

Lecture Notes on Optimal Estimation and Control

Moritz Diehl

July 16, 2014

Contents

1	Introduction	5
1.1	Dynamic System Classes	5
1.2	Continuous Time Systems	8
1.3	Discrete Time Systems	12
1.4	Optimization Problem Classes	14
1.5	Overview and Notation	16
2	Nonlinear Optimization	19
2.1	Important Special Classes	20
2.2	First Order Optimality Conditions	21
2.3	Second Order Optimality Conditions	23
3	Newton-Type Optimization Algorithms	25
3.1	Equality Constrained Optimization	25
3.2	Local Convergence of Newton-Type Methods	28
3.3	Inequality Constrained Optimization	29
3.4	Globalisation Strategies	31
4	Calculating Derivatives	33
4.1	Algorithmic Differentiation (AD)	33
4.2	The Forward Mode of AD	35
4.3	The Backward Mode of AD	36
4.4	Algorithmic Differentiation Software	39
5	Discrete Time Optimal Control Formulations	41
5.1	Optimal Control Problem (OCP) Formulations	41
5.2	Analysis of a Simplified Optimal Control Problem	43
6	Sparsity Structure of the Optimal Control Problem	45
6.1	Partial Separability of the Lagrangian	45
6.2	The Sparse QP Subproblem	46
6.3	Sparsity Exploitation in QP Solvers	47
7	Dynamic Programming	51
7.1	Dynamic Programming in Discrete State Space	51
7.2	Linear Quadratic Problems	52
7.3	Infinite Horizon Problems	54
7.4	The Linear Quadratic Regulator	54
7.5	Robust and Stochastic Dynamic Programming	55
7.6	Interesting Properties of the DP Operator	55
7.7	The Gradient of the Value Function	56
7.8	A Discrete Time Minimum Principle	57

8	Continuous Time Optimal Control Problems	59
8.1	Formulation of Continuous Time Optimal Control Problems	59
8.2	Overview of Numerical Approaches	60
9	Direct Approaches to Continuous Optimal Control	63
9.1	Direct Single Shooting	63
9.2	Direct Multiple Shooting	64
9.3	Direct Collocation	65
9.4	Sensitivity Computation in Shooting Methods	67
9.5	A Classification of Direct Optimal Control Methods	68
10	Nonlinear Model Predictive Control	69
10.1	NMPC Optimization Problem	70
10.2	Nominal Stability of NMPC	70
10.3	Online Initialization via Shift	71
10.4	Outline of Real-Time Optimization Strategies	72
11	Moving Horizon Estimation	73
11.1	State and Parameter Estimation Problem Formulation	73
11.2	The Trajectory Estimation Problem	74
11.3	Dynamic Programming for the Trajectory Estimation Problem	76
11.4	Linear Quadratic Trajectory Estimation	76
11.5	Recursive Bayesian Estimation of the Last State	78

Chapter 1

Introduction

Optimal control regards the *optimization* of *dynamic systems*. In this lecture we identify dynamic systems with processes that are evolving with time and that can be characterized by *states* x that allow us to predict the future behavior of the system. If the state is not known, we first need to *estimate* it based on the available measurement information. The estimation process is very often optimization-based and uses the same optimization methods that are used for optimal control. This is the reason why this lecture bundles both optimal control and estimation in one single course. Often, a dynamic system can be controlled by a suitable choice of inputs that we denote as *controls* u in this script. In optimal control, these controls shall be chosen optimally in order to optimize some *objective function* and respect some *constraints*.

For optimal control, we might think of an electric train where the state x consists of the current position and velocity, and where the control u is the engine power that the train driver can choose at each moment. We might regard the motion of the train on a time interval $[t_{\text{init}}, t_{\text{fin}}]$, and the objective could be to minimize the consumed energy to drive from Station A to Station B, and one of the constraints would be that the train should arrive in Station B at the fixed final time, t_{fin} .

On the other hand, in optimization-based estimation, we treat unknown disturbances as inputs, and the objective function is the misfit between the actual measurements and their model predictions. The resulting optimization problems are mathematically of the same form as the problems of optimal control, with the disturbances as controls. For this reason, we focus the larger part of the course on the topic “optimal control”. At several occasions we specialize to estimation problems as a special case.

A typical property of a dynamic system is that knowledge of an *initial state* x_{init} and a *control input trajectory* $u(t)$ for all $t \in [t_{\text{init}}, t_{\text{fin}}]$ allows one to determine the whole *state trajectory* $x(t)$ for $t \in [t_{\text{init}}, t_{\text{fin}}]$. As the motion of a train can very well be modelled by Newton’s laws of motion, the usual description of this dynamic system is deterministic and in continuous time and with continuous states.

But dynamic systems and their mathematical models can come in many variants, and it is useful to properly define the names given commonly to different dynamic system classes, which we do in the next section. Afterwards, we will discuss two important classes, continuous time and discrete time systems, in more mathematical detail, before we give an overview of optimization problem classes and finally outline the contents of the lecture chapter by chapter.

1.1 Dynamic System Classes

In this section, let us go, one by one, through the many dividing lines in the field of dynamic systems.

Continuous vs Discrete Time Systems

Any dynamic system evolves over time, but time can come in two variants: while the physical time is continuous and forms the natural setting for most technical and biological systems, other dynamic systems can best be modelled in discrete time, such as digitally controlled sampled-data systems, or games.

We call a system a *discrete time system* whenever the time in which the system evolves only takes values on a predefined time grid, usually assumed to be integers. If we have an interval of real numbers, like for the physical time, we call it a *continuous time system*. In this lecture, we usually denote the continuous time by the variable

$t \in \mathbb{R}$ and write for example $x(t)$. In case of discrete time systems, we use an index, usually $k \in \mathbb{N}$, and write x_k for the state at time point k .

Continuous vs Discrete State Spaces

Another crucial element of a dynamic system is its state x , which often lives in a continuous state space, like the position of the train, but can also be discrete, like the position of the figures on a chess game. We define the *state space* \mathbb{X} to be the set of all values that the state vector x may take. If \mathbb{X} is a subset of a real vector space such as \mathbb{R}^{n_x} or another differentiable manifold, we speak of a *continuous state space*. If \mathbb{X} is a finite or a countable set, we speak of a *discrete state space*. If the state of a system is described by a combination of discrete and continuous variables we speak of a *hybrid state space*.

A *multi-stage system* is the special case of a system with hybrid state space that develops through a sequence of stages and where the state space on each stage is continuous. An example for a multi-stage system is walking, where consecutive stages are characterized by the number of feet that are on the ground at a given moment. For multi-stage systems, the time instant when one stage ends and the next one starts can often be described by a *switching function*. This function is positive on one and negative on the other stage, and assumes the value zero at the time instant that separates the stages.

Another special case are systems that develop in a continuous state space and in continuous time, but are sometimes subject to discontinuous jumps, such as bouncing billiard balls. These can often be modelled as multi-stage systems with switching functions, plus so called *jump conditions* that describe the discontinuous state evolution at the time instant between the stages.

Finite vs Infinite Dimensional State Spaces

The class of continuous state spaces can be further subdivided into the finite dimensional ones, whose state can be characterized by a finite set of real numbers, and the infinite dimensional ones, which have a state that lives in function spaces. The evolution of finite dimensional systems in continuous time is usually described by *ordinary differential equations (ODE)* or their generalizations, such as *differential algebraic equations (DAE)*.

Infinite dimensional systems are sometimes also called *distributed parameter systems*, and in the continuous time case, their behaviour is typically described by *partial differential equations (PDE)*. An example for a controlled infinite dimensional system is the evolution of the airflow and temperature distribution in a building that is controlled by an air-conditioning system.

Continuous vs Discrete Control Sets

We denote by \mathbb{U} the set in which the controls u live, and exactly as for the states, we can divide the possible control sets into *continuous control sets* and *discrete control sets*. A mixture of both is a *hybrid control set*. An example for a discrete control set is the set of gear choices for a car, or any switch that we can choose to be either on or off, but nothing in between.

In the systems and control community, the term *hybrid system* denotes a dynamic system which has either a hybrid state or hybrid control space, or both. Generally speaking, hybrid systems are more difficult to optimize than systems with continuous control and state spaces.

However, an interesting and relevant class are hybrid systems that have continuous time and continuous states, but discrete controls. They might be called hybrid systems with *external switches* or *integer controls* and turn out to be tremendously easier to optimize than other forms of hybrid systems, if treated with the right numerical methods.

Time-Variant vs Time-Invariant Systems

A system whose dynamics depend on time is called a *time-variant system*, while a dynamic system is called *time-invariant* if its evolution does not depend on the time and date when it is happening. As the laws of physics are time-invariant, most technical systems belong to the latter class, but for example the temperature evolution of a house with hot days and cold nights might best be described by a time-variant system model. While the class of time-variant systems trivially comprises all time-invariant systems, it is an important observation that also the other direction holds: each time-variant system can be modelled by a nonlinear time-invariant system if the state space is augmented by an extra state that takes account of the advancement of time, and which we might call the “clock state”.

Linear vs Nonlinear Systems

If the state trajectory of a system depends linearly on the initial value and the control inputs, it is called a *linear system*. If the dependence is affine, one should ideally speak of an *affine system*, but often the term linear is used here as well. In all other cases, we speak of a *nonlinear system*.

A particularly important class of linear systems are *linear time invariant (LTI)* systems. An LTI system can be completely characterized in at least three equivalent ways: first, by two matrices that are typically called A and B ; second, by its *step response function*; and third, by its *frequency response function*. A large part of the research in the control community is devoted to the study of LTI systems.

Controlled vs Uncontrolled Dynamic Systems

While we are in this lecture mostly interested in *controlled dynamic systems*, i.e. systems that have a control input that we can choose, it is good to remember that there exist many systems that cannot be influenced at all, but that only evolve according to their intrinsic laws of motion. These *uncontrolled systems* have an empty control set, $\mathbb{U} = \emptyset$. If a dynamic system is both uncontrolled and time-invariant it is also called an *autonomous system*.

Note that an autonomous system with discrete state space that also lives in discrete time is often called an *automaton*.

Within the class of controlled dynamic systems, of special interest are the so called *controllable systems*, which have the desirable property that their state vector x can be steered from any initial state x_{init} to any final state x_{fin} in a finite time with suitably chosen control input trajectories. Many controlled systems of interest are not completely controllable because some parts of their state space cannot be influenced by the control inputs. If these parts are stable, the system is called *stabilizable*.

Stable vs Unstable Dynamic Systems

A dynamic system whose state trajectory remains bounded for bounded initial values and controls is called a *stable system*, and an *unstable system* otherwise. For autonomous systems, *stability* of the system around a fixed point can be defined rigorously: for any arbitrarily small neighborhood \mathcal{N} around the fixed point there exists a region so that all trajectories that start in this region remain in \mathcal{N} . *Asymptotic stability* is stronger and additionally requires that all considered trajectories eventually converge to the fixed point. For autonomous LTI systems, stability can be computationally characterized by the eigenvalues of the system matrix.

Deterministic vs Stochastic Systems

If the evolution of a system can be predicted when its initial state and the control inputs are known, it is called a *deterministic system*. When its evolution involves some random behaviour, we call it a *stochastic system*.

The movements of assets on the stockmarket are an example for a stochastic system, whereas the motion of planets in the solar system can usually be assumed to be deterministic. An interesting special case of deterministic systems with continuous state space are *chaotic systems*. These systems are so sensitive to their initial values that even knowing these to arbitrarily high, but finite, precisions does not allow one to predict the complete future of the system: only the near future can be predicted. The partial differential equations used in weather forecast models have this property, and one well-known chaotic system of ODE, the *Lorenz attractor*, was inspired by these.

Note that also games like chess can be interpreted as dynamic systems. Here the evolution is neither deterministic nor stochastic, but determined by the actions of an adverse player. If we assume that the adversary always chooses the worst possible control action against us, we enter the field of *game theory*, which in continuous state spaces and engineering applications is often denoted by *robust optimal control*.

Open-Loop vs Closed-Loop Controlled Systems

When choosing the inputs of a controlled dynamic system, one first way is decide in advance, before the process starts, which control action we want to apply at which time instant. This is called *open-loop control* in the systems and control community, and has the important property that the control u is a function of time only and does not depend on the current system state.

A second way to choose the controls incorporates our most recent knowledge about the system state which we might observe with the help of measurements. This knowledge allows us to apply feedback to the system by adapting the control action according to the measurements. In the systems and control community, this is called

closed-loop control, but also the more intuitive term *feedback control* is used. It has the important property that the control action does depend on the current state. The map from the state to the control action is called a *feedback control policy*. In case this policy optimizes our optimization objective, it is called the *optimal feedback control policy*.

Open-loop control can be compared to a cooking instruction that says: cook the potatoes for 25 minutes in boiling water. A closed-loop, or feedback control of the same process would for example say: cook the potatoes in boiling water until they are so soft that they do not attach anymore to a fork that you push into them. The feedback control approach promises the better result, but requires more work as we have to take the measurements.

This lecture is mainly concerned with numerical methods of how to compute optimal open-loop controls for given objective and constraints. But the last part of the lecture is concerned with a powerful method to approximate the optimal feedback control policy: *model predictive control*, a feedback control technique that is based on the repeated solution of open-loop optimal control problems.

Focus of This Script

In this script we have a strong focus on deterministic systems with continuous state and control spaces. Mostly, we consider discrete time systems, while in a follow up lecture on numerical optimal control we discuss the treatment of continuous time systems in much more detail.

The main reason for our focus on continuous state and control spaces is that the resulting optimal control problems can efficiently be treated by derivative-based optimization methods. They are thus tremendously easier to solve than most other classes, both in terms of the solvable system sizes and of computational speed. Also, these continuous optimal control problems comprise the important class of convex optimal control problems, which allow us to find a global solution reliably and fast. Convex optimal control problems are important in their own right, but also serve as an approximation of nonconvex optimal control problems within Newton-type optimization methods.

1.2 Continuous Time Systems

Most systems of interest in science and engineering are described in form of differential equations which live in continuous time. On the other hand, all numerical simulation methods have to discretize the time interval of interest in some form or the other and thus effectively generate discrete time systems. We will thus only briefly sketch some relevant properties of continuous time systems in this section, and sketch how they can be transformed into discrete time systems. Throughout the lecture, we will mainly be concerned with discrete time systems, while we occasionally come back to the continuous time case.

Ordinary Differential Equations

A controlled dynamic system in continuous time can in the simplest case be described by an ordinary differential equation (ODE) on a time interval $[t_{\text{init}}, t_{\text{fin}}]$ by

$$\dot{x}(t) = f(x(t), u(t), t), \quad t \in [t_{\text{init}}, t_{\text{fin}}] \quad (1.1)$$

where $t \in \mathbb{R}$ is the time, $u(t) \in \mathbb{R}^{n_u}$ are the controls, and $x(t) \in \mathbb{R}^{n_x}$ is the state. The function f is a map from states, controls, and time to the rate of change of the state, i.e. $f : \mathbb{R}^{n_x} \times \mathbb{R}^{n_u} \times [t_{\text{init}}, t_{\text{fin}}] \rightarrow \mathbb{R}^{n_x}$. Due to the explicit time dependence of the function f , this is a time-variant system.

We are first interested in the question if this differential equation has a solution if the initial value $x(t_{\text{init}})$ is fixed and also the controls $u(t)$ are fixed for all $t \in [t_{\text{init}}, t_{\text{fin}}]$. In this context, the dependence of f on the fixed controls $u(t)$ is equivalent to a further time-dependence of f , and we can redefine the ODE as $\dot{x} = \tilde{f}(x, t)$ with $\tilde{f}(x, t) := f(x, u(t), t)$. Thus, let us first leave away the dependence of f on the controls, and just regard the time-dependent uncontrolled ODE:

$$\dot{x}(t) = f(x(t), t), \quad t \in [t_{\text{init}}, t_{\text{fin}}]. \quad (1.2)$$

Initial Value Problems

An initial value problem (IVP) is given by (1.2) and the initial value constraint $x(t_{\text{init}}) = x_{\text{init}}$ with some fixed parameter x_{init} . Existence of a solution to an IVP is guaranteed under continuity of f with respect to x and t

according to a theorem from 1886 that is due to Giuseppe Peano. But existence alone is of limited interest as the solutions might be non-unique.

Example 1 (Non-Unique ODE Solution) *The scalar ODE with $f(x) = \sqrt{|x(t)|}$ can stay for an undetermined duration in the point $x = 0$ before leaving it at an arbitrary time t_0 . It then follows a trajectory $x(t) = (t - t_0)^2/4$ that can be easily shown to satisfy the ODE (1.2). We note that the ODE function f is continuous, and thus existence of the solution is guaranteed mathematically. However, at the origin, the derivative of f approaches infinity. It turns out that this is the reason which causes the non-uniqueness of the solution.*

As we are only interested in systems with well-defined and deterministic solutions, we would like to formulate only ODE with unique solutions. Here helps the following theorem by Charles Émile Picard (1890) and Ernst Leonard Lindelöf (1894).

Theorem 1 (Existence and Uniqueness of IVP) *Regard the initial value problem (1.2) with $x(t_{\text{init}}) = x_{\text{init}}$, and assume that $f : \mathbb{R}^{n_x} \times [t_{\text{init}}, t_{\text{fin}}] \rightarrow \mathbb{R}^{n_x}$ is continuous with respect to x and t . Furthermore, assume that f is Lipschitz continuous with respect to x , i.e., that there exists a constant L such that for all $x, y \in \mathbb{R}^{n_x}$ and all $t \in [t_{\text{init}}, t_{\text{fin}}]$*

$$\|f(x, t) - f(y, t)\| \leq L\|x - y\|. \quad (1.3)$$

Then there exists a unique solution $x : [t_{\text{init}}, t_{\text{fin}}] \rightarrow \mathbb{R}^{n_x}$ of the IVP.

Lipschitz continuity of f with respect to x is not easy to check. It is much easier to verify if a function is differentiable. It is therefore a helpful fact that every function f that is differentiable with respect to x is also locally Lipschitz continuous, and one can prove the following corollary to the Theorem of Picard-Lindelöf.

