbf{x}}(\cdot|\mathbf{u}^N)$  be the corresponding delayed state. For convenience,  $\mathbf{x}(\cdot|\mathbf{u}^N)$  and  $\tilde{\mathbf{x}}(\cdot|\mathbf{u}^N)$  are written as  $\mathbf{x}(\cdot|\boldsymbol{\sigma}^N)$  and  $\tilde{\mathbf{x}}(\cdot|\boldsymbol{\sigma}^N)$ , respectively.

From (4.8), we have

$$u_l^N(t - \beta_j) = \sum_{q=1}^N \sigma_l^{N,q} \chi_{I_q}(t - \beta_j), \quad t \in I_k, \quad j = 1, \dots, p, \quad l = 1, \dots, r. \quad (4.12)$$

For brevity, define

$$\tilde{\mathbf{u}}^N(t|\boldsymbol{\sigma}^N) = [\mathbf{u}^N(t - \beta_1|\boldsymbol{\sigma}^N)^\top, \dots, \mathbf{u}^N(t - \beta_j|\boldsymbol{\sigma}^N)^\top]^\top,$$

where  $\mathbf{u}^N(t - \beta_j|\boldsymbol{\sigma}^N) = [u_1^N(t - \beta_j|\boldsymbol{\sigma}^N), \dots, u_r^N(t - \beta_j|\boldsymbol{\sigma}^N)]^\top$ ,  $j = 1, \dots, p$ . Thus, by applying the control parametrization technique, system (4.1)-(4.3) becomes

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N), \tilde{\mathbf{x}}(t|\boldsymbol{\sigma}^N), \mathbf{u}^N(t|\boldsymbol{\sigma}^N), \tilde{\mathbf{u}}^N(t|\boldsymbol{\sigma}^N)), \quad t \in [0, t_f], \quad (4.13)$$

with initial conditions

$$\mathbf{x}(t) = \boldsymbol{\phi}(t), \quad t \in [-\tau_M, 0), \quad (4.14)$$

$$\mathbf{x}(0) = \mathbf{x}^0. \quad (4.15)$$

With  $\mathbf{u} \in \mathcal{U}_N$ , the continuous state inequality constraints (4.6) become:

$$h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) \leq 0, \quad i = 1, \dots, N_c. \quad (4.16)$$

We may now state an approximate problem of Problem (Q) as follows:

**Problem ( $Q_N$ ).** *Given system (4.13) with initial conditions (4.14) and (4.15), find a control vector  $\boldsymbol{\sigma}^N \in \Xi_N$  such that the following cost function*

$$\tilde{J}_0(\boldsymbol{\sigma}^N) = \Phi_0(\mathbf{x}(t_f|\boldsymbol{\sigma}^N)) + \int_0^{t_f} \mathcal{L}_0(t, \mathbf{x}(t|\boldsymbol{\sigma}^N), \tilde{\mathbf{x}}(t|\boldsymbol{\sigma}^N), \mathbf{u}^N(t|\boldsymbol{\sigma}^N), \tilde{\mathbf{u}}^N(t|\boldsymbol{\sigma}^N)) dt \quad (4.17)$$

is minimized over  $\Xi_N$  subject to constraints (4.16).

### 4.3.2 Constraints transformation

For each of the continuous inequality constraints, it contains infinite number of constraints and hence these continuous inequality constraints are difficult to handle directly. Note that each of these continuous inequality constraints (4.16) is equivalent to the following equality constraint in integral form:

$$g_i(\boldsymbol{\sigma}^N) = \int_0^{t_f} \max \{h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)), 0\} = 0, \quad i = 1, \dots, N_c. \quad (4.18)$$

Let  $\mathcal{B}_N$  be the feasible region defined by

$$\begin{aligned} \mathcal{B}_N &= \{\boldsymbol{\sigma}^N \in \Xi_N : h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) \leq 0, \quad \forall t \in [0, t_f], \quad i = 1, \dots, N_c\} \\ &= \{\boldsymbol{\sigma}^N \in \Xi_N : g_i(\boldsymbol{\sigma}^N) = 0, \quad i = 1, \dots, N_c\}. \end{aligned}$$

Furthermore, let  $\mathcal{B}_N^0$  denote the interior of  $\mathcal{B}_N$ , i.e.,

$$\mathcal{B}_N^0 = \{\boldsymbol{\sigma}^N \in \Xi_N : h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) < 0, \quad \forall t \in [0, t_f], \quad i = 1, \dots, N_c\}.$$

Note that the sets  $\mathcal{B}_N$  and  $\mathcal{B}_N^0$  are defined for the control parameter vectors in  $\Xi_N$ . The corresponding subsets in  $\mathcal{U}_N$  are denoted by  $\mathcal{F}_N$  and  $\mathcal{F}_N^0$ , respectively.

We assume that the following assumptions are satisfied.

**(4.A.6).**  $\mathcal{B}_N^0 \neq \emptyset$ .

**(4.A.7).** Suppose that  $\boldsymbol{\sigma}^N \in \mathcal{B}_N$ . Then, there exists a control vector  $\bar{\boldsymbol{\sigma}} \in \mathcal{B}_N^0$  such that  $\alpha \bar{\boldsymbol{\sigma}} + (1 - \alpha) \boldsymbol{\sigma}^N \in \mathcal{B}_N^0$  for all  $\alpha \in (0, 1]$ .

Since  $\max \{h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)), 0\}$ ,  $i = 1, \dots, N_c$ , are non-smooth, they are approximated

by the following smooth functions (see Chapter 8 of [80])

$$\mathcal{L}_{i,\varepsilon}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) = \begin{cases} h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)), & \text{if } h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) > \varepsilon, \\ \frac{\{h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) - \varepsilon\}^2}{4\varepsilon}, & \text{if } -\varepsilon \leq h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) \leq \varepsilon, \\ 0, & \text{if } h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) < -\varepsilon, \end{cases} \quad (4.19)$$

where  $\varepsilon > 0$  is a smoothing parameter controlling the accuracy of the approximation. Thus, constraints (4.18) are approximated by

$$g_{i,\varepsilon}(\boldsymbol{\sigma}^N) = \int_0^{t_f} \mathcal{L}_{i,\varepsilon}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) dt = 0, \quad i = 1, \dots, N_c, \quad (4.20)$$

Let  $\mathcal{B}_{N,\varepsilon}$  be the feasible region defined by

$$\begin{aligned} \mathcal{B}_{N,\varepsilon} &= \{\boldsymbol{\sigma}^N \in \Xi_N : g_{i,\varepsilon}(\boldsymbol{\sigma}^N) = 0, \quad i = 1, \dots, N_c\} \\ &= \{\boldsymbol{\sigma}^N \in \Xi_N : h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) \leq -\varepsilon, \quad t \in [0, t_f], \quad i = 1, \dots, N_c\}. \end{aligned}$$

Clearly,  $\mathcal{B}_{N,\varepsilon} \subset \mathcal{B}_N$  for each  $\varepsilon > 0$ .

However, it can be shown that constraints (4.20) fail to satisfy the usual constraint qualification (see [82, 89]) because when any of the constraints is satisfied as an equality, its gradient is a zero vector (see Chapter 6 of [80]). It is well-known in the optimization literature that any gradient-based optimization technique will not work well for optimization problems with constraints which fail to satisfy the usual constraint qualification. Here, the concept of the penalty function approach is used to append the approximate constraints to the cost function, forming an augmented cost function given below:

$$\begin{aligned} \tilde{J}_{\varepsilon,\gamma}(\boldsymbol{\sigma}^N) &= \Phi_0(\mathbf{x}(t_f|\boldsymbol{\sigma}^N)) + \int_0^{t_f} \mathcal{L}_0(t, \mathbf{x}(t|\boldsymbol{\sigma}^N), \tilde{\mathbf{x}}(t|\boldsymbol{\sigma}^N), \mathbf{u}^N(t|\boldsymbol{\sigma}^N), \tilde{\mathbf{u}}^N(t|\boldsymbol{\sigma}^N)) dt \\ &\quad + \gamma \sum_{i=1}^{N_c} \int_0^{t_f} \mathcal{L}_{i,\varepsilon}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) dt, \end{aligned} \quad (4.21)$$

where  $\gamma > 0$  is the penalty factor of the continuous state inequality constraints.

In this way, Problem  $(Q_N)$  is approximated by a sequence of approximate problems given below:

**Problem  $(Q_{N,\varepsilon,\gamma})$ .** *Given system (4.13) with initial conditions (4.14) and (4.15), find a control vector  $\boldsymbol{\sigma}^N \in \Xi_N$  such that the augmented cost function (4.21) is minimized.*

Problem  $(Q_{N,\varepsilon,\gamma})$  can be regarded as a nonlinear optimization problem subject to simple bounds on the decision variables specified by (4.5). It can be solved by a gradient-based optimization technique, such as the sequential quadratic approximate scheme with an active set strategy.

### 4.3.3 Convergence analysis

In this subsection, we shall show that the optimal cost of Problem  $(Q_{N,\varepsilon,\gamma})$  converges to the optimal cost of the original problem  $(Q)$ .

For each  $N \in \{2, 3, \dots\}$ , let  $S^N$  denote the set containing the partition points  $t_q$ ,  $q = 1, \dots, N$ . It is chosen such that  $S^{N+1} \supset S^N$  and  $\lim_{N \rightarrow \infty} S^N$  is dense in  $[0, t_f]$ . Let  $\boldsymbol{\sigma}^{N,*}$  be an optimal control vector to Problem  $(Q_N)$  and let  $\mathbf{u}^{N,*}$  be the corresponding piecewise-constant control in  $\mathcal{U}$ , (in fact, it is in  $\mathcal{U}_N$ ). Thus, for each integer  $N \geq 2$ ,  $\mathbf{u}^{N,*}$  is a suboptimal control to Problem  $(Q)$  such that  $J_0(\mathbf{u}^{N+1,*}) \leq J_0(\mathbf{u}^{N,*})$  for all  $N \geq 2$ .

We have the following result. Its proof is similar to that given for Theorem 6.4.2 in Chapter 6 of [80].

**Lemma 4.1.** *Let  $\{\mathbf{u}^N\}_{N=1}^\infty$  be a bounded sequence of functions in  $L_\infty^r$ . Then, the sequence  $\{\mathbf{x}(\cdot|\mathbf{u}^N)\}_{N=1}^\infty$  of the corresponding solutions of system (4.13) with initial conditions (4.14)-(4.15) is also bounded in  $L_\infty^r$ .*

We now relate the solutions of Problem  $(Q)$  and Problem  $(Q_{N,\varepsilon,\gamma})$  in the following theorems.

**Theorem 4.1.** *For any  $\varepsilon > 0$ , there exists a  $\rho(\varepsilon) > 0$  such that for any  $\rho$ ,  $0 < \rho < \rho(\varepsilon)$ , if  $g_{i,\varepsilon}(\boldsymbol{\sigma}^N) \leq \rho$ ,  $i = 1, \dots, N_c$ , then  $\boldsymbol{\sigma}^N \in \mathcal{B}_N$ .*

*Proof.* Since  $h_i$ ,  $i = 1, \dots, N_c$ , are continuously differentiable in  $[0, t_f] \times \mathbb{R}^n$  (recall assumption (4.A.4)), it follows that for each  $i$ ,  $i = 1, \dots, N_c$ , and any  $\boldsymbol{\sigma}^N \in \Xi_N$ , we have

$$\frac{dh_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N))}{dt} = \frac{\partial h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N))}{\partial \mathbf{x}} \dot{\mathbf{x}}(t) + \frac{\partial h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N))}{\partial t}.$$

By (4.A.3) and the fact that  $\mathbf{x}(t|\boldsymbol{\sigma}^N) \in \mathbf{X}$  all  $t \in [0, t_f]$ , where  $\mathbf{X}$  is a bounded set (see Lemma 4.1), there exists a positive constant  $m_i$  such that

$$\left| \frac{dh_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N))}{dt} \right| \leq m_i, \quad \forall t \in [0, t_f]. \quad (4.22)$$

For any  $\varepsilon > 0$ , define

$$\kappa_{i,\varepsilon} = \frac{\varepsilon}{16} \min \left\{ t_f, \frac{\varepsilon}{2m_i} \right\}.$$

Let

$$\mathcal{B}_{N,i} = \{\boldsymbol{\sigma}^N \in \Xi_N : g_i(\boldsymbol{\sigma}^N) = 0\},$$

and

$$\mathcal{B}_{N,i,\varepsilon,\rho} = \{\boldsymbol{\sigma}^N \in \Xi_N : g_{i,\varepsilon}(\boldsymbol{\sigma}^N) \leq \rho\},$$



where  $\rho$  is a positive real number. It suffices to show that  $\mathcal{B}_{N,i,\varepsilon,\rho} \subset \mathcal{B}_{N,i}$  for any  $\rho$  such that  $0 < \rho < \kappa_{i,\varepsilon}$ . We assume the contrary. Then there exists a  $\boldsymbol{\sigma}^N \in \Xi_N$  such that

$$g_{i,\varepsilon}(\boldsymbol{\sigma}^N) \leq \rho < \kappa_{i,\varepsilon}, \quad (4.23)$$

but

$$g_i(\boldsymbol{\sigma}^N) > 0. \quad (4.24)$$

Again by the continuity of  $h_i$ , it follows from (4.24) that there exists a  $\iota \in [0, t_f]$  such that

$$h_i(\iota, \mathbf{x}(\iota|\boldsymbol{\sigma}^N)) > 0,$$

and there exists an interval  $\mathcal{I}_i \subset [0, t_f]$  such that

$$h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) > -\frac{\varepsilon}{2}, \quad \forall t \in \mathcal{I}_i. \quad (4.25)$$

Using (4.22), it is clear from (4.25) that

$$|\mathcal{I}_i| \geq \min \left\{ t_f, \frac{\varepsilon}{2m_i} \right\},$$

where  $|\mathcal{I}_i|$  denotes the length of the interval  $\mathcal{I}_i$ . From the definition of  $g_{i,\varepsilon}$ , we have

$$\begin{aligned} g_{i,\varepsilon}(\boldsymbol{\sigma}^N) &= \int_0^{t_f} h_{i,\varepsilon}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) dt \geq \int_{\mathcal{I}_i} h_{i,\varepsilon}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) dt \geq \int_{\mathcal{I}_i} \min_{t \in \mathcal{I}_i} h_{i,\varepsilon}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) dt \\ &\geq \min_{t \in \mathcal{I}_i} \{ (h_{i,\varepsilon}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) + \varepsilon)^2 / 4\varepsilon \} |\mathcal{I}_i| \geq \frac{\varepsilon}{16} \min \left\{ t_f, \frac{\varepsilon}{2m_i} \right\} = \kappa_{i,\varepsilon}. \end{aligned}$$

This is a contradiction to (4.23). Thus, the proof is completed.  $\square$

**Theorem 4.2.** *For any  $\varepsilon > 0$ , there exists a  $\gamma(\varepsilon) > 0$  such that for any  $\gamma > \gamma(\varepsilon)$ , if  $\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}$  is an optimal control vector of Problem  $(Q_{N,\varepsilon,\gamma})$ , then it satisfies the continuous inequalities constraints (4.16) of Problem  $(Q_N)$ .*

*Proof.* Let  $\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}$  be an optimal control vector of Problem  $(Q_{N,\varepsilon,\gamma})$ . Then, for any  $\boldsymbol{\sigma}^N \in \Xi_N$ ,

$$\tilde{J}_{\varepsilon,\eta}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) = \tilde{J}_0(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) + \gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) \leq \tilde{J}_0(\boldsymbol{\sigma}^N) + \gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}^N).$$

Let  $\boldsymbol{\sigma}_\varepsilon^N \in \mathcal{B}_{N,\varepsilon}$  be fixed. Then, by the definition of  $g_{i,\varepsilon}$  given in (4.20), we have

$$\gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_\varepsilon^N) = 0.$$

Thus,

$$\tilde{J}_0(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) + \gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) \leq \tilde{J}_0(\boldsymbol{\sigma}_\varepsilon^N) + \gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_\varepsilon^N) = \tilde{J}_0(\boldsymbol{\sigma}_\varepsilon^N). \quad (4.26)$$

Since  $\Xi_N$  is a compact set, there exists, by (4.A.1)-(4.A.3) and (4.A.4)-(4.A.5), a  $\bar{\boldsymbol{\sigma}}^N \in \Xi_N$  such that  $\tilde{J}_0(\bar{\boldsymbol{\sigma}}^N) \leq \tilde{J}_0(\boldsymbol{\sigma}^N)$  for all  $\boldsymbol{\sigma}^N \in \Xi_N$ . Clearly,

$$\tilde{J}_0(\bar{\boldsymbol{\sigma}}^N) \leq \tilde{J}_0(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}). \quad (4.27)$$

Then, adding the penalty term  $\gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*})$  to each side of (4.27), we have

$$\tilde{J}_0(\bar{\boldsymbol{\sigma}}^N) + \gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) \leq \tilde{J}_0(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) + \gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}).$$

Using (4.26) gives

$$\tilde{J}_0(\bar{\boldsymbol{\sigma}}^N) + \gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) \leq \tilde{J}_0(\boldsymbol{\sigma}_\varepsilon^N). \quad (4.28)$$

Rearranging (4.28),

$$\gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) \leq \tilde{J}_0(\boldsymbol{\sigma}_\varepsilon^N) - \tilde{J}_0(\bar{\boldsymbol{\sigma}}^N). \quad (4.29)$$

Letting  $z = \tilde{J}_0(\boldsymbol{\sigma}_\varepsilon^N) - \tilde{J}_0(\bar{\boldsymbol{\sigma}}^N)$ , we obtain

$$\sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) \leq \frac{z}{\gamma}.$$

By choosing  $\rho(\varepsilon) = \frac{z}{\gamma(\varepsilon)}$ , it follows that for any  $\gamma > \gamma(\varepsilon)$ ,

$$0 < \rho < \rho(\varepsilon)$$

and

$$\sum_{i=1}^{N_c} g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) \leq \rho.$$

Consequently,  $g_{i,\varepsilon}(\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}) \leq \rho$ ,  $i = 1, \dots, N_c$ . Hence, it follows from Theorem 4.1 that  $\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*} \in \mathcal{B}_N$ . This completes the proof.  $\square$

**Theorem 4.3.** *Let  $\boldsymbol{\sigma}^{N,*}$  be an optimal control vector of Problem  $(Q_N)$  and let  $\boldsymbol{\sigma}_{\varepsilon,\gamma}^{N,*}$  be an optimal control vector of Problem  $(Q_{N,\varepsilon,\gamma})$ , where  $\gamma(\varepsilon)$  is chosen such that  $\boldsymbol{\sigma}_{\varepsilon,\gamma(\varepsilon)}^{N,*} \in \mathcal{B}_N$ . Then,*

$$\lim_{\varepsilon \rightarrow 0} \tilde{J}_{\varepsilon,\gamma(\varepsilon)}(\boldsymbol{\sigma}_{\varepsilon,\gamma(\varepsilon)}^{N,*}) = \tilde{J}_0(\boldsymbol{\sigma}^{N,*}).$$

*Proof.* By (4.A.7), there exists a  $\bar{\sigma}^N \in \mathcal{B}_N^0$  such that

$$\sigma_\alpha^N = (1 - \alpha)\sigma^{N,*} + \alpha\bar{\sigma}^N \in \mathcal{B}_N^0, \quad \alpha \in (0, 1].$$

Now, for any  $\delta_1 > 0$ , there exists an  $\alpha_1 \in (0, 1]$  such that

$$\tilde{J}_0(\sigma_\alpha^N) \leq \tilde{J}_0(\sigma^{N,*}) + \delta_1, \quad \forall \alpha \in (0, \alpha_1). \quad (4.30)$$

Choose  $\alpha_2 = \alpha_1/2$ . Then it is clear that  $\sigma_{\alpha_2}^N \in \mathcal{B}_N^0$ . Thus, there exists a  $\delta_2 > 0$  such that

$$\max (h_i(t, \mathbf{x}(t|\sigma_{\alpha_2}^N))) < -\delta_2, \quad i = 1, \dots, N_c.$$

Choosing  $\varepsilon = \delta_2$ , thus  $\sigma_{\alpha_2}^N$  satisfies (4.20). Since,  $\sigma_{\varepsilon,\gamma}^{N,*}$  is an optimal control vector of Problem  $(Q_{N,\varepsilon,\gamma})$ , it follows that

$$\tilde{J}_{\varepsilon,\gamma}(\sigma_{\varepsilon,\gamma}^{N,*}) = \tilde{J}_0(\sigma_{\varepsilon,\gamma}^{N,*}) + \gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\sigma_{\varepsilon,\gamma}^{N,*}) \leq \tilde{J}_0(\sigma_{\alpha_2}^N) + \gamma \sum_{i=1}^{N_c} g_{i,\varepsilon}(\sigma_{\alpha_2}^N) = \tilde{J}_0(\sigma_{\alpha_2}^N).$$

Noting that the penalty term is non-negative, we have

$$\tilde{J}_{\varepsilon,\gamma}(\sigma_{\varepsilon,\gamma}^{N,*}) \leq \tilde{J}_0(\sigma_{\alpha_2}^N). \quad (4.31)$$

Since, by Theorem 4.2,  $\sigma_{\varepsilon,\gamma}^{N,*}$  is a feasible point of Problem  $(Q_N)$ , it follows from (4.30) and (4.31) that

$$\tilde{J}_0(\sigma^{N,*}) \leq \tilde{J}_{\varepsilon,\gamma}(\sigma_{\varepsilon,\gamma}^{N,*}) \leq \tilde{J}_0(\sigma_{\alpha_2}^N) + \delta_1.$$

Since  $\delta_1 > 0$  is arbitrary, letting  $\varepsilon \rightarrow 0$ , the results follows.  $\square$

To proceed, we need the following lemmas. The first is quoted from Lemma 6.4.1 in [80].

**Lemma 4.2.** *For each  $\mathbf{u} \in \mathcal{U}$ , let*

$$\mathbf{u}^N(t) = \sum_{q=1}^N \sigma^{N,q} \chi_{I_q}(t), \quad t \in [0, t_f], \quad (4.32)$$

where  $I_q = [t_{q-1}, t_q)$ ,  $q = 1, \dots, N$ ,

$$\sigma^q = \frac{1}{|I_q|} \int_{I_q} \mathbf{u}(s) ds, \quad (4.33)$$

and  $|I_q| = t_q - t_{q-1}$ . Then,  $\mathbf{u}^N \rightarrow \mathbf{u}$  almost everywhere in  $[0, t_f]$  as  $N \rightarrow \infty$ ; and furthermore,

$$\lim_{N \rightarrow \infty} \int_0^{t_f} |\mathbf{u}^N(t) - \mathbf{u}(t)| dt = 0,$$

where,  $|\cdot|$  denotes the usual Euclidean norm.

**Lemma 4.3.** Let  $\{\mathbf{u}^N\}_{N=2}^\infty$  be a bounded sequence of functions in  $\mathcal{L}_\infty^r$ . Suppose that  $\mathbf{u}^N \rightarrow \mathbf{u}$  almost everywhere in  $[0, t_f]$  as  $N \rightarrow \infty$ . Then,

$$\lim_{N \rightarrow \infty} |\mathbf{x}(t|\mathbf{u}^N) - \mathbf{x}(t|\mathbf{u})| = 0,$$

$$\lim_{N \rightarrow \infty} J_0(\mathbf{u}^N) = J_0(\mathbf{u}).$$

*Proof.* The proof is similar to that given for Lemma 6.4.3 and Lemma 6.4.4 in [80].  $\square$

Define

$$\Omega = \{\mathbf{u} \in \mathcal{U} : h_i(t, \mathbf{x}(t|\mathbf{u})) \leq 0, \quad t \in [0, t_f], \quad i = 1, \dots, N_c\}.$$

Furthermore, let  $\Omega^0$  denote the interior of  $\Omega$ , i.e.,

$$\Omega^0 = \{\mathbf{u} \in \mathcal{U} : h_i(t, \mathbf{x}(t|\mathbf{u})) < 0, \quad t \in [0, t_f], \quad i = 1, \dots, N_c\}.$$

We assume that the following assumptions are satisfied.

(4.A.8).  $\Omega^0 \neq \emptyset$ .

(4.A.9). Suppose that  $\mathbf{u} \in \Omega$ . Then, there exists a control  $\bar{\mathbf{u}} \in \Omega^0$  such that

$$\alpha \bar{\mathbf{u}} + (1 - \alpha)\mathbf{u} \in \Omega^0, \quad \forall \alpha \in (0, 1].$$

**Theorem 4.4.** Suppose  $\mathbf{u}^*$  be an optimal control of the original problem (Q). Furthermore, for each  $N \geq 2$ , let  $\mathbf{u}^{N,*}$  be an optimal piecewise-constant control of Problem (Q<sub>N</sub>). Then,

$$\lim_{N \rightarrow \infty} J_0(\mathbf{u}^{N,*}) = J_0(\mathbf{u}^*).$$

*Proof.* By (4.A.9), there exists a  $\bar{\mathbf{u}} \in \Omega^0$  such that

$$\bar{\mathbf{u}}_\alpha = \alpha \bar{\mathbf{u}} + (1 - \alpha)\mathbf{u}^* \in \Omega^0, \quad \alpha \in (0, 1]. \quad (4.34)$$

Equation (4.34) implies that  $\bar{\mathbf{u}}_\alpha \rightarrow \mathbf{u}^*$  as  $\alpha \rightarrow 0$  for almost all  $t \in [0, t_f]$ . Hence, for any real number  $\delta > 0$ , there exists, by Lemma 4.3, an  $\alpha_1 \in (0, 1)$  such that

$$|J_0(\bar{\mathbf{u}}_\alpha) - J_0(\mathbf{u}^*)| < \frac{\delta}{2}, \quad \forall \quad 0 < \alpha \leq \alpha_1. \quad (4.35)$$

Since  $\bar{\mathbf{u}}_\alpha \in \Omega^0$ , it is clear that for any  $\alpha$ , there is a corresponding real number  $\rho > 0$  such that

$$h_i(t, \mathbf{x}(t|\bar{\mathbf{u}}_\alpha)) < -\rho, \quad i = 1, \dots, N_c, \quad t \in [0, t_f]. \quad (4.36)$$

We now fix  $\alpha$ . Let  $\{\bar{\mathbf{u}}_\alpha^N\}_{N=2}^\infty$  denote the sequence of piecewise-constant controls constructed from  $\bar{\mathbf{u}}_\alpha$  according to (4.32)-(4.33). Thus, by Lemma 4.2,  $\{\bar{\mathbf{u}}_\alpha^N\}_{N=2}^\infty \rightarrow \bar{\mathbf{u}}_\alpha$  almost everywhere in  $[0, t_f]$  as  $N \rightarrow \infty$ . By using (4.A.1) and (4.A.2) and Lemma 4.3, there exists an integer  $N_0 \geq 2$ , such that for all  $N > N_0$

$$|h_i(t, \mathbf{x}(t|\bar{\mathbf{u}}_\alpha^N)) - h_i(t, \mathbf{x}(t|\bar{\mathbf{u}}_\alpha))| < \frac{\rho}{2}, \quad i = 1, \dots, N_c, \quad t \in [0, t_f]. \quad (4.37)$$

From (4.36)-(4.37), it follows that

$$h_i(t, \mathbf{x}(t|\bar{\mathbf{u}}_\alpha^N)) < -\frac{\rho}{2}, \quad i = 1, \dots, N_c, \quad t \in [0, t_f],$$

for all  $N \geq N_0$ . This implies that  $\bar{\mathbf{u}}_\alpha^N \in \mathcal{B}_N^0$ . Furthermore, by Lemma 4.3, there exists an  $N_1 \geq 2$  such that for all  $N > N_1$ ,

$$|J_0(\bar{\mathbf{u}}_\alpha^N) - J_0(\bar{\mathbf{u}}_\alpha)| < \frac{\delta}{2}. \quad (4.38)$$

Set  $N_2 = \max\{N_0, N_1\}$ . Then, it follows from (4.35) and (4.38) that

$$|J_0(\bar{\mathbf{u}}_\alpha^N) - J_0(\mathbf{u}^*)| \leq |J_0(\bar{\mathbf{u}}_\alpha^N) - J_0(\bar{\mathbf{u}}_\alpha)| + |J_0(\bar{\mathbf{u}}_\alpha) - J_0(\mathbf{u}^*)| < \delta, \quad (4.39)$$

for all  $N > N_2$ . Since  $\mathbf{u}^*$  is the optimal control of Problem (Q), it follows that

$$J_0(\mathbf{u}^*) \leq J_0(\mathbf{u}^{N,*}).$$

Furthermore,  $\mathbf{u}^{N,*}$  is an optimal control of Problem (Q<sub>N</sub>). Thus,

$$J_0(\mathbf{u}^{N,*}) \leq J_0(\bar{\mathbf{u}}_\alpha^N).$$

On the other hand,  $\mathbf{u}^{N,*}$  is a suboptimal control of Problem (Q). It follows that, for all  $N > N_2$ ,

$$J_0(\mathbf{u}^*) \leq J_0(\mathbf{u}^{N,*}) \leq J_0(\bar{\mathbf{u}}_\alpha^N). \quad (4.40)$$

Combining (4.39) and (4.40) gives

$$J_0(\mathbf{u}^*) \leq J_0(\mathbf{u}^{N,*}) \leq J_0(\bar{\mathbf{u}}_\alpha^N) \leq J_0(\mathbf{u}^*) + \delta.$$

Since  $\delta > 0$  can be chosen arbitrarily, and  $\bar{\mathbf{u}}_\alpha^N \in \mathcal{B}_N^0$ , it is clear that  $J_0(\mathbf{u}^{N,*}) \rightarrow J_0(\mathbf{u}^*)$  as  $N \rightarrow \infty$ .  $\square$

Theorem 4.4 indicates that the optimal cost of Problem  $(Q_N)$  will converge to the optimal cost of Problem  $(Q)$  as  $N \rightarrow \infty$ . However, there is no guarantee that the optimal control of Problem  $(Q_N)$  itself will converge to the optimal control of Problem  $(Q)$ . However, we have the following result.