Corollary 1 (Local Existence and Uniqueness) *Regard the same initial value problem as in Theorem 1, but instead of global Lipschitz continuity, assume that f is continuously differentiable with respect to x for all $t \in [t_{\text{init}}, t_{\text{fin}}]$. Then there exists a possibly shortened, but non-empty interval $[t_{\text{init}}, t'_{\text{fin}}]$ with $t'_{\text{fin}} \in (t_{\text{init}}, t_{\text{fin}}]$ on which the IVP has a unique solution.*

Note that for nonlinear continuous time systems – in contrast to discrete time systems – it is very easily possible to obtain an “explosion”, i.e., a solution that tends to infinity for finite times, even with innocently looking and smooth functions f .

Example 2 (Explosion of an ODE) *Regard the scalar example $f(x) = x^2$ with $t_{\text{init}} = 0$ and $x_{\text{init}} = 1$, and let us regard the interval $[t_{\text{init}}, t_{\text{fin}}]$ with $t_{\text{fin}} = 10$. The IVP has the explicit solution $x(t) = 1/(1 - t)$, which is only defined on the half open interval $[0, 1)$, because it tends to infinity for $t \rightarrow 1$. Thus, we need to choose some $t'_{\text{fin}} < 1$ in order to have a unique and finite solution to the IVP on the shortened interval $[t_{\text{init}}, t'_{\text{fin}}]$. The existence of this local solution is guaranteed by the above corollary. Note that the explosion in finite time is due to the fact that the function f is not globally Lipschitz continuous, so Theorem 1 is not applicable.*

Discontinuities with Respect to Time

It is important to note that the above theorem and corollary can be extended to the case that there are finitely many discontinuities of f with respect to t . In this case the ODE solution can only be defined on each of the continuous time intervals separately, while the derivative of x is not defined at the time points at which the discontinuities of f occur, at least not in the strong sense. But the transition from one interval to the next can be determined by continuity of the state trajectory, i.e. we require that the end state of one continuous initial value problem is the starting value of the next one.

The fact that unique solutions still exist in the case of discontinuities is important because, first, many optimal control problems have discontinuous control trajectories $u(t)$ in their solution, and, second, many algorithms discretize the controls as piecewise constant functions which have jumps at the interval boundaries. Fortunately, this does not cause difficulties for existence and uniqueness of the IVPs.

Linear Time Invariant (LTI) Systems

A special class of tremendous importance are the linear time invariant (LTI) systems. These are described by an ODE of the form

$$\dot{x} = Ax + Bu \quad (1.4)$$

with fixed matrices $A \in \mathbb{R}^{n_x \times n_x}$ and $B \in \mathbb{R}^{n_x \times n_u}$. LTI systems are one of the principal interests in the field of automatic control and a vast literature exists on LTI systems. Note that the function $f(x, u) = Ax + Bu$ is Lipschitz continuous with respect to x with Lipschitz constant $L = \|A\|$, so that the global solution to any initial value problem with a piecewise continuous control input can be guaranteed.

Many important notions such as *controllability* or *stabilizability*, and computational results such as the *step response* or *frequency response function* can be defined in terms of the matrices A and B alone. Note that in the field of linear system analysis and control, usually also output equations $y = Cx$ are present, where the outputs y may be the only physically relevant quantities. Only the linear operator from u to y - the input-output-behaviour - is of interest, while the state x is just an intermediate quantity. In that context, the states are not even unique, because different state space realizations of the same input-output behavior exist. In this lecture, however, we are not interested in input-outputs-behaviours, but assume that the state is the principal quantity of interest. Output equations are not part of the models in this lecture. If one wants to make the connection to the LTI literature, one might set $C = \mathbb{I}$.

Zero Order Hold and Solution Map

In the age of digital control, the inputs u are often generated by a computer and implemented at the physical system as piecewise constant between two sampling instants. This is called *zero order hold*. The grid size is typically constant, say of fixed length $\Delta t > 0$, so that the sampling instants are given by $t_k = k \cdot \Delta t$. If our original model is a differentiable ODE model, but we have piecewise constant control inputs with fixed values $u(t) = u_k$ with $u_k \in \mathbb{R}^{n_u}$ on each interval $t \in [t_k, t_{k+1}]$, we might want to regard the transition from the state $x(t_k)$ to the state $x(t_{k+1})$ as a discrete time system. This is indeed possible, as the ODE solution exists and is unique on the interval $[t_k, t_{k+1}]$ for each initial value $x(t_k) = x_{\text{init}}$.

If the original ODE system is time-invariant, it is enough to regard one initial value problem with constant control $u(t) = u_{\text{const}}$

$$\dot{x}(t) = f(x(t), u_{\text{const}}), \quad t \in [0, \Delta t], \quad \text{with } x(0) = x_{\text{init}}. \quad (1.5)$$

The unique solution $x : [0, \Delta t] \rightarrow \mathbb{R}^{n_x}$ to this problem is a function of both, the initial value x_{init} and the control u_{const} , so we might denote the solution by

$$x(t; x_{\text{init}}, u_{\text{const}}), \quad \text{for } t \in [0, \Delta t]. \quad (1.6)$$

This map from $(x_{\text{init}}, u_{\text{const}})$ to the state trajectory is called the *solution map*. The final value of this short trajectory piece, $x(\Delta t; x_{\text{init}}, u_{\text{const}})$, is of major interest, as it is the point where the next sampling interval starts. We might define the transition function $f_{\text{dis}} : \mathbb{R}^{n_x} \times \mathbb{R}^{n_u} \rightarrow \mathbb{R}^{n_x}$ by $f_{\text{dis}}(x_{\text{init}}, u_{\text{const}}) = x(\Delta t; x_{\text{init}}, u_{\text{const}})$. This function allows us to define a discrete time system that uniquely describes the evolution of the system state at the sampling instants t_k :

$$x(t_{k+1}) = f_{\text{dis}}(x(t_k), u_k). \quad (1.7)$$

Solution Map of Linear Time Invariant Systems

Let us regard a simple and important example: for linear continuous time systems

$$\dot{x} = Ax + Bu$$

with initial value x_{init} at $t_{\text{init}} = 0$, and constant control input u_{const} , the solution map $x(t; x_{\text{init}}, u_{\text{const}})$ is explicitly given as

$$x(t; x_{\text{init}}, u_{\text{const}}) = \exp(At)x_{\text{init}} + \int_0^t \exp(A(t-\tau))Bu_{\text{const}}d\tau,$$

where $\exp(A)$ is the matrix exponential. It is interesting to note that this map is well defined for all times $t \in \mathbb{R}$, as linear systems cannot explode. The corresponding discrete time system with sampling time Δt is again a linear time invariant system, and is given by

$$f_{\text{dis}}(x_k, u_k) = A_{\text{dis}}x_k + B_{\text{dis}}u_k \quad (1.8)$$

with

$$A_{\text{dis}} = \exp(A\Delta t) \quad \text{and} \quad B_{\text{dis}} = \int_0^{\Delta t} \exp(A(\Delta t - \tau))Bd\tau. \quad (1.9)$$

One interesting observation is that the discrete time system matrix A_{dis} resulting from the solution of an LTI system in continuous time is by construction an invertible matrix, with inverse $A_{\text{dis}}^{-1} = \exp(-A\Delta t)$. For systems with strongly decaying dynamics, however, the matrix A_{dis} might have some very small eigenvalues and will thus be nearly singular.

Sensitivities

In the context of optimal control, derivatives of the dynamic system simulation are needed for nearly all numerical algorithms. Following Theorem 1 and Corollary 1 we know that the solution map to the IVP (1.5) exists on an interval $[0, \Delta t]$ and is unique under mild conditions even for general nonlinear systems. But is it also differentiable with respect to the initial value and control input?

In order to discuss the issue of derivatives, which in the dynamic system context are often called *sensitivities*, let us first ask what happens if we call the solution map with different inputs. For small perturbations of the values $(x_{\text{init}}, u_{\text{const}})$, we still have a unique solution $x(t; x_{\text{init}}, u_{\text{const}})$ on the whole interval $t \in [0, \Delta t]$. Let us restrict ourselves to a neighborhood \mathcal{N} of fixed values $(x_{\text{init}}, u_{\text{const}})$. For each fixed $t \in [0, \Delta t]$, we can now regard the well defined and unique solution map $x(t; \cdot) : \mathcal{N} \rightarrow \mathbb{R}^{n_x}$, $(x_{\text{init}}, u_{\text{const}}) \mapsto x(t; x_{\text{init}}, u_{\text{const}})$. A natural question to ask is if this map is differentiable. Fortunately, it is possible to show that if f is m -times continuously differentiable with respect to both x and u , then the solution map $x(t; \cdot)$, for each $t \in [0, \Delta t]$, is also m -times continuously differentiable with respect to $(x_{\text{init}}, u_{\text{const}})$.

In the general nonlinear case, the solution map $x(t; x_{\text{init}}, u_{\text{const}})$ can only be generated by a numerical simulation routine. The computation of derivatives of this numerically generated map is a delicate issue that we discuss in detail in a follow up course on numerical optimal control. To mention already the main difficulty, note that most numerical integration routines are adaptive, i.e., might choose to do different numbers of integration steps for different IVPs. This renders the numerical approximation of the map $x(t; x_{\text{init}}, u_{\text{const}})$ typically non-differentiable in the inputs $x_{\text{init}}, u_{\text{const}}$. Thus, multiple calls of a black-box integrator and application of finite differences might result in very wrong derivative approximations.

Numerical Integration Methods

A numerical simulation routine that approximates the solution map is often called an *integrator*. A simple but very crude way to generate an approximation for $x(t; x_{\text{init}}, u_{\text{const}})$ for $t \in [0, \Delta t]$ is to perform a linear extrapolation based on the time derivative $\dot{x} = f(x, u)$ at the initial time point:

$$\tilde{x}(t; x_{\text{init}}, u_{\text{const}}) = x_{\text{init}} + tf(x_{\text{init}}, u_{\text{const}}), \quad t \in [0, \Delta t]. \quad (1.10)$$

This is called one *Euler integration step*. For very small Δt , this approximation becomes very good. In fact, the error $\tilde{x}(\Delta t; x_{\text{init}}, u_{\text{const}}) - x(\Delta t; x_{\text{init}}, u_{\text{const}})$ is of second order in Δt . This motivated Leonhard Euler to perform several steps of smaller size, and propose what is now called the *Euler integration method*. We subdivide the interval $[0, \Delta t]$ into M subintervals each of length $h = \Delta t/M$, and perform M such linear extrapolation steps consecutively, starting at $\tilde{x}_0 = x_{\text{init}}$:

$$\tilde{x}_{j+1} = \tilde{x}_j + hf(\tilde{x}_j, u_{\text{const}}), \quad j = 0, \dots, M-1. \quad (1.11)$$

It can be proven that the Euler integration method is *stable*, i.e. that the propagation of local errors is bounded with a constant that is independent of the step size h . Therefore, the approximation becomes better and better when we decrease the step size h : since the *consistency* error in each step is of order h^2 , and the total number of steps is of order $\Delta t/h$, the accumulated error in the final step is of order $h\Delta t$. As this is linear in the step size h , we say that the Euler method has the *order one*. Taking more steps is more accurate, but also needs more computation time. One measure for the computational effort of an integration method is the number of evaluations of f , which for the Euler method grows linearly with the desired accuracy.

In practice, the Euler integrator is rarely competitive, because other methods exist that deliver the desired accuracy levels at much lower computational cost. We discuss several numerical simulation methods later, but present here already one of the most widespread integrators, the *Runge-Kutta Method of Order Four*, which we will often abbreviate as *RK4*. One step of the RK4 method needs four evaluations of f and stores the results in

four intermediate quantities $k_i \in \mathbb{R}^{n_x}$, $i = 1, \dots, 4$. Like the Euler integration method, the RK4 also generates a sequence of values \tilde{x}_j , $j = 0, \dots, M$, with $\tilde{x}_0 = x_{\text{init}}$. At \tilde{x}_j , and using the constant control input u_{const} , one step of the RK4 method proceeds as follows:

$$k_1 = f(\tilde{x}_j, u_{\text{const}}) \quad (1.12a)$$

$$k_2 = f(\tilde{x}_j + \frac{h}{2} k_1, u_{\text{const}}) \quad (1.12b)$$

$$k_3 = f(\tilde{x}_j + \frac{h}{2} k_2, u_{\text{const}}) \quad (1.12c)$$

$$k_4 = f(\tilde{x}_j + h k_3, u_{\text{const}}) \quad (1.12d)$$

$$\tilde{x}_{j+1} = \tilde{x}_j + \frac{h}{6} (k_1 + 2k_2 + 2k_3 + k_4) \quad (1.12e)$$

One step of RK4 is thus as expensive as four steps of the Euler method. But it can be shown that the accuracy of the final approximation \tilde{x}_M is of order $h^4 \Delta t$. In practice, this means that the RK4 method usually needs tremendously fewer function evaluations than the Euler method to obtain the same accuracy level.

From here on, and throughout the major part of the lecture, we will leave the field of continuous time systems, and directly assume that we control a discrete time system $x_{k+1} = f_{\text{dis}}(x_k, u_k)$. Let us keep in mind, however, that the transition map $f_{\text{dis}}(x_k, u_k)$ is usually not given as an explicit expression but can instead be a relatively involved computer code with several intermediate quantities. In the exercises of this lecture, we will usually discretize the occurring ODE systems by using only one Euler or RK4 step per control interval, i.e. use $M = 1$ and $h = \Delta t$. The RK4 step often gives already a sufficient approximation at relatively low cost.

1.3 Discrete Time Systems

Let us now discuss in more detail the discrete time systems that are at the basis of the control problems in the first part of this lecture. In the general time-variant case, these systems are characterized by the dynamics

$$x_{k+1} = f_k(x_k, u_k), \quad k = 0, 1, \dots, N-1 \quad (1.13)$$

on a time horizon of length N , with N control input vectors $u_0, \dots, u_{N-1} \in \mathbb{R}^{n_u}$ and $(N+1)$ state vectors $x_0, \dots, x_N \in \mathbb{R}^{n_x}$.

If we know the initial state x_0 and the controls u_0, \dots, u_{N-1} we could recursively call the functions f_k in order to obtain all other states, x_1, \dots, x_N . We call this a *forward simulation* of the system dynamics.

Definition 1 (Forward simulation) *The forward simulation is the map*

$$\begin{aligned} f_{\text{sim}} : \quad & \mathbb{R}^{n_x + N n_u} \rightarrow \mathbb{R}^{(N+1)n_x} \\ & (x_0; u_0, u_1, \dots, u_{N-1}) \mapsto (x_0, x_1, x_2, \dots, x_N) \end{aligned} \quad (1.14)$$

that is defined by solving Equation (1.13) recursively for all $k = 0, 1, \dots, N-1$.

The inputs of the forward simulation routine are the initial value x_0 and the controls u_k for $k = 0, \dots, N-1$. In many practical problems we can only choose the controls while the initial value is fixed. Though this is a very natural assumption, it is not the only possible one. In optimization, we might have very different requirements: We might, for example, have a free initial value that we want to choose in an optimal way. Or we might have both a fixed initial state and a fixed terminal state that we want to reach. We might also look for periodic sequences with $x_0 = x_N$, but do not know x_0 beforehand. All these desires on the initial and the terminal state can be expressed by suitable constraints. For the purpose of this manuscript it is important to note that the fundamental equation that is characterizing a dynamic optimization problem are the system dynamics stated in Equation (1.13), but no initial value constraint, which is optional.

Linear Time Invariant (LTI) Systems

As discussed already for the continuous time case, linear time invariant (LTI) systems are not only one of the simplest possible dynamic system classes, but also have a rich and beautiful history. In the discrete time case, they are determined by the system equation

$$x_{k+1} = Ax_k + Bu_k, \quad k = 0, 1, \dots, N-1. \quad (1.15)$$

with fixed matrices $A \in \mathbb{R}^{n_x \times n_x}$ and $B \in \mathbb{R}^{n_x \times n_u}$. An LTI system is asymptotically stable if all eigenvalues of the matrix A are strictly inside the unit disc of the complex plane, i.e. have a modulus smaller than one. It is easy to show that the forward simulation map for an LTI system on a horizon with length N is given by

$$f_{\text{sim}}(x_0; u_0, \dots, u_{N-1}) = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} x_0 \\ Ax_0 + Bu_0 \\ A^2x_0 + ABu_0 + Bu_1 \\ \vdots \\ A^N x_0 + \sum_{k=0}^{N-1} A^{N-1-k} Bu_k \end{bmatrix}$$

In order to check controllability, due to linearity, one might ask the question if after N steps any terminal state x_N can be reached from $x_0 = 0$ by a suitable choice of control inputs. Because of

$$x_N = \underbrace{\begin{bmatrix} A^{N-1}B & A^{N-2}B & \dots & B \end{bmatrix}}_{=C_N} \begin{bmatrix} u_0 \\ u_1 \\ \vdots \\ u_{N-1} \end{bmatrix}$$

this is possible if and only if the matrix $C_N \in \mathbb{R}^{n_x \times N n_u}$ has the rank n_x . Increasing N can only increase the rank, but one can show that the maximum possible rank is already reached for $N = n_x$, so it is enough to check if the so called *controllability matrix* C_{n_x} has the rank n_x .

Eigenvalues and Eigenvectors of LTI Systems

Every square matrix $A \in \mathbb{R}^{n_x \times n_x}$ can be brought into the Jordan canonical form $A = QJQ^{-1}$ with non-singular $Q \in \mathbb{C}^{n_x \times n_x}$ and J block diagonal, consisting of m -Jordan blocks J_i . Thus, it holds that

$$J = \begin{bmatrix} J_1 & & \\ & J_2 & \\ & & \ddots \\ & & & J_m \end{bmatrix} \quad \text{with} \quad J_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \lambda_i & 1 & \\ & & \ddots & \ddots \\ & & & \lambda_i \end{bmatrix}.$$

Many of the Jordan blocks might just have size one, i.e. $J_i = [\lambda_i]$. To better understand the uncontrolled system evolution with dynamics $x_{k+1} = Ax_k$ and initial condition $x_0 = x_{\text{init}}$, one can regard the solution map $x_N = A^N x_0$ in the eigenbasis, which yields the expression

$$x_N = Q J^N (Q^{-1} x_0)$$

First, it is seen that all Jordan blocks evolve independently, after the initial condition is represented in the eigenbasis. Second, a simple Jordan block J_i will just result in the corresponding component being multiplied by a factor λ_i^N . Third, for nontrivial Jordan blocks, one obtains more complex expressions with N upper diagonals of the form

$$J_i^N = \begin{bmatrix} \lambda_i^N & N\lambda_i^{N-1} & \dots & 1 \\ & \lambda_i^N & N\lambda_i^{N-1} & \ddots \\ & & \ddots & \\ & & & \lambda_i^N \end{bmatrix}.$$

If one eigenvalue has a larger modulus $|\lambda_i|$ than all others, the Jordan block J_i^N will grow faster (or shrink slower) than the others for increasing N . The result is that the corresponding eigenvector(s) will dominate the final state x_N for large N , while all others “die out”. Here, the second largest eigenvalues will result in the most slowly decaying components, and their corresponding eigenvectors will keep a visible contribution in x_N the longest.

Interestingly, complex eigenvalues as well as eigenvectors appear in complex conjugate pairs. If an eigenvalue λ_i is complex, the (real part of) the corresponding eigenvector will perform oscillatory motion. To understand the behaviour of complex eigenvectors, let us regard a complex conjugate pair of simple eigenvalues λ_i and $\lambda_j = \bar{\lambda}_i$,

and their eigenvectors $v_i, v_j \in \mathbb{C}^{n_x}$, i.e. $Av_i = \lambda_i v_i$ and $Av_j = \bar{\lambda}_i v_j$. It is easy to see that, because A is real, $v_j = \bar{v}_i$ is a possible choice for the eigenvector corresponding to $\bar{\lambda}_i$. Then holds that $\text{Re}\{v_i\} = \frac{1}{2}(v_i + v_j)$. Therefore,

$$A^N \text{Re}\{v_i\} = \frac{1}{2}(\lambda_i^N v_i + \lambda_j^N v_j) = \frac{1}{2}(\lambda_i^N v_i + \bar{\lambda}_i^N \bar{v}_i) = \text{Re}\{\lambda_i^N v_i\}.$$

If we represent λ_i as $\lambda_i = re^{\phi i}$ (where the i in the exponent is the imaginary unit while the other i remains just an integer index), then $\lambda_i^N = r^N e^{N\phi i}$. If ϕ is a fraction of 2π , there is an N such that $N\phi = 2\pi$, and after N iterations we will obtain the same real part as in the original eigenvector, but multiplied with r^N . We can conclude that the real part of the eigenvector to a complex eigenvalue $re^{\phi i}$ performs a form of damped or growing oscillatory motion with period duration $N = 2\pi/\phi$ and growth constant r .

Affine Systems and Linearizations along Trajectories

An important generalization of linear systems are affine time-varying systems of the form

$$x_{k+1} = A_k x_k + B_k u_k + c_k, \quad k = 0, 1, \dots, N-1. \quad (1.16)$$

These often appear as linearizations of nonlinear dynamic systems along a given reference trajectory. To see this, let us regard a nonlinear dynamic system and some given reference trajectory values $\bar{x}_0, \dots, \bar{x}_{N-1}$ as well as $\bar{u}_0, \dots, \bar{u}_{N-1}$. Then the Taylor expansion of each function f_k at the reference value (\bar{x}_k, \bar{u}_k) is given by

$$(x_{k+1} - \bar{x}_{k+1}) \approx \frac{\partial f_k}{\partial x}(\bar{x}_k, \bar{u}_k)(x_k - \bar{x}_k) + \frac{\partial f_k}{\partial u}(\bar{x}_k, \bar{u}_k)(u_k - \bar{u}_k) + (f_k(\bar{x}_k, \bar{u}_k) - \bar{x}_{k+1})$$

thus resulting in affine time-varying dynamics of the form (1.16). Note that even for a time-invariant nonlinear system the linearized dynamics becomes time-variant due to the different linearization points on the reference trajectory.

It is an important fact that the forward simulation map of an affine system (1.16) is again an affine function of the initial value and the controls. More specifically, this affine map is for any $N \in \mathbb{N}$ given by:

$$x_N = (A_{N-1} \cdots A_0) x_0 + \sum_{k=0}^{N-1} \left(\prod_{j=k+1}^{N-1} A_j \right) (B_k u_k + c_k).$$

1.4 Optimization Problem Classes

Mathematical optimization refers to finding the best, or *optimal* solution among a set of possible decisions, where optimality is defined with the help of an *objective function*. Some solution candidates are *feasible*, others not, and it is assumed that *feasibility* of a solution candidate can be checked by evaluation of some *constraint functions* that need for example be equal to zero. Like the field of dynamic systems, the field of mathematical optimization comprises many different problem classes, which we will briefly try to classify in this section.

Historically, optimization has been identified with programming, where a program was understood as a deterministic plan, e.g., in logistics. For this reason, many of the optimization problem classes have been given names that contain the words *program* or *programming*. In this script we will often use these names and their abbreviations, because they are still widely used. Thus, we use e.g. the term *linear program (LP)* as a synonym for a *linear optimization problem*. It is interesting to note that the major society for mathematical optimization, which had for decades the name *Mathematical Programming Society (MPS)*, changed its name in 2011 to *Mathematical Optimization Society (MOS)*, while it decided not to change the name of its major journal, that still is called *Mathematical Programming*. In this script we chose a similarly pragmatic approach to the naming conventions.

Finite vs Infinite Dimensional Optimization

An important dividing line in the field of optimization regards the dimension of the space in which the decision variable, say x , is chosen. If x can be represented by finitely many numbers, e.g. $x \in \mathbb{R}^n$ with some $n \in \mathbb{N}$, we speak of a *finite dimensional optimization problem*, otherwise, of an *infinite dimensional optimization problem*. The second might also be referred to as *optimization in function spaces*. Discrete time optimal control problems fall into the first, continuous time optimal control problems into the second class.

Besides the dimension of the decision variable, also the dimension of the constraint functions can be finite or infinite. If an infinite number of inequality constraints is present while the decision variable is finite dimensional, one speaks of a *semi-infinite optimization problem*. This class naturally arises in the context of *robust optimization*, where one wants to find the best choice of the decision variable that satisfies the constraints for all possible values of an unknown but bounded disturbance.

Continuous vs Integer Optimization

A second dividing line concerns the type of decision variables. These can be either *continuous*, like for example real valued vectors $x \in \mathbb{R}^n$, or any other elements of a smooth manifold. On the other hand, the decision variable can be *discrete*, or *integer valued*, i.e. we have $z \in \mathbb{Z}^n$, or, when a set of binary choices has to be made, $z \in \{0, 1\}^n$. In this case one often also speaks of *combinatorial optimization*. If an optimization problem has both, continuous and integer variables, it is called a *mixed-integer optimization problem*.

An important class of continuous optimization problems are the so called *nonlinear programs (NLP)*. They can be stated in the form

$$\begin{aligned} &\underset{x \in \mathbb{R}^n}{\text{minimize}} && f(x) \end{aligned} \tag{1.17a}$$

$$\text{subject to} \quad g(x) = 0, \tag{1.17b}$$

$$h(x) \leq 0, \tag{1.17c}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $g : \mathbb{R}^n \rightarrow \mathbb{R}^{n_g}$, and $h : \mathbb{R}^n \rightarrow \mathbb{R}^{n_h}$ are assumed to be at least once continuously differentiable. Note that we use function and variable names such as f and x with a very different meaning than before in the context of dynamic systems. In the first part of the lecture we discuss algorithms to solve this kind of optimization problems, and the discrete time optimal control problems treated in this lecture can also be regarded as a specially structured form of NLPs. Two important subclasses of NLPs are the *linear programs (LP)*, which have affine problem functions f, g, h , and the *quadratic programs (QP)*, which have affine constraint functions g, h and a more general linear quadratic objective $f(x) = c^T x + \frac{1}{2} x^T H x$ with a symmetric matrix $H \in \mathbb{R}^{n \times n}$.

A large class of mixed-integer optimization problems are the so called *mixed integer nonlinear programs (MINLP)*, which can be stated as

$$\begin{aligned} &\underset{\substack{x \in \mathbb{R}^n \\ z \in \mathbb{Z}^m}}{\text{minimize}} && f(x, z) \end{aligned} \tag{1.18a}$$

$$\text{subject to} \quad g(x, z) = 0, \tag{1.18b}$$

$$h(x, z) \leq 0. \tag{1.18c}$$

Among the MINLPs, an important special case arises if the problem functions f, g, h are affine in both variables, x and z , which is called a *mixed integer linear program (MILP)*. If the objective is allowed to be linear quadratic, one speaks of a *mixed integer quadratic program (MIQP)*. If in an MILP only integer variables are present, one usually just calls it an *integer program (IP)*. The field of (linear) integer programming is huge and has powerful algorithms available. Most problems in logistics fall into this class, a famous example being the *travelling salesman problem*, which concerns the shortest closed path that one can travel through a given number of towns, visiting each town exactly once.

An interesting class of mixed-integer optimization problems arises in the context of optimal control of hybrid dynamic systems, which in the discrete time case can be regarded a special case of MINLP. In continuous time, we enter the field of infinite dimensional mixed-integer optimization, often also called *Mixed-integer optimal control problems (MIOCP)*.

Convex vs Nonconvex Optimization

Arguably the most important dividing line in the world of optimization is between convex and nonconvex optimization problems. Convex optimization problems are a subclass of the continuous optimization problems and arise if the objective function is a convex function and the set of feasible points a convex set. In this case one can show that any *local solution*, i.e. values for the decision variables that lead to the best possible objective value in a neighborhood, is also a *global solution*, i.e. has the best possible objective value among all feasible points.

Practically very important is the fact that convexity of a function or a set can be checked just by checking convexity of its building blocks and if they are constructed in a way that preserves convexity.

Several important subclasses of NLPs are convex, such as LPs. Also QPs are convex if they have a convex objective f . Another example are *Quadratically Constrained Quadratic Programs (QCQP)* which have quadratic inequalities and whose feasible set is the intersection of ellipsoids. Some other optimization problems are convex but do not form part of the NLP family. Two widely used classes are *second-order cone programs (SOCP)* and *semi-definite programs (SDP)* which have linear objective functions but more involved convex feasible sets: for SOCP, it is the set of vectors which have one component that is larger than the Euclidean norm of all the other components and which it is called the *second order cone*, and for SDP it is the set of symmetric matrices that are positive semi-definite, i.e. have all eigenvalues larger than zero. SDPs are often used when designing linear feedback control laws. Also infinite dimensional optimization problems such as optimal control problems in continuous time can be convex under fortunate circumstances.

In this context, it is interesting to note that a sufficient condition for convexity of an optimal control problem is that the underlying dynamic system is linear and that the objective and constraints are convex in controls and states. On the other hand, optimal control problems with underlying nonlinear dynamic systems, which are the focus of this lecture, are usually nonconvex.

Optimization problems with integer variables can never be convex due to the nonconvexity of the set of integers. However, it is of great algorithmic advantage if mixed-integer problems have a convex substructure in the sense that convex problems arise when the integer variables are allowed to also take real values. These so called *convex relaxations* are at the basis of nearly all competitive algorithms for mixed-integer optimization. For example, linear integer programs can be solved very efficiently because their convex relaxations are just linear programs, which are convex and can be solved very efficiently.

1.5 Overview and Notation

The chapters of these lecture notes can roughly be divided into six major areas.

- Introduction
- Optimization Background
- Discrete Time Optimal Control
- Continuous Time Optimal Control
- Model Predictive Control and Moving Horizon Estimation

Notation

Within this lecture we use \mathbb{R} for the set of real numbers, \mathbb{R}_+ for the non-negative ones and \mathbb{R}_{++} for the positive ones, \mathbb{Z} for the set of integers, and \mathbb{N} for the set of natural numbers including zero, i.e. we identify $\mathbb{N} = \mathbb{Z}_+$. The set of real-valued vectors of dimension n is denoted by \mathbb{R}^n , and $\mathbb{R}^{n \times m}$ denotes the set of matrices with n rows and m columns. By default, all vectors are assumed to be column vectors, i.e. we identify $\mathbb{R}^n = \mathbb{R}^{n \times 1}$. We usually use square brackets when presenting vectors and matrices elementwise. Because we will often deal with concatenations of several vectors, say $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^m$, yielding a vector in \mathbb{R}^{n+m} , we abbreviate this concatenation sometimes as (x, y) in the text, instead of the correct but more clumsy equivalent notations $[x^\top, y^\top]^\top$ or

$$\begin{bmatrix} x \\ y \end{bmatrix}.$$

Square and round brackets are also used in a very different context, namely for intervals in \mathbb{R} , where for two real numbers $a < b$ the expression $[a, b] \subset \mathbb{R}$ denotes the closed interval containing both boundaries a and b , while an open boundary is denoted by a round bracket, e.g. (a, b) denotes the open interval and $[a, b)$ the half open interval containing a but not b .

When dealing with norms of vectors $x \in \mathbb{R}^n$, we denote by $\|x\|$ an arbitrary norm, and by $\|x\|_2$ the Euclidean norm, i.e. we have $\|x\|_2^2 = x^\top x$. We denote a weighted Euclidean norm with a positive definite weighting matrix $Q \in \mathbb{R}^{n \times n}$ by $\|x\|_Q$, i.e. we have $\|x\|_Q^2 = x^\top Q x$. The L_1 and L_∞ norms are defined by $\|x\|_1 = \sum_{i=1}^n |x_i|$

and $\|x\|_\infty = \max\{|x_1|, \dots, |x_n|\}$. Matrix norms are the induced operator norms, if not stated otherwise, and the Frobenius norm $\|A\|_F$ of a matrix $A \in \mathbb{R}^{n \times m}$ is defined by $\|A\|_F^2 = \text{trace}(AA^\top) = \sum_{i=1}^n \sum_{j=1}^m A_{ij}A_{ij}$.

When we deal with derivatives of functions f with several real inputs and several real outputs, i.e. functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m, x \mapsto f(x)$, we define the Jacobian matrix $\frac{\partial f}{\partial x}(x)$ as a matrix in $\mathbb{R}^{m \times n}$, following standard conventions. For scalar functions with $m = 1$, we denote the gradient vector as $\nabla f(x) \in \mathbb{R}^n$, a column vector, also following standard conventions. Slightly less standard, we generalize the gradient symbol to all functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ even with $m > 1$, i.e. we generally define in this lecture

$$\nabla f(x) = \frac{\partial f}{\partial x}(x)^\top \in \mathbb{R}^{n \times m}.$$

Using this notation, the first order Taylor series is e.g. written as

$$f(x) = f(\bar{x}) + \nabla f(\bar{x})^\top (x - \bar{x}) + o(\|x - \bar{x}\|)$$

The second derivative, or Hessian matrix will only be defined for scalar functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and be denoted by $\nabla^2 f(x)$.

For square symmetric matrices of dimension n we sometimes use the symbol \mathbb{S}_n , i.e. $\mathbb{S}_n = \{A \in \mathbb{R}^{n \times n} | A = A^\top\}$. For any symmetric matrix $A \in \mathbb{S}_n$ we write $A \succcurlyeq 0$ if it is a positive semi-definite matrix, i.e. all its eigenvalues are larger or equal to zero, and $A \succ 0$ if it is positive definite, i.e. all its eigenvalues are positive. This notation is also used for *matrix inequalities* that allow us to compare two symmetric matrices $A, B \in \mathbb{S}_n$, where we define for example $A \succcurlyeq B$ by $A - B \succcurlyeq 0$.

When using logical symbols, $A \Rightarrow B$ is used when a proposition A implies a proposition B . In words the same is expressed by “If A then B ”. We write $A \Leftrightarrow B$ for “ A if and only if B ”, and we sometimes shorten this to “ A iff B ”, with a double “f”, following standard practice.

Chapter 2

Nonlinear Optimization

”The great watershed in optimization is not between linearity and nonlinearity, but convexity and nonconvexity.”
R. Tyrrell Rockafellar

In this first part of the book we discuss several concepts from the field of mathematical optimization that are important for optimal control. Our focus is on quickly arriving at a point where the specific optimization methods for dynamic systems can be treated, while the same material can be found in much greater detail in many excellent textbooks on numerical optimization such as [NW06].

The reason for keeping this part on optimization self-contained and without explicit reference to optimal control is that this allows us to separate between the general concepts of optimization and those specific to optimal control. For this reason, we use in this part the language and notation that is customary in mathematical optimization. The optimization problem with which we are concerned in this part is the standard *Nonlinear Program (NLP)* that was already stated in the introduction:

$$\begin{aligned} \text{minimize} \quad & f(x) \\ & x \in \mathbb{R}^n \end{aligned} \tag{2.1a}$$

$$\text{subject to} \quad g(x) = 0, \tag{2.1b}$$

$$h(x) \leq 0, \tag{2.1c}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $g : \mathbb{R}^n \rightarrow \mathbb{R}^{n_g}$, and $h : \mathbb{R}^n \rightarrow \mathbb{R}^{n_h}$ are assumed to be twice continuously differentiable. Function f is called the *objective function*, function g is the vector of *equality constraints*, and h the vector of *inequality constraints*. We start with some fundamental definitions. First, we collect all points that satisfy the constraints in one set.

Definition 2 (Feasible set) *The feasible set Ω is the set*

$$\Omega := \{x \in \mathbb{R}^n \mid g(x) = 0, h(x) \leq 0\}.$$

The points of interest in optimization are those feasible points that minimize the objective, and they come in two different variants.

Definition 3 (Global minimum) *The point $x^* \in \mathbb{R}^n$ is a global minimizer if and only if (iff) $x^* \in \Omega$ and $\forall x \in \Omega : f(x) \geq f(x^*)$. The value $f(x^*)$ is the global minimum.*

Unfortunately, the global minimum is usually difficult to find, and most algorithms allow us to only find *local minimizers*, and to verify optimality only locally.

Definition 4 (Local minimum) *The point $x^* \in \mathbb{R}^n$ is a local minimizer iff $x^* \in \Omega$ and there exists a neighborhood \mathcal{N} of x^* (e.g., an open ball around x^*) so that $\forall x \in \Omega \cap \mathcal{N} : f(x) \geq f(x^*)$. The value $f(x^*)$ is a local minimum.*

In order to be able to state the optimality conditions that allow us to check if a candidate point x^* is a local minimizer or not, we need to describe the feasible set in the neighborhood of x^* . It turns out that not all inequality constraints need to be considered locally, but only the *active* ones.