**Theorem 4.5.** *let  $\mathbf{u}^{N,*}$  be a piecewise-constant control constructed from the optimal control vector  $\boldsymbol{\sigma}^{N,*}$  of Problem  $(Q_N)$ . Let  $\mathbf{u}^*$  be an optimal control of the original problem  $(Q)$ . If  $\{\mathbf{u}^{N,*}\}_{N=2}^{\infty}$  converges to  $\bar{\mathbf{u}} \in \mathcal{U}$  almost everywhere on  $[0, t_f]$  as  $N \rightarrow \infty$ , then  $\bar{\mathbf{u}}$  is also an optimal control of Problem  $(Q)$ .*

*Proof.* First, by Lemma 4.3, we have

$$\lim_{N \rightarrow \infty} J_0(\mathbf{u}^{N,*}) \rightarrow J_0(\bar{\mathbf{u}}).$$

From Theorem 4.4, we recall that

$$\lim_{N \rightarrow \infty} J_0(\mathbf{u}^{N,*}) \rightarrow J_0(\mathbf{u}^*).$$

Since the limit of a convergence sequence is unique, we have

$$J_0(\bar{\mathbf{u}}) = J_0(\mathbf{u}^*).$$

It remains to show that  $\bar{\mathbf{u}}$  is a feasible control of Problem  $(Q)$ . On the contrary, suppose that it is not true. Then, there exists an integer  $i \in \{1, \dots, N_c\}$  and a non-zero interval  $\mathcal{I} \subset [0, t_f]$  such that

$$h_i(\bar{t}, \mathbf{x}(\bar{t}|\bar{\mathbf{u}})) > 0, \quad \forall t \in \mathcal{I}. \quad (4.41)$$

Since, by (4.A.4),  $h_i$  is continuous, it follows from Lemma 4.3 that

$$\lim_{N \rightarrow \infty} |\mathbf{x}(t|\mathbf{u}^{N,*}) - \mathbf{x}(t|\bar{\mathbf{u}})| = 0,$$

for each  $t \in [0, t_f]$ . Furthermore,  $\mathbf{x}(t|\mathbf{u}^{N,*}) \in \mathbf{X}$  for all  $t \in [0, t_f]$ , where  $\mathbf{X}$  is a bounded set (see Lemma 4.1). Thus,

$$\lim_{N \rightarrow \infty} \int_{\mathcal{I}} |h_i(t, \mathbf{x}(t|\mathbf{u}^{N,*})) - h_i(t, \mathbf{x}(t|\bar{\mathbf{u}}))| dt = 0.$$

Therefore, if (4.41) is valid, then

$$\int_{\mathcal{I}} h_i(t, \mathbf{x}(t|\bar{\mathbf{u}})) dt = \int_{\mathcal{I}} \{h_i(t, \mathbf{x}(t|\bar{\mathbf{u}})) - h_i(t, \mathbf{x}(t|\mathbf{u}^{N,*}))\} dt + \int_{\mathcal{I}} h_i(t, \mathbf{x}(t|\mathbf{u}^{N,*})) dt > 0.$$

Since  $\mathbf{u}^{N,*}$  is an optimal control of Problem  $(Q_N)$ , we have

$$h_i(t, \mathbf{x}(t|\mathbf{u}^{N,*})) \leq 0, \quad \forall t \in [0, t_f].$$

Thus,

$$0 < \int_{\mathcal{I}} h_i(t, \mathbf{x}(t|\bar{\mathbf{u}})) dt = \lim_{N \rightarrow \infty} \int_{\mathcal{I}} \{h_i(t, \mathbf{x}(t|\bar{\mathbf{u}}) - h_i(t, \mathbf{x}(t|\mathbf{u}^{N,*}))\} dt = 0. \quad (4.42)$$

This is a contradiction. Thus,  $\bar{\mathbf{u}}$  is feasible as required.  $\square$

#### 4.3.4 Computational method

To solve Problem  $(Q_{N,\varepsilon,\gamma})$ , the gradient formula of the augmented cost function (4.21) with respect to the control vector  $\boldsymbol{\sigma}^N$  is needed. It is derived, although rather involved, via variational formulae given below. Let  $H_{\varepsilon,\gamma}$  be the Hamiltonian function defined by

$$\begin{aligned} H_{\varepsilon,\gamma} &= \mathcal{L}_0(t) + \boldsymbol{\lambda}(t)^\top \mathbf{f}(t) + \gamma \sum_{i=1}^{N_c} \mathcal{L}_{i,\varepsilon}(t, \mathbf{x}(t|\mathbf{u})) + \sum_{k=1}^M \mathcal{L}_0(t + \tau_k) \\ &\quad + \sum_{k=1}^M (\bar{\boldsymbol{\lambda}}^k)^\top(t) \bar{\mathbf{f}}_k(t) e(t_f - t - \tau_k), \end{aligned} \quad (4.43)$$

where the following abbreviations are used

$$\begin{aligned} \mathbf{f}(t) &= \mathbf{f}(t, \mathbf{x}(t|\mathbf{u}), \tilde{\mathbf{x}}(t|\mathbf{u}), \mathbf{u}(t), \tilde{\mathbf{u}}(t)) \\ \mathcal{L}_0(t) &= \mathcal{L}_0(t, \mathbf{x}(t|\mathbf{u}), \tilde{\mathbf{x}}(t|\mathbf{u}), \mathbf{u}(t), \tilde{\mathbf{u}}(t)), \end{aligned}$$

and for each  $k$ ,  $k = 1, \dots, M$ ,

$$\bar{\boldsymbol{\lambda}}^k(t) = \boldsymbol{\lambda}(t + \tau_k), \quad (4.44)$$

$$\bar{\mathbf{f}}_k(t) = \mathbf{f}(t + \tau_k, \mathbf{x}(t + \tau_k|\mathbf{u}), \tilde{\mathbf{x}}(t + \tau_k|\mathbf{u}), \mathbf{u}(t + \tau_k), \tilde{\mathbf{u}}(t + \tau_k)), \quad (4.45)$$

$$\mathcal{L}_0(t + \tau_k) = \mathcal{L}_0(t + \tau_k, \mathbf{x}(t + \tau_k|\mathbf{u}), \tilde{\mathbf{x}}(t + \tau_k|\mathbf{u}), \mathbf{u}(t + \tau_k), \tilde{\mathbf{u}}(t + \tau_k)),$$

while  $e(\cdot)$  is the unit step function defined by

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

and

$$\mathcal{L}_{i,\varepsilon}(t, \mathbf{x}(t|\mathbf{u})) = \begin{cases} h_i(t, \mathbf{x}(t|\mathbf{u})), & \text{if } h_i(t, \mathbf{x}(t|\mathbf{u})) > \varepsilon, \\ \frac{\{h_i(t, \mathbf{x}(t|\mathbf{u})) - \varepsilon\}^2}{4\varepsilon}, & \text{if } -\varepsilon \leq h_i(t, \mathbf{x}(t|\mathbf{u})) \leq \varepsilon, \\ 0, & \text{if } h_i(t, \mathbf{x}(t|\mathbf{u})) < -\varepsilon. \end{cases} \quad (4.46)$$

Let  $\boldsymbol{\lambda}(t)$  be the corresponding solution of the co-state system defined by

$$\dot{\boldsymbol{\lambda}}(t) = -\frac{\partial H_{\varepsilon,\gamma}}{\partial \mathbf{x}}, \quad (4.47)$$

with boundary conditions

$$\boldsymbol{\lambda}(t_f) = \frac{\partial \Phi_0(\mathbf{x}(t_f|\mathbf{u}))}{\partial \mathbf{x}}, \quad (4.48)$$

$$\boldsymbol{\lambda}(t) = \mathbf{0}, \quad t > t_f. \quad (4.49)$$

Then, for each pair of  $\varepsilon$  and  $\gamma$ , we have the following result.

**Theorem 4.6.** *Let  $\mathbf{u}$  be any control in  $\mathcal{U}$  and let  $\Delta\mathbf{u}(t) \in \mathbb{R}^r$  be any bounded measurable function defined in  $[-\tau_M, t_f]$  with  $\Delta\mathbf{u}(t) = \mathbf{0}$  for all  $t \in [-\tau_M, 0]$ . Then, the directional derivative of the function  $\tilde{J}_{\varepsilon,\gamma}$  given by (4.21) is:*

$$\Delta J_{\varepsilon,\gamma}(\mathbf{u}) = \int_0^{t_f} \frac{\partial H_{\varepsilon,\gamma}}{\partial \mathbf{u}} \Delta\mathbf{u}(t) dt,$$

where

$$\begin{aligned} J_{\varepsilon,\gamma}(\mathbf{u}) &= \Phi_0(\mathbf{x}(t_f|\mathbf{u})) + \int_0^{t_f} \mathcal{L}_0(t, \mathbf{x}(t|\mathbf{u}), \tilde{\mathbf{x}}(t|\mathbf{u}), \mathbf{u}(t), \tilde{\mathbf{u}}(t)) dt \\ &+ \gamma \sum_{i=1}^{N_c} \int_0^{t_f} \mathcal{L}_{i,\varepsilon}(t, \mathbf{x}(t|\mathbf{u})) dt. \end{aligned} \quad (4.50)$$

*Proof.* Let  $\mathbf{u}(t) \in \mathcal{U}$  be arbitrary but fixed. Let the control vector  $\mathbf{u}(t)$  be perturbed by  $\varepsilon\Delta\mathbf{u}(t)$ , where  $\varepsilon > 0$  is a small real number and  $\Delta\mathbf{u}(t)$  is an arbitrary but fixed perturbation of  $\mathbf{u}(t)$  given by

$$\begin{aligned} \Delta\mathbf{u}(t) &= [\Delta u_1(t), \Delta u_2(t), \dots, \Delta u_r(t)]^\top, \quad t \in [0, t_f] \\ \Delta\mathbf{u}(t) &= \mathbf{0}, \quad t < 0. \end{aligned}$$

where  $\Delta u_j(t)$ ,  $j = 1, \dots, r$ , are arbitrary but given functions. Let,

$$\mathbf{u}_\varepsilon(t) = \mathbf{u}(t) + \varepsilon\Delta\mathbf{u}(t). \quad (4.51)$$

Furthermore, let

$$\tilde{\mathbf{u}}_\varepsilon^k(t) = \mathbf{u}_\varepsilon(t - \tau_k) = \mathbf{u}(t - \tau_k) + \varepsilon\Delta\mathbf{u}(t - \tau_k), \quad k = 1, \dots, M,$$

and

$$\tilde{\mathbf{u}}_\varepsilon(t) = [\tilde{\mathbf{u}}_\varepsilon^1(t)^\top, \dots, \tilde{\mathbf{u}}_\varepsilon^M(t)^\top]^\top.$$



For brevity, let  $\mathbf{x}(\cdot)$  denote the solution of system (4.1)-(4.3) with the control  $\mathbf{u}$ , and  $\mathbf{x}_\epsilon(\cdot)$  denote the solution of system (4.1)-(4.3) with the control  $\mathbf{u}_\epsilon$ . Clearly,

$$\begin{aligned}\mathbf{x}(t) &= \mathbf{x}(0) + \int_0^t \mathbf{f}(s, \mathbf{x}(s|\mathbf{u}), \tilde{\mathbf{x}}(s|\mathbf{u}), \mathbf{u}(s), \tilde{\mathbf{u}}(t)) ds, \\ \mathbf{x}_\epsilon(t) &= \mathbf{x}(0) + \int_0^t \mathbf{f}(s, \mathbf{x}(s|\mathbf{u}_\epsilon), \tilde{\mathbf{x}}(s|\mathbf{u}_\epsilon), \mathbf{u}_\epsilon, \tilde{\mathbf{u}}_\epsilon(s)) ds.\end{aligned}$$

We will use the notation  $\frac{\partial}{\partial \tilde{\mathbf{x}}^k}$  to denote the partial differentiation with respect to the  $k$ th delayed state in  $\tilde{\mathbf{x}}(t)$  (i.e. the partial differentiation with respect to  $\mathbf{x}(t - \tau_k)$ ), and the notation  $\frac{\partial}{\partial \tilde{\mathbf{u}}^k}$  to denote the partial differentiation with respect to the  $k$ th delayed control in  $\tilde{\mathbf{u}}(t)$  (i.e. the partial differentiation with respect to  $\mathbf{u}(t - \tau_k)$ ). Then, by the chain rule, we have

$$\begin{aligned}\Delta \mathbf{x}(t) &= \left. \frac{d\mathbf{x}_\epsilon(t)}{d\epsilon} \right|_{\epsilon=0} \\ &= \int_0^t \left\{ \frac{\partial \mathbf{f}(s)}{\partial \mathbf{x}} \Delta \mathbf{x}(s) + \frac{\partial \mathbf{f}(s)}{\partial \mathbf{u}} \Delta \mathbf{u}(s) + \sum_{k=1}^M \frac{\partial \mathbf{f}(s)}{\partial \tilde{\mathbf{x}}^k} \Delta \mathbf{x}(s - \tau_k) \right\} ds \\ &\quad + \sum_{k=1}^M \int_0^t \frac{\partial \mathbf{f}(s)}{\partial \tilde{\mathbf{u}}^k} \Delta \mathbf{u}(s - \tau_k) ds.\end{aligned}$$

Clearly,

$$\begin{aligned}\frac{d(\Delta \mathbf{x}(t))}{dt} &= \frac{\partial \mathbf{f}(t)}{\partial \mathbf{x}} \Delta \mathbf{x}(t) + \frac{\partial \mathbf{f}(t)}{\partial \mathbf{u}} \Delta \mathbf{u}(t) + \sum_{k=1}^M \frac{\partial \mathbf{f}(t)}{\partial \tilde{\mathbf{x}}^k} \Delta \mathbf{x}(t - \tau_k) \\ &\quad + \sum_{k=1}^M \frac{\partial \mathbf{f}(t)}{\partial \tilde{\mathbf{u}}^k} \Delta \mathbf{u}(t - \tau_k).\end{aligned}\tag{4.52}$$

Then, by the chain rule,

$$\begin{aligned}\Delta J_{\epsilon, \gamma}(\mathbf{u}) &= \left. \frac{dJ_{\epsilon, \gamma}(\mathbf{u}_\epsilon)}{d\epsilon} \right|_{\epsilon=0} \\ &= \frac{\partial \Phi_0(\mathbf{x}(t_f|\mathbf{u}))}{\partial \mathbf{x}(t_f)} \Delta \mathbf{x}(t_f) + \int_0^{t_f} \left\{ \frac{\partial \mathcal{L}_0(t)}{\partial \mathbf{x}} \Delta \mathbf{x}(t) + \frac{\partial \mathcal{L}_0(t)}{\partial \mathbf{u}} \Delta \mathbf{u}(t) \right\} dt \\ &\quad + \int_0^{t_f} \left\{ \sum_{k=1}^M \frac{\partial \mathcal{L}_0(t)}{\partial \tilde{\mathbf{x}}^k} \Delta \mathbf{x}(t - \tau_k) + \sum_{k=1}^M \frac{\partial \mathcal{L}_0(t)}{\partial \tilde{\mathbf{u}}^k} \Delta \mathbf{u}(t - \tau_k) \right\} dt \\ &\quad + \gamma \sum_{i=1}^{N_c} \int_0^{t_f} \frac{\partial \mathcal{L}_{i, \epsilon}(t, \mathbf{x}(t|\mathbf{u}))}{\partial \mathbf{x}} \Delta \mathbf{x}(t) dt.\end{aligned}\tag{4.53}$$

Consider the definition of  $\bar{\mathbf{f}}_k(t)$  in (4.45), we have

$$\sum_{k=1}^M \frac{\partial \mathcal{L}_0(t)}{\partial \bar{\mathbf{x}}^k} \Delta \mathbf{x}(t - \tau_k) = \sum_{k=1}^M e(t_f - t - \tau_k) \frac{\partial \mathcal{L}_0(t + \tau_k)}{\partial \mathbf{x}} \Delta \mathbf{x}(t), \quad (4.54)$$

$$\sum_{k=1}^M \frac{\partial \mathcal{L}_0(t)}{\partial \bar{\mathbf{u}}^k} \Delta \mathbf{u}(t - \tau_k) = \sum_{k=1}^M e(t_f - t - \tau_k) \frac{\partial \mathcal{L}_0(t + \tau_k)}{\partial \mathbf{u}} \Delta \mathbf{u}(t). \quad (4.55)$$

Substituting (4.54)-(4.55) into (4.56) yields,

$$\begin{aligned} \Delta J_{\varepsilon, \gamma}(\mathbf{u}) &= \frac{\partial \Phi_0(\mathbf{x}(t_f | \mathbf{u}))}{\partial \mathbf{x}(t_f)} \Delta \mathbf{x}(t_f) + \int_0^{t_f} \frac{\partial \mathcal{L}_0(t)}{\partial \mathbf{u}} \Delta \mathbf{u}(t) dt \\ &+ \sum_{k=1}^M \int_0^{t_f} \frac{\partial \mathcal{L}_0(t + \tau_k)}{\partial \mathbf{u}} e(t_f - t - \tau_k) \Delta \mathbf{u}(t) dt + \int_0^{t_f} \frac{\partial \mathcal{L}_0(t)}{\partial \mathbf{x}} \Delta \mathbf{x}(t) dt \\ &+ \int_0^{t_f} \left\{ \sum_{k=1}^M \frac{\partial \mathcal{L}_0(t + \tau_k)}{\partial \mathbf{x}} e(t_f - t - \tau_k) + \gamma \sum_{i=1}^{N_c} \frac{\partial \mathcal{L}_{i, \varepsilon}(t, \mathbf{x}(t | \mathbf{u}))}{\partial \mathbf{x}} \right\} \Delta \mathbf{x}(t) dt. \end{aligned} \quad (4.56)$$

By using the definition of the Hamiltonian function  $H_{\varepsilon, \gamma}$ , we have

$$\begin{aligned} \mathcal{L}_0(t) + \sum_{k=1}^M \mathcal{L}_0(t + \tau_k) e(t_f - t - \tau_k) + \gamma \sum_{i=1}^{N_c} \mathcal{L}_{i, \varepsilon}(t, \mathbf{x}(t | \mathbf{u})) \\ = H_{\varepsilon, \gamma} - \boldsymbol{\lambda}(t)^\top \mathbf{f}(t) + \sum_{k=1}^M (\bar{\boldsymbol{\lambda}}^k)^\top(t) \bar{\mathbf{f}}_k(t) e(t_f - t - \tau_k). \end{aligned}$$

Thus,

$$\begin{aligned} \Delta J_{\varepsilon, \gamma}(\mathbf{u}) &= \frac{\partial \Phi_0(\mathbf{x}(t_f | \mathbf{u}))}{\partial \mathbf{x}(t_f)} \Delta \mathbf{x}(t_f) + \int_0^{t_f} \left\{ \frac{\partial H_{\varepsilon, \gamma}}{\partial \mathbf{x}} \Delta \mathbf{x}(t) + \frac{\partial H_{\varepsilon, \gamma}}{\partial \mathbf{u}} \Delta \mathbf{u}(t) \right\} dt \\ &- \int_0^{t_f} \boldsymbol{\lambda}(t)^\top \left\{ \frac{\partial \mathbf{f}(t)}{\partial \mathbf{x}} \Delta \mathbf{x}(t) + \frac{\partial \mathbf{f}(t)}{\partial \mathbf{u}} \Delta \mathbf{u}(t) \right\} dt \\ &- \sum_{k=1}^M \int_0^{t_f} (\bar{\boldsymbol{\lambda}}^k)^\top(t) e(t_f - t - \tau_k) \left\{ \frac{\partial \bar{\mathbf{f}}_k(t)}{\partial \mathbf{x}} \Delta \mathbf{x}(t) + \frac{\partial \bar{\mathbf{f}}_k(t)}{\partial \mathbf{u}} \Delta \mathbf{u}(t) \right\} dt. \end{aligned} \quad (4.57)$$

By (4.43)-(4.45), and the definition of  $e(\cdot)$ , it follows that

$$\begin{aligned} \sum_{k=1}^M \int_0^{t_f} (\bar{\boldsymbol{\lambda}}^k)^\top(t) e(t_f - t - \tau_k) \left\{ \frac{\partial \bar{\mathbf{f}}_k(t)}{\partial \mathbf{x}} \Delta \mathbf{x}(t) + \frac{\partial \bar{\mathbf{f}}_k(t)}{\partial \mathbf{u}} \Delta \mathbf{u}(t) \right\} \\ = \sum_{k=1}^M \int_0^{t_f - \tau_k} (\bar{\boldsymbol{\lambda}}^k)^\top(t) e(t_f - t - \tau_k) \left\{ \frac{\partial \bar{\mathbf{f}}_k(t)}{\partial \mathbf{x}} \Delta \mathbf{x}(t) + \frac{\partial \bar{\mathbf{f}}_k(t)}{\partial \mathbf{u}} \Delta \mathbf{u}(t) \right\} \end{aligned}$$

$$= \sum_{k=1}^M \int_{\tau_k}^{t_f} (\boldsymbol{\lambda}^k)^\top(t) \left\{ \frac{\partial \mathbf{f}(t)}{\partial \tilde{\mathbf{x}}} \Delta \mathbf{x}(t - \tau_k) + \frac{\partial \hat{\mathbf{f}}_k(t)}{\partial \tilde{\mathbf{u}}} \Delta \mathbf{u}(t - \tau_k) \right\}. \quad (4.58)$$

Since  $\Delta \mathbf{x}(t - \tau_k) = \mathbf{0}$ , for  $0 \leq t \leq \tau_k$ , and  $\Delta \mathbf{u}(t - \tau_k) = \mathbf{0}$ , for  $0 \leq t < \tau_k$ , we have

$$\begin{aligned} & \sum_{k=1}^M \int_{\tau_k}^{t_f} (\boldsymbol{\lambda}^k)^\top(t) \left\{ \frac{\partial \mathbf{f}(t)}{\partial \tilde{\mathbf{x}}} \Delta \mathbf{x}(t - \tau_k) + \frac{\partial \mathbf{f}(t)}{\partial \tilde{\mathbf{u}}} \Delta \mathbf{u}(t - \tau_k) \right\} \\ &= \sum_{k=1}^M \int_0^{t_f} (\boldsymbol{\lambda}^k)^\top(t) \left\{ \frac{\partial \mathbf{f}(t)}{\partial \tilde{\mathbf{x}}} \Delta \mathbf{x}(t - \tau_k) + \frac{\partial \mathbf{f}(t)}{\partial \tilde{\mathbf{u}}} \Delta \mathbf{u}(t - \tau_k) \right\}. \end{aligned} \quad (4.59)$$

Substituting (4.59) and (4.58) into (4.57), and then combining with (4.58), it gives

$$\begin{aligned} \Delta J_{\varepsilon, \gamma}(\mathbf{u}) &= \frac{\partial \Phi_0(\mathbf{x}(t_f | \mathbf{u}))}{\partial \mathbf{x}(t_f)} \Delta \mathbf{x}(t_f) + \int_0^{t_f} \frac{\partial H_{\varepsilon, \gamma}}{\partial \mathbf{u}} \Delta \mathbf{u}(t) dt \\ & \quad \int_0^{t_f} \left\{ \frac{\partial H_{\varepsilon, \gamma}}{\partial \mathbf{x}} \Delta \mathbf{x}(t) - \boldsymbol{\lambda}(t)^\top \frac{d(\Delta \mathbf{x}(t))}{dt} \right\} dt. \end{aligned} \quad (4.60)$$

Integrating the last term of (4.60) by parts gives

$$\begin{aligned} \Delta J_{\varepsilon, \gamma}(\mathbf{u}) &= \frac{\partial \Phi_0(\mathbf{x}(t_f | \mathbf{u}))}{\partial \mathbf{x}(t_f)} \Delta \mathbf{x}(t_f) - \boldsymbol{\lambda}(t)^\top \Delta \mathbf{x}(t) \Big|_0^{t_f} \\ & \quad + \int_0^{t_f} \left\{ \frac{\partial H_{\varepsilon, \gamma}}{\partial \mathbf{x}} \Delta \mathbf{x}(t) + \frac{\partial H_{\varepsilon, \gamma}}{\partial \mathbf{u}} \Delta \mathbf{u}(t) + \frac{d(\boldsymbol{\lambda}(t))^\top}{dt} \Delta \mathbf{x}(t) \right\} dt. \end{aligned} \quad (4.61)$$

Since  $\mathbf{x}(0)$  is a given constant, it is clear that  $\Delta \mathbf{x}(0) = 0$ . Hence,  $\boldsymbol{\lambda}(0)^\top \Delta \mathbf{x}(0) = 0$ . Since  $\Delta \mathbf{u}(t)$  is arbitrary on  $[0, t_f]$ , substituting (4.47), (4.48), and (4.49) into (4.61) yields

$$\Delta J_{\varepsilon, \gamma}(\mathbf{u}) = \int_0^{t_f} \frac{\partial H_{\varepsilon, \gamma}}{\partial \mathbf{u}} \Delta \mathbf{u}(t) dt. \quad (4.62)$$

This completes the proof.  $\square$

**Theorem 4.7.** For each  $q = 1, \dots, N$ , the gradient of the augmented cost function  $\tilde{J}_{\varepsilon, \gamma}$  with respect to  $\boldsymbol{\sigma}^{N, q}$  is

$$\frac{\partial \tilde{J}_{\varepsilon, \gamma}(\boldsymbol{\sigma}^N)}{\partial \boldsymbol{\sigma}^{N, q}} = \int_{t_{q-1}}^{t_q} \frac{\partial H_{\varepsilon, \gamma}}{\partial \mathbf{u}} dt. \quad (4.63)$$

*Proof.* Let  $\mathbf{u}_\varepsilon(t)$  be defined as (4.51), where

$$\Delta \mathbf{u}(t) = \Delta \boldsymbol{\sigma}^N \chi_{[t_{q-1}, t_q)}(t), \quad (4.64)$$

$$\Delta \mathbf{u}(t) = \mathbf{0}, \quad t < 0, \quad (4.65)$$

and

$$\Delta\boldsymbol{\sigma}^N = [\mathbf{0}^\top, \dots, \mathbf{0}^\top, (\Delta\boldsymbol{\sigma}^{N,q})^\top, \mathbf{0}^\top, \dots, \mathbf{0}^\top]^\top.$$

Substitute (4.64) into (4.62), it is clear that

$$\Delta J_{\varepsilon,\gamma}(\mathbf{u}) = \int_{t_{q-1}}^{t_q} \frac{\partial H_{\varepsilon,\gamma}}{\partial \mathbf{u}} \Delta\boldsymbol{\sigma}^{N,q} dt. \quad (4.66)$$

Restricting the controls to  $\mathcal{U}^N$  yields

$$J_{\varepsilon,\gamma}(\mathbf{u}^N) = \tilde{J}_{\varepsilon,\gamma}(\boldsymbol{\sigma}^N).$$

Using Theorem 4.6, it is clear that

$$\begin{aligned} \Delta J_{\varepsilon,\gamma}(\mathbf{u}^N) &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{\partial J_{\varepsilon,\gamma}(\mathbf{u}^N + \epsilon \Delta \mathbf{u}) - J_{\varepsilon,\gamma}(\mathbf{u}^N)}{\partial \epsilon} \right\} \\ &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{\partial \tilde{J}_{\varepsilon,\gamma}(\boldsymbol{\sigma}^N + \epsilon \Delta \boldsymbol{\sigma}^N) - \tilde{J}_{\varepsilon,\gamma}(\boldsymbol{\sigma}^N)}{\partial \epsilon} \right\} \\ &= \left\langle \frac{\partial \tilde{J}_{\varepsilon,\gamma}(\boldsymbol{\sigma}^N)}{\partial \boldsymbol{\sigma}^N}, \Delta \boldsymbol{\sigma}^N \right\rangle \\ &= \left\langle \frac{\partial \tilde{J}_{\varepsilon,\gamma}(\boldsymbol{\sigma}^N)}{\partial \boldsymbol{\sigma}^{N,q}}, \Delta \boldsymbol{\sigma}^{N,q} \right\rangle. \end{aligned} \quad (4.67)$$

Combining (4.66) and (4.67), and noting that  $\Delta\boldsymbol{\sigma}^{N,q}$  is arbitrary, the theorem follows readily.  $\square$

With the gradient formula given in Theorem 4.7, Problem  $(Q_{N,\varepsilon,\gamma})$ , for each  $\varepsilon > 0$  and  $\gamma > 0$ , can be solved by using a gradient-based optimization technique, such as the sequential quadratic programming approximation method. We propose the following algorithm.