Definition 5 (Active Constraints and Active Set) An inequality constraint $h_i(x) \leq 0$ is called active at $x^* \in \Omega$ iff $h_i(x^*) = 0$ and otherwise inactive. The index set $\mathcal{A}(x^*) \subset \{1, \dots, n_h\}$ of active inequality constraint indices is called the "active set".

Often, the name *active set* also comprises all equality constraint indices, as equalities could be considered to be always active.

Problem (2.1) is very generic. In Section 2.1 we review some special cases, which still yield large classes of optimization problems. In order to choose the right algorithm for a practical problem, we should know how to classify it and which mathematical structures can be exploited. Replacing an inadequate algorithm by a suitable one can reduce solution times by orders of magnitude. E.g., an important structure is convexity. It allows us to find global minima by searching for local minima only.

For the general case we review the first and second order conditions of optimality in Sections 2.2 and 2.3, respectively.

2.1 Important Special Classes

Linear Optimization

An obvious special case occurs when the functions f , g , and h in (2.1) are linear, resulting in a linear optimization problem (or Linear Program, LP)

$$\begin{aligned} &\text{minimize} && c^\top x \\ &x \in \mathbb{R}^n \end{aligned} \tag{2.2a}$$

$$\text{subject to} \quad Ax - b = 0, \tag{2.2b}$$

$$Cx - d \leq 0. \tag{2.2c}$$

Here, the problem data are $c \in \mathbb{R}^n$, $A \in \mathbb{R}^{n_g \times n}$, $b \in \mathbb{R}^{n_g}$, $C \in \mathbb{R}^{n_h \times n}$, and $d \in \mathbb{R}^{n_h}$.

It is easy to show that one optimal solution of any LP – if the LP does have a solution and is not unbounded – has to be a vertex of the polytope of feasible points. Vertices can be represented and calculated by means of basis solution vectors, with a basis of *active inequality constraints*. Thus, there are only finitely many vertices, giving rise to Simplex algorithms that compare all possible solutions in a clever way. However, naturally also the optimality conditions of Section 2.2 are valid and can be used for algorithms, in particular interior point methods.

Quadratic Optimization

If in the general NLP formulation (2.1) the constraints g, h are affine, and the objective is a linear-quadratic function, we call the resulting problem a Quadratic Optimization Problem or Quadratic Program (QP). A general QP can be formulated as follows.

$$\begin{aligned} &\text{minimize} && c^\top x + \frac{1}{2} x^\top B x \\ &x \in \mathbb{R}^n \end{aligned} \tag{2.3a}$$

$$\text{subject to} \quad Ax - b = 0, \tag{2.3b}$$

$$Cx - d \leq 0. \tag{2.3c}$$

Here, the problem data are $c \in \mathbb{R}^n$, $A \in \mathbb{R}^{n_g \times n}$, $b \in \mathbb{R}^{n_g}$, $C \in \mathbb{R}^{n_h \times n}$, $d \in \mathbb{R}^{n_h}$, as well as the "Hessian matrix" $B \in \mathbb{R}^{n \times n}$. Its name stems from the fact that $\nabla^2 f(x) = B$ for $f(x) = c^\top x + \frac{1}{2} x^\top B x$.

The eigenvalues of B decide on convexity or non-convexity of a QP, i.e., the possibility to solve it in polynomial time to global optimality, or not. If $B \succcurlyeq 0$ we speak of a convex QP, and if $B \succ 0$ we speak of a strictly convex QP. The latter class has the property that it always has unique minimizers.

Convex Optimization

Roughly speaking, a set is convex, if all connecting lines lie inside the set:

Definition 6 (Convex Set) A set $\Omega \subset \mathbb{R}^n$ is convex if

$$\forall x, y \in \Omega, t \in [0, 1] : x + t(y - x) \in \Omega. \tag{2.4}$$

A function is convex, if all secants are above the graph:

Definition 7 (Convex Function) A function $f : \Omega \rightarrow \mathbb{R}$ is convex, if Ω is convex and if

$$\forall x, y \in \Omega, t \in [0, 1] : f(x + t(y - x)) \leq f(x) + t(f(y) - f(x)). \quad (2.5)$$

Note that this definition is equivalent to saying that the Epigraph of f , i.e., the set $\{(x, s) \in \mathbb{R}^n \times \mathbb{R} \mid x \in \Omega, s \geq f(x)\}$, is a convex set.

Definition 8 (Concave Function) A function $f : \Omega \rightarrow \mathbb{R}$ is called “concave” if $(-f)$ is convex.

Note that the feasible set Ω of an optimization problem (2.1) is convex if the function g is affine and the functions h_i are convex, as supported by the following theorem.

Theorem 2 (Convexity of Sublevel Sets) The sublevel set $\{x \in \Omega \mid h(x) \leq 0\}$ of a convex function $h : \Omega \rightarrow \mathbb{R}$ is convex.

Definition 9 (Convex Optimization Problem) An optimization problem with convex feasible set Ω and convex objective function $f : \Omega \rightarrow \mathbb{R}$ is called a convex optimization problem.

Theorem 3 (Local Implies Global Optimality for Convex Problems) For a convex optimization problem, every local minimum is also a global one.

We leave the proofs of Theorems 2 and 3 as an exercise.

There exists a whole algebra of operations that preserve convexity of functions and sets, which is excellently explained in the text books on convex optimization [BTN01, BV04]. Here we only mention an important fact that is related to the positive curvature of a function.

Theorem 4 (Convexity for C^2 Functions) Assume that $f : \Omega \rightarrow \mathbb{R}$ is twice continuously differentiable and Ω convex and open. Then f is convex if and only if for all $x \in \Omega$ the Hessian is positive semi-definite, i.e.,

$$\forall x \in \Omega : \nabla^2 f(x) \succeq 0. \quad (2.6)$$

Again, we leave the proof as an exercise. As an example, the quadratic objective function $f(x) = c^\top x + \frac{1}{2}x^\top Bx$ of (2.3) is convex if and only if $B \succeq 0$, because $\forall x \in \mathbb{R}^n : \nabla^2 f(x) = B$.

2.2 First Order Optimality Conditions

An important question in continuous optimization is if a feasible point $x^* \in \Omega$ satisfies necessary first order optimality conditions. If it does not satisfy these conditions, x^* cannot be a local minimizer. If it does satisfy these conditions, it is a hot candidate for a local minimizer. If the problem is convex, these conditions are even *sufficient* to guarantee that it is a global optimizer. Thus, most algorithms for nonlinear optimization search for such points. The first order condition can only be formulated if a technical “constraint qualification” is satisfied, which in its simplest and numerically most attractive variant comes in the following form.

Definition 10 (LICQ) The linear independence constraint qualification (LICQ) holds at $x^* \in \Omega$ iff all vectors $\nabla g_i(x^*)$ for $i \in \{1, \dots, n_g\}$ and $\nabla h_i(x^*)$ for $i \in \mathcal{A}(x^*)$ are linearly independent.

To give further meaning to the LICQ condition, let us combine all active inequalities with all equalities in a map \tilde{g} defined by stacking all functions on top of each other in a column vector as follows:

$$\tilde{g}(x) = \begin{bmatrix} g(x) \\ h_i(x) (i \in \mathcal{A}(x^*)) \end{bmatrix}. \quad (2.7)$$

LICQ is then equivalent to full row rank of the Jacobian matrix $\frac{\partial \tilde{g}}{\partial x}(x^*)$.

The Karush-Kuhn-Tucker Optimality Conditions

This condition allows us to formulate the famous KKT conditions that are due to Karush [Kar39] and Kuhn and Tucker [KT51].

Theorem 5 (KKT Conditions) *If x^* is a local minimizer of the NLP (2.1) and LICQ holds at x^* then there exist so called multiplier vectors $\lambda \in \mathbb{R}^{n_g}$ and $\mu \in \mathbb{R}^{n_h}$ with*

$$\nabla f(x^*) + \nabla g(x^*)\lambda^* + \nabla h(x^*)\mu^* = 0 \quad (2.8a)$$

$$g(x^*) = 0 \quad (2.8b)$$

$$h(x^*) \leq 0 \quad (2.8c)$$

$$\mu^* \geq 0 \quad (2.8d)$$

$$\mu_i^* h_i(x^*) = 0, \quad i = 1, \dots, n_h. \quad (2.8e)$$

Regarding the notation used in the first line above, please observe that in this script we use the gradient symbol ∇ also for functions g, h with multiple outputs, not only for scalar functions like f . While ∇f is a column vector, in ∇g we collect the gradient vectors of all output components in a matrix which is the transpose of the Jacobian, i.e., $\nabla g(x) := \frac{\partial g}{\partial x}(x)^\top$.

Definition 11 (KKT Point) *We call a triple (x^*, λ^*, μ^*) a “KKT Point” if it satisfies LICQ and the KKT conditions Eqs. (2.8a)-(2.8e).*

Note: The KKT conditions are the First order necessary conditions for optimality (FONC) for constrained optimization, and are thus the equivalent to $\nabla f(x^*) = 0$ in unconstrained optimization. In the special case of convex problems, the KKT conditions are not only *necessary* for a *local* minimizer, but even *sufficient* for a *global* minimizer. In fact, the following extremely important statement holds.

Theorem 6 *Regard a convex NLP and a point x^* at which LICQ holds. Then:*

$$x^* \text{ is a global minimizer} \iff \exists \lambda, \mu \text{ so that the KKT conditions hold.}$$

The Lagrangian Function

Definition 12 (Lagrangian Function) *We define the so called “Lagrangian function” to be*

$$\mathcal{L}(x, \lambda, \mu) = f(x) + \lambda^\top g(x) + \mu^\top h(x). \quad (2.9)$$

Here, we have used again the so called “Lagrange multipliers” or “dual variables” $\lambda \in \mathbb{R}^{n_g}$ and $\mu \in \mathbb{R}^{n_h}$. The Lagrangian function plays a crucial role in both convex and general nonlinear optimization, not only as a practical shorthand within the KKT conditions: using the definition of the Lagrangian, we have (2.8a) $\iff \nabla_x \mathcal{L}(x^*, \lambda^*, \mu^*) = 0$.

Remark 1: In the absence of inequalities, the KKT conditions simplify to $\nabla_x \mathcal{L}(x, \lambda) = 0$, $g(x) = 0$, a formulation that is due to Lagrange and was much earlier known than the KKT conditions.

Remark 2: The KKT conditions require the inequality multipliers μ to be positive, $\mu \geq 0$, while the sign of the equality multipliers λ is arbitrary. An interesting observation is that for a convex problem with f and all h_i convex and g affine, and for $\mu \geq 0$, the Lagrangian function is a convex function in x . This often allows us to explicitly find the unconstrained minimum of the Lagrangian for any given λ and $\mu \geq 0$, which is called the Lagrange dual function, and which can be shown to be an underestimator of the minimum. Maximizing this underestimator over all λ and $\mu \geq 0$ leads to the concepts of weak and strong duality, which we omit here for brevity.

Complementarity

The last three KKT conditions (2.8c)-(2.8e) are called the *complementarity* conditions. For each index i , they define an L-shaped set in the (h_i, μ_i) space. This set is not a smooth manifold but has a non-differentiability at the origin, i.e., if $h_i(x^*) = 0$ and also $\mu_i^* = 0$. This case is called a *weakly active constraint*. Often we want to exclude this case. On the other hand, an active constraint with $\mu_i^* > 0$ is called strictly active.

Definition 13 *Regard a KKT point (x^*, λ^*, μ^*) . We say that strict complementarity holds at this KKT point iff all active constraints are strictly active.*

Strict complementarity is a favourable condition because, together with a second order condition, it implies that the active set is stable against small perturbations. It also makes many theorems easier to formulate and to prove, and is also required to prove convergence of some numerical methods.

2.3 Second Order Optimality Conditions

In case of strict complementarity at a KKT point (x^*, λ^*, μ^*) , the optimization problem can locally be regarded to be a problem with equality constraints only, namely those within the function \tilde{g} defined in Equation (2.7). Though more complex second order conditions can be formulated that are applicable even when strict complementarity does not hold, we restrict ourselves here to this special case.

Theorem 7 (Second Order Optimality Conditions) *Let us regard a point x^* at which LICQ holds together with multipliers λ^*, μ^* so that the KKT conditions (2.8a)-(2.8e) are satisfied and let strict complementarity hold. Regard a basis matrix $Z \in \mathbb{R}^{n \times (n-n_{\tilde{g}})}$ of the null space of $\frac{\partial \tilde{g}}{\partial x}(x^*) \in \mathbb{R}^{n_{\tilde{g}} \times n}$, i.e., Z has full column rank and $\frac{\partial \tilde{g}}{\partial x}(x^*)Z = 0$.*

Then the following two statements hold:

(a) *If x^* is a local minimizer, then $Z^\top \nabla_x^2 \mathcal{L}(x^*, \lambda^*, \mu^*)Z \succcurlyeq 0$.
(Second Order Necessary Condition, short : SONC)*

(b) *If $Z^\top \nabla_x^2 \mathcal{L}(x^*, \lambda^*, \mu^*)Z \succ 0$, then x^* is a local minimizer.
This minimizer is unique in its neighborhood, i.e., a strict local minimizer, and stable against small differentiable perturbations of the problem data. (Second Order Sufficient Condition, short: SOSC)*

The matrix $\nabla_x^2 \mathcal{L}(x^*, \lambda^*, \mu^*)$ plays an important role in optimization algorithms and is called the *Hessian of the Lagrangian*, while its projection on the null space of the Jacobian, $Z^\top \nabla_x^2 \mathcal{L}(x^*, \lambda^*, \mu^*)Z$, is called the *reduced Hessian*.

Quadratic Problems with Equality Constraints

To illustrate the above optimality conditions, let us regard a QP with equality constraints only.

$$\begin{aligned} \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad & c^\top x + \frac{1}{2}x^\top Bx \end{aligned} \quad (2.10a)$$

$$\text{subject to} \quad Ax + b = 0. \quad (2.10b)$$

We assume that A has full row rank i.e., LICQ holds. The Lagrangian is $\mathcal{L}(x, \lambda) = c^\top x + \frac{1}{2}x^\top Bx + \lambda^\top (Ax + b)$ and the KKT conditions have the explicit form

$$c + Bx + A^\top \lambda = 0 \quad (2.11a)$$

$$b + Ax = 0. \quad (2.11b)$$

This is a linear equation system in the variable (x, λ) and can be solved if the so called *KKT matrix*

$$\begin{bmatrix} B & A^\top \\ A & 0 \end{bmatrix}$$

is invertible. In order to assess if the unique solution (x^*, λ^*) of this linear system is a minimizer, we need first to construct a basis Z of the null space of A , e.g., by a full QR factorization of $A^\top = QR$ with $Q = (Y|Z)$ square orthonormal and $R = (\bar{R}^\top | 0)^\top$. Then we can check if the reduced Hessian matrix $Z^\top BZ$ is positive semidefinite. If it is not, the objective function has negative curvature in at least one of the feasible directions and x^* cannot be a minimizer. If on the other hand $Z^\top BZ \succ 0$ then x^* is a strict local minimizer. Due to convexity this would also be the global solution of the QP.

Invertibility of the KKT Matrix and Stability under Perturbations

An important fact is the following. If the second order sufficient conditions for optimality of Theorem 7 (b) hold, then it can be shown that the KKT-matrix

$$\begin{bmatrix} \nabla_x^2 \mathcal{L}(x^*, \lambda^*, \mu^*) & \frac{\partial \tilde{g}}{\partial x}(x^*)^\top \\ \frac{\partial \tilde{g}}{\partial x}(x^*) & \end{bmatrix}$$

is invertible. This implies that the solution is stable against perturbations. To see why, let us regard a perturbed variant of the optimization problem (2.1)

$$\begin{aligned} \text{minimize} \quad & f(x) + \delta_f^\top x \\ \text{subject to} \quad & x \in \mathbb{R}^n \end{aligned} \tag{2.12a}$$

$$g(x) + \delta_g = 0, \tag{2.12b}$$

$$h(x) + \delta_h \leq 0, \tag{2.12c}$$

with small vectors $\delta_f, \delta_g, \delta_h$ of appropriate dimensions that we summarize as $\delta = (\delta_f, \delta_g, \delta_h)$. If a solution exists for $\delta = 0$, the question arises if a solution exists also for small $\delta \neq 0$, and how this solution depends on the perturbation δ . This is answered by the following theorem.

Theorem 8 (SOSC implies Stability of Solutions) *Regard the family of perturbed optimization problems (2.12) and assume that for $\delta = 0$ exists a local solution $(x^*(0), \lambda^*(0), \mu^*(0))$ that satisfies LICQ, the KKT condition, strict complementarity, and the second order sufficient condition of Theorem 7 (b). Then there exists an $\epsilon > 0$ so that for all $\|\delta\| \leq \epsilon$ exists a unique local solution $(x^*(\delta), \lambda^*(\delta), \mu^*(\delta))$ that depends differentiably on δ . This local solution has the same active set as the nominal one, i.e., its inactive constraint multipliers remain zero and the active constraint multipliers remain positive. The solution does not depend on the inactive constraint perturbations. If \tilde{g} is the combined vector of equalities and active inequalities, and $\tilde{\lambda}$ and $\tilde{\delta}_2$ the corresponding vectors of multipliers and constraint perturbations, then the derivative of the solution $(x^*(\delta), \tilde{\lambda}^*(\delta))$ with respect to $(\delta_1, \tilde{\delta}_2)$ is given by*

$$\frac{d}{d(\delta_1, \tilde{\delta}_2)} \begin{bmatrix} x^*(\delta) \\ \tilde{\lambda}^*(\delta) \end{bmatrix} \bigg|_{\delta=0} = - \begin{bmatrix} \nabla_x^2 \mathcal{L}(x^*, \lambda^*, \mu^*) & \frac{\partial \tilde{g}}{\partial x}(x^*)^\top \\ \frac{\partial \tilde{g}}{\partial x}(x^*) & \end{bmatrix}^{-1} \tag{2.13}$$

This differentiability formula follows from differentiation of the necessary optimality conditions of the parametrized optimization problems with respect to $(\delta_1, \tilde{\delta}_2)$

$$\nabla f(x^*(\delta)) + \frac{\partial \tilde{g}}{\partial x}(x^*)^\top \tilde{\lambda} + \delta_1 = 0 \tag{2.14}$$

$$\tilde{g}(x^*(\delta)) + \tilde{\delta}_2 = 0 \tag{2.15}$$

Invertibility of the KKT matrix and stability of the solution under perturbations are very useful facts for the applicability of Newton-type optimization methods that are discussed in the next chapter.