**Algorithm 4.1.**

*Step 1.* Set  $\varepsilon = 0.01$  and  $\gamma = 10$ .

*Step 2.* Solve the state differential equation (4.1) with the initial conditions (4.3) forward in time from  $t = 0$  to  $t = t_f$ . Let the solution obtained be denoted by  $\mathbf{x}(\cdot | \boldsymbol{\sigma}^N)$ .

*Step 3.* Compute the value of the augmented cost function given by (4.21).

*Step 4.* Solve the co-state system (4.47) backward in time from  $t = t_f$  to  $t = 0$  with the boundary condition (4.48) and (4.49), where  $\mathbf{x}(\cdot | \boldsymbol{\sigma}^N)$  is from Step 2. Let the solution obtained be referred to as  $\boldsymbol{\lambda}(\cdot | \boldsymbol{\sigma}^N)$ .

- Step 5.* Compute the gradient of the augmented cost function (4.21) with respect to  $\sigma^N$  according to (4.63).
- Step 6.* Solve the approximate Problem  $(Q_{N,\varepsilon,\gamma})$  by using the sequential quadratic programming approximation scheme with active set strategy. Let the optimal control vector obtained be denoted as  $\sigma_{\varepsilon,\gamma}^{N,*}$ .
- Step 7.* Check the feasibility of the continuous state inequality constraints for all  $t \in [0, t_f]$ . If  $\sigma_{\varepsilon,\gamma}^{N,*}$  is feasible, go to Step 8; otherwise, set  $\gamma := 10 \times \gamma$  and go to Step 2.
- Step 8.* Set  $\varepsilon := \varepsilon/10$ . If  $\varepsilon > \varepsilon_{min}$ , go to Step 2, else successfully exist.

**Remark 4.1.** Problem  $(Q_{N,\varepsilon,\gamma})$  with  $\varepsilon > 0$  and  $\eta_z > 0$ ,  $z = 1, \dots, N_g$ , chosen as detailed in Step 7 of Algorithm 4.1 is an approximate problem of Problem (Q). By using arguments similar to that given for Theorem 4.4, it can be shown that the approximate optimal cost will converge to the true optimal cost as  $N \rightarrow \infty$ . In practice, we could start with a small integer  $N$ , and obtain the optimal control vector of the corresponding Problem  $(Q_{N,\varepsilon,\gamma})$ . We then double the value of  $N$  and re-calculate the optimal control vector of the corresponding Problem  $(Q_{N,\varepsilon,\gamma})$  with the previous optimal control vector taken as the initial guess in the optimization process. We repeat this process until the reduction in the cost value is negligible. From extensive simulation studies, it is observed that  $N$  does not need to be very large, certainly the one used in our simulation study is more than sufficient. In fact, the approximation of the control by a piecewise-constant function should also be followed in some real applications. The extension to the case where the control is approximated by piecewise linear or piecewise smooth function is straightforward. The switching times will affect the cost value but is insignificant. The main advantage of taking the switching times as decision variables is that the number of switching times could be reduced for achieving the same cost value. However, the price to pay is a significant increase in the computational burden.

**Remark 4.2.** Theorem 4.2 ensures that for each  $\varepsilon > 0$ , a suitable  $\gamma(\varepsilon)$  can be obtained in finite number of iterations. This is due to the special structure of the penalty function used.

## 4.4 Application: optimal control of an evaporation process

We now demonstrate the applicability of our approach to a realistic optimal control problem arising in evaporation process. Specifically, we consider the industrial evaporation process described in [119]. The alumina production process mainly includes aluminium

hydroxide solution preparing process, clarifying process, dissolving process, decomposing process, evaporation process, and roasting process. The main contents of the mother liquor discharged from the decomposing process are sodium hydroxide and aluminum oxide which are valuable materials needed in the recycling. However, the concentration of the mother liquor is lower than the required concentration needed for the leaching process or the grinding process. Thus, it cannot be used directly, and hence the evaporation process is needed to improve the concentration of the mother liquor, such that the acid and caustic materials can be re-used.

The returned lye discharged from the evaporation process is one of the main raw materials needed for converting the bauxite to aluminium hydroxide solution (which is also called the raw slurry). It is known [157] that the quality of raw slurry has a direct influence on the quality of the final product. Unacceptable fluctuation in the composition of the returned lye can lead to instability of the blending process during the preparation of the raw slurry. Consequently, the quality of the product obtained cannot be guaranteed. Usually, only the solution concentration at the outlet of the evaporation process is measured at every two-hour interval, but it takes about one hour for the feed flowing through the evaporation process. Clearly, inappropriate control of the evaporation process will lead to unacceptable output solution, yet with high steam consumption. Hence, optimal control is needed.

In this section, the control method proposed in Section 4.3 is applied to study the optimal control of a practical alumina evaporation process, in which the objective is to find a control such that the specific quality of the sodium aluminate solution control is met with the least energy usage, while the constraints on the state and the control are satisfied.

#### 4.4.1 The evaporation system

In a typical alumina production factory in China, the objective of the evaporation process is to increase the concentration of the industrial sodium aluminate solution (the sodium hydroxide content of the solution is to be increased to about 160~170g/L for Bayer process). The industrial sodium aluminate solution is highly viscous. It contains many impurities, such as sodium carbonate and sodium sulfate, which can easily emit from the solution due to the increase of caustic alkali concentration or the decrease of temperature. Crystallization of the impurities will cause serious pipe plug problem. To avoid the forming of high viscosity fluid at low temperature, the alumina production factory employs multiple falling film evaporators for the evaporation process as shown in Figure 4.1. It consists of four falling film tube evaporators, three direct pre-heaters, three flash evaporators, four flash tanks and a condenser.

The four falling film tube evaporators are connected in series with reference to their

vapor and liquor lines. Heat is supplied to the first evaporator by live steam generated in the power plant after the pressure is reduced to about 0.5Mpa. The vapor, produced by each of the first three evaporators, is used as the heating source for the next evaporators in series. The vapor produced by the fourth evaporator is condensed and then discharged from the process through the condenser installed in the water circuit.

Between each of the two adjacent evaporators there is a preheater, which is used to preheat the solution fed into the previous (with respect to the vapor lines) evaporator. The heating source of each preheater comes from the vapor produced by the previous (with respect to the vapor lines) evaporator and flash evaporator.

The feed enters the system at the third and the fourth evaporators, and flows backward. Finally the solution leaves from the first evaporator and is fed into the flash evaporators where the final product is drained. The structure of a falling tube evaporator is as shown

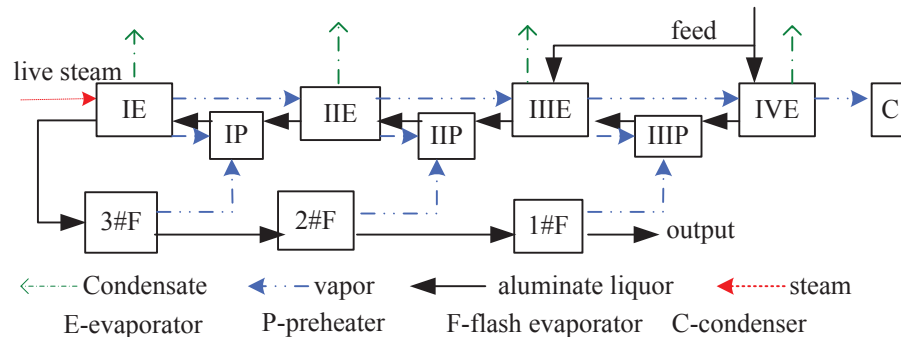


Figure 4.1: Flow sheet of alumina evaporation process

in Figure 4.2. It consists of an evaporation vessel (A), a heat exchanger (B), a number of pumps that realizes the transfer of solution in and out the evaporator, and valves that control the steam in and out of the evaporator.

The liquor, i.e. the industrial sodium aluminate solution, is injected at the bottom of the evaporation vessel through Pump 1. It is then pumped to the top of the evaporator, where a distributor is used to provide a uniformly distributed falling liquid film inside each heating tube. The effect of the distribution depends on the viscosity of the solution and the cycling rate. Because the broken up of the falling film will cause serious scaring problem [109], the change of the cycling rate should only be slight such that the change in the surface tension force of the film is mild.

The liquor flows through the inside of the heating tubes, and the heating steam fills the outside of the heating tubes. It leads to a heat transfer from the steam to the liquor. The heating steam is governed by Valve 1. As the liquor goes through the heating tube, it reaches its boiling temperature, causing water within the liquor to evaporate. A separator is used to split the vapor from the solution at the top of the evaporator vessel. The vapor is drained through Valve 2 which is used to connect the two adjacent evaporators, and

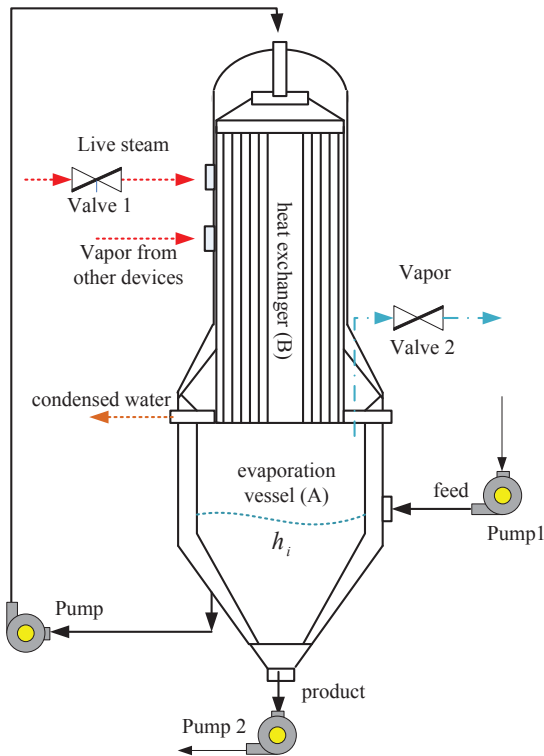


Figure 4.2: Structure of a vertical tube falling film evaporator

hence cannot be adjusted. Finally, the liquor is drained through Pump 2.

In practice, once the control variables are set, they must be used for at least 5 minutes. In other words, the control variables (i.e. flow rate of live steam) are to be adjusted in a piecewise-constant manner. Thus, the control variables can be approximated by piecewise-constant functions with possible discontinuities at the preset switching points. Hence, the optimal control approach introduced in Section 4.3 can be used to deal with the optimal control for the evaporation process.

#### 4.4.2 Mathematical model of the evaporation system

##### *A Dynamic system for the evaporation process*

Three types of dynamical process models have been developed for a falling-film evaporator [46,145]. In this chapter, a pilot-scale evaporation process is considered. As the direct preheaters are used to mix the solution and vapor, there is little change in its level. Thus, the dynamics of the preheaters are omitted; and the evaporator and the corresponding preheater are taken as a whole. Also, we neglect the dynamics of the distributor and the separator. Only the dynamics involved in the evaporation vessel is considered.

A fundamental model that describes the multi-effect falling film evaporation process can be derived under the following assumptions:



- Perfect mixing in each phase in each vessel;
- The absence of any non-condensable gases and the process is adiabatic;
- No transportation lags associated with the movement of steam;
- The cross-sectional area of the evaporator and the specific heat capacity of water are constant;
- All pipes are full.

According to experienced engineers in the factory, the changes of the input solution, for the evaporator, take about 15 minutes to cause an effect on the changes in the evaporation vessel. This is due to the hysteresis during the solution flowing to the evaporation vessel, such as in the distributor and the place where the input solution is injected. Similarly, the time-delays for the flash evaporators are about 5 minutes. From simulation studies, it appears that the values for these delays are acceptable in practice. The variables that are of interest are product temperature in the effect, solution level in the evaporation vessel and product concentration of each effect. In particular, sodium carbonate, sodium hydroxide, and alumina are the three components of the industrial sodium aluminate solution which are measured at every two-hour interval. Thus, under the aforementioned assumptions, the alumina evaporation process model is built based on the principles of heat balance and material balance in unit operations (evaporator and flash evaporator). They are described by the following differential equations with multiple time-delays [119]:

(a) The time variations of the solution levels are:

$$\frac{dh_i(t)}{dt} = \frac{1}{A_i \rho_i(t)} \Delta M_i(t), \quad i = 1, \dots, 7. \quad (4.68)$$

(b) The concentration of the solution in the evaporation vessel is assumed to be equal to that of the output solution. For each  $i$ ,  $i = 1, \dots, 7$ , the three ingredients of the solution, i.e. sodium carbonate, sodium hydroxide, and alumina, are denoted by  $C_i^j$ ,  $j = 1, 2, 3$ . For each  $j$ ,  $j = 1, 2, 3$ , the time variations of the concentrations can be expressed as:

$$\frac{dC_i^j(t)}{dt} = \frac{1}{A_i h_i(t)} [F_{i+1}(t - \tau_1) C_{i+1}^j(t - \tau_1) - F_i(t) C_i^j(t) - \frac{dh_i(t)}{dt} A_i C_i^j(t)], \quad i = 1, 2, 3, \quad (4.69a)$$

$$\frac{dC_i^j(t)}{dt} = \frac{1}{A_i h_i(t)} [F_{i+1}(t - \tau_2) C_{i+1}^j(t - \tau_2) - F_i(t) C_i^j(t) - \frac{dh_i(t)}{dt} A_i C_i^j(t)], \quad i = 4, 5, \quad (4.69b)$$

$$\frac{dC_i^j(t)}{dt} = \frac{1}{A_i h_i(t)} [F_{i+1}(t - \tau_2) C_{i+1}^j(t - \tau_2) + F_{01}(t - \tau_2) C_0^j(t - \tau_2) - F_i(t) C_i^j(t) - \frac{dh_i(t)}{dt} A_i C_i^j(t)], \quad i = 6, \quad (4.69c)$$

$$\frac{dC_i^j(t)}{dt} = \frac{1}{A_i h_i(t)} \left[ F_0(t - \tau_2) C_0^j(t - \tau_2) - F_i(t) C_i^j(t) - \frac{dh_i(t)}{dt} A_i C_i^j(t) \right], \quad i = 7. \quad (4.69d)$$

(c) The time variations of the solution temperatures are:

$$\frac{dT_i(t)}{dt} = \frac{\Delta Q_i(t)}{dt} \frac{1}{A_i h_i(t) cp_i(t)} - \frac{T_i(t)}{cp_i(t)} \frac{dcp_i}{dt} - \frac{T_i(t) \rho_i(t)}{h_i(t)} \frac{dh_i(t)}{dt}, \quad i = 1, \dots, 7. \quad (4.70)$$

Here,  $C_i^j$ , for each  $i = 1, 2, 3$ ;  $j = 1, 2, 3$ , denotes the concentration of the  $j$ th ingredient of the  $i$ th flash evaporator;  $C_i^j$ , for each  $i = 4, \dots, 7$ ;  $j = 1, 2, 3$ , denotes the concentration of the  $j$ th ingredient of the  $(i - 3)$ th evaporator;  $h$  is the solution level of the evaporation vessel;  $A$  is the cross-sectional area of the evaporation vessel;  $cp$  and  $\rho$  are, respectively, the specific heat capacity and the density of the output solution;  $\Delta Q_i$  and  $\Delta M_i$ ,  $i = 1, \dots, 7$ , are the heat changes and the mass changes in the evaporation vessel, respectively;  $T$  is the product temperature;  $F$  is the flow rate of the output solution;  $F_0$  and  $F_{01}$  are, respectively, the feed inputs into the third and fourth evaporators;  $C_0$  is the condensation of the feed injected to the process;  $\tau_1$  and  $\tau_2$  are the time-delays, where  $\tau_1 = 5$  minutes and  $\tau_2 = 15$  minutes.

The mass changes in the evaporation vessels can be calculated by using the following formulas:

$$\Delta M_i(t) = F_{i+1}(t - \tau_1) \rho_{i+1}(t - \tau_1) - V_i(t) - F_i(t) \rho_i(t), \quad i = 1, 2, 3, \quad (4.71)$$

$$\Delta M_4(t) = F_5(t - \tau_2) \rho_5(t - \tau_2) - V_4(t) - F_4(t) \rho_4(t) + V_3(t), \quad (4.72)$$

$$\Delta M_5(t) = F_6(t - \tau_2) \rho_6(t - \tau_2) - V_5(t) - F_5(t) \rho_5(t) + V_2(t), \quad (4.73)$$

$$\Delta M_6(t) = F_7(t - \tau_2) \rho_7(t - \tau_2) - V_6(t) - F_6(t) \rho_6(t) + V_1(t) + F_{01}(t - \tau_2) \rho_0(t - \tau_2), \quad (4.74)$$

$$\Delta M_7(t) = F_0(t - \tau_2) \rho_0(t - \tau_2) - V_7(t) - F_7(t) \rho_7(t), \quad (4.75)$$

where  $V_i$ , for each  $i = 1, 2, 3$ , is the vapor produced by the  $i$ th flash evaporator;  $V_i$ , for each  $i = 4, \dots, 7$ , is the vapor produced by the  $(i - 3)$ th evaporator. They are calculated as follows:

$$V_i(t) = \frac{F_{i+1}(t - \tau_1) \rho_{i+1}(t - \tau_1) [cp_{i+1}(t - \tau_1) T_{i+1}(t - \tau_1) - cp_i(t) T_i(t)]}{H_i(t) - cp_i(t) T_i(t)}, \quad (4.76a)$$

$$V_4(t) = \frac{V_0(t) r_0(t) + V_3(t) (H_3(t) - cp_4(t) T_4(t)) - cp_4(t) T_4(t)}{H_4(t) - cp_4(t) T_4(t)} \quad (4.76b)$$

$$+ \frac{F_5(t - \tau_2) \rho_5(t - \tau_2) [cp_5(t - \tau_2) T_5(t - \tau_2)]}{H_4(t) - cp_4(t) T_4(t)}, \quad (4.76c)$$

$$V_5(t) = \frac{V_4(t) r_4(t) + V_2(t) (H_2(t) - cp_5(t) T_5(t))}{H_5(t) - cp_5(t) T_5(t)} \quad (4.76d)$$

$$+ \frac{F_6(t - \tau_2) \rho_6(t - \tau_2) [cp_6(t - \tau_2) T_6(t - \tau_2) - cp_5(t) T_5(t)]}{H_5(t) - cp_5(t) T_5(t)}, \quad (4.76e)$$

$$\begin{aligned}
V_6(t) = & \frac{V_1(t)(H_1(t) - cp_6(t)T_6(t)) + V_5(t)r_5(t)}{H_6(t) - cp_6(t)T_6(t)} \\
& + \frac{F_7(t - \tau_2)\rho_7(t - \tau_2)[cp_7(t - \tau_2)T_7(t - \tau_2) - cp_6(t)T_6(t)]}{H_6(t) - cp_6(t)T_6(t)} \\
& + \frac{V_5(t)r_5(t) + F_{01}(t - \tau_2)\rho_0(t - \tau_2)[cp_0(t - \tau_2)T_0(t - \tau_2) - cp_6(t)T_6(t)]}{H_6(t) - cp_6(t)T_6(t)}, \quad (4.76f)
\end{aligned}$$

$$V_7(t) = \frac{V_6(t)r_6(t) + F_0(t - \tau_2)\rho_0(t - \tau_2)[cp_0(t - \tau_2)T_0(t - \tau_2) - cp_7(t)T_7(t)]}{H_7(t) - cp_7(t)T_7(t)}, \quad (4.76g)$$

where  $H$  and  $r$  are, respectively, the enthalpy and the latent heat of the output vapor;  $V_0$  is the flow rate of the live steam;  $r_0$  is the latent heat of the live steam; and  $\rho_0$ ,  $cp_0$ ,  $T_0$  are, respectively, the density, specific heat capacity, and temperature of the feed injected to the process. It is assumed that the vapor is saturated. The relationships between the latent heat and the enthalpy are obtained in [21] as given below:

$$r_i = 2495.0 - 2.219T_{vi} - 0.002128(T_{vi})^2, \quad i = 1, \dots, 7, \quad (4.77)$$

$$H_i = r_i + 4.18T_{vi}, \quad i = 1, \dots, 7, \quad (4.78)$$

where

$$T_{vi} = T_i - \Delta T_i, \quad i = 1, \dots, 7. \quad (4.79)$$

The density and the boiling point rise are important properties that must be specified in a multiple-effect evaporator [46], especially, when the soluble solid concentration is high. In order to obtain their relationships, several thermal balance tests were carried out during the whole acid cycle of the evaporation process. Correlations describing the relations between the boiling point rise and the density are determined by the regression method using the experimental data. For each  $i$ ,  $i = 1, \dots, 7$ , they are obtained as follows:

$$\rho_i = 1045 + 1.2C_i^1 + C_i^2 + 0.8C_i^3, \quad (4.80)$$

$$\begin{aligned}
\Delta T_i = & \frac{0.0162(T_i + 273)^2(75.77C_i^2\rho^{-1} - 3.608)}{r_i} \\
& - 0.23C_i^1 - 0.073C_i^2 - 0.1094C_i^3 + 0.3206T_i, \quad (4.81)
\end{aligned}$$

where  $\Delta T_i$ , for each  $i = 1, 2, 3$ , is the boiling point rise of the solution output from the  $i$ th flash evaporator; and  $\Delta T_i$ , for each  $i = 4, \dots, 7$ , is the boiling point rise of the solution output from the  $(i - 3)$ th evaporator.

The specific heat of the solution at each of these evaporators can be calculated from its component concentration as follows:

$$cp_i = 4.18 - \rho^{-1}(2.994C_i^1 + 2.923C_i^2 + 3.266C_i^3), \quad i = 1, \dots, 7 \quad (4.82)$$

Furthermore, the heat changes in the evaporation vessels are:

$$\begin{aligned} \Delta Q_i(t) &= F_{i+1}(t - \tau_1)\rho_{i+1}(t - \tau_1)cp_{i+1}(t - \tau_1)T_{i+1}(t - \tau_1) \\ &\quad - V_i(t)H_i(t) - F_i(t)\rho_i(t)cp_i(t)T_i(t), \quad i = 1, 2, 3, \end{aligned} \quad (4.83a)$$

$$\begin{aligned} \Delta Q_4(t) &= V_0(t)r_0(t) + V_3(t)H_3(t) + F_5(t - \tau_2)\rho_5(t - \tau_2)cp_5(t - \tau_2)T_5(t - \tau_2) \\ &\quad - V_4(t)H_4(t) - F_4(t)\rho_4(t)cp_4(t)T_4(t), \end{aligned} \quad (4.83b)$$

$$\begin{aligned} \Delta Q_5(t) &= V_4(t)r_4(t) + V_2(t)H_2(t) + F_6(t - \tau_2)\rho_6(t - \tau_2)cp_6(t - \tau_2)T_6(t - \tau_2) \\ &\quad - V_5(t)H_5(t) - F_5(t)\rho_5(t)cp_5(t)T_5(t), \end{aligned} \quad (4.83c)$$

$$\begin{aligned} \Delta Q_6(t) &= V_5(t)r_5(t) + V_1(t)H_1(t) + F_7(t - \tau_2)\rho_7(t - \tau_2)cp_7(t - \tau_2)T_7(t - \tau_2) \\ &\quad - V_6(t)H_6(t) - F_6(t)\rho_6(t)cp_6(t)T_6(t) \\ &\quad + F_{01}(t - \tau_2)\rho_0(t - \tau_2)cp_0(t - \tau_2)T_0(t - \tau_2), \end{aligned} \quad (4.83d)$$

$$\begin{aligned} \Delta Q_7(t) &= V_6(t)r_6(t) + F_0(t - \tau_2)\rho_0(t - \tau_2)cp_0(t - \tau_2)T_0(t - \tau_2) \\ &\quad - V_7(t)H_7(t) - F_7(t)\rho_7(t)cp_7(t)T_7(t). \end{aligned} \quad (4.83e)$$

Let

$$\mathbf{x} = [x_1, \dots, x_{35}]^\top = [h_1, \dots, h_7, C_1^1, \dots, C_7^1, C_1^2, \dots, C_7^2, C_1^3, \dots, C_7^3, T_1, \dots, T_7]^\top \in \mathbb{R}^{35},$$

denote the state with 35 variables. Taking into account the thermophysical properties of the solution and vapor, and substituting (4.71)-(4.83e) into (4.68) to (4.70), the evaporation system model can be expressed as a system of 35 differential equations. The split flow rate of feed, product flow rates of each of the evaporators and flash evaporators, and the live steam flow rate are the control variables denoted as  $\mathbf{u} = [u_1, \dots, u_9]^\top = [F_1, \dots, F_7, F_{01}, V_0]^\top \in \mathbb{R}^9$ . Let  $\mathbf{f} = [f_1, \dots, f_{35}]^\top$ , where  $f_i$ ,  $i = 1, \dots, 35$ , denote the functions appeared on the right-hand sides of (4.68)-(4.70).