Software: An excellent tool to formulate and solve convex optimization problems in a MATLAB environment is CVX, which is available as open-source code and easy to install.

Software for solving a QP Problem: MATLAB: quadprog. Commercial: CPLEX, MOSEK. Open-source: CVX, qpOASES.

For anyone not really familiar with the concepts of nonlinear optimization that are only very briefly outlined here, it is highly recommended to have a look at the excellent Springer text book “Numerical Optimization” by Jorge Nocedal and Steve Wright [NW06]. Who likes to know more about convex optimization than the much too brief outline given in this script is recommended to have a look at the equally excellent Cambridge University Press text book “Convex Optimization” by Stephen Boyd and Lieven Vandenberghe [BV04], whose PDF is freely available.

Chapter 3

Newton-Type Optimization Algorithms

”Nature and nature’s laws lay hid in night;
God said “Let Newton be” and all was light.”
Alexander Pope

3.1 Equality Constrained Optimization

Let us first regard an optimization problem with only equality constraints,

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && f(x) \end{aligned} \tag{3.1a}$$

$$\text{subject to} \quad g(x) = 0 \tag{3.1b}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}^{n_g}$ are both smooth functions. The idea of the Newton-type optimization methods is to apply a variant of Newton’s method to solve the nonlinear KKT conditions

$$\nabla_x \mathcal{L}(x, \lambda) = 0 \tag{3.2a}$$

$$g(x) = 0 \tag{3.2b}$$

In order to simplify notation, we define

$$w := \begin{bmatrix} x \\ \lambda \end{bmatrix} \text{ and } F(w) := \begin{bmatrix} \nabla_x \mathcal{L}(x, \lambda) \\ g(x) \end{bmatrix} \tag{3.3}$$

with $w \in \mathbb{R}^{n+n_g}$, $F : \mathbb{R}^{n+n_g} \rightarrow \mathbb{R}^{n+n_g}$, so that we can compactly formulate the above nonlinear root finding problem as

$$F(w) = 0. \tag{3.4}$$

Starting from an initial guess w_0 , Newton’s method generates a sequence of iterates $\{w_k\}_{k=0}^{\infty}$ by linearizing the nonlinear equation at the current iterate

$$F(w_k) + \frac{\partial F}{\partial w_k}(w_k)(w - w_k) = 0 \tag{3.5}$$

and obtaining the next iterate as its solution, i.e.

$$w_{k+1} = w_k - \frac{\partial F}{\partial w_k}(w_k)^{-1} F(w_k) \tag{3.6}$$

For equality constrained optimization, the linear system (3.5) has the specific form¹

$$\begin{bmatrix} \nabla_x \mathcal{L}(x_k, \lambda_k) \\ g(x_k) \end{bmatrix} + \underbrace{\begin{bmatrix} \nabla_x^2 \mathcal{L}(x_k, \lambda_k) & \nabla g(x_k) \\ \nabla g(x_k)^T & 0 \end{bmatrix}}_{\text{KKT-matrix}} \begin{bmatrix} x - x_k \\ \lambda - \lambda_k \end{bmatrix} = 0 \quad (3.7)$$

Using the definition

$$\nabla_x \mathcal{L}(x_k, \lambda_k) = \nabla f(x_k) + \nabla g(x_k) \lambda_k \quad (3.8)$$

we see that the contributions depending on the old multiplier λ_k cancel each other, so that the above system is equivalent to

$$\begin{bmatrix} \nabla f(x_k) \\ g(x_k) \end{bmatrix} + \begin{bmatrix} \nabla_x^2 \mathcal{L}(x_k, \lambda_k) & \nabla g(x_k) \\ \nabla g(x_k)^T & 0 \end{bmatrix} \begin{bmatrix} x - x_k \\ \lambda \end{bmatrix} = 0. \quad (3.9)$$

This formulation shows that the data of the linear system only depend on λ_k via the Hessian matrix. We need not use the exact Hessian matrix, but can approximate it with different methods. This leads to the more general class of Newton-type optimization methods. Using any such approximation $B_k \approx \nabla_x^2 \mathcal{L}(x_k, \lambda_k)$, we finally obtain the Newton-type iteration as

$$\begin{bmatrix} x_{k+1} \\ \lambda_{k+1} \end{bmatrix} = \begin{bmatrix} x_k \\ 0 \end{bmatrix} - \begin{bmatrix} B_k & \nabla g(x_k) \\ \nabla g^T(x_k) & 0 \end{bmatrix}^{-1} \begin{bmatrix} \nabla f(x_k) \\ g(x_k) \end{bmatrix} \quad (3.10)$$

The general *Newton-type method* is summarized in Algorithm 1. If we use $B_k = \nabla_x^2 \mathcal{L}(x_k, \lambda_k)$, we recover the *exact Newton method*.

Algorithm 1 Equality constrained full step Newton-type method

Choose: initial guesses x_0, λ_0 , and a tolerance ϵ

Set: $k = 0$

while $\|\nabla \mathcal{L}(x_k, \lambda_k)\| \geq \epsilon$ or $\|g(x_k)\| \geq \epsilon$ **do**

 obtain a Hessian approximation B_k

 get x_{k+1}, λ_{k+1} from (3.10)

$k = k + 1$

end while

3.1.1 Quadratic Model Interpretation

It is easy to show that x_{k+1} and λ_{k+1} from (3.10) can equivalently be obtained from the solution of a QP:

$$\begin{aligned} \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad & \nabla f(x_k)^T (x - x_k) + \frac{1}{2} (x - x_k)^T B_k (x - x_k) \end{aligned} \quad (3.11a)$$

$$\text{subject to} \quad g(x_k) + \nabla g(x_k)^T (x - x_k) = 0 \quad (3.11b)$$

So we can interpret the Newton-type optimization method as a “Sequential Quadratic Programming” (SQP) method, where we find in each iteration the solution x^{QP} and λ^{QP} of the above QP and take it as the next NLP solution guess and linearization point x_{k+1} and λ_{k+1} . This interpretation will turn out to be crucial when we treat inequality constraints. But let us first discuss what methods exist for the choice of the Hessian approximation B_k .

¹Recall that in this script we use the convention $\nabla g(x) := \frac{\partial g}{\partial x}(x)^T$ that is consistent with the definition of the gradient $\nabla f(x)$ of a scalar function f being a column vector.

3.1.2 The Exact Newton Method

The first and obvious way to obtain B_k is to use the exact Newton method and just set

$$B_k := \nabla_x^2 \mathcal{L}(x_k, \lambda_k)$$

But how can this matrix be computed? Many different ways for computing this second derivative exist. The most straightforward way is a finite difference approximation where we perturb the evaluation of $\nabla \mathcal{L}$ in the direction of all unit vectors $\{e_i\}_{i=1}^n$ by a small quantity $\delta > 0$. This yields each time one column of the Hessian matrix, as

$$\nabla_x^2 \mathcal{L}(x_k, \lambda_k) e_i = \frac{\nabla_x \mathcal{L}(x_k + \delta e_i, \lambda_k) - \nabla_x \mathcal{L}(x_k, \lambda_k)}{\delta} + O(\delta) \quad (3.12)$$

Unfortunately, the evaluation of the numerator of this quotient suffers from numerical cancellation, so that δ cannot be chosen arbitrarily small, and the maximum attainable accuracy for the derivative is $\sqrt{\epsilon}$ if ϵ is the accuracy with which the gradient $\nabla_x \mathcal{L}$ can be obtained. Thus, we loose half the valid digits. If $\nabla_x \mathcal{L}$ was itself already approximated by finite differences, this means that we have lost three quarters of the originally valid digits. More accurate and also faster ways to obtain derivatives of arbitrary order will be presented in the chapter on algorithmic differentiation.

Local convergence rate: The exact Newton method has a *quadratic convergence rate*, i.e. $\|w_{k+1} - w^*\| \leq \frac{\omega}{2} \|w_k - w^*\|^2$. This means that the number of accurate digits doubles in each iteration. As a rule of thumb, once a Newton method is in its area of quadratic convergence, it needs at maximum 6 iterations to reach the highest possible precision.

3.1.3 The Constrained Gauss-Newton Method

Let us regard the special case that the objective $f(x)$ has a nonlinear least-squares form, i.e. $f(x) = \frac{1}{2} \|R(x)\|_2^2$ with some function $R : \mathbb{R}^n \rightarrow \mathbb{R}^{n_R}$. In this case we can use a very powerful Newton-type method which approximates the Hessian B_k using only first order derivatives. It is called the *Gauss-Newton method*. To see how it works, let us thus regard the nonlinear least-squares problem

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && \frac{1}{2} \|R(x)\|_2^2 \end{aligned} \quad (3.13a)$$

$$\text{subject to} \quad g(x) = 0 \quad (3.13b)$$

The idea of the Gauss-Newton method is to linearize at a given iterate x_k both problem functions R and g , in order to obtain the following approximation of the original problem.

$$\begin{aligned} & \underset{x \in \mathbb{R}^n}{\text{minimize}} && \frac{1}{2} \|R(x_k) + \nabla R(x_k)^T (x - x_k)\|_2^2 \end{aligned} \quad (3.14a)$$

$$\text{subject to} \quad g(x_k) + \nabla g(x_k)^T (x - x_k) = 0 \quad (3.14b)$$

This is a convex QP which can easily be seen by noting that the objective (3.14a) is equal to

$$\frac{1}{2} R(x_k)^T R(x_k) + (x - x_k)^T \underbrace{\nabla R(x_k) R(x_k)}_{=\nabla f(x_k)} + \frac{1}{2} (x - x_k)^T \underbrace{\nabla R(x_k) \nabla R(x_k)^T}_{=: B_k} (x - x_k)$$

which is convex because $B_k \succcurlyeq 0$. Note that the constant term does not influence the solution and can be dropped. Thus, the Gauss-Newton subproblem (3.14) is identical to the SQP subproblem (3.11) with a special choice of the Hessian approximation, namely

$$B_k := \nabla R(x_k) \nabla R(x_k)^T = \sum_{i=1}^{n_R} \nabla R_i(x_k) \nabla R_i(x_k)^T$$

Note that no multipliers λ_k are needed in order to compute B_k . In order to assess the quality of the Gauss-Newton Hessian approximation, let us compare it with the exact Hessian, that is given by

$$\begin{aligned} \nabla_x^2 \mathcal{L}(x, \lambda) &= \sum_{i=1}^{n_R} \nabla R_i(x_k) \nabla R_i(x_k)^T + \sum_{i=1}^{n_F} R_i(x) \nabla^2 R_i(x) + \sum_{i=1}^{n_g} \lambda_i \nabla^2 g_i(x) \end{aligned} \quad (3.15)$$

$$= B_k + O(\|R(x_k)\|) + O(\|\lambda\|). \quad (3.16)$$

One can show that in the solution of a problem holds $\|\lambda^*\| = O(\|R(x^*)\|)$. Thus, in the vicinity of the solution, the difference between the exact Hessian and the Gauss-Newton approximation B_k is of order $O(\|R(x^*)\|)$.

Local convergence rate: The Gauss-Newton method converges *linearly*, $\|w_{k+1} - w^*\| \leq \kappa \|w_k - w^*\|$ with a contraction rate $\kappa = O(\|R(x^*)\|)$. Thus, it converges fast if the residuals $R_i(x^*)$ are small, or equivalently, if the objective is close to zero, which is our desire in least-squares problems. In estimation problems, a low objective corresponds to a “good fit”. Thus the Gauss-Newton method is only attracted by local minima with a small function value, a favourable feature in practice.

3.1.4 Hessian Approximation by Quasi-Newton BFGS Updates

Besides the exact Hessian and the Gauss-Newton Hessian approximation, there is another widely used way to obtain a Hessian approximation B_k within the Newton-type framework. It is based on the observation that the evaluation of $\nabla_x \mathcal{L}$ at different points can deliver curvature information that can help us to estimate $\nabla_x^2 \mathcal{L}$, similar as it can be done by finite differences, cf. Equation (3.12), but without any extra effort per iteration besides the evaluation of $\nabla f(x_k)$ and $\nabla g(x_k)$ that we need anyway in order to compute the next step. Quasi-Newton Hessian update methods use the previous Hessian approximation B_k , the step $s_k := x_{k+1} - x_k$ and the gradient difference $y_k := \nabla_x \mathcal{L}(x_{k+1}, \lambda_{k+1}) - \nabla_x \mathcal{L}(x_k, \lambda_{k+1})$ in order to obtain the next Hessian approximation B_{k+1} . As in the finite difference formula (3.12), this approximation shall satisfy the *secant condition*

$$B_{k+1} s_k = y_k \quad (3.17)$$

but because we only have one single direction s_k , this condition does not uniquely determine B_{k+1} . Thus, among all matrices that satisfy the secant condition, we search for the ones that minimize the distance to the old B_k , measured in some suitable norm. The most widely used Quasi-Newton update formula is the Broyden-Fletcher-Goldfarb-Shanno (BFGS) update that can be shown to minimize a weighted Frobenius norm. It is given by the explicit formula:

$$B_{k+1} = B_k - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k} + \frac{y_k y_k^T}{s_k^T y_k}. \quad (3.18)$$

Local convergence rate: It can be shown that $B_k \rightarrow \nabla_x^2 \mathcal{L}(x^*, \lambda^*)$ in the relevant directions, so that *superlinear convergence* is obtained with the BFGS method, i.e. $\|w_{k+1} - w^*\| \leq \kappa_k \|w_k - w^*\|$ with $\kappa_k \rightarrow 0$.

3.2 Local Convergence of Newton-Type Methods

We have seen three examples for Newton-type optimization methods which have different rates of local convergence if they are started close to a solution. They are all covered by the following theorem that exactly states the conditions that are necessary in order to obtain local convergence.

Theorem 9 (Newton-Type Convergence) *Regard the root finding problem*

$$F(w) = 0, \quad F : \mathbb{R}^n \rightarrow \mathbb{R}^n \quad (3.19)$$

with w^* a local solution satisfying $F(w^*) = 0$ and a Newton-type iteration $w_{k+1} = w_k - M_k^{-1} F(w_k)$ with $M_k \in \mathbb{R}^{n \times n}$ invertible for all k . Let us assume a Lipschitz condition on the Jacobian $J(w) := \frac{\partial F}{\partial w}(w)$ as follows:

$$\|M_k^{-1}(J(w_k) - J(w))\| \leq \omega \|w_k - w^*\| \quad (3.20)$$

Let us also assume a bound on the distance of approximation M_k from the true Jacobian $J(w_k)$:

$$\|M_k^{-1}(J(w_k) - M_k)\| \leq \kappa_k \quad (3.21)$$

where $\kappa_k \leq \kappa$ with $\kappa < 1$. Finally, we assume that the initial guess w_0 is sufficiently close to the solution w^* ,

$$\|w_0 - w^*\| \leq \frac{2}{\omega}(1 - \kappa). \quad (3.22)$$

Then $w_k \rightarrow w^*$ with the following linear contraction in each iteration:

$$\|w_{k+1} - w^*\| \leq \left(\kappa_k + \frac{\omega}{2} \|w_k - w^*\| \right) \cdot \|w_k - w^*\|. \quad (3.23)$$

If $\kappa_k \rightarrow 0$, this results in a superlinear convergence rate, and if $\kappa = 0$ quadratic convergence results.

Noting that in Newton-type optimization we have

$$J(w_k) = \begin{bmatrix} \nabla_x^2 \mathcal{L}(x_k, \lambda_k) & \frac{\partial g}{\partial x}(x_k)^T \\ \frac{\partial g}{\partial x}(x_k) & 0 \end{bmatrix} \quad (3.24)$$

$$M_k = \begin{bmatrix} B_k & \frac{\partial g}{\partial x}(x_k)^T \\ \frac{\partial g}{\partial x}(x_k) & 0 \end{bmatrix} \quad (3.25)$$

$$J(w_k) - M_k = \begin{bmatrix} \nabla_x^2 \mathcal{L}(\cdot) - B_k & 0 \\ 0 & 0 \end{bmatrix} \quad (3.26)$$

the above theorem directly implies the three convergence rates that we had already mentioned.

Corollary 2 *Newton-type optimization methods converge*

- quadratically if $B_k = \nabla_x^2 \mathcal{L}(x_k, \lambda_k)$ (exact Newton),
- superlinearly if $B_k \rightarrow \nabla_x^2 \mathcal{L}(x_k, \lambda_k)$ (BFGS),
- linearly if $\|B_k - \nabla_x^2 \mathcal{L}(x_k, \lambda_k)\|$ is small (Gauss-Newton).

3.3 Inequality Constrained Optimization

When a nonlinear optimization problem with inequality constraints shall be solved, two big families of methods exist, first, nonlinear interior point (IP), and second, sequential quadratic programming (SQP) methods. Both aim at solving the KKT conditions (2.8) which include the non-smooth complementarity conditions, but have different ways to deal with this non-smoothness.

3.3.1 Interior Point Methods

The basic idea of an interior point method is to replace the non-smooth L-shaped set resulting from the complementarity conditions with a smooth approximation, typically a hyperbola. Thus, a smoothing constant $\tau > 0$ is introduced and the KKT conditions are replaced by the smooth equation system

$$\nabla f(x^*) + \nabla g(x^*)\lambda^* + \nabla h(x^*)\mu^* = 0 \quad (3.27a)$$

$$g(x^*) = 0 \quad (3.27b)$$

$$\mu_i^* h_i(x^*) + \tau = 0, \quad i = 1, \dots, n_h. \quad (3.27c)$$

Note that the last equation ensures that $-h_i(x^*)$ and μ_i^* are both strictly positive and on a hyperbola.² For τ very small, the L-shaped set is very closely approximated by the hyperbola, but the nonlinearity is increased. Within an interior point method, we usually start with a large value of τ and solve the resulting nonlinear equation system by a Newton method, and then iteratively decrease τ , always using the previously obtained solution as initialization for the next one.