Let this system model be referred to as System (S1).

#### B Initial conditions for the evaporation system

For System (S1), the initial values for the state and the control variables at and prior to  $t = 0$  are obtained from a real-life evaporation process of an alumina production factory in China. Specifically, for each of the temperature  $T_i$ ,  $i = 1, \dots, 7$ , it has relatively little fluctuation. However, the values of the seven temperatures  $T_i$ ,  $i = 1, \dots, 7$ , are quit different. Thus, we shall minus each of the temperature  $T_i$ ,  $i = 1, \dots, 7$ , by 97.7, 106.15, 116.1, 131.2, 105.6, 75, 54.5, respectively. Then, mark the values after treatments in Figure 4.3(a). The values of the levels are marked in Figure 4.3(b). Let  $\{\zeta_k\}_{k=1}^K$  with  $\zeta_k < \zeta_{k+1}$ ,  $k = 1, \dots, K - 1$ , be the set of the observed time points on the interval  $[-15, 0]$ . The value of the  $i$ th state variable  $x_i$  at the observed time point  $\zeta_k$  is denoted as  $\phi_i^k$ . The

temperatures and the liquor levels at and prior to  $t = 0$  are calculated according to

$$x_i(t) = \phi_i^k + \frac{\phi_i^k - \phi_i^{k+1}}{\zeta_k - \zeta_{k+1}}(t - \zeta_k), \quad (4.84)$$

$$t \in [\zeta_k, \zeta_{k+1}], \quad k = 1, \dots, K - 1; \quad i = 1, \dots, 7, 29, \dots, 35.$$

In addition, the concentrations of the solution  $x_i$ ,  $i = 8, \dots, 28$ , can only be accessed every two hours through analyzing sample solution collected from the practical evaporation process in the factory. It is assumed that the concentrations do not change during the time interval  $[-15, 0]$ . They are list below.

$$[x_8(t), \dots, x_{14}(t)]^\top = [74.85, 73.27, 71.75, 66.39, 54.86, 47.39, 44.83]^\top, \quad t \in [-15, 0],$$

$$[x_{15}(t), \dots, x_{21}(t)]^\top = [163.94, 161.07, 157.70, 145.94, 120.57, 104.17, 98.54]^\top, \quad t \in [-15, 0],$$

$$[x_{22}(t), \dots, x_{28}(t)]^\top = [75.61, 74.29, 72.73, 67.31, 55.61, 48.04, 45.45]^\top, \quad t \in [-15, 0].$$

The control in time horizon  $t \in [-15, 0]$  for System (S1) are given by

$$\boldsymbol{\psi}(t) = [\psi_1(t), \dots, \psi_9(t)]^\top = [2.54, 2.59, 2.64, 2.76, 3.44, 3.99, 4.22, 0.165, 986.5]^\top. \quad (4.85)$$

Let these initial state and control conditions obtained over the time interval  $[-15, 0]$  be referred to as Initial Condition (IC).

### 4.4.3 Optimal control problem formulation

The optimal control of the evaporation process is to find a control such that the tasks listed below are accomplished.

- (i) The energy usage is minimized.
- (ii) The specific requirements of the industrial sodium aluminate solution are met.
- (iii) The solution level for each effect is maintained to within its operation limits.

The energy usage is measured in terms of the mass units of live steam used for evaporating one mass unit of water. It is defined by the following equation:

$$J_0 = V_0(t)/W(t) = u_9(t)/W(t),$$

where  $W(t)$  is the total water evaporated for the evaporation process at time  $t$ . It is the difference between the mass flow rate of the feed and the final product, and is calculated

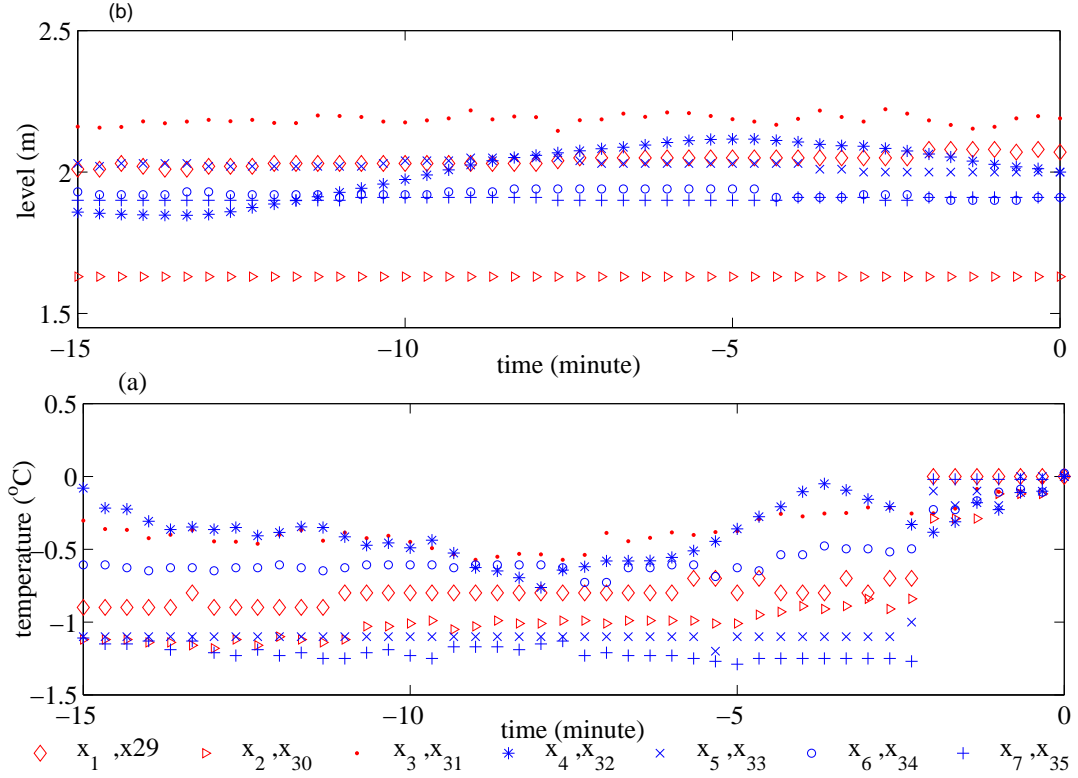


Figure 4.3: The variations of temperatures and liquor levels at observed time points

as given below:

$$\begin{aligned} W(t) &= (F_0(t) + F_{01}(t))\rho_0(t) - F_1(t)\rho_1(t) \\ &= (F_0(t) + F_{01}(t))\rho_0(t) - u_1(t)(1045 + 1.2x_8(t) + x_{15}(t) + 0.8x_{22}(t)). \end{aligned}$$

Thus, the cost function to be minimized is

$$\begin{aligned} J &= \Phi_0(\mathbf{x}(t_f|\mathbf{u})) \\ &+ \int_0^{t_f} \left\{ \frac{u_9(t)}{(F_0(t) + F_{01}(t))\rho_0(t) - u_1(t)(1045 + 1.2x_8(t) + x_{15}(t) + 0.8x_{22}(t))} \right\}^2 dt, \end{aligned} \quad (4.86)$$

where  $t_f$  is the final time of the time horizon  $[0, t_f]$ . Note that the change of the cost function value caused by the change of the product concentration takes about 75 minutes to be evaluated. Thus, the final time  $t_f$  should be much larger than 75.  $\Phi_0$  is the terminal cost given by

$$\Phi_0(\mathbf{x}(t_f|\mathbf{u})) = \sum_{i=1}^7 (\hat{x}_i - x_i(t_f|\mathbf{u}))^2, \quad (4.87)$$

where  $\hat{x}_i$ ,  $i = 1, \dots, 7$ , are specified desired solution levels, which can be determined according to experience. One way to specify these desired solution levels is to solve the optimal control problem with  $\Phi_0$  deleted initially. Then, the time when the concentration of the product solution has become stable is identified. The values of the solution levels at this particular time instant are chosen as these specified desired solution levels. We then re-solve the optimal problem with  $\Phi_0$  being included. The inclusion of  $\Phi_0$  is for regulating the solution levels toward the end of the time horizon.

To proceed further, the bounds on the state and control variables are specified through analyzing the production data as well as utilizing the experience of the operators from the evaporation process of an alumina production factory in China. It is found that the solution levels of the flash evaporators are less important than those of the evaporators. Moreover, according to the production data, the solution levels of the flash evaporators must be in the range of 1.5~2.5 m; the solution level of the first evaporator is limited to 1.8~2.3 m; and the solution levels of the last three evaporators are constrained to lie between 1.9 m and 2.1 m; and the sodium hydroxide concentration of the final product must reach 160~170 g/L. For the control variables, the flow rate of the solution should be operated within 1.6~5.3 m<sup>3</sup>/min; the split flow rate of the feed input into the third evaporator is allowed to vary between 0 m<sup>3</sup>/min to 0.6 m<sup>3</sup>/min. The live steam supplied to the process can be adjusted within 566 kg/min to 1230 kg/min. Let us write down explicitly the bounds for these variables as follows:

$$\mathbf{a} = [a_1, \dots, a_8]^\top = [2.5, 2.5, 2.5, 2.3, 2.1, 2.1, 2.1, 170]^\top, \quad (4.88a)$$

$$\mathbf{b} = [b_1, \dots, b_8]^\top = [1.5, 1.5, 1.5, 1.8, 1.9, 1.9, 1.9, 160]^\top, \quad (4.88b)$$

$$\mathbf{c} = [c_1, \dots, c_9]^\top = [5.3, 5.3, 5.3, 5.3, 5.3, 5.3, 5.3, 0.6, 1230]^\top, \quad (4.88c)$$

$$\mathbf{d} = [d_1, \dots, d_8]^\top = [1.6, 1.6, 1.6, 1.6, 1.6, 1.6, 1.6, 0, 566]^\top, \quad (4.88d)$$

where  $a_i$  and  $b_i$ , for each  $i = 1, \dots, 7$ , are the upper and lower bounds for the level of the  $i$ th equipment, respectively;  $a_8$  and  $b_8$  are the upper and lower bounds for the product concentration, respectively;  $d_l$  and  $c_l$ , for each  $l = 1, \dots, 9$ , are the lower and upper bounds for the level of the  $l$ th control, respectively. The continuous inequality constraints on the states and controls may now be stated explicitly as follows:

$$b_i \leq x_i(t) \leq a_i, \quad i = 1, \dots, 7, \quad t \in [0, t_f], \quad (4.89a)$$

$$b_8 \leq x_{15}(t) \leq a_8, \quad t \in [0, t_f], \quad (4.89b)$$

$$d_l \leq u_l(t) \leq c_l, \quad l = 1, \dots, 9, \quad t \in [0, t_f]. \quad (4.89c)$$

Any measurable function  $\mathbf{u} = [u_1, \dots, u_9]^\top : [0, t_f] \rightarrow \mathbb{R}^9$  such that the constraints (4.89c) are satisfied is called an admissible control. Let  $\mathcal{U}$  be the set which consists of all such admissible controls.

The optimal control problem may now be stated formally below.

**Problem (Q1).** *Given System (S1) with Initial Conditions (IC), find a control  $\mathbf{u} \in \mathcal{U}$  such that the cost functional (4.86) is minimized subject to the continuous inequality constraints on the states given by (4.89a)-(4.89b).*

For Problem (Q1), the dimension of the state variables is 35 and the differential equations of the dynamics are nonlinear with multiple delays. Furthermore, there are eight continuous inequality constraints on the state variables which are not allowed to be violated at any time point during the time horizon. The final time  $t_f$  is taken as 8 hours, which is rather long. It does not appear that the nonlinear model predictive control (NMPC) technique could be applied directly due to the complexity of Problem (Q1). Nonetheless, Problem (Q1) will be solved by using the NMPC. As expected, the computational time is much too long for it to be used in real operation. Thus, we shall make use of the control parameterization method and apply the proposed optimal control method to solve Problem (Q1). We shall also check the robustness of the optimal control obtained. Furthermore, the results will be compared with the data collected from the real plant to ensure proper operation of the process and those obtained by NMPC. Under normal operation, the problem will be re-solved 3 hours before the end of the 8 hours period so that a new optimal control can be obtained and used for the next 8 hours period.

#### 4.4.4 Numerical results

##### A Numerical calculation

Consider Problem (Q1), i.e., the optimal control problem with its dynamical system, initial condition, and the cost function described by System (S1), Initial Condition (IC) and (4.86), respectively. Clearly,

$$\mathcal{L}_0(t, \mathbf{x}(t), \mathbf{u}) = \left[ \frac{u_9(t)}{(F_0(t) + F_{01}(t))\rho_0(t) - u_1(t)(1045 + 0.8x_{22}(t) + 1.2x_8(t) + x_{15}(t))} \right]^2,$$

$$\Phi_0(\mathbf{x}(t_f|\mathbf{u})) = \sum_{i=1}^7 (\hat{x}_i - x_i(t_f|\mathbf{u}))^2.$$

By using the control parameterization technique and the constraints transformation method described in Sections 4.3.1 and 4.3.2, the augmented cost function is:

$$\tilde{J}_{\varepsilon, \gamma}(\boldsymbol{\sigma}^N) = \Phi_0(\mathbf{x}(t_f|\boldsymbol{\sigma}^N)) + \int_0^{t_f} \left\{ \mathcal{L}_0(t) + \gamma \sum_{i=1}^{16} [\mathcal{L}_{i, \varepsilon}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N))] \right\} dt \quad (4.90)$$



Here,  $\mathcal{L}_{i,\varepsilon}(t, \mathbf{x}(t|\boldsymbol{\sigma}^N))$ ,  $i = 1, \dots, 16$ , are specified in (4.19), where  $h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N))$  are:

$$\begin{aligned} h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) &= x_i(t|\boldsymbol{\sigma}^N) - a_i, \quad i = 1, \dots, 7, \\ h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) &= b_{i-7} - x_{i-7}(t|\boldsymbol{\sigma}^N), \quad i = 8, \dots, 14, \\ h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) &= x_{15}(t|\boldsymbol{\sigma}^N) - a_{i-7}, \quad i = 15, \\ h_i(t, \mathbf{x}(t|\boldsymbol{\sigma}^N)) &= b_{i-8} - x_{15}(t|\boldsymbol{\sigma}^N), \quad i = 16, \end{aligned}$$

where  $a_i$  and  $b_i$ ,  $i = 1, \dots, 8$ , are given in (4.88a)-(4.88b).

Define the Hamiltonian function as:

$$\begin{aligned} H_{\varepsilon,\gamma} &= \mathcal{L}_0(t, \mathbf{x}(t), \mathbf{u}^N(t)) \\ &+ \sum_{k=1}^2 [(\bar{\boldsymbol{\lambda}}^k(t))^\top \bar{\mathbf{f}}_k e(t_f - t + \tau_k)] + \gamma \sum_{i=1}^{16} [\mathcal{L}_{i,\varepsilon}(t, \mathbf{x}(t|\mathbf{u}^N))] + (\boldsymbol{\lambda}(t))^\top \mathbf{f}(t), \end{aligned}$$

where  $\bar{\boldsymbol{\lambda}}^k$  and  $\bar{\mathbf{f}}_k$ ,  $k = 1, 2$ , are defined by (4.44) and (4.45), and  $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_{35}]^\top \in \mathbb{R}^{35}$  is the solution of the co-state system defined by

$$\frac{d\boldsymbol{\lambda}(t)^\top}{dt} = -\frac{\partial\{\mathcal{L}_0(t) + \gamma \sum_{i=1}^{16} \mathcal{L}_{i,\varepsilon}(t, \mathbf{x}(t|\mathbf{u}^N))\}}{\partial \mathbf{x}} + (\boldsymbol{\lambda}(t))^\top \frac{\partial \mathbf{f}}{\partial \mathbf{x}} + \sum_{k=1}^2 (\bar{\boldsymbol{\lambda}}^k(t))^\top \frac{\partial \bar{\mathbf{f}}_k}{\partial \mathbf{x}} e(t_f - t + \tau_k),$$

with terminal conditions

$$\begin{aligned} \boldsymbol{\lambda}(t_f) &= \frac{\partial \Phi_0(\mathbf{x}(t_f|\boldsymbol{\sigma}^N))^\top}{\partial \mathbf{x}}, \\ \boldsymbol{\lambda}(t) &= \mathbf{0}, \quad t > t_f. \end{aligned}$$

The gradient formula of the augmented cost function with respect to each component of the control parameter vector can be calculated by using Theorem 4.7, where

$$\begin{aligned} \frac{\partial H_{\varepsilon,\gamma}}{\partial u_1} &= \frac{2u_9(t)^2(1045 + 0.8x_{22}(t) + 1.2x_8(t) + x_{15}(t))}{[(F_0(t) + F_{01}(t))\rho_0(t) - u_1(t)(1045 + 0.8x_{22}(t) + 1.2x_8(t) + x_{15}(t))]^3} \\ &+ (\boldsymbol{\lambda}(t))^\top \frac{\partial \mathbf{f}}{\partial u_1} + \sum_{k=1}^2 (\bar{\boldsymbol{\lambda}}^k(t))^\top \frac{\partial \bar{\mathbf{f}}_k}{\partial u_1} e(t_f - t + \tau_k), \\ \frac{\partial H_{\varepsilon,\gamma}}{\partial u_l} &= (\boldsymbol{\lambda}(t))^\top \frac{\partial \mathbf{f}}{\partial u_l} + \sum_{k=1}^2 (\bar{\boldsymbol{\lambda}}^k(t))^\top \frac{\partial \bar{\mathbf{f}}_k}{\partial u_l} e(t_f - t + \tau_k), \quad l = 2, \dots, 8, \\ \frac{\partial H_{\varepsilon,\gamma}}{\partial u_9} &= \frac{2u_9(t)}{[(F_0(t) + F_{01}(t))\rho_0(t) - u_1(t)(1045 + 0.8x_{22}(t) + 1.2x_8(t) + x_{15}(t))]^2} \\ &+ (\boldsymbol{\lambda}(t))^\top \frac{\partial \mathbf{f}}{\partial u_9} + \sum_{k=1}^2 (\bar{\boldsymbol{\lambda}}^k(t))^\top \frac{\partial \bar{\mathbf{f}}_k}{\partial u_9} e(t_f - t + \tau_k), \end{aligned}$$

and the partial derivatives of the functions  $H_{\varepsilon,\gamma}$ ,  $\mathbf{f}$ ,  $\bar{\mathbf{f}}_1$ , and  $\bar{\mathbf{f}}_2$ , with respect to  $x_i$ ,  $i = 1, \dots, 35$ , are calculated by using Maple software. Now, by using Algorithm 4.1, the optimal control problem can be solved using any gradient-based optimization algorithm, where the gradient of  $\tilde{J}_{\varepsilon,\gamma}(\boldsymbol{\sigma}^N)$  with respect to  $\boldsymbol{\sigma}^N$  are normalized at each iteration of the optimization process.

### B Result and discussion

Simulation studies are performed using MATLAB on a computer with Intel Core 2 Quad Q9400 processor, where the final time  $t_f$  is taken as 480 minutes.

The optimal control problem is solved with the penalty factor taken as  $\gamma = 10^5$ . According to Remark 4.1, we choose  $N = 80$ , which means that the control is allowed to switch its value at every 6 minutes. Furthermore, the desired solution levels  $\hat{x}_i$ ,  $i = 1, \dots, 7$ , are chosen as 1.73, 2.25, 2.22, 2.14, 1.98, 2.03, and 1.98, respectively.

The disturbances that commonly affect the evaporation process are: disturbances due to the changes of the concentrations of feed; and the disturbances on the flow rate of the live steam. We now consider the same optimal control problem under the optimal control obtained. However, we assume that the feed concentration and the flow rate of the live steam are perturbed by a Gaussian noise with standard deviation of  $\pm 5\%$ .

For comparison, we shall use the model predictive control (MPC) [61,136] to construct the controller for the problem considered, where the objective function for MPC is given below:

$$J_M = 100 \sum_{k=1}^M \left\{ \frac{u_9(t_k)}{(F_0(t_k) + F_{01}(t_k))\rho_0(t_k) - u_1(t_k)(1045 + 0.8x_{22}(t_k) + 1.2x_8(t_k) + x_{15}(t_k))} \right\}^2 + 10 \sum_{k=1}^P (x_{15}(t_k) - x_{ref})^2 + \sum_{k=1}^M \Delta \mathbf{u}(t_k)^\top \mathbf{R} \Delta \mathbf{u}(t_k)$$

where  $\mathbf{R} = \text{diag}[1, 1, 1, 1, 1, 1, 1, 0.1, 1]$ .  $t_k$ ,  $k = 1, \dots, P$ , are  $k$ th sampling time. It is 0.1 hours.  $\Delta \mathbf{u}(t_k)$  is the change rate of the control at time  $t_k$ .  $M = 10$  is the control horizon,  $P = 13$  is the predict horizon.  $x_{ref}$  is the reference trajectory, which is taken as 162.5—the concentration of the sodium hydroxide achieved by the optimal control at the end of the simulation time. Furthermore, the states and the controls are required to satisfy constraints (4.89a)-(4.89b).

Consider the problem with disturbance as described above. Figure 4.4 shows the results for this problem under the optimal control, the level controller used in the current practice, and the MPC. These controls are depicted in Figure 4.5. It takes about 3 hours to accomplish the optimal control calculation. The computational time of MPC for each predicted horizon is about 20 minute, which is significantly larger than the sampling time 6 minutes. The total computational time of MPC is over 5 times longer when compared with the optimal control method proposed in this chapter. Detailed comparisons between

the results obtained by MPC, the level controller and those obtained by the optimal control method are as follows.

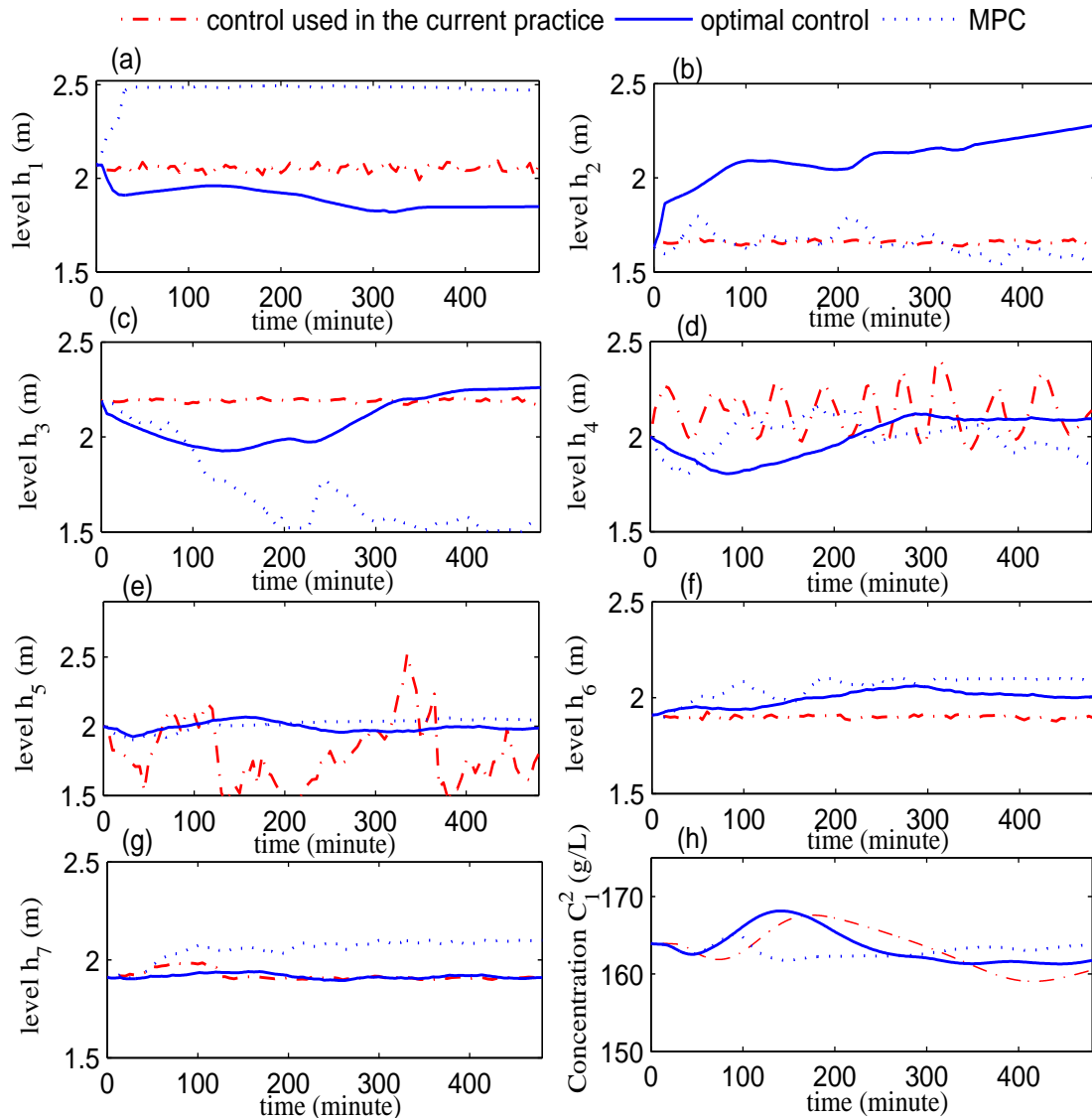


Figure 4.4: State of the evaporation process under the disturbances of the feed concentration and the live steam flow rate

The energy usage value (the mass units of live steam used for evaporating one mass unit of water) obtained by the proposed optimal control is 0.368. With disturbances, the energy usage is 0.369. The energy usage obtained by the MPC controller is 0.371. It is 0.38 under the level controller used in current actual operation.

Figures 4.4(a)-4.4(g) show the changes of the solution levels. Although the solution levels touch the permitted bounds at certain time points, both the optimal control and the MPC drive the solution levels towards inside of the permitted ranges as the simulation time increases. The results obtained by the optimal control have much less oscillations in the solution levels when compared with those obtained using the level controller in the

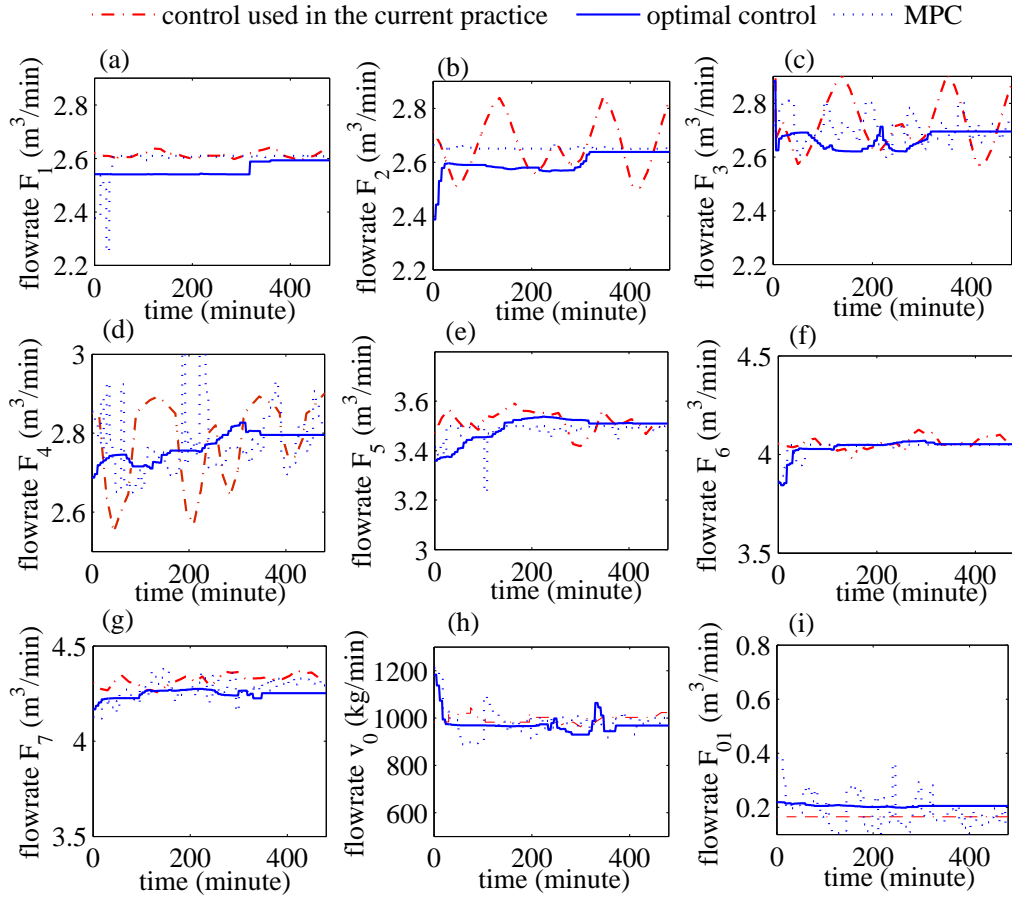


Figure 4.5: Control of the evaporation process

current practice.

In Figure 4.4(h), the plot of the product concentrations shows that the solution concentration is above 160g/L and below 170g/L during the whole time horizon. It means that the optimal control achieves disturbance rejection while maintaining the quality of product concentrations when disturbances occur in the feed concentration and live steam. Furthermore, we can see that only slight fluctuation in the concentration of the final product is observed after 250 minutes. The concentration obtained by using the MPC takes more than 8 hours to approach the desired value.

The live steam flow rate is as shown in Figure 4.5(h). By using the optimal control, the total live steam consumption is 466.96T. Under disturbances, it is 468.50T. On the other hand, by using the level controller in the current practice, the total live steam consumption is 481.76T. By using MPC, the live steam consumption is 471.5T. This represents a significant reduction of live steam consumption being achieved by using the optimal control, even in the presence of disturbances. The two main reasons are: (i) The optimal control improves the performance of the flash evaporator, and (ii) the live steam consumption and the final concentration are involved in the objective function which

is being minimized. Thus, unnecessary live steam usage is reduced while achieving the quality of the final product concentration.

As shown in Figure 4.5, the changes of all the variables of the optimal control and MPC vary strictly inside their bounds. This is due to the imposed continuous inequality constraints on the levers. It is a useful feature in practice, as there are rooms for adjusting the pumps or valves with dead-zone of the alumina evaporation process. The variations of the optimal control variables are much less when compared with those of the MPC.

## 4.5 Conclusion

In this chapter, we consider a time-delayed optimal control problem subject to continuous state inequality constraints and control constraints. This optimal control problem arises from practical production processes, where the systems are time-delayed dynamic systems. The objective is to find an admissible control such that the energy consumption or the material consumption are minimized, while practical limitations and engineering specifications, which are expressed as continuous inequality constraints on the state variables and the control constants, are satisfied. An efficient numerical algorithm is developed based on the control parameterization technique for solving this constrained time-delayed optimal problem. From solving an optimal control problem in a practical evaporation process, it is observed that the results obtained by the optimal control are superior to those obtained by MPC controller and the controller used in the current practice.

---

---

# CHAPTER 5

---

## A max-min control problem arising in gradient elution chromatography

### 5.1 Introduction

Chromatography plays an important role as a separation and purification process in many industrial settings, especially in the preparation of biochemical and pharmaceutical products. A typical chromatographic system shown in Figure 5.1 consists of a column containing an absorbent called the stationary phase, and a liquid that flows through the column called the mobile phase. The absorbent is fixed in the column. The mixture to be separated is injected into the mobile phase and flows through the column. Because different components in the mixture are attracted to the stationary phase at varying degrees, they travel through the column at different speeds, and thus they exit the column at different times (called peak times or retention times). Therefore, the mixture is gradually separated while moving through the column. The separated components are analyzed by a detector at the outlet of the column. The chromatography signal is shown in Figure 5.2.

In practice, chromatography is time-consuming, and the purity of the final product must satisfy strict conditions. Thus, it is essential that the chromatographic process be controlled judiciously by varying mobile phase conditions such as pH value, ionic strength, and flow rate. To achieve a high-quality separation and improve productivity, the minimum duration between successive retention times should be maximized.

In this chapter, we consider an optimal control problem in which manipulative variables in the chromatographic process need to be determined to maximize separation efficiency and achieve optimum separation. This problem has been formulated as a max-min optimal control problem in [90]. It has two non-standard characteristics: (i) The objective function is non-smooth; and (ii) each state variable is defined over a different time horizon. The final time for each state variable, the so-called retention time, is not fixed and actually depends on the control variables. A computational method for solving this problem, based on the control parameterization technique is proposed in [80]. This method involves reformulating the max-min objective function into a more convenient form, then

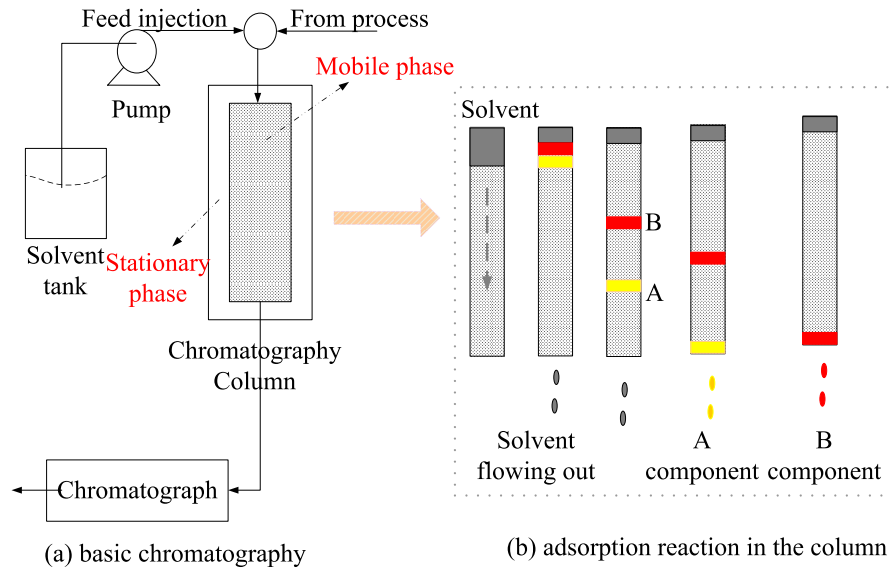


Figure 5.1: Basic chromatography and the reaction principle

approximating the control by a piecewise-constant function, before finally transforming the time horizon into the fixed interval  $[0, 1]$ . This yields an approximate mathematical programming problem that can be solved using existing optimization algorithms.

The time transformation method used in [90] is based on the substitution  $t = st_f$ , where  $t \in [0, t_f]$  is the original time variable,  $s \in [0, 1]$  is the new time variable, and  $t_f$  is the final time of the process. The same transformation has been successfully applied to solve time-optimal control problems in [13, 67]. When applied to the chromatography optimal control problem, this transformation does not map the retention times to fixed points, only the final time. Furthermore, the retention times do not necessarily coincide with the control switching times in the new time horizon. In fact, the retention times remain variable under this transformation, which makes them very difficult to compute numerically. This is a major disadvantage, as studies show that the productivity of a chromatographic process is highly sensitive to the retention times [11, 98]. However, although difficult, accurate determination of the retention times is crucial for industrial applications [152].

In the chromatography optimal control problem, an equality state constraint is imposed at each retention time. Such constraints are called characteristic - time constraints in the optimal control literature [120, 121]. Computational methods for solving optimal control problems with characteristic-time constraints is developed in [120, 121]. These methods, however, assume that the ordering of the characteristic times is fixed and known. In the chromatography optimal control problem, the characteristic times are the retention times, and their ordering depends on the control variables. One approach that can be

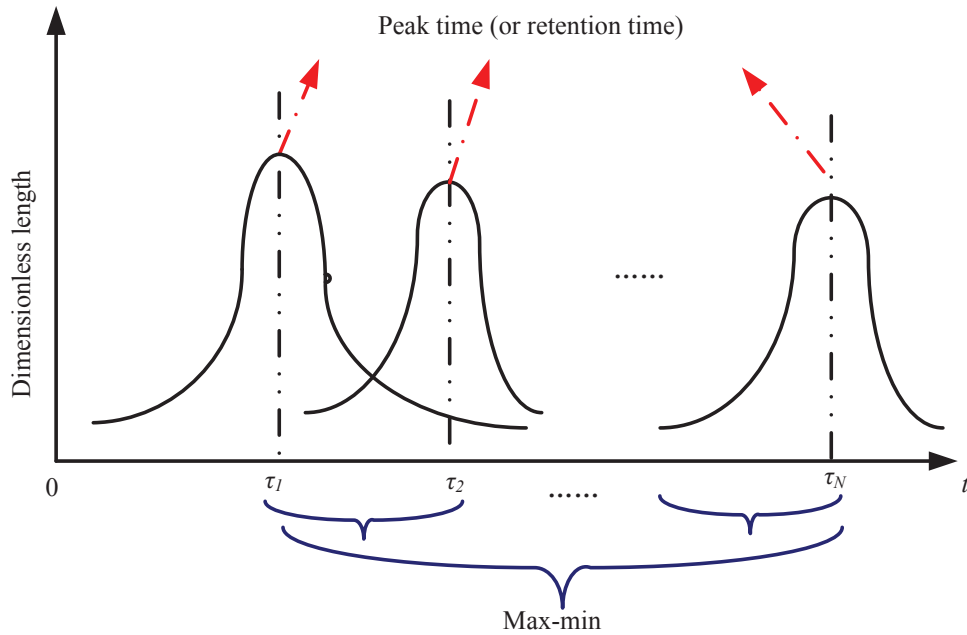


Figure 5.2: The chromatography signal

used to deal with this situation is the bilevel optimization approach [14], whereby the ordering is optimized in the outer level, and the control variables and retention times are optimized in the inner level. Another possible approach is the grid search algorithm [152], whereby the ordering of the retention times is determined after the controls are chosen. However, both the bilevel approach and the grid search algorithm are inefficient because they involve solving a computationally-intensive discrete optimization problem.

In this chapter, we consider the same chromatography optimal control problem formulated in [90]. We propose a new method for reformulating this problem that facilitates accurate determination of the retention times. First, a set of auxiliary decision variables are introduced to govern the ordering of the retention times. Then, after approximating the control variables by piecewise-constant functions, a novel time-scaling transformation is used to map the retention times to fixed points in a new time horizon. We then show that the max-min optimal control problem under consideration is equivalent to a minimization problem subject to additional inequality constraints. This minimization problem can be solved using an exact penalty method [41]. This method is then applied to solve two real world problems. The results show that the approach is both accurate and efficient.



## 5.2 Problem statement

A gradient elution chromatographic process with  $N$  components can be described by the following dynamical system of differential equations [90]:

$$\dot{x}_k(t) = f_k(t, x_k(t), \mathbf{u}(t)), \quad t > 0, \quad k = 1, \dots, N, \quad (5.1)$$

with initial conditions

$$x_k(0) = 0, \quad k = 1, \dots, N, \quad (5.2)$$

where  $x_k$  is the chromatography signal corresponding to the  $k$ th component,  $f_k$  is the signal velocity corresponding to the  $k$ th component, and  $\mathbf{u}$  is a vector representing mobile phase conditions such as pH value, ionic strength, temperature, and flow rate. In the language of control theory,  $\mathbf{x} = [x_1, \dots, x_N]^\top \in \mathbb{R}^N$  is called the state, and  $\mathbf{u} = [u_1, \dots, u_p]^\top \in \mathbb{R}^p$  is called the control. Each  $f_k : \mathbb{R} \times \mathbb{R} \times \mathbb{R}^p \rightarrow \mathbb{R}$  is assumed to be a given continuously differentiable function.

For each  $k = 1, \dots, N$ , the retention time  $\tau_k$  for the  $k$ th component is defined by the following equality constraint:

$$x_k(\tau_k) = \int_0^{\tau_k} f_k(t, x_k(t), \mathbf{u}(t)) dt = L_k, \quad k = 1, \dots, N, \quad (5.3)$$

where  $L_k$  is the peak height of the chromatography signal  $x_k$ . Thus, the state variable  $x_k$  is defined on the time horizon  $[0, \tau_k]$ .

We assume that the  $p$  control variables in the chromatographic process are bounded:

$$a_j \leq u_j(t) \leq b_j, \quad t \geq 0, \quad j = 1, \dots, p, \quad (5.4)$$

where  $a_j$  and  $b_j$  are the lower and upper bounds of the  $j$ th control variable, respectively. Given a piecewise continuous function  $\mathbf{u} : [0, \infty) \rightarrow \mathbb{R}^p$  satisfying (5.4), we can solve the system defined by (5.1) and (5.2) to yield a corresponding state trajectory.

To ensure that the gradient elution process is effective, the control variables should be chosen so that the minimum duration between successive retention times is maximized. Hence, we introduce the following objective function:

$$J(\mathbf{u}) = \min_{i \neq j} \left\{ \frac{(\tau_j - \tau_i)^2}{\tau_f} \right\}, \quad (5.5)$$

where  $\tau_f$  is the terminal time for the chromatographic process defined by

$$\tau_f = \max_{k=1, \dots, N} \{\tau_k\}.$$

We now introduce the following optimal control problem.

**Problem (P).** *Given the system defined by (5.1) and (5.2), choose a control  $\mathbf{u} : [0, \infty) \rightarrow \mathbb{R}^p$  and the corresponding retention times  $\tau_k, k = 1, \dots, N$ , such that the objective function (5.5) is maximized subject to the terminal constraints (5.3) and the control constraints (5.4).*

Compared with standard optimal control problems, Problem (P) has several unusual characteristics:

- (i) The max-min objective function is non-smooth.
- (ii) Each state variable is defined over a different time horizon.
- (iii) The retention times and the ordering of the retention times are not fixed, but instead depend on the control through (5.3).

Thus, Problem (P) is a highly non-standard optimal control problem and cannot be solved directly using standard methods such as the Pontryagin minimum principle [92] or state discretization [38, 148].

### 5.3 Problem transformation

Let  $\boldsymbol{\tau} = [\tau_1, \dots, \tau_N]^\top$  be a vector containing the retention times. Furthermore, let  $v_{ij}, i = 1, \dots, N; j = 1, \dots, N$ , be a set of auxiliary decision variables controlling the order of the retention times, where

$$v_{ij} = \begin{cases} 1, & \text{if component } j \text{ has the } i\text{th earliest retention time,} \\ 0, & \text{otherwise.} \end{cases} \quad (5.6)$$

Clearly,

$$\sum_{j=1}^N v_{ij} = 1, \quad i = 1, \dots, N, \quad (5.7)$$

and

$$\sum_{i=1}^N v_{ij} = 1, \quad j = 1, \dots, N. \quad (5.8)$$

We collect the auxiliary variables into a vector  $\mathbf{v} = [\mathbf{v}_1^\top, \dots, \mathbf{v}_N^\top]^\top \in \mathbb{R}^{NN}$ , where  $\mathbf{v}_i = [v_{i1}, \dots, v_{iN}]^\top$ . As an example, consider a 3-component mixture in which component 2 is the first component to exit the chromatography column, component 1 is the second, and component 3 is the last. Then  $\tau_2 < \tau_1 < \tau_3$ , and thus  $\mathbf{v}_1 = [0, 1, 0]^\top$ ,  $\mathbf{v}_2 = [1, 0, 0]^\top$ , and  $\mathbf{v}_3 = [0, 0, 1]^\top$ .

Standard gradient-based optimization techniques cannot handle binary constraints such as (5.6). Thus, we replace the 0-1 constraints on  $v_{ij}$  by the following constraints:

$$\sum_{j=1}^N v_{ij}(j^2 - j + \frac{1}{3}) - \left\{ \sum_{j=1}^N v_{ij}(j - \frac{1}{2}) \right\}^2 = \frac{1}{12}, \quad i = 1, \dots, N, \quad (5.9)$$

and

$$0 \leq v_{ij} \leq 1, \quad i = 1, \dots, N, \quad j = 1, \dots, N. \quad (5.10)$$

The following result shows that (5.9) and (5.10) imply  $v_{ij} \in \{0, 1\}$ .

**Theorem 5.1.** *Suppose that  $v_{ij}$ ,  $i = 1, \dots, N$ ;  $j = 1, \dots, N$ , satisfy (5.7) and (5.10). Then, for each  $i = 1, \dots, N$ , (5.9) holds if and only if there exists a  $k \in \{1, \dots, N\}$  such that  $v_{ik} = 1$  and  $v_{ij} = 0$  for all  $j \neq k$ .*

*Proof.* Let  $i \in \{1, \dots, N\}$  be fixed but arbitrary, and assume that  $v_{ik} = 1$  and  $v_{ij} = 0$  for all  $j \neq k$ . Then

$$\sum_{j=1}^N v_{ij}(j^2 - j + \frac{1}{3}) - \left\{ \sum_{j=1}^N v_{ij}(j - \frac{1}{2}) \right\}^2 = k^2 - k + \frac{1}{3} - (k - \frac{1}{2})^2 = \frac{1}{12}. \quad (5.11)$$

To prove the opposite implication, we use similar arguments to those used in the proof of Lemma 3.1 in [65]. First, consider the following optimization problem, which we call Problem (Q):

$$\min_{v_{i1}, \dots, v_{iN}} \sum_{j=1}^N v_{ij}(j^2 - j + \frac{1}{3}) - \left\{ \sum_{j=1}^N v_{ij}(j - \frac{1}{2}) \right\}^2 \quad (5.12)$$

$$\text{s.t.} \quad \sum_{j=1}^N v_{ij} = 1, \quad (5.13)$$

$$v_{ij} \geq 0, \quad j = 1, \dots, N. \quad (5.14)$$

Note that Problem (Q) has a continuous objective function and a compact feasible region. Hence, it admits at least one optimal solution. The Lagrangian for Problem (Q) is defined by

$$\mathcal{L} = \sum_{j=1}^N v_{ij}(j^2 - j + \frac{1}{3}) - \left\{ \sum_{j=1}^N v_{ij}(j - \frac{1}{2}) \right\}^2 - \rho \left\{ \sum_{j=1}^N v_{ij} - 1 \right\} - \sum_{j=1}^N \eta_j v_{ij},$$

where  $\rho$  is the Lagrange multiplier for the equality constraint, and  $\eta_j \geq 0$ ,  $j = 1, \dots, N$ , are the Lagrange multipliers for the inequality constraints.

Since the linear independence constraint qualification is clearly satisfied in Prob-

lem (Q), any optimal solution must satisfy the following Kuhn-Tucker conditions [73]:

$$\frac{\partial \mathcal{L}}{\partial v_{ik}} = k^2 - k + \frac{1}{3} - 2 \left\{ \sum_{j=1}^N v_{ij} (j - \frac{1}{2}) \right\} (k - \frac{1}{2}) - \rho - \eta_k = 0, \quad (5.15)$$

$$k = 1, \dots, N,$$

and

$$\eta_k v_{ik} = 0, \quad k = 1, \dots, N. \quad (5.16)$$

From (5.15), we obtain

$$g(k) - \eta_k = 0, \quad k = 1, \dots, N,$$

where  $g : \mathbb{R} \rightarrow \mathbb{R}$  is a quadratic function defined by

$$\begin{aligned} g(y) &= y^2 - y + \frac{1}{3} - 2 \left\{ \sum_{j=1}^N v_{ij} (j - \frac{1}{2}) \right\} (y - \frac{1}{2}) - \rho \\ &= y^2 - \left[ 1 + 2 \left\{ \sum_{j=1}^N v_{ij} (j - \frac{1}{2}) \right\} \right] y + \left[ \frac{1}{3} + \sum_{j=1}^N v_{ij} (j - \frac{1}{2}) - \rho \right]. \end{aligned}$$

Clearly  $g(k) = 0$  if and only if  $\eta_k = 0$ . Since  $g(y)$  is a quadratic function with at most two real roots, no more than two of the multipliers for the inequality constraints in Problem (Q) are zero.

Suppose that  $\eta_k > 0$  for all  $k = 1, \dots, N$ . Then, by (5.16),  $v_{ik} = 0$  for all  $k = 1, \dots, N$ , contradicting (5.13) in Problem (Q). Thus, it suffices to consider the following two cases:

- (a) There exist an integer  $k_1$  such that  $\eta_{k_1} = 0$  and  $\eta_k > 0$  for all  $k \neq k_1$  (exactly one of the multipliers is zero).
- (b) There exists integers  $k_1$  and  $k_2$  such that  $\eta_{k_1} = \eta_{k_2} = 0$  and  $\eta_k > 0$  for all  $k \neq k_1, k_2$  (exactly two of the multipliers are zero).

Consider Case (a). In this case, (5.16) and the equality constraint (5.13) in Problem (Q) imply that  $v_{ik_1} = 1$  and  $v_{ik} = 0$  for all  $k \neq k_1$ . Thus, as shown in (5.11), the optimal cost of Problem (Q) is equal to  $\frac{1}{12}$ .

We now consider Case (b). In this case, (5.16) implies that  $v_{ik} = 0$  for all  $k \neq k_1, k_2$ . Hence, from the equality constraint (5.13) in Problem (Q), we have  $v_{ik_2} = 1 - v_{ik_1}$ . Thus,

the optimal cost of Problem (Q) is

$$\begin{aligned}
& \sum_{j=1}^N v_{ij} \left( j^2 - j + \frac{1}{3} \right) - \left\{ \sum_{j=1}^N v_{ij} \left( j - \frac{1}{2} \right) \right\}^2 \\
&= v_{ik_1} \left( k_1^2 - k_1 + \frac{1}{3} \right) + v_{ik_2} \left( k_2^2 - k_2 + \frac{1}{3} \right) - \left\{ v_{ik_1} \left( k_1 - \frac{1}{2} \right) + v_{ik_2} \left( k_2 - \frac{1}{2} \right) \right\}^2 \\
&= v_{ik_1} \left( k_1^2 - k_1 + \frac{1}{3} \right) + (1 - v_{ik_1}) \left( k_2^2 - k_2 + \frac{1}{3} \right) \\
&\quad - \left\{ v_{ik_1} \left( k_1 - \frac{1}{2} \right) + (1 - v_{ik_1}) \left( k_2 - \frac{1}{2} \right) \right\}^2 \\
&= \frac{1}{12} + (v_{ik_1} - v_{ik_1}^2) (k_1 - k_2)^2 \geq \frac{1}{12},
\end{aligned}$$

where the last inequality follows from  $0 \leq v_{ij} \leq 1$ . Clearly, from Cases (a) and (b), the minimum value of the cost function in Problem (Q) is  $\frac{1}{12}$ . Furthermore, this minimum is achieved only when there exists a  $k \in \{1, \dots, N\}$  such that  $v_{ik} = 1$  and  $v_{ij} = 0$ ,  $j \neq k$ . This completes the proof.  $\square$

On the basis of Theorem 5.1, we can replace the binary constraints (5.6) by the non-discrete constraints (5.7), (5.9), and (5.10). As we will see later, this reformulation enables us to determine the optimal retention time ordering by using an exact penalty function method.

To proceed, we now use the control parameterization technique [80] to approximate Problem (P) by a finite-dimensional optimization problem. This is done by approximating the control  $\mathbf{u}$  by a piecewise-constant function that switches value at each retention time and at  $q-1$  times between each pair of successive retention times. Thus, the time horizon is divided into  $qN$  subintervals, with  $q$  subintervals between each pair of successive retention times, and the control is approximated by a constant value on each subinterval.

For each  $j = 1, \dots, p$ , the control  $u_j$  is approximated as follows:

$$u_j(t) = \sum_{i=1}^{qN} \sigma_j^i \chi_{[t_{i-1}, t_i]}(t), \quad t \in [0, \tau_f], \quad (5.17)$$

where  $t_i$ ,  $i = 0, \dots, qN$ , are the control switching times;  $\sigma_j^i$ ,  $i = 1, \dots, qN$ ;  $j = 1, \dots, p$ , are the control heights; and  $\chi_{[t_{i-1}, t_i]}$  is the indicator function for the subinterval  $[t_{i-1}, t_i]$  defined by

$$\chi_{[t_{i-1}, t_i]}(t) = \begin{cases} 1, & \text{if } t \in [t_{i-1}, t_i), \\ 0, & \text{otherwise.} \end{cases} \quad (5.18)$$

Here,  $0 \leq t_0 \leq t_1 \leq \dots \leq t_{qN} = \tau_f$ . Furthermore, for each  $k = 1, \dots, N$ , the switching time  $t_{kq}$  (the right end-point of subinterval  $[t_{kq-1}, t_{kq}]$ ) coincides with one of the retention times. Our objective is to choose the control heights and control switching times in (5.17) appropriately so that the objective function (5.5) is maximized. Note that the control

approximation scheme used in (5.17) is more flexible than the one used in [90], which does not allow the switching times to be determined optimally (instead they are pre-fixed).

It is well-known that treating the control switching times as decision variables causes major problems in numerical computation, see [67, 121]. Hence, we will use the so-called time-scaling transformation [81] to map the switching times to fixed points in a new time horizon. This involves introducing a new time variable  $s \in [0, qN]$ , and then relating  $s$  to  $t$  through the following differential equation:

$$\begin{aligned} \frac{dt(s)}{ds} &= \omega(s), \\ t(0) &= 0, \end{aligned} \quad (5.19)$$

where  $\omega : [0, qN] \rightarrow [0, \infty)$  is a piecewise-constant function with switching points at the fixed locations  $s = i$ ,  $i = 1, \dots, qN - 1$ . We express  $\omega$  mathematically as follows:

$$\omega(s) = \sum_{i=1}^{qN} \theta_i \chi_{[i-1, i)}(s), \quad (5.20)$$

where

$$\theta_i = t_i - t_{i-1} \geq 0, \quad i = 1, \dots, qN. \quad (5.21)$$

It follows from (5.19) that for  $s \in [l - 1, l]$ , we have

$$t(s) = \int_0^s \omega(\eta) d\eta = \sum_{j=1}^{l-1} \theta_j + \theta_l (s - l + 1).$$

Furthermore, for each  $l = 1, \dots, qN$ ,

$$t(l) = \sum_{j=1}^l \theta_j = \sum_{j=1}^l (t_j - t_{j-1}) = t_l.$$

This shows that the time-scaling transformation defined by (5.19) and (5.20) maps the control switching times to fixed integers.