One way to interpret the above smoothened KKT-conditions is to use the last condition to eliminate $\mu_i^* = -\frac{\tau}{h_i(x^*)}$ and to insert this expression into the first equation, and to note that $\nabla_x (\log(-h_i(x))) = \frac{1}{h_i(x)} \nabla h_i(x)$. Thus, the above smooth form of the KKT conditions is nothing else than the optimality conditions of a *barrier problem*

$$\begin{aligned} \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad & f(x) - \tau \sum_{i=1}^{n_h} \log(-h_i(x)) \end{aligned} \quad (3.28a)$$

$$\text{subject to} \quad g(x) = 0. \quad (3.28b)$$

Note that the objective function of this problem tends to infinity when $h_i(x) \rightarrow 0$. Thus, even for very small $\tau > 0$, the barrier term in the objective function will prevent the inequalities to be violated. The *primal barrier method* just solves the above barrier problem with a Newton-type optimization method for equality constrained

²In the numerical solution algorithms for this system, we have to ensure that the iterates do not jump to a second hyperbola of infeasible shadow solutions, by shortening steps if necessary to keep the iterates in the correct quadrant.

optimization for each value of τ . Though easy to implement and to interpret, it is not necessarily the best for numerical treatment, among other because its KKT matrices become very ill-conditioned for small τ . This is not the case for the *primal-dual IP method* that solves the full nonlinear equation system (3.27) including the dual variables μ .

For convex problems, very strong complexity results exist that are based on *self-concordance* of the barrier functions and give upper bounds on the total number of Newton iterations that are needed in order to obtain a numerical approximation of the global solution with a given precision. When an IP method is applied to a general NLP that might be non-convex, we can of course only expect to find a local solution, but convergence to KKT points can still be proven, and these *nonlinear IP methods* perform very well in practice.

Software: A very widespread and successful implementation of the nonlinear IP method is the open-source code IPOPT [WB06, WB09]. Though IPOPT can be applied to convex problems and will yield the global solution, dedicated IP methods for different classes of convex optimization problems can exploit more problem structure and will solve these problems faster and more reliably. Most commercial LP and QP solution packages such as CPLEX or MOSEK make use of IP methods, as well as many open-source implementations such as the sparsity exploiting QP solver OOQP.

3.3.2 Sequential Quadratic Programming (SQP) Methods

Another approach to address NLPs with inequalities is inspired by the quadratic model interpretation that we gave before for Newton-type methods. It is called *Sequential Quadratic Programming (SQP)* and solves in each iteration an inequality constrained QP that is obtained by linearizing the objective and constraint functions:

$$\begin{aligned} \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad & \nabla f(x_k)^T(x - x_k) + \frac{1}{2}(x - x_k)^T B_k(x - x_k) \end{aligned} \quad (3.29a)$$

$$\text{subject to} \quad g(x_k) + \nabla g(x_k)^T(x - x_k) = 0 \quad (3.29b)$$

$$h(x_k) + \nabla h(x_k)^T(x - x_k) \leq 0 \quad (3.29c)$$

Note that the active set is automatically discovered by the QP solver and can change from iteration to iteration. However, under strict complementarity, it will be the same as in the true NLP solution x^* once the SQP iterates x_k are in the neighborhood of x^* .

As before for equality constrained problems, the Hessian B_k can be chosen in different ways. First, in the *exact Hessian SQP method* we use $B_k = \nabla_x^2 \mathcal{L}(x_k, \lambda_k, \mu_k)$, and it can be shown that under the second order sufficient conditions (SOSC) of Theorem 7 (b), this method has locally quadratic convergence. Second, in the case of a least-squares objective $f(x) = \frac{1}{2}\|R(x)\|_2^2$, we can use the Gauss-Newton Hessian approximation $B_k = \nabla R(x_k)\nabla R(x_k)^T$, yielding linear convergence with a contraction rate $\kappa = O(\|R(x^*)\|)$. Third, quasi-Newton updates such as BFGS can directly be applied, using the Lagrange gradient difference $y_k := \nabla_x \mathcal{L}(x_{k+1}, \lambda_{k+1}, \mu^{k+1}) - \nabla_x \mathcal{L}(x_k, \lambda_{k+1}, \mu^{k+1})$ in formula (3.18).

Note that in each iteration of an SQP method, an inequality constrained QP needs to be solved, but that we did not mention yet how this should be done. One way would be to apply an IP method tailored to QP problems. This is indeed done, in particular within SQP methods for large sparse problems. Another way is to use a QP solver that is based on an *active set method*, as sketched in the next subsection.

Software: A successful and sparsity exploiting SQP code is SNOPT [GMS97]. Many optimal control packages such as MUSCOD-II [LSBS03] or the open-source package ACADO [HFD11] contain at their basis structure exploiting SQP methods. Also the MATLAB solver `fmincon` is based on an SQP algorithm.

3.3.3 Active Set Methods

Another class of algorithms to address optimization problems with inequalities, the *active set methods*, are based on the following observation: if we would know the active set, then we could solve directly an equality constrained optimization problem and obtain the correct solution. The main task is thus to find the correct active set, and an active set method iteratively refines a guess for the active set that is often called the *working set*, and solves in each iteration an equality constrained problem. This equality constrained problem is particularly easy to solve in the case of linear inequality constraints, for example in LPs and QPs. Many very successful LP solvers are based on

an active set method which is called the *simplex algorithm*, whose invention by Dantzig [Dan63] was one of the great breakthroughs in the field of optimization. Also many successful QP solvers are based on active set methods. A major advantage of active set strategies is that they can very efficiently be warm-started under circumstances where a series of related problems have to be solved, e.g. within an SQP method, within codes for mixed integer programming, or in the context of model predictive control (MPC) [FBD08].

3.4 Globalisation Strategies

In all convergence results for the Newton-type algorithms stated so far, we had to assume that the initialization was sufficiently close to the true solution in order to make the algorithm converge, which is not always the case. An approach often used to overcome this problem is to use a *homotopy* between a problem we have already solved and the problem we want to solve: in this procedure, we start with the known solution and then proceed slowly, step by step modifying the relevant problem parameters, towards the problem we want to solve, each time converging the Newton-type algorithm and using the obtained solution as initial guess for the next problem. Applying a homotopy requires more user input than just the specification of the problem, so most available Newton-type optimization algorithms have so called *globalisation strategies*. Most of these strategies can be interpreted as automatically generated homotopies.

In the ideal case, a globalisation strategy ensures *global convergence*, i.e. the Newton-type iterations converge to a local minimum from arbitrary initial guesses. Note that the terms *global convergence* and *globalisation strategies* have nothing to do with *global optimization*, which is concerned with finding global minima for non-convex problems.

Here, we only touch the topic of globalisation strategies very superficially, and for all details we refer to textbooks on nonlinear optimization and recommend in particular [NW06].

Two ingredients characterize a globalization strategy: first, a measure of progress, and second, a way to ensure that progress is made in each iteration.

3.4.1 Measuring Progress: Merit Functions and Filters

When two consecutive iterations of a Newton-type algorithm for solution of a constrained optimization problem shall be compared with each other it is not trivial to judge if progress is made by the step. The objective function might be improved, while the constraints might be violated more, or conversely. A *merit function* introduces a scalar measure of progress with the property that each local minimum of the NLP is also a local minimum of the merit function. Then, during the optimization routine, it can be monitored if the next Newton-type iteration gives a better merit function than the iterate before. If this is not the case, the step can be rejected or modified.

A widely used merit function is the *exact L1 merit function*

$$T_1(x) = f(x) + \sigma(\|g(x)\|_1 + \|h^+(x)\|_1)$$

with $f(x)$ the objective, $g(x)$ the residual vector of the equality constraints, and $h^+(x)$ the violations of the inequality constraints, i.e. $h_i^+(x) = \max(0, h_i(x))$ for $i = 1, \dots, n_h$. Note that the L1 penalty function is non-smooth. If the penalty parameter σ is larger than the largest modulus of any Lagrange multiplier at a local minimum and KKT point (x^*, λ^*, μ^*) , i.e. if $\sigma > \max(\|\lambda^*\|_\infty, \|\mu^*\|_\infty)$, then the L1 penalty is exact in the sense that x^* also is a local minimum of $T_1(x)$. Thus, in a standard procedure we require that in each iteration a descent is achieved, i.e. $T_1(x_{k+1}) < T_1(x_k)$, and if it is not the case, the step is rejected or modified, e.g. by a line search or a trust region method.

A disadvantage of requiring a descent in the merit function in each iteration is that the full Newton-type steps might be too often rejected, which can slow down the speed of convergence. Remedies to are e.g. a “watchdog technique” that starting at some iterate x_k allows up to $M - 1$ full Newton-type steps without merit function improvement if the M th iterate is better, i.e. if at the end holds $T_1(x_{k+M}) < T_1(x_k)$, so that the generosity was justified. If this is not the case, the algorithm jumps back to x_k and enforces strict descent for a few iterations.

A different approach that avoids the arbitrary weighting of objective function and constraint violations within a merit function and often allows to accept more full Newton-steps comes in the form of *filter methods*. They regard the pursuit of a low objective function and low constraint violations as two equally important aims, and accept each step that leads to an improvement in at least one of the two, compared to all previous iterations. To ensure this, a so called *filter* keeps track of the best objective and constraint violation pairs that have been achieved so far, and the method rejects only those steps that are *dominated by the filter* i.e., for which one of the previous iterates

had both, a better objective and a lower constraint violation. Otherwise the new iterate is accepted and added to the filter, possibly dominating some other pairs in the filter that can then be removed from the filter. Filter methods are popular because of the fact that they often allow the full Newton-step and still have a global convergence guarantee.

3.4.2 Ensuring Progress: Line Search and Trust-Region Methods

If a full Newton-type step does not lead to progress in the chosen measure, it needs to be rejected. But how can a step be generated that is acceptable? Two very popular ways for this exist, one called *line search*, the other *trust region*.

A line search method takes the result of the QP subproblem as a trial step only, and shortens the step if necessary. If $(x_k^{\text{QP}}, \lambda_k^{\text{QP}}, \mu_k^{\text{QP}})$ is the solution of the QP at an SQP iterate x_k , it can be shown (if the QP multipliers are smaller than σ) that the step vector or *search direction* $(x_k^{\text{QP}} - x_k)$ is a descent direction for the L1 merit function T_1 , i.e. descent in T_1 can be enforced by performing, instead of the full SQP step $x_{k+1} = x_k^{\text{QP}}$, a shorter step

$$x_{k+1} = x_k + t(x_k^{\text{QP}} - x_k)$$

with a damping factor or *step length* $t \in (0, 1]$. One popular way to ensure global convergence with help of a merit function is to require in each step the so called *Armijo condition*, a tightened descent condition, and to perform a *backtracking* line search procedure that starts by trying the full step ($t = 1$) first and iteratively shortens the step by a constant factor ($t \leftarrow t/\beta$ with $\beta > 1$) until this descent condition is satisfied. As said, the L1 penalty function has the desirable property that the search direction is a descent direction so that the Armijo condition will eventually be satisfied if the step is short enough. Line-search methods can also be combined with a filter as a measure of progress, instead of the merit function.

An alternative way to ensure progress is to modify the QP subproblem by adding extra constraints that enforce the QP solution to be in a small region around the previous iterate, the *trust region*. If this region is small enough, the QP solution shall eventually lead to an improvement of the merit function, or be acceptable by the filter. The underlying philosophy is that the linearization is only valid in a region around the linearization point and only here we can expect our QP approximation to be a good model of the original NLP. Similar as for line search methods with the L1 merit function, it can be shown for suitable combinations that the measure of progress can always be improved when the trust region is made small enough. Thus, a trust region algorithm checks in each iteration if enough progress was made to accept the step and adapts the size of the trust region if necessary.

As said above, a more detailed description of different globalisation strategies is given in [NW06].

Chapter 4

Calculating Derivatives

”Progress is measured by the degree of differentiation within a society.”
Herbert Read

Derivatives of computer coded functions are needed everywhere in optimization. In order to just check optimality of a point, we need already to compute the gradient of the Lagrangian function. Within Newton-type optimization methods, we need the full Jacobian of the constraint functions. If we want to use an exact Hessian method, we even need second order derivatives of the Lagrangian.

There are many ways to compute derivatives: Doing it by hand is error prone and nearly impossible for longer evaluation codes. Computer algebra packages like Mathematica or Maple can help us, but require that the function is formulated in their specific language. More annoyingly, the resulting derivative code can become extremely long and slow to evaluate.

On the other hand, *finite differences* can always be applied, even if the functions are only available as black-box codes. They are easy to implement and relatively fast, but they necessarily lead to a loss of precision of half the valid digits, as they have to balance the numerical errors that originate from Taylor series truncation and from finite precision arithmetic. Second derivatives obtained by recursive application of finite differences are even more inaccurate. The best perturbation sizes are difficult to find in practice. Note that the computational cost to compute the gradient $\nabla f(x)$ of a scalar function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is $(n + 1)$ times the cost of one function evaluation.

We will see that a more efficient way exists to evaluate the gradient of a scalar function, which is also more accurate. The technology is called *algorithmic differentiation (AD)* and requires in principle nothing more than that the function is available in the form of source code in a standard programming language such as C, C++ or FORTRAN.

4.1 Algorithmic Differentiation (AD)

Algorithmic differentiation uses the fact that each differentiable function $F : \mathbb{R}^n \rightarrow \mathbb{R}^{n_F}$ is composed of several *elementary operations*, like multiplication, division, addition, subtraction, sine-functions, exp-functions, etc. If the function is written in a programming language like e.g. C, C++ or FORTRAN, special AD-tools can have access to all these elementary operations. They can process the code in order to generate new code that does not only deliver the function value, but also desired derivative information. Algorithmic differentiation was traditionally called *automatic differentiation*, but as this might lead to confusion with symbolic differentiation, most AD people now prefer the term *algorithmic differentiation*, which fortunately has the same abbreviation. A good and authoritative textbook on AD is [GW08].

In order to see how AD works, let us regard a function $F : \mathbb{R}^n \rightarrow \mathbb{R}^{n_F}$ that is composed of a sequence of m elementary operations. While the inputs x_1, \dots, x_n are given before, each elementary operation ϕ_i , $i = 0, \dots, m - 1$ generates another intermediate variable, x_{n+i+1} . Some of these intermediate variables are used as output of the code, but in principle we can regard all variables as possible outputs, which we do here. This way to regard a function evaluation is stated in Algorithm 4.1 and illustrated in Example 3 below.

Example 3 (Function Evaluation via Elementary Operations) *Let us regard the simple scalar function*

$$f(x_1, x_2, x_3) = \sin(x_1 x_2) + \exp(x_1 x_2 x_3)$$

Input: x_1, \dots, x_n

Output: x_1, \dots, x_{n+m}

for $i = 0$ to $m - 1$ **do**

$x_{n+i+1} \leftarrow \phi_i(x_1, \dots, x_{n+i})$

end for

Note: each ϕ_i depends on only one or two out of $\{x_1, \dots, x_{n+i}\}$.

with $n = 3$. We can decompose this function into $m = 5$ elementary operations, namely

$$\begin{aligned} x_4 &= x_1 x_2 \\ x_5 &= \sin(x_4) \\ x_6 &= x_4 x_3 \\ x_7 &= \exp(x_6) \\ x_8 &= x_5 + x_7 \end{aligned}$$

Thus, if the $n = 3$ inputs x_1, x_2, x_3 are given, the $m = 5$ elementary operations ϕ_0, \dots, ϕ_4 compute the $m = 5$ intermediate quantities, x_4, \dots, x_8 , the last of which is our desired scalar output, x_{n+m} .

The idea of AD is to use the chain rule and differentiate each of the elementary operations ϕ_i separately. There are two modes of AD, on the one hand the “forward” mode of AD, and on the other hand the “backward”, “reverse”, or “adjoint” mode of AD. In order to present both of them in a consistent form, we first introduce an alternative formulation of the original user function, that uses augmented elementary functions, as follows¹: we introduce new augmented states

$$\tilde{x}_0 = x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad \tilde{x}_1 = \begin{bmatrix} x_1 \\ \vdots \\ x_{n+1} \end{bmatrix}, \quad \dots, \quad \tilde{x}_m = \begin{bmatrix} x_1 \\ \vdots \\ x_{n+m} \end{bmatrix} \quad (4.1)$$

as well as new augmented elementary functions $\tilde{\phi}_i : \mathbb{R}^{n+i} \rightarrow \mathbb{R}^{n+i+1}$, $\tilde{x}_i \mapsto \tilde{x}_{i+1} = \tilde{\phi}_i(\tilde{x}_i)$ with

$$\tilde{\phi}_i(\tilde{x}_i) = \begin{bmatrix} x_1 \\ \vdots \\ x_{n+i} \\ \phi_i(x_1, \dots, x_{n+i}) \end{bmatrix}, \quad i = 0, \dots, m-1. \quad (4.2)$$

Thus, the whole evaluation tree of the function can be summarized as a concatenation of these augmented functions followed by a multiplication with a “selection matrix” C that selects from \tilde{x}_m the final outputs of the computer code.

$$F(x) = C \cdot \tilde{\phi}_{m-1}(\tilde{\phi}_{m-2}(\dots \tilde{\phi}_1(\tilde{\phi}_0(x))))).$$

The full Jacobian of F , that we denote by $J_F = \frac{\partial F}{\partial x}$ is given by the chain rule as the product of the Jacobians of the augmented elementary functions $\tilde{J}_i = \frac{\partial \tilde{\phi}_i}{\partial \tilde{x}_i}$, as follows:

$$J_F = C \cdot \tilde{J}_{m-1} \cdot \tilde{J}_{m-2} \cdots \tilde{J}_1 \cdot \tilde{J}_0. \quad (4.3)$$

Note that each elementary Jacobian is given as a unit matrix plus one extra row. Also note that the extra row that is here marked with stars $*$ has at maximum two non-zero entries.

$$\tilde{J}_i = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \\ * & * & * & * \end{bmatrix}.$$

¹MD thanks Carlo Savorgnan for having outlined to him this way of presenting forward and backward AD

Input: $\dot{x}_1, \dots, \dot{x}_n$ and all partial derivatives $\frac{\partial \phi_i}{\partial x_j}$

Output: $\dot{x}_1, \dots, \dot{x}_{n+m}$

for $i = 0$ to $m - 1$ **do**
 $\dot{x}_{n+i+1} \leftarrow \sum_{j=1}^{n+i} \frac{\partial \phi_i}{\partial x_j} \dot{x}_j$
end for

Note: each sum consist of only one or two non-zero entries.