After applying the time-scaling transformation, the control variables defined in (5.17) are written as:

$$\tilde{u}_j(s) = u_j(t(s)) = \sum_{i=1}^{qN} \sigma_j^i \chi_{[i-1, i)}(s), \quad j = 1, \dots, p, \quad (5.22)$$

where the control heights satisfy the following constraints:

$$a_j \leq \sigma_j^i \leq b_j, \quad j = 1, \dots, p, \quad i = 1, \dots, qN. \quad (5.23)$$

Define  $\boldsymbol{\sigma} = [(\boldsymbol{\sigma}^1)^\top, \dots, (\boldsymbol{\sigma}^{qN})^\top]^\top \in \mathbb{R}^{qNp}$ , where  $\boldsymbol{\sigma}^i = [\sigma_1^i, \dots, \sigma_p^i]^\top \in \mathbb{R}^p$ ,  $i = 1, \dots, qN$ .

Furthermore, let  $\boldsymbol{\theta} = [\theta_1, \dots, \theta_{qN}]^\top$ .

Under the time-scaling transformation, the system defined by (5.1) and (5.2) becomes

$$\begin{aligned} \frac{d\tilde{x}_k(s)}{ds} &= \sum_{i=1}^{qN} \theta_i f_k(t(s), \tilde{x}_k(s), \boldsymbol{\sigma}^i) \chi_{[i-1, i)}(s), \quad k = 1, \dots, N, \\ \frac{dt(s)}{ds} &= \omega(s), \end{aligned} \quad (5.24)$$

with initial conditions

$$\begin{aligned} \tilde{x}_k(0) &= 0, \quad k = 1, \dots, N, \\ t(0) &= 0. \end{aligned} \quad (5.25)$$

The objective function (5.5) becomes

$$\begin{aligned} \tilde{J}(\boldsymbol{\sigma}, \boldsymbol{\theta}, \mathbf{v}) &= \min_{i=1, \dots, N-1} \left\{ \frac{(t_{(i+1)q} - t_{iq})^2}{t_{qN}} \right\} \\ &= \min_{i=1, \dots, N-1} \frac{(\theta_{iq+1} + \theta_{iq+2} + \dots + \theta_{iq+q})^2}{\theta_1 + \theta_2 + \dots + \theta_{qN}}. \end{aligned} \quad (5.26)$$

For each  $i = 1, \dots, N$ , exactly one component in the chromatography system will have its retention time at  $t = t_{iq}$ . Hence, we have the following interior point constraints:

$$\sum_{k=1}^N v_{ik} (\tilde{x}_k(iq) - L_k) = 0, \quad i = 1, \dots, N. \quad (5.27)$$

In view of our discussion above, Problem (P) can be approximated by the following optimization problem.

**Problem (P1).** *Given the system defined by (5.24) and (5.25), find vectors  $\boldsymbol{\sigma} \in \mathbb{R}^{qNp}$ ,  $\boldsymbol{\theta} \in \mathbb{R}^{qN}$ , and  $\mathbf{v} \in \mathbb{R}^{NN}$ , such that the objective function (5.26) is maximized subject to (5.7)-(5.10), (5.21), (5.23) and interior point constraints (5.27).*

Although Problem (P1) is a finite-dimensional optimization problem, it is still difficult to solve because the objective function (5.26) is non-smooth. Thus, we introduce a new decision parameter  $\xi$ , where

$$\xi = \min_{i=1, \dots, N-1} \frac{(\theta_{iq+1} + \theta_{iq+2} + \dots + \theta_{iq+q})^2}{\theta_1 + \theta_2 + \dots + \theta_{qN}}.$$

Clearly, for each  $i = 1, \dots, N - 1$ , the following inequality is satisfied:

$$\frac{(\theta_{iq+1} + \theta_{iq+2} + \dots + \theta_{iq+q})^2}{\theta_1 + \theta_2 + \dots + \theta_{qN}} \geq \xi.$$

Since  $\theta_1 + \cdots + \theta_{qN} > 0$ , this inequality can be rewritten as:

$$\xi \sum_{k=1}^{qN} \theta_k - \left\{ \sum_{l=iq+1}^{iq+q} \theta_l \right\}^2 \leq 0, \quad i = 1, \dots, N-1. \quad (5.28)$$

Thus, Problem (P1) is equivalent to the following problem.

**Problem (P2).** *Given the system defined by (5.24) and (5.25), find vectors  $\boldsymbol{\sigma} \in \mathbb{R}^{qNp}$ ,  $\boldsymbol{\theta} \in \mathbb{R}^{qN}$ , and  $\mathbf{v} \in \mathbb{R}^{N^2}$ , and the parameter  $\xi$ , to minimize the cost function*

$$\hat{J}(\boldsymbol{\sigma}, \boldsymbol{\theta}, \mathbf{v}, \xi) = -\xi \quad (5.29)$$

subject to (5.7)-(5.10), (5.21), (5.23), (5.27) and (5.28).

## 5.4 A computational method

Problem (P2) is a finite-dimensional optimization problem with nonlinear equality and nonlinear inequality constraints. The major difficulty with solving this problem is that the constraints (5.7) - (5.10) force  $v_{ij}$  to be binary decision variables, and thus the feasible region for Problem (P2) is disjoint. Standard gradient-based optimization methods usually fail miserably when applied to problems with disjoint feasible regions. In this section, we will show how to solve Problem (P2) using an exact penalty function method.

First, let

$$g_m(\boldsymbol{\theta}, \xi) = \xi \sum_{k=1}^{qN} \theta_k - \left\{ \sum_{l=mq+1}^{mq+q} \theta_l \right\}^2, \quad m = 1, \dots, N-1.$$

Furthermore, for each  $m = 1, \dots, N$ , let

$$\begin{aligned} h_m(\mathbf{v}) &= \sum_{j=1}^N v_{mj} - 1, \\ \hat{h}_m(\mathbf{v}) &= \sum_{i=1}^N v_{im} - 1, \\ \tilde{h}_m(\mathbf{v}) &= \sum_{j=1}^N v_{mj} (j^2 - j + \frac{1}{3}) - \left\{ \sum_{j=1}^N v_{mj} (j - \frac{1}{2}) \right\}^2 - \frac{1}{12}. \end{aligned}$$

We will use an exact penalty function approach, recently developed in [41], to handle the equality constraints (5.7)-(5.9), and the inequality constraints (5.28). This approach has been successively used to solve semi-infinite and discrete optimization problems [40-42].



Consider the following penalty function:

$$\hat{J}_\eta(\boldsymbol{\sigma}, \boldsymbol{\theta}, \mathbf{v}, \xi, \varepsilon) = -\xi + \varepsilon^{-\alpha}(\Delta_1(\boldsymbol{\theta}, \xi, \varepsilon) + \Delta_2(\mathbf{v})) + \eta\varepsilon^\beta, \quad (5.30)$$

where  $\varepsilon > 0$  is a new decision variable and

$$\begin{aligned} \Delta_1(\boldsymbol{\theta}, \xi, \varepsilon) &= \sum_{m=1}^{N-1} \left\{ \max\{0, g_m(\boldsymbol{\theta}, \xi) - \varepsilon^\gamma \kappa_m\} \right\}^2, \\ \Delta_2(\mathbf{v}) &= \sum_{m=1}^N \left\{ (h_m(\mathbf{v}))^2 + (\hat{h}_m(\mathbf{v}))^2 + (\tilde{h}_m(\mathbf{v}))^2 \right\}. \end{aligned}$$

Here,  $\alpha > 0$ ,  $\beta > 2$ ,  $\gamma > 0$ , and  $\kappa_m \in (0, 1)$ ,  $m = 1, \dots, N-1$ , are fixed constants, and  $\eta > 0$  is the penalty parameter. Note that  $\Delta_1(\boldsymbol{\theta}, \xi, \varepsilon)$  measures violations in (5.28), while  $\Delta_2(\mathbf{v})$  measures violations in (5.7) - (5.9). The idea is that when the penalty parameter  $\eta$  is large, the final term in (5.30) forces  $\varepsilon$  to be small, which in turn causes the middle term in (5.30) to penalize constraint violations very severely. Hence, minimizing the penalty function will lead to a feasible point of Problem (P2). With this in mind, we introduce the following exact penalty function problem for Problem (P2).

**Problem (P3).** *Given the system defined by (5.24) and (5.25), find vectors  $\boldsymbol{\sigma} \in \mathbb{R}^{qNp}$ ,  $\boldsymbol{\theta} \in \mathbb{R}^{qN}$ ,  $\mathbf{v} \in \mathbb{R}^{NN}$ , and the parameters  $\xi$  and  $\varepsilon$ , such that the penalty function (5.30) is minimized subject to the bound constraints (5.10), (5.21), and (5.23), and the interior point constraints (5.27).*

Problem (P3) is an optimal parameter selection problem with bound constraints on decision parameters  $\boldsymbol{\sigma}$ ,  $\boldsymbol{\theta}$ ,  $\mathbf{v}$ ,  $\xi$ ,  $\varepsilon$  and interior point constraints (5.27). Unlike the non-smooth approximate problem derived in [90], Problem (P3) can be solved using standard computational methods such as those developed in Chapter 5 of [80]. These computational methods have been implemented in the optimal control software package MISER [91].

To solve Problem (P3), MISER requires gradient formulae for both the objective function and the constraint functions. The gradients of  $\hat{J}_\eta$  are given below:

$$\begin{aligned} \frac{\partial \hat{J}_\eta}{\partial \theta_i} &= 2\varepsilon^{-\alpha} \sum_{m=1}^{N-1} \left[ \max\{0, g_m(\boldsymbol{\theta}, \xi) - \varepsilon^\gamma \kappa_m\} \right] \left[ \xi - 2 \left\{ \sum_{k=mq+1}^{mq+q} \theta_k \right\} \chi_{[mq+1, mq+q]}(i) \right], \\ \frac{\partial \hat{J}_\eta}{\partial v_{ij}} &= 2\varepsilon^{-\alpha} \sum_{m=1}^N \left\{ h_m(\mathbf{v}) \frac{\partial h_m(\mathbf{v})}{\partial v_{ij}} + \hat{h}_m(\mathbf{v}) \frac{\partial \hat{h}_m(\mathbf{v})}{\partial v_{ij}} + \tilde{h}_m(\mathbf{v}) \frac{\partial \tilde{h}_m(\mathbf{v})}{\partial v_{ij}} \right\}, \\ \frac{\partial \hat{J}_\eta}{\partial \xi} &= -1 + 2\varepsilon^{-\alpha} \sum_{m=1}^{N-1} \left\{ \left[ \max\{0, g_m(\boldsymbol{\theta}, \xi) - \varepsilon^\gamma \kappa_m\} \right] \sum_{k=1}^{qN} \theta_k \right\}, \end{aligned}$$

$$\begin{aligned}\frac{\partial \hat{J}_\eta}{\partial \varepsilon} &= 2\varepsilon^{-\alpha} \sum_{m=1}^{N-1} \left\{ -\gamma \kappa_m \varepsilon^{\gamma-1} \left[ \max\{0, g_m(\boldsymbol{\theta}, \xi) - \varepsilon^\gamma \kappa_m\} \right] \right\} \\ &\quad - \alpha \varepsilon^{-\alpha-1} (\Delta_1(\boldsymbol{\theta}, \xi, \varepsilon) + \Delta_2(\mathbf{v})) + \beta \eta \varepsilon^{\beta-1}, \\ \frac{\partial \hat{J}_\eta}{\partial \sigma_j^i} &= 0,\end{aligned}$$

where

$$\frac{\partial \tilde{h}_m(\mathbf{v})}{\partial v_{ij}} = \begin{cases} (j^2 - j + \frac{1}{3}) - 2(j - \frac{1}{2}) \left\{ \sum_{j=1}^N v_{ij} (j - \frac{1}{2}) \right\}, & \text{if } i = m, \\ 0, & \text{if } i \neq m, \end{cases}$$

and

$$\frac{\partial h_m(\mathbf{v})}{\partial v_{ij}} = \begin{cases} 1, & \text{if } i = m, \\ 0, & \text{if } i \neq m, \end{cases} \quad \frac{\partial \hat{h}_m(\mathbf{v})}{\partial v_{ij}} = \begin{cases} 1, & \text{if } j = m, \\ 0, & \text{if } j \neq m. \end{cases}$$

We now consider the gradient of the interior point constraints (5.27). First, define

$$\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v}) = \sum_{k=1}^N v_{mk} (\tilde{x}_k(mq) - L_k), \quad m = 1, \dots, N. \quad (5.31)$$

That is,  $\Phi_m$  is the left-hand side of the interior point constraints (5.27) (with index  $i$  replaced by  $m$ ). The partial derivatives of  $\Phi_m$  with respect to  $\xi$ ,  $\varepsilon$ , and  $v_{ij}$  can be computed directly:

$$\frac{\partial \Phi_m}{\partial \varepsilon} = \frac{\partial \Phi_m}{\partial \xi} = 0,$$

and

$$\frac{\partial \Phi_m}{\partial v_{ij}} = \begin{cases} \tilde{x}_j(mq) - L_j, & \text{if } i = m, \\ 0, & \text{if } i \neq m. \end{cases}$$

The partial derivatives with respect to  $\boldsymbol{\sigma}$  and  $\boldsymbol{\theta}$ , however, are more difficult because  $\boldsymbol{\sigma}$  and  $\boldsymbol{\theta}$  influence  $\Phi_m$  implicitly through the dynamic system defined by (5.24) and (5.25). They are derived via formulas given below. Define the Hamiltonian function

$$H_m = \sum_{k=1}^N \lambda_k^m(s) \omega(s) f_k(t(s), \tilde{x}_k(s), \tilde{\mathbf{u}}(s)) + \lambda_{N+1}^m(s) \omega(s), \quad (5.32)$$

where the costate functions  $\lambda_k^m : \mathbb{R} \rightarrow \mathbb{R}$  satisfy the dynamic systems

$$\dot{\lambda}_k^m(s) = -\frac{\partial H_m}{\partial \tilde{x}_k}, \quad s \in [0, mq], \quad k = 1, \dots, N, \quad (5.33)$$

$$\lambda_k^m(mq) = v_{mk}, \quad (5.34)$$

and

$$\dot{\lambda}_{N+1}^m(s) = -\frac{\partial H_m}{\partial t}, \quad s \in [0, mq], \quad (5.35)$$

$$\lambda_{N+1}^m(mq) = 0. \quad (5.36)$$

For each  $m$ ,  $m = 1, \dots, N$ , we have the following results.

**Theorem 5.2.** *The partial derivatives of  $\Phi_m$  with respect to  $\sigma$  is*

$$\frac{\partial \Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial \sigma_j^i} = \begin{cases} \int_{i-1}^i \frac{\partial H_m}{\partial \sigma_j} ds, & \text{if } i \leq mq, \\ 0, & \text{if } i > mq, \end{cases} \quad (5.37)$$

*Proof.* Let the control parameter vector  $\sigma$  be perturbed by  $\epsilon \rho$ , where  $\epsilon > 0$  is a small real number and  $\rho$  is an arbitrary fixed perturbation of  $\sigma$ . Then, we have

$$\sigma(\epsilon) = \sigma + \epsilon \rho,$$

where  $\rho = [(\rho^1)^\top, \dots, (\rho^{qN})^\top]^\top$ , and  $\sigma(\epsilon) = [(\sigma^1(\epsilon))^\top, \dots, (\sigma^{qN}(\epsilon))^\top]^\top$ . Consequently, the system state  $\tilde{\mathbf{x}}$  as well as the function  $\Phi_m$  will be perturbed. Let

$$\tilde{x}_k(s, \epsilon) = \tilde{x}_k(s | \sigma(\epsilon)), \quad k = 1, \dots, N.$$

Clearly,

$$\tilde{x}_k(s, \epsilon) = \tilde{x}_k(0) + \sum_{i=1}^{qN} \int_0^s \theta_i f_k(t(l), \tilde{x}_k(l), \sigma^i(\epsilon))(l) dl.$$

Let  $f_k(t(\cdot), \tilde{x}_k(\cdot), \sigma^i(\epsilon))(\cdot)$  be written as  $f_k(\cdot)$  for brevity. Then, by the chain rule, we have

$$\Delta \tilde{x}_k(s) = \left. \frac{d\tilde{x}_k(s, \epsilon)}{d\epsilon} \right|_{\epsilon=0} = \sum_{i=1}^{qN} \int_0^s \theta_i \left\{ \frac{\partial f_k(l)}{\partial \tilde{x}_k} \Delta \tilde{x}_k(l) + \frac{\partial f_k(l)}{\partial \sigma^i} \rho^i \right\} dl.$$

Clearly,

$$\frac{d\Delta \tilde{x}_k(s)}{ds} = \sum_{i=1}^{qN} \theta_i \left\{ \frac{\partial f_k(s)}{\partial \tilde{x}_k} \Delta \tilde{x}_k(s) + \frac{\partial f_k(s)}{\partial \sigma^i} \rho^i \right\}. \quad (5.38)$$

By using (5.32), we have

$$\Delta H_m = \frac{\partial H_m}{\partial \sigma} \rho + \sum_{k=1}^N \frac{\partial H_m}{\partial \tilde{x}_k} \Delta \tilde{x}_k(s) = \sum_{k=1}^N \lambda_k^m(s) \omega(s) \left\{ \frac{\partial f_k(s)}{\partial \tilde{x}_k} \Delta \tilde{x}_k(s) + \frac{\partial f_k(s)}{\partial \sigma} \rho \right\}.$$

Using the definition of  $\omega(s)$  in (5.20), it is clear that

$$\frac{\partial H_m}{\partial \boldsymbol{\sigma}} \boldsymbol{\rho} + \sum_{k=1}^N \frac{\partial H_m}{\partial \tilde{x}_k} \Delta \tilde{x}_k(s) = \sum_{i=1}^{qN} \sum_{k=1}^N \lambda_k^m(s) \theta_i \left\{ \frac{\partial f_k(s)}{\partial \tilde{x}_k} \Delta \tilde{x}_k(s) + \frac{\partial f_k(s)}{\partial \boldsymbol{\sigma}^i} \boldsymbol{\rho}^i \right\}. \quad (5.39)$$

Substituting (5.38) into (5.39) yields

$$\frac{\partial H_m}{\partial \boldsymbol{\sigma}} \boldsymbol{\rho} + \sum_{k=1}^N \frac{\partial H_m}{\partial \tilde{x}_k} \Delta \tilde{x}_k(s) - \sum_{k=1}^N \lambda_k^m(s) \frac{d\Delta \tilde{x}_k(s)}{ds} = 0. \quad (5.40)$$

Then, by the chain rule, we obtain

$$\Delta \Phi_m(\tilde{\boldsymbol{x}}(mq), \boldsymbol{v}) = \left. \frac{d\Phi_m(\tilde{\boldsymbol{x}}(mq), \boldsymbol{v})}{d\epsilon} \right|_{\epsilon=0} = \frac{\partial \Phi_m(\tilde{\boldsymbol{x}}(mq), \boldsymbol{v})}{\partial \tilde{x}_k} \Delta \tilde{x}_k(mq). \quad (5.41)$$

Adding (5.40) to the right hand side of (5.41), gives

$$\begin{aligned} \Delta \Phi_m(\tilde{\boldsymbol{x}}(mq), \boldsymbol{v}) &= \sum_{k=1}^N \frac{\partial \Phi_m(\tilde{\boldsymbol{x}}(mq), \boldsymbol{v})}{\partial \tilde{x}_k} \Delta \tilde{x}_k(mq) + \int_0^{mq} \frac{\partial H_m}{\partial \boldsymbol{\sigma}} \boldsymbol{\rho} dl \\ &\quad + \sum_{k=1}^N \int_0^{mq} \left\{ \frac{\partial H_m}{\partial \tilde{x}_k} \Delta \tilde{x}_k(l) - \lambda_k^m(l) \frac{d\Delta \tilde{x}_k(l)}{dl} \right\} dl \end{aligned} \quad (5.42)$$

Integrating the last term of (5.42) by parts gives

$$\begin{aligned} \Delta \Phi_m(\tilde{\boldsymbol{x}}(mq), \boldsymbol{v}) &= \sum_{k=1}^N \frac{\partial \Phi_m(\tilde{\boldsymbol{x}}(mq), \boldsymbol{v})}{\partial \tilde{x}_k} \Delta \tilde{x}_k(mq) - \sum_{k=1}^N \lambda_k^m(l) \Delta \tilde{x}_k(l) \Big|_0^{mq} \\ &\quad + \int_0^{mq} \frac{\partial H_m}{\partial \boldsymbol{\sigma}} \boldsymbol{\rho} dl + \sum_{k=1}^N \int_0^{mq} \left\{ \frac{\partial H_m}{\partial \tilde{x}_k} \Delta \tilde{x}_k(l) + \frac{d\lambda_k^m(l)}{dl} \Delta \tilde{x}_k(l) \right\} dl. \end{aligned} \quad (5.43)$$

Since  $\tilde{x}_k(0) = 0$ ,  $k = 1, \dots, N$ , it follows that  $\lambda_k^m(0) \Delta \tilde{x}_k(0) = 0$ ,  $k = 1, \dots, N$ . Substituting (5.33)-(5.34) into (5.43) yields

$$\Delta \Phi_m(\tilde{\boldsymbol{x}}(mq), \boldsymbol{v}) = \int_0^{mq} \frac{\partial H_m}{\partial \boldsymbol{\sigma}} \boldsymbol{\rho} dl.$$

Since  $\boldsymbol{\rho}$  is chosen arbitrary, using the definition of  $\boldsymbol{\sigma}$ , the result follows directly.  $\square$

**Theorem 5.3.** For each  $m$ ,  $m = 1, \dots, N$ , the partial derivatives of  $\Phi_m$  with respect to  $\boldsymbol{\theta}$  is

$$\frac{\partial \Phi_m(\tilde{\boldsymbol{x}}(mq), \boldsymbol{v})}{\partial \theta_i} = \begin{cases} \int_{i-1}^i \frac{\partial H_m}{\partial \theta_i} ds, & \text{if } i \leq mq, \\ 0, & \text{if } i > mq. \end{cases} \quad (5.44)$$

*Proof.* Let the time scaling parameter vector  $\boldsymbol{\theta}$  be perturbed by  $\epsilon \boldsymbol{\rho}$ , where  $\epsilon > 0$  is a small

real number and  $\boldsymbol{\rho}$  is an arbitrary fixed perturbation of  $\boldsymbol{\theta}$ . Then, we have

$$\boldsymbol{\theta}(\epsilon) = \boldsymbol{\theta} + \epsilon\boldsymbol{\rho}.$$

where  $\boldsymbol{\rho} = [\rho_1, \dots, \rho_{qN}]^\top$ , and  $\boldsymbol{\theta}(\epsilon) = [\theta_1(\epsilon), \dots, \theta_{qN}(\epsilon)]^\top$ . Consequently, the system state  $\tilde{\boldsymbol{x}}$  as well as the function  $\Phi_m$  will be perturbed. Let

$$\tilde{x}_k(s, \epsilon) = \tilde{x}_k(s|\boldsymbol{\theta}(\epsilon)), \quad k = 1, \dots, N.$$

Clearly,

$$\tilde{x}_k(s, \epsilon) = \tilde{x}_k(0) + \sum_{i=1}^{qN} \int_0^s \theta_i(\epsilon) f_k(t(l), \tilde{x}_k(l), \boldsymbol{\sigma}^i)(l) dl.$$

Let  $f_k(t(\cdot), \tilde{x}_k(\cdot), \boldsymbol{\sigma}^i)(\cdot)$  be written as  $f_k(\cdot)$  for brevity. Then, by the chain rule, we have

$$\Delta\tilde{x}_k(s) = \left. \frac{d\tilde{x}_k(s, \epsilon)}{d\epsilon} \right|_{\epsilon=0} = \sum_{i=1}^{qN} \int_0^s \left\{ \rho_i f_k(l) + \theta_i \frac{\partial f_k(l)}{\partial \tilde{x}_k} \Delta\tilde{x}_k(l) + \theta_i \frac{\partial f_k(l)}{\partial t} \Delta t(l) \right\} dl.$$

Clearly,

$$\frac{d\Delta\tilde{x}_k(s)}{ds} = \sum_{i=1}^{qN} \left\{ \rho_i f_k(s) + \theta_i \frac{\partial f_k(s)}{\partial \tilde{x}_k} \Delta\tilde{x}_k(s) + \theta_i \frac{\partial f_k(s)}{\partial t} \Delta t(s) \right\}. \quad (5.45)$$

By (5.32), we have

$$\begin{aligned} \Delta H_m &= \frac{\partial H_m}{\partial \boldsymbol{\theta}} \boldsymbol{\rho} + \sum_{k=1}^N \frac{\partial H_m}{\partial \tilde{x}_k} \Delta\tilde{x}_k(s) + \frac{\partial H_m}{\partial t} \Delta t(s) \\ &= \sum_{k=1}^N \lambda_k^m(s) \left\{ \boldsymbol{\rho} f_k(s) + \omega(s) \frac{\partial f_k(s)}{\partial \tilde{x}_k} \Delta\tilde{x}_k(s) + \omega(s) \frac{\partial f_k(s)}{\partial t} \Delta t(s) \right\} + \lambda_{N+1}^m(s) \frac{d\omega(s)}{ds}. \end{aligned}$$

Using the definition of  $\omega(s)$  in (5.20), it is clear that

$$\begin{aligned} &\frac{\partial H_m}{\partial \boldsymbol{\theta}} \boldsymbol{\rho} + \sum_{k=1}^N \frac{\partial H_m}{\partial \tilde{x}_k} \Delta\tilde{x}_k(s) + \frac{\partial H_m}{\partial t} \Delta t(s) \\ &= \sum_{i=1}^{qN} \sum_{k=1}^N \lambda_k^m(s) \left\{ \rho_i f_k(s) + \theta_i \frac{\partial f_k(s)}{\partial \tilde{x}_k} \Delta\tilde{x}_k(s) + \theta_i \frac{\partial f_k(s)}{\partial t} \Delta t(s) \right\} + \lambda_{N+1}^m(s) \frac{d\omega(s)}{ds}. \end{aligned} \quad (5.46)$$

Substituting (5.45) into (5.46) yields

$$\frac{\partial H_m}{\partial \boldsymbol{\theta}} \boldsymbol{\rho} + \frac{\partial H_m}{\partial t} \Delta t(s) - \lambda_{N+1}^m(s) \frac{d\omega(s)}{ds} + \sum_{k=1}^N \frac{\partial H_m}{\partial \tilde{x}_k} \Delta\tilde{x}_k(s) - \sum_{k=1}^N \lambda_k^m(s) \frac{d\Delta\tilde{x}_k(s)}{ds} = 0. \quad (5.47)$$

Then, by the chain rule, we obtain

$$\begin{aligned}\Delta\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v}) &= \left. \frac{d\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{d\epsilon} \right|_{\epsilon=0} \\ &= \frac{\partial\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial t} \Delta t(mq) + \frac{\partial\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial \tilde{x}_k} \Delta \tilde{x}_k(mq).\end{aligned}\quad (5.48)$$

Adding (5.47) to the right hand side of (5.48), we obtain

$$\begin{aligned}\Delta\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v}) &= \frac{\partial\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial t} \Delta t(mq) + \sum_{k=1}^N \frac{\partial\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial \tilde{x}_k} \Delta \tilde{x}_k(mq) \\ &+ \int_0^{mq} \frac{\partial H_m}{\partial \boldsymbol{\theta}} \boldsymbol{\rho} dl + \int_0^{mq} \frac{\partial H_m}{\partial t} \Delta t(l) dl - \int_0^{mq} \lambda_{N+1}^m(l) \frac{d\omega(l)}{dl} dl \\ &+ \sum_{k=1}^N \int_0^{mq} \left\{ \frac{\partial H_m}{\partial \tilde{x}_k} \Delta \tilde{x}_k(l) - \lambda_k^m(l) \frac{d\Delta \tilde{x}_k(l)}{dl} \right\} dl\end{aligned}\quad (5.49)$$

Integrating the last term of (5.49) by parts gives

$$\begin{aligned}\Delta\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v}) &= \frac{\partial\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial t} \Delta t(mq) + \sum_{k=1}^N \frac{\partial\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial \tilde{x}_k} \Delta \tilde{x}_k(mq) \\ &+ \int_0^{mq} \frac{\partial H_m}{\partial \boldsymbol{\theta}} \boldsymbol{\rho} dl + \int_0^{mq} \frac{\partial H_m}{\partial t} \Delta t(l) dl - \int_0^{mq} \lambda_{N+1}^m(l) \frac{d\omega(l)}{dl} dl \\ &- \sum_{k=1}^N \lambda_k^m(l) \Delta \tilde{x}_k(l) \Big|_0^{mq} + \sum_{k=1}^N \int_0^{mq} \left\{ \frac{\partial H_m}{\partial \tilde{x}_k} - \frac{d\lambda_k^m(l)}{dl} \right\} \Delta \tilde{x}_k(l) dl\end{aligned}\quad (5.50)$$

Since  $\tilde{x}_k(0) = 0$ ,  $k = 1, \dots, N$ , it follows that  $\lambda_k^m(0) \Delta \tilde{x}_k(0) = 0$ ,  $k = 1, \dots, N$ . Substituting (5.33)-(5.34) into (5.50) yields

$$\begin{aligned}\Delta\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v}) &= \frac{\partial\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial t} \Delta t(mq) + \int_0^{mq} \frac{\partial H_m}{\partial \boldsymbol{\theta}} \boldsymbol{\rho} dl \\ &+ \int_0^{mq} \left\{ \frac{\partial H_m}{\partial t} \Delta t(l) - \lambda_{N+1}^m(s) \frac{d\omega(l)}{dl} \right\} dl\end{aligned}\quad (5.51)$$

Integrating the last term of (5.51) by parts gives

$$\begin{aligned}\Delta\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v}) &= \frac{\partial\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial t} \Delta t(mq) - \lambda_{N+1}^m(s) \omega(l) \Big|_0^{mq} + \int_0^{mq} \frac{\partial H_m}{\partial \boldsymbol{\theta}} \boldsymbol{\rho} dl \\ &+ \int_0^{mq} \left\{ \frac{\partial H_m}{\partial t} \Delta t(l) + \frac{d\lambda_{N+1}^m(s)}{ds} \omega(l) \right\} dl\end{aligned}\quad (5.52)$$

By (5.19), clearly,  $\omega(l) = \Delta t(l)$ . By (5.31), we have

$$\frac{\partial\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v})}{\partial t} = 0.\quad (5.53)$$

Substituting (5.35)-(5.36) and (5.53) into (5.52) yields

$$\Delta\Phi_m(\tilde{\mathbf{x}}(mq), \mathbf{v}) = \int_0^{mq} \frac{\partial H_m}{\partial \boldsymbol{\theta}} \boldsymbol{\rho} dl.$$

Since  $\boldsymbol{\rho}$  is chosen arbitrary, using the definition of  $\boldsymbol{\theta}$ , the result follows directly.  $\square$

The gradient formula given above for the penalty and constraint functions can be combined with a standard optimization algorithm—for example, a conjugate gradient method or sequential quadratic programming [48]—to solve Problem (P3) as a nonlinear programming problem. The optimal control software MISER does this automatically. It can be shown that a local solution of Problem (P3) converges to a local solution of Problem (P2) as the penalty parameter becomes sufficiently large [41, 42]. Hence, we can obtain an approximate solution of Problem (P2) by solving Problem (P3) for large  $\eta$ . A corresponding suboptimal control for Problem (P) can then be constructed according to (5.17). In the next section, we use this approach to solve two examples.

## 5.5 Numerical examples

### 5.5.1 Example 5.1

In [90], a linear gradient elution chromatographic process for separating four protein solutes is described. The retention time for each protein solute is controlled by adjusting the ionic strength of the mobile phase composition. The rate of change of each chromatography signal is given by the following system functions:

$$\begin{aligned} f_1(u(t)) &= \frac{(u(t))^2}{1.4(u(t))^2 + 3.6}, \\ f_2(u(t)) &= \frac{u(t)}{1.4u(t) + 5.4}, \\ f_3(u(t)) &= \frac{u(t)}{1.4u(t) + 7.2}, \\ f_4(u(t)) &= \frac{(u(t))^2}{1.4(u(t))^2 + 9}, \end{aligned}$$

where the control is subject to the bound constraints  $0.5 \leq u(t) \leq 2$ . A chromatography signal of 1 indicates that the protein solute has left the chromatography column. Thus, the state variables must satisfy the following interior point constraints at the retention times:

$$x_k(\tau_k) = 1, \quad k = 1, 2, 3, 4. \quad (5.54)$$

We consider the problem of maximizing the objective function (5.5) with  $N = 4$  subject

to the interior point constraints (5.54). This problem is in the form of Problem (P). We use the method outlined in Sections 5.3 and 5.4 with  $q = 1$  to approximate Problem (P) by Problem (P3). Problem (P3), the exact penalty function problem, can be solved by using the software package MISER 3. The parameters in the penalty function for Problem (P3) are set as follows:

$$\alpha = 1.5; \quad \beta = 2.2; \quad \gamma = 3; \quad \kappa_1 = \kappa_2 = \kappa_3 = 0.3.$$

We start by solving Problem (P3) using MISER 3 with  $\eta = 10$ . Then we increase  $\eta$  and re-solve Problem (P3), using the previous solution as the initial guess. We continue to increase  $\eta$  until  $\eta = 30$ , at which point the objective function value is 0.8205 with  $\varepsilon = 0.0024$ . The objective function value is slightly better than the result reported in [90]. The largest violation of the interior point constraints (5.27) is  $6 \times 10^{-4}$ , which is smaller than the violation of  $1.07 \times 10^{-3}$  reported in [90].

We next consider Problem (P3) with  $q = 3$ . Then, the time horizon is divided into 12 sub-intervals. We again use MISER 3 to solve Problem (P3). The optimal trajectories and optimal control computed by MISER 3 are shown in Figure 5.3. The optimal control values and the corresponding time-scaling transformation parameters are given in Tables 5.1 and 5.2, respectively. The optimal objective function value is 0.8276 when  $\varepsilon = 0.001$ ,  $\eta = 300$ , and  $\boldsymbol{\tau} = [2.7507, 5.9795, 9.1807, 12.3820]$ . As expected, the optimal value for  $q = 3$  is higher than the optimal value for  $q = 1$ .

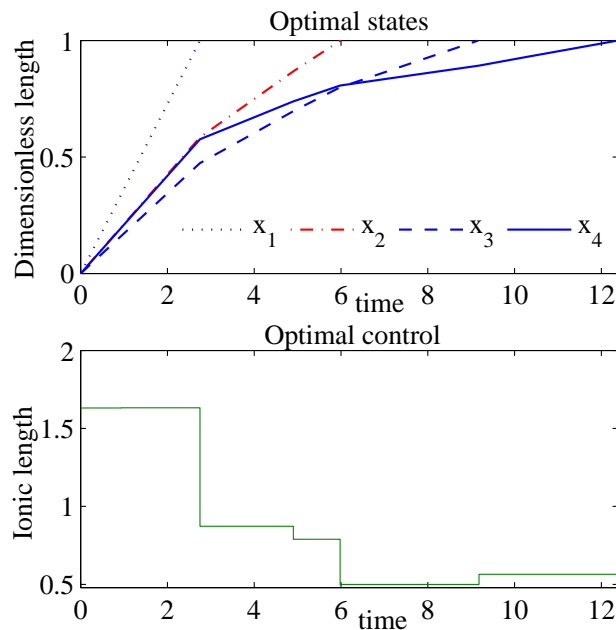


Figure 5.3: Optimal controls and states for Example 5.1.

The equality and inequality constraints (5.7)-(5.10), (5.21), (5.23), and 5.28 in Problem (P2) are all satisfied. Compared with the results obtained in [90], our optimal ob-



Table 5.1: Optimal control values for Example 5.1.

$\sigma_1$	$\sigma_2$	$\sigma_3$	$\sigma_4$	$\sigma_5$	$\sigma_6$
1.6314	1.6324	1.6326	0.8729	0.8730	0.7901
$\sigma_7$	$\sigma_8$	$\sigma_9$	$\sigma_{10}$	$\sigma_{11}$	$\sigma_{12}$
0.5000	0.5000	0.5000	0.5644	0.5645	0.5645

Table 5.2: Optimal interval durations for Example 5.1.

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$	$\theta_5$	$\theta_6$
0.9347	0.9153	0.9007	1.0740	1.0777	1.0770
$\sigma_7$	$\sigma_8$	$\sigma_9$	$\sigma_{10}$	$\sigma_{11}$	$\sigma_{12}$
1.0801	1.0608	1.0603	1.0541	1.0728	1.0743

jective function value is better. Moreover, the largest violation of the interior point constraints (5.27) is  $8 \times 10^{-4}$  which is slightly smaller than the violation achieved in [90].

### 5.5.2 Example 5.2

We now apply our optimal control method to a more complicated system with multiple inputs. Consider the ion-exchange chromatographic process for separating three protein species described in [112]. The peak time for each protein solute is controlled by adjusting the ionic strength ( $I$ ) and the pH value ( $z$ ) of the mobile phase composition. The rate of change of each chromatography signal is given by the following system functions:

$$f_k(I(t), z(t)) = \left[ 1 + \frac{(1-\varepsilon_b)\varepsilon_m}{\varepsilon_b} + \frac{(1-\varepsilon_b)(1-\varepsilon_m)}{\varepsilon_b} K_k(I(t), z(t)) \right]^{-1}, \quad k = 1, 2, 3$$

where  $\varepsilon_b = 0.25$ ,  $\varepsilon_m = 0.8$ , and the affinity distribution coefficients are given by

$$\begin{aligned} K_1(I(t), z(t)) &= 0.49 \times 10^6 (0.1867I^{-1})^{-18.5+9.9z-0.77z^2}, \\ K_2(I(t), z(t)) &= 5.48 \times 10^6 (0.1867I^{-1})^{-26.9+13.8z-1.15z^2}, \\ K_3(I(t), z(t)) &= 2.21 \times 10^6 (0.1867I^{-1})^{-24.7+12.9z-1.1z^2}. \end{aligned}$$

The controls are bounded by

$$0.405 \leq I(t) \leq 0.495, \quad 4.8 \leq z(t) \leq 5.2.$$

A protein species is considered to be separated when its corresponding chromatography signal reaches 1. That is, the state variables must satisfy the following interior point constraints:

$$x_k(\tau_k) = 1, \quad k = 1, 2, 3.$$

Our aim is to solve Problem (P), which involves maximizing the objective function (5.5) for the system with  $N = 3$  subject to the interior point constraints given above. Using the process outlined in Section 5.3 and Section 5.4 with  $q = 3$ , we approximate Problem (P) by Problem (P3). The parameters for the exact penalty function are the same as in Example 5.1. Using MISER 3 to solve Problem (P3), we obtain the optimal

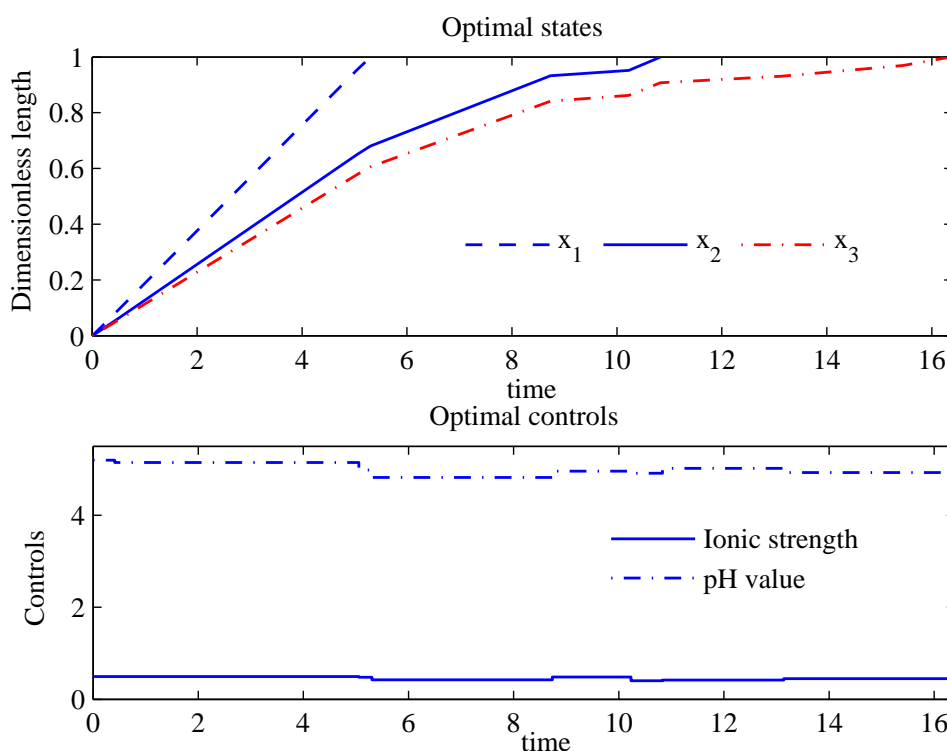


Figure 5.4: Optimal results for Example 5.2

trajectories and optimal controls shown in Figure 5.4. The optimal objective function value is 1.8673, when  $\eta = 200$ ,  $\varepsilon = 2.24895 \times 10^{-4}$ , and  $\boldsymbol{\tau} = [5.3029, 10.8293, 16.3557]$ . The largest violation of the interior point constraints (5.27) is  $1 \times 10^{-5}$ . All of the other constraints in Problem (P2) are satisfied.

## 5.6 Conclusions

In this chapter, we have presented a new computational method for solving a challenging max-min optimal control problem arising in gradient elution chromatography. The state

variables in this problem represent the chromatography signals of the mixture components, and each signal is required to reach its desired target at a different retention time. Our method is based on an approximation scheme whereby the control is approximated by a piecewise-constant function. We use a novel time-scaling transformation to map the retention times to fixed points, and another transformation to convert the max-min objective function into a smooth function. An exact penalty method is then used to solve the resulting approximate optimization problem. The numerical simulations in Section 5.5 show that our new method gives better results than those obtained by using the previous method proposed in [90]. More importantly, our new method is capable of accurately computing the retention times, and it can also be applied to more general problems of larger dimension.

---

---

# CHAPTER 6

---

## Summary and future research directions

### 6.1 Main contributions of the thesis

In this thesis, we considered two parameter identification problems and two nonstandard optimal control problems. We developed new methods for solving these problems as nonlinear programming problems, which are based on gradient-based optimization algorithms. This involved combining a variety of novel techniques, including the time-scaling transformation, constraint transcription, control parameterization, as well as complicated derivation of gradient formulas and detailed convergence analysis. We summarize our main contributions below.

In Chapter 2, we considered the problem of identifying unknown delays and unknown system parameters in a general nonlinear delay-differential system with smooth inputs. This problem was formulated as an optimization problem in which the decision variables are the time-delays and the system parameters. We showed that the gradient of the cost function in this optimization problem can be computed by solving a set of auxiliary delay-differential systems. On this basis, the optimization problem can be solved as a nonlinear programming problem using a gradient-based optimization algorithms. Most other methods for delay identification are only applicable to systems containing only a single delay, or to linear systems. The novelties of our approach are: (i) It is applicable to general nonlinear delay-differential systems; and (ii) it makes use of efficient nonlinear programming techniques, and hence it has excellent potential for real-time implementation.

In Chapter 3, we considered a general nonlinear time-delay system in which the delay is involved in a piecewise constant input. We formulated the problem of identifying the unknown state delay and control delay as a nonlinear optimization problem in which the cost function measures the least-squares error between predicted output and observed system output. Since the input function is discontinuous, the dynamics are clearly discontinuous with respect to  $\beta$ . Thus, the gradient of the cost function with respect to input delay does not exist at those time points which lie in the set of the control switching time points. We then showed that, except at those special time points, the gradient of the cost function with respect to the input delay can be computed by integrating

an auxiliary impulsive system with instantaneous jumps forward in time. To solve the identification problem, we propose a heuristic strategy which can be combined with our gradient computation procedure. Then, a computation method is developed based on any standard nonlinear programming algorithm. Two industrial examples are solved by using the proposed method. The results obtained show that this approach is highly effective.

In Chapter 4, we developed a computational method for solving optimal control problems with multiple time-delay systems and subject to continuous state inequality constraints. In this method, the control is approximated by a piecewise constant function whose heights are decision variables to be determined optimally. The approximate control is allowed to change its value at  $N - 1$  switching times in the time horizon. By novel application of the transformation method used in conjunction with a local smoothing technique, the continuous state inequality constraints are approximated by smooth canonical constraints, which are then appended to the cost function forming an augmented cost function. On this basis, the optimal control problem is approximated by a sequence of optimal parameter selection problems involving time-delay systems and subject to bound constraints on the control vector. Furthermore, we proved under some mild assumptions that the cost of the optimal control vector of the approximate problem converges to the optimal cost of the original problem. Finally, a computational algorithm is developed to solve each of the optimal parameter selection problems. This algorithm involves integrating the time-delay system forward in time and a set of time-delay costate systems backward in time. Then, the gradient of the augmented cost function is calculated. This approach is then applied to solve a control problem arising in a practical evaporation process. The results obtained are highly satisfactory.

In Chapter 5, we considered the problem of determining an optimal operating schedule for a chromatography process. The optimal control problem is formulated as a max-min optimal control problem with multiple characteristic-time equality constraints, which is similar to the one considered in [90]—it involves choosing the retention time for each component so that the minimum time interval between each two successive components is maximized. However, the problem considered in [90] assumes that the order of the retention times is known and fixed, whereas we considered a general problem where the ordering of the retention times is not fixed, but instead needs to be determined optimally. We then developed a new computational method for solving this max-min optimal control problem. The main idea of this method is to approximate the control by a piecewise constant function whose values and switching times are decision variables to be determined optimally. The approximate control is allowed to change its value at each switching time, and also at  $N - 1$  times between consecutive characteristic times ( $N$  is a fixed integer). Then a new time scaling transformation method is used to map these switching times to fixed points in a new time horizon. Note that the method used in [90] only maps the terminal time to a fixed point, but the intermediate retention times remain

as variables which change with the control. Thus, in the numerical computation, it is very difficult to determine accurately the values for these retention times. In our new method, we introduce a superior time-scaling transformation that maps both the control switching times and the retention times to fixed points. The retention times and the control switching times in the original time horizon become optimization variables in the new time horizon. This allows us to simultaneously determine accurate values for the retention times and the optimal control switching times. In Chapter 5, we use a set of auxiliary decision variables to explicitly keep track of the retention time ordering which involves approximating the integer constraints by a set of linear constraints and continuous quadratic constraints. The final approximate problem is solved by a recently developed exact penalty function method.

## 6.2 Future research directions

The work in this thesis has opened several interesting new areas for future research. We discuss some of them below.

The computational methods developed in Chapter 2 and Chapter 3 are only applicable to the delay system with time-invariant (constant) delays. In some practical applications, the delays are time varying. For example, the chemical reactor recycle system [20], cooling system [2], and the anesthesia control during intensive care [27]. For these systems with time varying delays, the methods developed in Chapters 2 and 3 are not applicable. Thus, considerable effort is needed to extend these methods to delay systems with time varying delays. It is a mathematically challenging and practically significant task to do the parameter identification on-line for systems with uncertainties, where multiple time delays and system parameters are to be identified.

The computational method developed in Chapter 4 involves integrating the time-delay system forward in time and integrating the corresponding time-delay co-state system backward in time. The computational burden is heavy. For example, it takes more than 3 hours to obtain an optimal control for eight hours of simulation time for the optimal control of the evaporation process considered in Section 4.4. In practice, the computation time is required to be light. Since many real practical processes involve complex dynamics, an interesting future research direction is to develop computational methods for which the computational load is light so that they are suitable for online implementation.

In the field of chromatography, the simulated moving bed (SMB) chromatography is becoming more and more popular. The model of the dynamic system for SMB chromatography is usually described by partial differential equations with position and time as the independent variables. Clearly, our method is not directly applicable to SMB chromatography. Thus, it presents an interesting and challenging task to extend our approach to SMB chromatography.

---

---

# Bibliography

---

- [1] A. Ghasemi, “Application of linear model predictive control and input-output linearization to constrained control of 3D cable robots,” *Modern Mechanical Engineering*, vol. 1, no. 2, pp. 69–76, 2011.
- [2] A. H. Tan and C. L. Cham, “Continuous-time model identification of a cooling system with variable delay,” *IET Control Theory and Applications*, vol. 5, no. 7, pp. 913–922, 2011.
- [3] A. Jamali, A. Hajiloo, and N. Nariman-zadeh, “Reliability-based robust pareto design of linear state feedback controllers using a multi-objective uniform-diversity genetic algorithm (muga),” *Expert Systems with Applications*, vol. 37, no. 1, pp. 401–413, 2010.
- [4] A. M. Cramer, S. D. Sudhoff, and E. L. Zivi, “Evolutionary algorithms for minimax problems in robust design,” *IEEE Transactions on Evolutionary Computation*, vol. 13, no. 2, pp. 444–453, 2009.
- [5] A. Miele, “Recent advances in gradient algorithms for optimal control problems,” *Journal of Optimization Theory and Applications*, vol. 17, no. 5, pp. 361–430, 1975.
- [6] A. Miele and R. E. Pritchard, “Gradient methods in control theory, part 2, sequential gradient-restoration algorithm,” Rice University, Aero-Astronautics Report 62, 1969.
- [7] A. Miele and T. Wang, “Primal-dual properties of sequential gradient-restoration algorithms for optimal control problems 2. general problem,” *Journal of Mathematical Analysis and Applications*, vol. 119, no. 1-2, pp. 21–54, 1986.
- [8] A. Miele, T. Wang, C. S. Chao, and J. B. Dabney, “Optimal control of a ship for collision avoidance maneuvers,” *Journal of optimization theory and applications*, vol. 103, no. 3, pp. 495–519, 1999.
- [9] A. Pot, S. Bhulai, and G. Koole, “A simple staffing method for multiskill call centers,” *Manufacturing and Service Operations Management*, vol. 10, no. 3, pp. 421–428, 2008.

- [10] A. Sbihi, “A cooperative local search-based algorithm for the multiple-scenario max-min knapsack problem,” *European Journal of Operational Research*, vol. 202, no. 2, pp. 339–346, 2010.
- [11] A. Toumi, F. Hanisch, and S. Engell, “Optimal operation of continuous chromatographic processes: mathematical optimization of the VARICOL process,” *Industrial and Engineering Chemistry Research*, vol. 41, pp. 4328–4337, 2002.
- [12] A. Woinaroschy, “Time-optimal control of startup distillation columns by iterative dynamic programming,” *Industrial and Engineering Chemistry Research*, vol. 47, no. 12, pp. 4158–4169, 2008.
- [13] A. Woinaroschy and R. Isopescu, “Time-optimal control of dividing-wall distillation columns,” *Industrial and Engineering Chemistry Research*, vol. 49, pp. 9195–9208, 2010.
- [14] B. Farhadinia, K. L. Teo, and R. C. Loxton, “A computational method for a class of non-standard time optimal control problems involving multiple time horizons,” *Mathematical and Computer Modelling*, vol. 49, pp. 1682–1691, 2009.
- [15] B. Gu and P. G. Yash, “Control of nonlinear processes by using linear model predictive control algorithms,” *ISA Transactions*, vol. 47, no. 2, pp. 211–216, 2008.
- [16] B. Li, K. L. Teo, G. H. Zhao, and G. R. Duan, “An efficient computational approach to a class of minmax optimal control problems with applications,” *The ANZIAM Journal*, vol. 51, pp. 162–177, 2009.
- [17] B. Pluymers, L. Roobrouck, J. Buijs, J. A. K. Suykens, and M. B. De, “Constrained linear MPC with time-varying terminal cost using convex combinations,” *Automatica*, vol. 41, no. 5, pp. 831–837, 2005.
- [18] B. Rustem, R. G. Becker, and W. Marty, “Robust min-max portfolio strategies for rival forecast and risk scenarios,” *Journal of Economic Dynamics and Control*, vol. 24, no. 112, pp. 1591–1621, 2000.
- [19] C. Buskens, “Optimierungs methoden und sensitivitätsanalyse für optimale steuerprozesse mit steuer- und zustands-beschränkungen,” Ph.D. dissertation, Institut für Numerische und instrumentelle Mathematik, Universität Münster, Germany, 1998.
- [20] C. C. Hua, J. Leng, and X. P. Guan, “Decentralized mrac for large-scale interconnected systems with time-varying delays and applications to chemical reactor systems,” *Journal of Process Control*, 2012.



- [21] C. Charalambous and A. R. Conn, “An efficient method to solve the minimax problem directly,” *SIAM Journal on Numerical Analysis*, vol. 15, no. 1, pp. 162–187, 1978.
- [22] C. F. Wen, “Continuous-time generalized fractional programming problems. part i: basic theory,” *Journal of Optimization Theory and Applications*, pp. 1–35, 2012.
- [23] C. Feller, T. A. Johansen, and S. Olaru, “An improved algorithm for combinatorial multi-parametric quadratic programming,” *Automatica*, 2012.
- [24] C. Jiang, K. L. Teo, R. Loxton, and G. R. Duan, “A neighboring extremal solution for an optimal switched impulsive control problem,” *Journal of Industrial and Management Optimization*, vol. 8, no. 3, pp. 591–609, 2012.
- [25] C. Lin, Z. Wang, and F. Yang, “Observer-based networked control for continuous-time systems with random sensor delays,” *Automatica*, vol. 45, no. 2, pp. 578–584, 2009.
- [26] C. Liu, Z. Gong, E. Feng, and H. Yin, “Modelling and optimal control for nonlinear multistage dynamical system of microbial fed-batch culture,” *Journal of Industrial and Management Optimization*, vol. 5, no. 4, pp. 835–850, 2009.
- [27] C. M. Ionescu, R. Hodrea, and R. De Keyser, “Variable time-delay estimation for anesthesia control during intensive care,” *IEEE Transactions on Biomedical Engineering*, vol. 58, no. 2, pp. 363–369, 2011.
- [28] C. Pignotti, “A note on stabilization of locally damped wave equations with time delay,” *Systems and Control Letters*, vol. 61, no. 1, pp. 92–97, 2012.
- [29] C. R. Hargraves and S. W. Paris, “Direct trajectory optimization using nonlinear programming and collocation,” in *Astrodynamics Conference*, vol. 3-12, Williamsburg, VA, August 18-20 1986.
- [30] —, “Direct trajectory optimization using nonlinear programming and collocation,” *AIAA J. Guidance*, vol. 10, pp. 338–342, 1987.
- [31] C. R. Porfrio and D. Odloak, “Optimizing model predictive control of an industrial distillation column,” *Control Engineering Practice*, vol. 19, no. 10, pp. 1137–1146, 2011.
- [32] C. Sonntag, W. J. Su, O. Stursberg, and S. Engell, “Optimized start-up control of an industrial-scale evaporation system with hybrid dynamics,” *Control Engineering Practice*, vol. 16, no. 8, pp. 976–990, 2008.