For the generation of first order derivatives, algorithmic differentiation uses two alternative ways to evaluate the product of these Jacobians, the *forward* and the *backward mode* as described in the next two sections.

4.2 The Forward Mode of AD

In forward AD we first define a *seed vector* $p \in \mathbb{R}^n$ and then evaluate the directional derivative $J_F p$ in the following way:

$$J_F p = C \cdot (\tilde{J}_{m-1} \cdot (\tilde{J}_{m-2} \cdots (\tilde{J}_1 \cdot (\tilde{J}_0 p))))). \quad (4.4)$$

In order to write down this long matrix product as an efficient algorithm where the multiplications of all the ones and zeros do not cause computational costs, it is customary in the field of AD to use a notation that uses “dot quantities” \dot{x}_i that we might think of as the velocity with which a certain variable changes, given that the input x changes with speed $\dot{x} = p$. We can interpret them as

$$\dot{x}_i \equiv \frac{dx_i}{dx} p.$$

In the augmented formulation, we can introduce dot quantities $\dot{\tilde{x}}_i$ for the augmented vectors \tilde{x}_i , for $i = 0, \dots, m - 1$, and the recursion of these dot quantities is just given by the initialization with the seed vector, $\dot{\tilde{x}}_i = p$, and then the recursion

$$\dot{\tilde{x}}_{i+1} = \tilde{J}_i(\tilde{x}_i) \dot{\tilde{x}}_i, \quad i = 0, 1, \dots, m - 1.$$

Given the special structure of the Jacobian matrices, most elements of $\dot{\tilde{x}}_i$ are only multiplied by one and nothing needs to be done, apart from the computation of the last component of the new vector $\dot{\tilde{x}}_{i+1}$. This last component is \dot{x}_{n+i+1} . Thus, in an efficient implementation, the forward AD algorithm works as the algorithm below. It first sets the seed $\dot{x} = p$ and then proceeds as follows.

In forward AD, the function evaluation and the derivative evaluation can be performed in parallel, which eliminates the need to store any internal information. This is best illustrated using an example.

Example 4 (Forward Automatic Differentiation) *We regard the same example as above, $f(x_1, x_2, x_3) = \sin(x_1 x_2) + \exp(x_1 x_2 x_3)$. First, each intermediate variable has to be computed, and then each line can be differentiated. For given x_1, x_2, x_3 and $\dot{x}_1, \dot{x}_2, \dot{x}_3$, the algorithm proceeds as follows:*

$$\begin{array}{ll} x_4 = x_1 x_2 & \dot{x}_4 = \dot{x}_1 x_2 + x_1 \dot{x}_2 \\ x_5 = \sin(x_4) & \dot{x}_5 = \cos(x_4) \dot{x}_4 \\ x_6 = x_4 x_3 & \dot{x}_6 = \dot{x}_4 x_3 + x_4 \dot{x}_3 \\ x_7 = \exp(x_6) & \dot{x}_7 = \exp(x_6) \dot{x}_6 \\ x_8 = x_5 + x_7 & \dot{x}_8 = \dot{x}_5 + \dot{x}_7 \end{array}$$

The result is $\dot{x}_8 = (\dot{x}_1, \dot{x}_2, \dot{x}_3) \nabla f(x_1, x_2, x_3)$.

It can be proven that the computational cost of Algorithm 4.2 is smaller than two times the cost of Algorithm 4.1, or short

$$\text{cost}(J_F p) \leq 2 \text{cost}(F).$$

If we want to obtain the full Jacobian of F , we need to call Algorithm 4.2 several times, each time with the seed vector corresponding to one of the n unit vectors in \mathbb{R}^n , i.e. we have

$$\text{cost}(J_F) \leq 2n \text{cost}(F).$$

AD in forward mode is slightly more expensive than numerical finite differences, but it is exact up to machine precision.

4.2.1 The “Imaginary trick” in MATLAB

An easy way to obtain high precision derivatives in MATLAB is closely related to AD in forward mode. It is based on the following observation: if $F : \mathbb{R}^n \rightarrow \mathbb{R}^{n_F}$ is analytic and can be extended to complex numbers as inputs and outputs, then for any $t > 0$ holds

$$J_F(x)p = \frac{j(F(x + itp))}{t} + O(t^2). \quad (4.5)$$

In contrast to finite differences, there is no subtraction in the numerator, so there is no danger of numerical cancellation errors, and t can be chosen extremely small, e.g. $t = 10^{-100}$, which means that we can compute the derivative up to machine precision. This “imaginary trick” can most easily be used in a programming language like MATLAB that does not declare the type of variables beforehand, so that real-valued variables can automatically be overloaded with complex-valued variables. This allows us to obtain high-precision derivatives of a given black-box MATLAB code. We only need to be sure that the code is analytic (which most codes are) and that matrix or vector transposes are not expressed by a prime $'$ (which conjugates a complex number), but by `transp`.

4.3 The Backward Mode of AD

In backward AD we evaluate the product in Eq. (4.3) in the reverse order compared with forward AD. Backward AD does not evaluate forward directional derivatives. Instead, it evaluates *adjoint directional derivatives*: when we define a *seed vector* $\lambda \in \mathbb{R}^{n_F}$ then backward AD is able to evaluate the product $\lambda^T J_F$. It does so in the following way:

$$\lambda^T J_F = (((\lambda^T C) \cdot \tilde{J}_{m-1}) \cdot \tilde{J}_{m-2}) \cdots \tilde{J}_1) \cdot \tilde{J}_0. \quad (4.6)$$

When writing this matrix product as an algorithm, we use “bar quantities” instead of the “dot quantities” that we used in the forward mode. These quantities can be interpreted as derivatives of the final output with respect to the respective intermediate quantity. We can interpret

$$\bar{x}_i \equiv \lambda^T \frac{dF}{dx_i}.$$

Each intermediate variable has a bar variable and at the start, we initialize all bar variables with the value that we obtain from $C^T \lambda$. Note that most of these seeds will usually be zero, depending on the output selection matrix C . Then, the backward AD algorithm modifies all bar variables. Backward AD gets most transparent in the augmented formulation, where we have bar quantities $\bar{\tilde{x}}_i$ for the augmented states \tilde{x}_i . We can transpose the above Equation (4.6) in order to obtain

$$J_F^T \lambda = \tilde{J}_0^T \cdot (\underbrace{\tilde{J}_1^T \cdots \tilde{J}_{m-1}^T}_{=\bar{\tilde{x}}_m} (C^T \lambda)).$$

$=\bar{\tilde{x}}_{m-1}$

In this formulation, the initialization of the backward seed is nothing else than setting $\bar{\tilde{x}}_m = C^T \lambda$ and then going in reverse order through the recursion

$$\bar{\tilde{x}}_i = \tilde{J}_i(\tilde{x}_i)^T \bar{\tilde{x}}_{i+1}, \quad i = m-1, m-2, \dots, 0.$$

Again, the multiplication with ones does not cause any computational cost, but an interesting feature of the reverse mode is that some of the bar quantities can get several times modified in very different stages of the algorithm.

Input: seed vector $\bar{x}_1, \dots, \bar{x}_{n+m}$ and all partial derivatives $\frac{\partial \phi_i}{\partial x_j}$

Output: $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$

```

for  $i = m - 1$  down to 0 do
  for all  $j = 1, \dots, n + i$  do
     $\bar{x}_j \leftarrow \bar{x}_j + \bar{x}_{n+i+1} \frac{\partial \phi_i}{\partial x_j}$ 
  end for
end for

```

Note: each inner loop will only update one or two bar quantities.

Note that the multiplication $\tilde{J}_i^T \bar{x}_{i+1}$ with the transposed Jacobian

$$\tilde{J}_i^T = \begin{bmatrix} 1 & & & * \\ & 1 & & * \\ & & \ddots & * \\ & & & 1 & * \end{bmatrix}.$$

modifies at maximum two elements of the vector \bar{x}_{i+1} by adding to them the partial derivative of the elementary operation multiplied with \bar{x}_{n+i+1} . In an efficient implementation, the backward AD algorithm looks as follows.

Example 5 (Reverse Automatic Differentiation) We regard the same example as before, and want to compute the gradient $\nabla f(x) = (\bar{x}_1, \bar{x}_2, \bar{x}_3)^T$ given (x_1, x_2, x_3) . We set $\lambda = 1$. Because the selection matrix C selects only the last intermediate variable as output, i.e. $C = (0, \dots, 0, 1)$, we initialize the seed vector with zeros apart from the last component, which is one. In the reverse mode, the algorithm first has to evaluate the function with all intermediate quantities, and only then it can compute the bar quantities, which it does in reverse order. At the end it obtains, among other, the desired quantities $(\bar{x}_1, \bar{x}_2, \bar{x}_3)$. The full algorithm is the following.

```

// *** forward evaluation of the function ***

```

```

 $x_4 = x_1 x_2$ 

```

```

 $x_5 = \sin(x_4)$ 

```

```

 $x_6 = x_4 x_3$ 

```

```

 $x_7 = \exp(x_6)$ 

```

```

 $x_8 = x_5 + x_7$ 

```

```

// *** initialization of the seed vector ***

```

```

 $\bar{x}_i = 0, \quad i = 1, \dots, 7$ 

```

```

 $\bar{x}_8 = 1$ 

```

```

// *** backwards sweep ***
// * differentiation of  $x_8 = x_5 + x_7$ 
 $\bar{x}_5 = \bar{x}_5 + 1 \bar{x}_8$ 
 $\bar{x}_7 = \bar{x}_7 + 1 \bar{x}_8$ 
// * differentiation of  $x_7 = \exp(x_6)$ 
 $\bar{x}_6 = \bar{x}_6 + \exp(x_6) \bar{x}_7$ 
// * differentiation of  $x_6 = x_4 x_3$ 
 $\bar{x}_4 = \bar{x}_4 + x_3 \bar{x}_6$ 
 $\bar{x}_3 = \bar{x}_3 + x_4 \bar{x}_6$ 
// * differentiation of  $x_5 = \sin(x_4)$ 
 $\bar{x}_4 = \bar{x}_4 + \cos(x_4) \bar{x}_5$ 
// differentiation of  $x_4 = x_1 x_2$ 
 $\bar{x}_1 = \bar{x}_1 + x_2 \bar{x}_4$ 
 $\bar{x}_2 = \bar{x}_2 + x_1 \bar{x}_4$ 

```

The desired output of the algorithm is $(\bar{x}_1, \bar{x}_2, \bar{x}_3)$, equal to the three components of the gradient $\nabla f(x)$. Note that all three are returned in only one reverse sweep.

It can be shown that the cost of Algorithm 4.3 is less than 3 times the cost of Algorithm 4.1, i.e.,

$$\text{cost}(\lambda^T J_F) \leq 3 \text{cost}(F).$$

If we want to obtain the full Jacobian of F , we need to call Algorithm 4.3 several times with the n_F seed vectors corresponding to the unit vectors in \mathbb{R}^{n_F} , i.e. we have

$$\text{cost}(J_F) \leq 3 n_F \text{cost}(F).$$

This is a remarkable fact: it means that the backward mode of AD can compute the full Jacobian at a cost that is independent of the state dimension n . This is particularly advantageous if $n_F \ll n$, e.g. if we compute the gradient of a scalar function like the objective or the Lagrangian. The reverse mode can be much faster than what we can obtain by finite differences, where we always need $(n+1)$ function evaluations. To give an example, if we want to compute the gradient of a scalar function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with $n = 1\,000\,000$ and each call of the function needs one second of CPU time, then the finite difference approximation of the gradient would take 1 000 001 seconds, while the computation of the same quantity with the backward mode of AD needs only 4 seconds (1 call of the function plus one backward sweep). Thus, besides being more accurate, backward AD can also be much faster than finite differences.

The only disadvantage of the backward mode of AD is that we have to store all intermediate variables and partial derivatives, in contrast to finite differences or forward AD. A partial remedy to this problem exist in form of *checkpointing* that trades-off computational speed and memory requirements. Instead of all intermediate variables, it only stores some “checkpoints” during the forward evaluation. During the backward sweep, starting at these checkpoints, it re-evaluates parts of the function to obtain those intermediate variables that have not been stored. The optimal number and location of checkpoints is a science of itself. Generally speaking, checkpointing reduces the memory requirements, but comes at the expense of runtime.

From a user perspective, the details of implementation are not too relevant, but it is most important to just know that the reverse mode of AD exists and that it allows in many cases a much more efficient derivative generation than any other technique.

4.3.1 Efficient Computation of the Hessian

A particularly important quantity in Newton-type optimization methods is the Hessian of the Lagrangian. It is the second derivative of the scalar function $\mathcal{L}(x, \lambda, \mu)$ with respect to x . As the multipliers are fixed for the purpose of differentiation, we can for notational simplicity just regard a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ of which we want to compute the Hessian $\nabla^2 f(x)$. With finite differences we would at least need $(n+2)(n+1)/2$ function evaluations in order to compute the Hessian, and due to round-off and truncation errors, the accuracy of a finite difference Hessian would be much lower than the accuracy of the function f : we loose three quarters of the valid digits.

In contrast to this, algorithmic differentiation can without problems be applied recursively, yielding a code that computes the Hessian matrix at the same precision as the function f itself, i.e. typically at machine precision. Moreover, if we use the reverse mode of AD at least once, e.g. by first generating an efficient code for $\nabla f(x)$ (using backward AD) and then using forward AD to obtain the Jacobian of it, we can reduce the CPU time considerably compared to finite differences. Using the above procedure, we would obtain the Hessian $\nabla^2 f$ at a cost of $2n$ times the cost of a gradient ∇f , which is about four times the cost of evaluating f alone. This means that we have the following runtime bound:

$$\text{cost}(\nabla^2 f) \leq 8n \text{cost}(f).$$

A compromise between accuracy and ease of implementation that is equally fast in terms of CPU time is to use backward AD only for computing the first order derivative $\nabla f(x)$, and then to use finite differences for the differentiation of $\nabla f(x)$.

4.4 Algorithmic Differentiation Software

Most algorithmic differentiation tools implement both forward and backward AD, and most are specific to one particular programming language. They come in two different variants: either they use *operator overloading* or *source-code transformation*.

The first class does not modify the code but changes the type of the variables and overloads the involved elementary operations. For the forward mode, each variable just gets an additional dot-quantity, i.e. the new variables are the pairs (x_i, \dot{x}_i) , and elementary operations just operate on these pairs, like e.g.

$$(x, \dot{x}) \cdot (y, \dot{y}) = (xy, x\dot{y} + y\dot{x}).$$

An interesting remark is that operator overloading is also at the basis of the imaginary trick in MATLAB where we use the overloading of real numbers by complex numbers and used the small imaginary part as dot quantity and exploited the fact that the extremely small higher order terms disappear by numerical cancellation.

A prominent and widely used AD tool for generic user supplied C++ code that uses operator overloading is ADOL-C. Though it is not the most efficient AD tool in terms of CPU time it is well documented and stable. Another popular tool in this class is CppAD.

The other class of AD tools is based on source-code transformation. They work like a text-processing tool that gets as input the user supplied source code and produces as output a new and very differently looking source code that implements the derivative generation. Often, these codes can be made extremely fast. Tools that implement source code transformations are ADIC for ANSI C, and ADIFOR and TAPENADE for FORTRAN codes.

In the context of ODE or DAE simulation, there exist good numerical integrators with forward and backward differentiation capabilities that are more efficient and reliable than a naive procedure that would consist of taking an integrator and processing it with an AD tool. Examples for integrators that use the principle of forward and backward AD are the code DAESOL-II or the open-source codes from the ACADO Integrators Collection or from the SUNDIALS Suite.

Chapter 5

Discrete Time Optimal Control Formulations

Throughout this part of the script we regard for notational simplicity time-invariant dynamical systems with dynamics

$$x_{k+1} = f(x_k, u_k), \quad k = 0, \dots, N-1. \quad (5.1)$$

Recall that u_k are the *controls* and x_k the *states*, with $x_k \in \mathbb{R}^{n_x}$ and $u_k \in \mathbb{R}^{n_u}$.

As discussed in the first chapter, if we know the initial state x_0 and the controls u_0, \dots, u_{N-1} , we could simulate the system to obtain all other states. But in optimization, we might have different requirements than just a fixed initial state. We might, for example, have both a fixed initial state and a fixed terminal state that we want to reach. Or we might just look for periodic sequences with $x_0 = x_N$. All these desires on the initial and the terminal state can be expressed by a boundary constraint function

$$r(x_0, x_N) = 0. \quad (5.2)$$

For the case of fixed initial value, this function would just be

$$r(x_0, x_N) = x_0 - \bar{x}_0 \quad (5.3)$$

where \bar{x}_0 is the fixed initial value and not an optimization variable. Another example would be to have both ends fixed, resulting in a function r of double the state dimension, namely:

$$r(x_0, x_N) = \begin{bmatrix} x_0 - \bar{x}_0 \\ x_N - \bar{x}_N \end{bmatrix}. \quad (5.4)$$

Finally, periodic boundary conditions can be imposed by setting

$$r(x_0, x_N) = x_0 - x_N. \quad (5.5)$$

Other constraints that are usually present are *path constraint* inequalities of the form

$$h(x_k, u_k) \leq 0, \quad k = 0, \dots, N-1. \quad (5.6)$$

In the case of upper and lower bounds on the controls, $u_{\min} \leq u_k \leq u_{\max}$, the function h would just be

$$h(x, u) = \begin{bmatrix} u - u_{\max} \\ u_{\min} - u \end{bmatrix}.$$

5.1 Optimal Control Problem (OCP) Formulations

Two major approaches can be distinguished to formulate and numerically solve a discrete time optimal control problem, the *simultaneous* and the *sequential* approach, which we will outline after having formulated the optimal control problem in its standard form.