- [33] C. Wu, K. L. Teo, R. Li, and Y. Zhao, “Optimal control of switched systems with time delay,” *Applied mathematics letters*, vol. 19, no. 10, pp. 1062–1067, 2006.
- [34] C. Y. Kaya, “The leap-frog algorithm and optimal control: background and demonstration,” in *Proceedings of ICOTA 98, Perth, Australia*, 1998.
- [35] C. Y. Kaya and J. L. Noakes, “Computations and time-optimal controls,” *Optimal Control Applications and Methods*, vol. 17, no. 3, pp. 171–185, 1996.
- [36] —, “Geodesics and an optimal control algorithm,” in *Proceedings of the 36th IEEE Conference on Decision and Control*, vol. 5, 1997, pp. 4918–4919.
- [37] —, “Leapfrog for optimal control,” *SIAM Journal on Numerical Analysis*, vol. 46, no. 6, pp. 2795–2817, 2008.
- [38] C. Y. Kaya and J. M. Martínez, “Euler discretization and inexact restoration for optimal control,” *Journal of Optimization Theory and Applications*, vol. 134, pp. 191–206, 2007.
- [39] C. Y. Kaya, S. K. Lucas, and S. T. Simakov, “Computations for bangbang constrained optimal control using a mathematical programming formulation,” *Optimal Control Applications and Methods*, vol. 25, no. 6, pp. 295–308, 2004.
- [40] C. Yu, K. L. Teo, and Y. Bai, “An exact penalty function method for nonlinear mixed discrete programming problems,” *Optimization Letters*, 2011.
- [41] C. Yu, K. L. Teo, L. Zhang, and Y. Bai, “A new exact penalty function method for continuous inequality constrained optimization problems,” *Journal of Industrial and Management Optimization*, vol. 6, pp. 895–910, 2010.
- [42] C. Yu, B. Li, R. Loxton, and K. L. Teo, “Optimal discrete-valued control computation,” *Journal of Global Optimization*, pp. 1–16, 2012.
- [43] D. Cao, M. Chen, H. Wang, and J. Ji, “Interval method for global solutions of a class of min-max problems,” *Applied Mathematics and Computation*, vol. 196, no. 2, pp. 594–602, 2008.
- [44] D. D. Morrison, J. D. Riley, and J. F. Zancanaro, “Multiple shooting method for two-point boundary value problems,” *Communications of the ACM*, vol. 5, no. 12, pp. 613–614, 1962.
- [45] D. DeHaan and M. Guay, “A real-time framework for model-predictive control of continuous-time nonlinear systems,” *IEEE Transactions on Automatic Control*, vol. 52, no. 11, pp. 2047–2057, 2007.

- [46] D. E. Hartley and W. Murgatroyd, “Criteria for the break-up of thin liquid layers flowing isothermally over solid surfaces,” *International Journal of Heat and Mass Transfer*, vol. 7, no. 9, pp. 1003–1015, 1964.
- [47] D. Eller, J. Aggarwal, and H. Banks, “Optimal control of linear time-delay systems,” *IEEE Transactions on Automatic Control*, vol. 14, no. 6, pp. 678–687, 1969.
- [48] D. G. Luenberger and Y. Ye, *Linear and nonlinear programming*, 3rd ed. Springer: New York, 2008.
- [49] D. Han, X. Yuan, W. Zhang, and X. Cai, “An adm-based splitting method for separable convex programming,” *Computational Optimization and Applications*, vol. 2, no. 54, pp. 343–369, 2013.
- [50] E. B. Lee and L. Markus, *Foundation of optimal control theory*. New York, John Wiley and Sons Inc., 1967.
- [51] E. Fridman, “Stability of linear descriptor systems with delay: a lyapunov-based approach,” *Journal of Mathematical Analysis and Applications*, vol. 273, no. 1, pp. 24–44, 2002.
- [52] E. Obasanjo, G. R. Tzallas, and B. Rustem, “An interior-point algorithm for nonlinear minimax problems,” *Journal of Optimization Theory and Applications*, vol. 144, no. 2, pp. 291–318, 2010.
- [53] E. Polak, J. O. Royset, and R. S. Womersley, “Algorithms with adaptive smoothing for finite minimax problems,” *Journal of Optimization Theory and Applications*, vol. 119, no. 3, pp. 459–484, 2003.
- [54] F. E. Curtis and M. L. Overton, “A sequential quadratic programming algorithm for nonconvex, nonsmooth constrained optimization,” *SIAM Journal on Optimization*, vol. 2, no. 22, pp. 474–500, 2012.
- [55] F. E. Haoussi, E. H. Tissir, F. Tadeo, and A. Hmamed, “Delay-dependent stabilisation of systems with time-delayed state and control: application to a quadruple-tank process,” *International Journal of Systems Science*, vol. 42, no. 1, pp. 41–49, 2011.
- [56] F. Fahroo and I. M. Ross, “Direct trajectory optimization by a chebyshev pseudospectral method,” in *Proceedings of the 2000 American Control Conference*, vol. 6, 2000, pp. 3860–3864.
- [57] F. Pan, R. C. Han, and D. M. Feng, “An identification method of time-varying delay based on genetic algorithm,” *Proceedings of the 2003 International Conference on Machine Learning and Cybernetics*, vol. 2, pp. 781–783, Nov. 2003.

- [58] F. S. Mjalli and M. A. Hussain, “Approximate predictive versus self-tuning adaptive control strategies of biodiesel reactors,” *Industrial and Engineering Chemistry Research*, vol. 48, no. 24, pp. 11 034–11 047, 2009.
- [59] G. A. Fraser, “A multiple-shooting technique for optimal control,” *Journal of Optimization Theory and Applications*, vol. 102, no. 2, pp. 299–313, 1999.
- [60] G. Fernholz, S. Engell, L. U. Kreul, and A. Gorak, “Optimal operation of a semi-batch reactive distillation column,” *Computers and Chemical Engineering*, vol. 24, no. 2-7, pp. 1569–1575, 2000.
- [61] G. P. Rangaiah, P. Saha, and M. O. Tad, “Nonlinear model predictive control of an industrial four-stage evaporator system via simulation,” *Chemical Engineering Journal*, vol. 87, no. 3, pp. 285–299, 2002.
- [62] H. Chung, E. Polak, and S. Sastry, “An external active-set strategy for solving optimal control problems,” *IEEE Transactions on Automatic Control*, vol. 54, no. 5, pp. 1129–1133, 2009.
- [63] H. Gao, T. Chen, and J. Lam, “A new delay system approach to network-based control,” *Automatica*, vol. 44, no. 1, pp. 39–52, 2008.
- [64] H. Kubin, M. Benesch, A. Dementjev, K. Kabitzsch, T. Unkelbach, F. H. Rogner, and C. Metzner, “Identification of process models and controller design for vacuum coating processes with a long dead time using an identification tool with advisory support,” in *IEEE Conference on Emerging Technologies Factory Automation*, Sept. 2009, pp. 1–8.
- [65] H. W. J. Lee, K. L. Teo, and X. Q. Cai, “An optimal control approach to nonlinear mixed integer programming problems,” *Computers and Mathematics with Applications*, vol. 36, pp. 87–105, 1998.
- [66] H. W. J. Lee, K. L. Teo, L. S. Jennings, and V. Rehbock, “Control parametrization enhancing technique for time optimal control problem,” *Dynamical Systems and Applications*, vol. 6, no. , pp. 243–261, 1997.
- [67] —, “Control parametrization enhancing technique for time optimal control problems,” *Dynamical Systems and Applications*, vol. 6, pp. 243–261, 1997.
- [68] I. Bonis, W. Xie, and C. Theodoropoulos, “A linear model predictive control algorithm for nonlinear large-scale distributed parameter systems,” *AIChE Journal*, vol. 58, no. 3, pp. 801–811, 2012.

- [69] I. Mallocci, J. Daafouz, and C. Iung, “Stability and stabilization of two time scale switched systems in discrete time,” *IEEE Transactions on Automatic Control*, vol. 55, no. 6, pp. 1434–1438, 2010.
- [70] I. V. Kolmanovskiy and T. L. Maizenberg, “Optimal control of continuous-time linear systems with a time-varying, random delay,” *Systems and Control Letters*, vol. 44, no. 2, pp. 119–126, 2001.
- [71] J. A. Tropp, “Just relax: Convex programming methods for identifying sparse signals in noise,” *IEEE Transactions on Information Theory*, vol. 52, no. 5, pp. 1030–1051, 2006.
- [72] J. Li and J. Huo, “Inexact smoothing method for large scale minimax optimization,” *Applied Mathematics and Computation*, vol. 218, no. 6, pp. 2750–2760, 2011.
- [73] J. Nocedal and S. Wright, *Numerical optimization*, 2nd ed. Springer: New York, 2006.
- [74] J. P. Richard, “Time-delay systems: an overview of some recent advances and open problems,” *Automatica*, vol. 39, no. 10, pp. 1667–1694, 2003.
- [75] J. Ruths, A. Zlotnik, and S. Li, “Convergence of a pseudospectral method for optimal control of complex dynamical systems,” in *The 50th IEEE Conference on Decision and Control and European Control Conference (CDC-ECC)*, Dec. 2011, pp. 5553–5558.
- [76] J. Zhou and A. L. Tits, “User’s guide for FFSQP version 3.7: A Fortran code for solving optimization programs, possibly minimax, with general inequality constraints and linear equality constraints, generating feasible iterates,” Institute for Systems Research, University of Maryland, College Park, MD, Tech. Rep. 20742, 1997.
- [77] K. Gu, V. L. Kharitonov, and J. Chen, *Stability of Time-Delay Systems*, 3 ed. Berlin, Germany: Birkhauser, 2003.
- [78] K. H. Wong, D. J. Clements, and K. L. Teo, “Optimal control computation for non-linear time-lag systems,” *Journal of Optimization Theory and Applications*, vol. 47, pp. 91–107, 1985.
- [79] K. H. Wong, L. S. Jennings, and F. Benyah, “The control parametrization enhancing transform for constrained time-delayed optimal control problems,” *The ANZIAM Journal*, vol. 43, no. E, pp. 154–185, 2002.
- [80] K. L. Teo, C. J. Goh, and K. H. Wong, *A unified computational approach to optimal control problems*. Longman Scientific and Technical: Essex, 1991.

- [81] K. L. Teo, L. S. Jennings, H. W. J. Lee, and V. Rehbock, “The control parameterization enhancing transform for constrained optimal control problems,” *The Journal of the Australian Mathematical Society. Series B. Applied Mathematics*, vol. 40, no. 3, pp. 314–335, 1999.
- [82] K. Teo and C. Goh, “A simple computational procedure for optimization problems with functional inequality constraints,” *IEEE Transactions on Automatic Control*, vol. 32, no. 10, pp. 940–941, 1987.
- [83] K. Uchida, K. Ikeda, T. Azuma, and A. Kojima, “Finite-dimensional characterizations of  $h_\infty$  control for linear systems with delays in input and output,” *International Journal of Robust and Nonlinear Control*, vol. 13, no. 9, pp. 833–843, 2003.
- [84] K. Watanabe and M. Ito, “A process-model control for linear systems with delay,” *IEEE Transactions on Automatic Control*, vol. 26, no. 6, pp. 1261–1269, 1981.
- [85] L. Belkoura, J. P. Richard, and M. Fliess, “Parameters estimation of systems with delayed and structured entries,” *Automatica*, vol. 45, no. 5, pp. 1117–1125, 2009.
- [86] L. Göllmann, D. Kern, and H. Maurer, “Optimal control problems with delays in state and control variables subject to mixed control—state constraints,” *Optimal Control Applications and Methods*, vol. 30, no. 4, pp. 341–365, 2009.
- [87] L. Jin, R. Kumar, and N. Elia, “Model predictive control-based real-time power system protection schemes,” *IEEE Transactions on Power Systems*, vol. 25, no. 2, pp. 988–998, 2010.
- [88] L. Pronzato and E. Walter, “Robust experiment design via maximin optimization,” *Mathematical Biosciences*, vol. 89, no. 2, pp. 161–176, 1988.
- [89] L. S. Jennings and K. L. Teo, “A computational algorithm for functional inequality constrained optimization problems,” *Automatica*, vol. 26, no. 2, pp. 371–375, 1990.
- [90] L. S. Jennings, K. L. Teo, F. Y. Wang, and Q. Yu, “Optimal protein separation,” *Computers and Chemical Engineering*, vol. 19, pp. 567–573, 1995.
- [91] L. S. Jennings, M. E. Fisher, K. L. Teo, and C. J. Goh, “MISER3: Solving optimal control problems an update,” *Advances in Engineering Software and Workstations*, vol. 13, no. 4, pp. 190–196, 1991.
- [92] L. S. Pontryagin, V. G. Boltyanskii, R. V. Gamkrelidze, and E. F. Mishchenko, *The mathematical theory of optimal processes*. Gordon and Breach Science Publishers: Montreux, 1986.

- [93] L. Y. Wang, W. H. Gui, K. L. Teo, R. C. Loxton, and C. H. Yang, “Time delayed optimal control problems with multiple characteristic time points: computation and industrial applications,” *Journal of Industrial and Management Optimization*, vol. 5, no. 4, pp. 705–718, 2009.
- [94] L. Y. Wang, W. H. Gui, K. L. Teo, R. Loxton, and C. H. Yang, “Optimal control problems arising in the zinc sulphate electrolyte purification process,” *Journal of Global Optimization*, vol. 54, no. 2, pp. 307–323, 2012.
- [95] L. Zunino, M. C. Soriano, I. Fischer, O. A. Rosso, and C. R. Mirasso, “Permutation-information-theory approach to unveil delay dynamics from time-series analysis,” *Physical Review E*, vol. 82, no. 4, p. ID 046212, 2010.
- [96] R. Luus, “Application of dynamic programming to differential-algebraic process systems,” *Computers and Chemical Engineering*, vol. 17, no. 4, pp. 373–377, 1993.
- [97] M. Baric, S. V. Rakovic, T. Besselmann, and M. Morari, “Max-min optimal control of constrained discrete-time systems,” in *World Congress*, vol. 17, no. 1, 2008, pp. 8803–8808.
- [98] M. G. García, C. E. Balsa, and J. R. Banga, “Dynamic optimization of a simulated moving bed (SMB) chromatographic separation process,” *Industrial and Engineering Chemistry Research*, vol. 45, pp. 9033–9041, 2006.
- [99] M. K. H. Fan, A. L. Tits, J. L. Zhou, L. S. Wang, and J. Koninckx, *CONSOL-OPTCAD user’s manual*. Systems Research Center, University of Maryland, 1992.
- [100] M. Kamgarpour and C. Tomlin, “On optimal control of non-autonomous switched systems with a fixed mode sequence,” *Automatica*, vol. 48, pp. 1177–1181, 2012.
- [101] M. L. A. Marquez and C. H. Moog, “Input-output feedback linearization of time-delay systems,” *IEEE Transactions on Automatic Control*, vol. 49, no. 5, pp. 781–785, 2004.
- [102] M. L. Kaplan and J. H. Heegaard, “Predictive algorithms for neuromuscular control of human locomotion,” *Journal of Biomechanics*, vol. 34, no. 6, pp. 1077–1083, 2001.
- [103] M. Liu, Q. G. Wang, B. Huang, and C. C. Hang, “Improved identification of continuous-time delay processes from piecewise step tests,” *Journal of Process Control*, vol. 17, no. 1, pp. 51–57, 2007.
- [104] M. S. Bazaraa, H. D. Sherali, and C. M. Shetty, *Nonlinear Programming: Theory and Algorithms*. New Jersey: John Wiley & Sons, 2006.

- [105] M. S. Charalambides, N. Shah, and C. C. Pantelides, “Synthesis of batch reaction/distillation processes using detailed dynamic models,” *Computers and Chemical Engineering*, vol. 19, pp. 167–174, 1995.
- [106] M. V. Kothare, V. Balakrishnan, and M. Morari, “Robust constrained model predictive control using linear matrix inequalities,” *Automatica*, vol. 32, no. 10, pp. 1361–1379, 1996.
- [107] M. Verstraeten, M. Pursch, P. Eckerle, J. Luong, and G. Desmet, “Modelling the thermal behaviour of the low-thermal mass liquid chromatography system,” *Journal of Chromatography A*, vol. 1218, no. 16, pp. 2252–2263, 2011.
- [108] N. Lyle, “A global algorithm for geodesics,” *Journal of the Australian Mathematical Society*, vol. 65, no. 1, pp. 37–50, 1998.
- [109] N. T. Russell, H. H. C. Bakker, and R. L. Chaplin, “A comparison of dynamic models for an evaporation process,” *Chemical Engineering Research and Design*, vol. 78, no. 8, pp. 1120–1128, 2000.
- [110] N. U. Ahmed, *Elements of finite-dimensional systems and control theory*. Essex: Longman Scientific and Technical, 1988.
- [111] ———, *Dynamic Systems and Control with Applications*. Singapore: World Scientific, 2006.
- [112] O. Kaltenbrunner, O. Giaverini, D. Woehle, and J. A. Asenjo, “Application of chromatographic theory for process characterization towards validation of an ion-exchange operation,” *Biotechnology and Bioengineering*, vol. 98, no. 1, pp. 201–210, 2007.
- [113] O. Stryk and R. Bulirsch, “Direct and indirect methods for trajectory optimization,” *Annals of Operations Research*, vol. 37, no. 1, pp. 357–373, 1992.
- [114] O. V. Stryk, “User’s guide for DIRCOL (Version 2.1): A direct collocation method for the numerical solution of optimal control problems,” Simulation and Systems Optimization, Technical University of Darmstadt, Tech. Rep., November 1999.
- [115] P. Berkmann and H. J. Pesch, “Abort landing in windshear: Optimal control problem with third-order state constraint and varied switching structure,” *Journal of Optimization Theory and Applications*, vol. 85, no. 1, pp. 21–57, 1995.
- [116] P. J. Gawthrop and M. T. Nihtila, “Identification of time delays using a polynomial identification method,” *System and Control Letters*, vol. 5, no. 4, pp. 267–271, 1985.



- [117] Q. Lin, R. Loxton, K. L. Teo, and Y. H. Wu, “A new computational method for a class of free terminal time optimal control problems,” *Pacific Journal of Optimization*, vol. 7, pp. 63–81, 2011.
- [118] —, “A new computational method for optimizing nonlinear impulsive systems,” *Dynamics of Continuous, Discrete and Impulsive Systems B*, vol. 18, no. 1, pp. 59–76, 2011.
- [119] Q. Q. Chai, C. H. Yang, K. L. Teo, and W. H. Gui, “Optimal Control of an Industrial-scale Evaporation Process: Sodium Aluminate Solution,” *Control Engineering Practice*, vol. 20, pp. 618–628, 2012.
- [120] R. B. Martin, “Optimal control drug scheduling of cancer chemotherapy,” *Automatica*, vol. 28, pp. 1113–1123, 1992.
- [121] R. C. Loxton, K. L. Teo, and V. Rehbock, “Optimal control problems with multiple characteristic time points in the objective and constraints,” *Automatica*, vol. 44, pp. 2923–2929, 2008.
- [122] R. C. Loxton, K. L. Teo, V. Rehbock, and K. F. C. Yiu, “Optimal control problems with a continuous inequality constraint on the state and the control,” *Automatica*, vol. 45, no. 10, pp. 2250–2257, 2009.
- [123] R. Datko, “Two examples of ill-posedness with respect to time delays revisited,” *IEEE Transactions on Automatic Control*, vol. 42, no. 4, pp. 511–515, 1997.
- [124] R. Loxton, K. L. Teo, and V. Rehbock, “Computational method for a class of switched system optimal control problems,” *IEEE Transactions on Automatic Control*, vol. 54, no. 10, pp. 2455–2460, 2009.
- [125] —, “An optimization approach to state-delay identification,” *IEEE Transactions on Automatic Control*, vol. 55, no. 9, pp. 2113–2119, 2010.
- [126] —, “Robust suboptimal control of nonlinear systems,” *Applied Mathematics and Computation*, vol. 217, no. 14, pp. 6566–6576, 2011.
- [127] R. Loxton, Q. Lin, V. Rehbock, and K. L. Teo, “Control parameterization for optimal control problems with continuous inequality constraints: New convergence results,” *Numerical Algebra, Control and Optimization*, vol. 2, no. 3, pp. 571–599, 2012.
- [128] R. Luus, “Application of dynamic programming to high-dimensional non-linear optimal control problems,” *International Journal of Control*, vol. 52, no. 1, pp. 239–250, 1990.

- [129] —, “Piecewise linear continuous optimal control by iterative dynamic programming,” *Industrial and Engineering Chemistry Research*, vol. 32, no. 5, pp. 859–865, 1993.
- [130] —, *Iterative dynamic programming*. Chapman and Hall/CRC, London, UK, 2000.
- [131] R. Song, H. Zhang, Y. Luo, and Q. Wei, “Optimal control laws for time-delay systems with saturating actuators based on heuristic dynamic programming,” *Neurocomputing*, vol. 73, no. 16-18, pp. 3020–3027, 2010.
- [132] T. Ruby and V. Rehbock, “Numerical solutions of optimal switching control problems,” *Optimization and control with applications*, vol. 96, pp. 447–459, 2005.
- [133] S. Dimartino, C. Boi, and G. C. Sarti, “A validated model for the simulation of protein purification through affinity membrane chromatography,” *Journal of Chromatography A*, vol. 1218, no. 13, pp. 1677–1690, 2011.
- [134] S. Diop, I. Kolmanovsky, P. E. Moraal, and M. V. Nieuwstadt, “Preserving stability/performance when facing an unknown time-delay,” *Control Engineering Practice*, vol. 9, no. 12, pp. 1319–1325, 2001.
- [135] S. Gonzalez and A. Miele, “Sequential gradient-restoration algorithm for optimal control problems with general boundary conditions,” *Journal of Optimization Theory and Applications*, vol. 26, pp. 395–425, 1978.
- [136] S. J. Qin and T. A. Badgwell, “A survey of industrial model predictive control technology,” *Control Engineering Practice*, vol. 11, no. 7, pp. 733–764, 2003.
- [137] S. K. Lucas and C. Y. Kaya, “Switching-time computation for bang-bang control laws,” in *Proceedings of the 2001 American Control Conference*, vol. 1, 2001, pp. 176–181.
- [138] S. Kameswaran and L. Biegler, “Convergence rates for direct transcription of optimal control problems using collocation at radau points,” *Computational Optimization and Applications*, vol. 41, pp. 81–126, 2008.
- [139] S. Trimboli, S. Di Cairano, A. Bemporad, and I. Kolmanovsky, “Model predictive control with delay compensation for air-to-fuel ratio control,” *Time Delay Systems: Methods, Applications and New Trends*, pp. 319–330, 2012.
- [140] S. Uma, M. Chidambaram, A. S. Rao, and C. K. Yoo, “Enhanced control of integrating cascade processes with time delays using modified Smith predictor,” *Chemical Engineering Science*, vol. 65, no. 3, pp. 1065–1075, 2010.

- [141] S. V. Drakunov, W. Perruquetti, J. P. Richard, and L. Belkoura, “Delay identification in time-delay systems using variable structure observers,” *Annual Reviews in Control*, vol. 30, no. 2, pp. 143–158, 2006.
- [142] R. W. H. Sargent, “Optimal control,” *Journal of Computational and Applied Mathematics*, vol. 124, no. 1-2, pp. 361–371, 2000.
- [143] V. L. Denis, C. Jauberthie, and B. G. Joly, “Identifiability of a nonlinear delayed-differential aerospace model,” *IEEE Transactions on Automatic Control*, vol. 51, no. 1, pp. 154–158, 2006.
- [144] V. M. Zavala and L. T. Biegler, “The advanced-step NMPC controller: optimality, stability and robustness,” *Automatica*, vol. 45, no. 1, pp. 86–93, 2009.
- [145] V. Miranda and R. Simpson, “Modelling and simulation of an industrial multiple effect evaporator: tomato concentrate,” *Journal of Food Engineering*, vol. 66, no. 2, pp. 203–210, 2005.
- [146] V. Rehbock, K. L. Teo, L. S. Jennings, and H. W. J. Lee, “A survey of the control parametrization and control parametrization enhancing methods for constrained optimal control problems,” *Progress in Optimization: Contributions from Australasia*, vol. 30, pp. 247–275, 1999.
- [147] W. Sun and Y. Yuan, *Optimization theory and methods: nonlinear programming*. Springer Verlag, 2006, vol. 1.
- [148] W. W. Hager, “Runge-Kutta methods in optimal control and the transformed adjoint system,” *Numerische Mathematik*, vol. 87, pp. 247–282, 2000.
- [149] X. L. Liu, G. R. Duan, and K. L. Teo, “Optimal soft landing control for moon lander,” *Automatica*, vol. 44, no. 4, pp. 1097–1103, 2008.
- [150] X. Liu, “Constrained control of positive systems with delays,” *IEEE Transactions on Automatic Control*, vol. 54, no. 7, pp. 1596–1600, 2009.
- [151] X. Xu and P. J. Antsaklis, “Optimal control of switched systems based on parameterization of the switching instants,” *IEEE Transactions on Automatic Control*, vol. 49, no. 1, pp. 2–16, 2004.
- [152] X. Y. Yin, W. F. Yao, and Y. Z. Hu, “Optimization of gradient separation for chromatographic fingerprint of herbal medicine: a predictive model combined with grid search,” *Journal of Chromatographic Science*, vol. 46, pp. 722–729, 2008.

- [153] Y. B. Kim and S. J. Kang, "Time delay control for fuel cells with bidirectional DC/DC converter and battery," *International Journal of Hydrogen Energy*, vol. 35, no. 16, pp. 8792–8803, 2010.
- [154] Y. H. Kim, C. K. Yoo, and I. B. Lee, "Optimization of biological nutrient removal in a sbr using simulation-based iterative dynamic programming," *Chemical Engineering Journal*, vol. 139, no. 1, pp. 11–19, 2008.
- [155] Y. Sun, B. van Wyk, and Z. Wang, "A new multi-swarm multi-objective particle swarm optimization based on pareto front set," *Advanced Intelligent Computing Theories and Applications. With Aspects of Artificial Intelligence*, pp. 203–210, 2012.
- [156] Y. Wang and S. Boyd, "Fast model predictive control using online optimization," *IEEE Transactions on Control Systems Technology*, vol. 18, no. 2, pp. 267–278, 2010.
- [157] Z. K. Chen and H. W. Chen, "New research on burden calculation for raw mix slurry in production of alumina with sintering process," *World Nonferrous Met*, pp. 41–45, 2004.
- [158] Z. Wan and M. V. Kothare, "Robust output feedback model predictive control using off-line linear matrix inequalities," *Journal of Process Control*, vol. 12, no. 7, pp. 763–774, 2002.

*Every reasonable effort has been made to acknowledge the owners of copyright material. I would be pleased to hear from any copyright owner who has been omitted or incorrectly acknowledged.*