5.1.1 Original Problem Formulation

Given the system model and constraints, a quite generic discrete time optimal control problem can be formulated as the following constrained NLP:

$$\begin{aligned} & \underset{x_0, u_0, x_1, \dots, u_{N-1}, x_N}{\text{minimize}} && \sum_{k=0}^{N-1} L(x_k, u_k) + E(x_N) \end{aligned} \quad (5.7a)$$

$$\text{subject to} \quad x_{k+1} - f(x_k, u_k) = 0, \quad \text{for } k = 0, \dots, N-1, \quad (5.7b)$$

$$h(x_k, u_k) \leq 0, \quad \text{for } k = 0, \dots, N-1, \quad (5.7c)$$

$$r(x_0, x_N) = 0. \quad (5.7d)$$

We remark that other optimization variables could be present as well, such as a free parameter p that can be chosen but is constant over time, like e.g. the size of a vessel in a chemical reactor or the length of a robot arm. Such parameters could be added to the optimisation formulation above by defining dummy states $\{p_k\}_{k=1}^N$ that satisfy the dummy dynamic model equations

$$p_{k+1} = p_k, \quad k = 0, \dots, N-1. \quad (5.8)$$

Note that the initial value of p_0 is not fixed by these constraints and thus we would have obtained our aim of having a time constant parameter vector that is free for optimization.

5.1.2 The Simultaneous Approach

The nonlinear program (5.7) is large and structured and can thus in principle be solved by any NLP solver. This is called the *simultaneous approach* to optimal control and requires the use of a structure exploiting NLP solver in order to be efficient. Note that in this approach, all original variables, i.e. u_k and x_k remain optimization variables of the NLP. Its name stems from the fact that the NLP solver has to simultaneously solve both, the simulation and the optimization problem. It is interesting to remark that the model equations (5.7b) will for most NLP solvers only be satisfied once the NLP iterations are converged. The simultaneous approach is therefore sometimes referred to as an *infeasible path* approach. The methods *direct multiple shooting* and *direct collocation* that we explain in the third part of this script are simultaneous approaches.

5.1.3 The Reduced Formulation and the Sequential Approach

On the other hand, we know that we could eliminate nearly all states by a forward simulation, and in this way we could reduce the variable space of the NLP. The idea is to keep only x_0 and $U = [u_0^T, \dots, u_{N-1}^T]^T$ as variables. The states x_1, \dots, x_N are eliminated recursively by

$$\bar{x}_0(x_0, U) = x_0 \quad (5.9a)$$

$$\bar{x}_{k+1}(x_0, U) = f(\bar{x}_k(x_0, U), u_k), \quad k = 0, \dots, N-1. \quad (5.9b)$$

Then the optimal control problem is equivalent to a *reduced problem* with much less variables, namely the following nonlinear program:

$$\begin{aligned} & \underset{x_0, U}{\text{minimize}} && \sum_{k=0}^{N-1} L(\bar{x}_k(x_0, U), u_k) + E(\bar{x}_N(x_0, U)) \end{aligned} \quad (5.10a)$$

$$\text{subject to} \quad h(\bar{x}_k(x_0, U), u_k) \leq 0, \quad \text{for } k = 0, \dots, N-1, \quad (5.10b)$$

$$r(x_0, \bar{x}_N(x_0, U)) = 0. \quad (5.10c)$$

Note that the model Equation (5.9b) is implicitly satisfied by definition, but is not anymore a constraint of the optimization problem. This reduced problem can now be addressed again by Newton-type methods, but the exploitation of sparsity in the problem is less important. This is called the *sequential* approach, because the simulation problem and optimization problem are solved sequentially, one after the other. Note that the user can observe during all iterations of the optimization procedure what is the resulting state trajectory for the current iterate, as the model equations are satisfied by definition.

If the initial value is fixed, i.e. if $r(x_0, x_N) = x_0 - \bar{x}_0$, one can also eliminate $x_0 \equiv \bar{x}_0$, which reduces the variables of the NLP further.

5.2 Analysis of a Simplified Optimal Control Problem

In order to learn more about the structure of optimal control problems and the relation between the simultaneous and the sequential approach, we regard in this section a simplified optimal control problem in discrete time:

$$\begin{aligned} & \underset{x_0, u_0, x_1, \dots, u_{N-1}, x_N}{\text{minimize}} && \sum_{k=0}^{N-1} L(x_k, u_k) + E(x_N) \end{aligned} \quad (5.11a)$$

$$\text{subject to} \quad f(x_k, u_k) - x_{k+1} = 0 \text{ for } k = 0, \dots, N-1 \quad (5.11b)$$

$$r(x_0, x_N) = 0 \quad (5.11c)$$

5.2.1 KKT Conditions of the Simplified Problem

We first summarize the variables as $w = (x_0, u_0, x_1, u_1, \dots, u_{N-1}, x_N)$ and summarize the multipliers as $\lambda = (\lambda_1, \dots, \lambda_N, \lambda_r)$. Then the above optimal control problem can be summarized as

$$\underset{w}{\text{minimize}} \quad F(w) \quad (5.12a)$$

$$\text{subject to} \quad G(w) = 0. \quad (5.12b)$$

Here, the objective $F(w)$ is just copied from (5.11a) while $G(w)$ collects all constraints:

$$G(w) = \begin{bmatrix} f(x_0, u_0) - x_1 \\ f(x_1, u_1) - x_2 \\ \vdots \\ f(x_{N-1}, u_{N-1}) - x_N \\ r(x_0, x_N) \end{bmatrix}. \quad (5.12c)$$

The Lagrangian function has the form

$$\begin{aligned} \mathcal{L}(w, \lambda) &= F(w) + \lambda^T G(w) \\ &= \sum_{k=0}^{N-1} L(x_k, u_k) + E(x_N) + \sum_{k=0}^{N-1} \lambda_{k+1}^T (f(x_k, u_k) - x_{k+1}) \\ &\quad + \lambda_r^T r(x_0, x_N), \end{aligned} \quad (5.13)$$

and the summarized KKT-conditions of the problem are

$$\nabla_w \mathcal{L}(w, \lambda) = 0 \quad (5.14a)$$

$$G(w) = 0. \quad (5.14b)$$

But let us look at these KKT-conditions in more detail. First, we evaluate the derivative of \mathcal{L} with respect to all state variables x_k , one after the other. We have to treat $k = 0$ and $k = N$ as special cases. For $k = 0$ we obtain:

$$\nabla_{x_0} \mathcal{L}(w, \lambda) = \nabla_{x_0} L(x_0, u_0) + \frac{\partial f}{\partial x_0}(x_0, u_0)^T \lambda_1 + \frac{\partial r}{\partial x_0}(x_0, x_N)^T \lambda_r = 0. \quad (5.15a)$$

Then the case for $k = 1, \dots, N-1$ is treated

$$\nabla_{x_k} \mathcal{L}(w, \lambda) = \nabla_{x_k} L(x_k, u_k) - \lambda_k + \frac{\partial f}{\partial x_k}(x_k, u_k)^T \lambda_{k+1} = 0. \quad (5.15b)$$

Last, the special case $k = N$

$$\nabla_{x_N} \mathcal{L}(w, \lambda) = \nabla_{x_N} E(x_N) - \lambda_N + \frac{\partial r}{\partial x_N}(x_0, x_N)^T \lambda_r = 0. \quad (5.15c)$$

Second, let us calculate the derivative of the Lagrangian with respect to all controls u_k , for $k = 0, \dots, N-1$. Here, no special cases need to be considered, and we obtain the general formula

$$\nabla_{u_k} \mathcal{L}(w, \lambda) = \nabla_{u_k} L(x_k, u_k) + \frac{\partial f}{\partial u_k}(x_k, u_k)^T \lambda_{k+1} = 0. \quad (5.15d)$$

Until now, we have computed in detail the components of the first part of the KKT-condition (5.14a), i.e. $\nabla_w \mathcal{L}(w, \lambda) = 0$. The other part of the KKT-condition, $G(w) = 0$, is trivially given by

$$f(x_k, u_k) - x_{k+1} = 0, \quad k = 0, \dots, N-1 \quad (5.15e)$$

$$r(x_0, x_N) = 0 \quad (5.15f)$$

Thus, collecting all equations (5.15a) to (5.15f), we have stated the KKT-conditions of the OCP. They can be treated by Newton-type methods in different ways. The *simultaneous approach* addresses equations (5.15a) to (5.15f) directly by a Newton-type method in the space of all variables (w, λ) . In contrast to this, the *sequential approach* eliminates all the states x_1, \dots, x_N in (5.15e) by a forward simulation, and if it is implemented efficiently, it also uses Eqs. (5.15c) and (5.15b) to eliminate all multipliers $\lambda_N, \dots, \lambda_1$ in a backward simulation, as discussed in the following subsection.

5.2.2 Computing Gradients in the Sequential Approach

A naive implementation of the sequential approach would start by coding routines that evaluate the objective and constraint functions, and then passing these routines as black-box codes to a generic NLP solver, like `fmincon` in MATLAB. But this would not be the most efficient way to implement the sequential approach. The reason is the generation of derivatives, which a generic NLP solver will compute by finite differences. On the other hand, many generic NLP solvers allow the user to deliver explicit functions for the derivatives as well. This allows us to compute the derivatives of the reduced problem functions more efficiently. The key technology here is algorithmic differentiation in the backward mode, as explained in Chapter 4.

To see how this relates to the optimality conditions (5.15a) to (5.15f) of the optimal control problem, let us simplify the setting even more by assuming a fixed initial value and no constraint on the terminal state, i.e. $r(x_0, x_N) = \bar{x}_0 - x_0$. In this case, the KKT conditions simplify to the following set of equations, which we bring already into a specific order:

$$x_0 = \bar{x}_0 \quad (5.16a)$$

$$x_{k+1} = f(x_k, u_k), \quad k = 0, \dots, N-1, \quad (5.16b)$$

$$\lambda_N = \nabla_{x_N} E(x_N) \quad (5.16c)$$

$$\lambda_k = \nabla_{x_k} L(x_k, u_k) + \frac{\partial f}{\partial x_k}(x_k, u_k)^T \lambda_{k+1}, \quad k = N-1, \dots, 1, \quad (5.16d)$$

$$\nabla_{u_k} L(x_k, u_k) + \frac{\partial f}{\partial u_k}(x_k, u_k)^T \lambda_{k+1} = 0, \quad k = 0, \dots, N-1. \quad (5.16e)$$

It can easily be seen that the first four equations can trivially be satisfied, by a forward sweep to obtain all x_k and a backward sweep to obtain all λ_k . Thus, x_k and λ_k can be made explicit functions of u_0, \dots, u_{N-1} . The only equation that is non-trivial to satisfy is the last one, the partial derivatives of the Lagrangian w.r.t. the controls u_0, \dots, u_{N-1} . Thus we could decide to eliminate x_k and λ_k and only search with a Newton-type scheme for the variables $U = (u_0, \dots, u_{N-1})$ such that these last equations are satisfied. It turns out that the left hand side residuals (5.16e) are nothing else than the derivative of the reduced problem's objective (5.10a), and the forward-backward sweep algorithm described above is nothing else than the reverse mode of algorithmic differentiation. It is much more efficient than the computation of the gradient by finite differences.

The forward-backward sweep is well known in the optimal control literature and often introduced without reference to the reverse mode of AD. On the other hand, it is good to know the general principles of AD in forward or backward mode, because AD can also be beneficial in other contexts, e.g. for the evaluation of derivatives of the other problem functions in (5.10a)-(5.10c). Also, when second order derivatives are needed, AD can be used and more structure can be exploited, but this is most easily derived in the context of the simultaneous approach, which we do in the following chapter.

Chapter 6

Sparsity Structure of the Optimal Control Problem

Let us in this chapter regard a very general optimal control problem in the original formulation, i.e. the NLP that would be treated by the simultaneous approach.

$$\begin{aligned} & \underset{x_0, u_0, x_1, \dots, u_{N-1}, x_N}{\text{minimize}} && \sum_{k=0}^{N-1} L_k(x_k, u_k) + E(x_N) \end{aligned} \quad (6.1a)$$

$$\text{subject to} \quad f_k(x_k, u_k) - x_{k+1} = 0, \quad \text{for } k = 0, \dots, N-1, \quad (6.1b)$$

$$\sum_{k=0}^{N-1} r_k(x_k, u_k) + r_N(x_N) = 0, \quad (6.1c)$$

$$h_k(x_k, u_k) \leq 0, \quad \text{for } k = 0, \dots, N-1, \quad (6.1d)$$

$$h_N(x_N) \leq 0. \quad (6.1e)$$

Compared to the OCP (5.7) in the previous chapter, we now allow indices on all problem functions making the system time dependent; also, we added terminal inequality constraints (6.1e), and as boundary conditions we now allow now very general coupled multipoint constraints (6.1c) that include the cases of fixed initial or terminal values or periodicity, but are much more general. Note that in these boundary constraints terms arising from different time points are only coupled by addition, because this allows us to maintain the sparsity structure we want to exploit in this chapter.

Collecting all variables in a vector w , the objective in a function $F(w)$, all equalities in a function $G(w)$ and all inequalities in a function $H(w)$, the optimal control problem could be summarized as

$$\underset{w}{\text{minimize}} \quad F(w) \quad (6.2a)$$

$$\text{subject to} \quad G(w) = 0, \quad (6.2b)$$

$$H(w) \leq 0. \quad (6.2c)$$

Its Lagrangian function is given by

$$\mathcal{L}(w, \lambda, \mu) = F(w) + \lambda^T G(w) + \mu^T H(w).$$

But this summarized form does not reveal any of the structure that is present in the problem.

6.1 Partial Separability of the Lagrangian

In fact, the above optimal control problem is a very sparse problem because each of its functions depends only on very few of its variables. This means for example that the Jacobian matrix of the equality constraints has many zero

entries. But not only first order derivatives are sparse, also the second order derivative that we need in Newton-type optimization algorithms, namely the Hessian of the Lagrangian, is a very sparse matrix. This is due to the fact that the Lagrangian is a *partially separable* function [GT82].

Definition 14 (Partial Separability) *A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is called partially separable if it can be decomposed as a sum of m functions $f_j : \mathbb{R}^{n_j} \rightarrow \mathbb{R}$ with $n_j < n$ for all $j = 1, \dots, m$. This means that for each j exists a subset I_j of indices from $\{1, \dots, n\}$ and subvectors x_{I_j} of x such that*

$$f(x) = \sum_{j=1}^m f_j(x_{I_j}).$$

The Lagrangian function of the above optimization problem can explicitly be decomposed into subfunctions that each depend on some of the multipliers and only on the variables (x_k, u_k) with the same index k . Let us collect again all variables in a vector w but decompose it as¹ $w = (w_1, \dots, w_N)$ with $w_k = (x_k, u_k)$ for $k = 0, \dots, N-1$ and $w_N = x_N$. Collecting all equality multipliers in a vector $\lambda = (\lambda_1, \dots, \lambda_N, \lambda_r)$ and the inequality multipliers in a vector $\mu = (\mu_0, \dots, \mu_N)$ we obtain for the Lagrangian

$$\mathcal{L}(w, \lambda, \mu) = \sum_{k=0}^N \mathcal{L}_k(w_k, \lambda, \mu)$$

with the local Lagrangian subfunctions defined as follows. The first subfunction is given as

$$\mathcal{L}_0(w_0, \lambda, \mu) = L_0(x_0, u_0) + \lambda_1^T f_0(x_0, u_0) + \mu_0^T h_0(x_0, u_0) + \lambda_r^T r_0(x_0, u_0)$$

and for $k = 1, \dots, N-1$ we have the subfunctions

$$\mathcal{L}_k(w_k, \lambda, \mu) = L_k(x_k, u_k) + \lambda_{k+1}^T f_k(x_k, u_k) - \lambda_k^T x_k + \mu_k^T h_k(x_k, u_k) + \lambda_r^T r_k(x_k, u_k)$$

while the last subfunction is given as

$$\mathcal{L}_N(w_N, \lambda, \mu) = E(x_N) - \lambda_N^T x_N + \mu_N^T h_N(x_N) + \lambda_r^T r_N(x_N).$$

In fact, while each of the equality multipliers appears in several $(\lambda_1, \dots, \lambda_N)$ or even all problem functions (λ_r) , the primal variables of the problem do not have any overlap in the subfunctions. This leads to the remarkable observation that the Hessian matrix $\nabla_w^2 \mathcal{L}$ is *block diagonal*, i.e. it consists only of small symmetric matrices that are located on its diagonal. All other second derivatives are zero, i.e.

$$\frac{\partial^2 \mathcal{L}}{\partial w_i \partial w_j}(w, \lambda, \mu) = 0, \quad \text{for any } i \neq j.$$

This block diagonality of the Hessian leads to several very favourable facts, namely that (i) the Hessian can be approximated by *high-rank* or *block updates* within a BFGS method [GT82, BP84], and (ii) that the QP subproblem in all Newton-type methods has the same decomposable objective function as the original optimal control problem itself.

6.2 The Sparse QP Subproblem

In order to analyse the sparsity structure of the optimal control problem, let us regard the quadratic subproblem that needs to be solved in one iteration of an exact Hessian SQP method. In order not to get lost in too many indices, we disregard the SQP iteration index completely. We regard the QP that is formulated at a current iterate (x, λ, μ) and use the SQP step $\Delta w = (\Delta x_0, \Delta u_0, \dots, \Delta x_N)$ as the QP variable. This means that in the summarized formulation we would have the QP subproblem

$$\begin{aligned} \text{minimize} \quad & \nabla F(w)^T \Delta w + \frac{1}{2} \Delta w^T \nabla_w^2 \mathcal{L}(w, \lambda, \mu) \Delta w \end{aligned} \tag{6.3a}$$

$$\text{subject to} \quad G(w) + \nabla G(w)^T \Delta w = 0, \tag{6.3b}$$

$$H(w) + \nabla H(w)^T \Delta w \leq 0. \tag{6.3c}$$

¹Note that for notational beauty we omit here and in many other occasions the transpose signs that would be necessary to make sure that the collection of column vectors is again a column vector, when this is clear from the context.

Let us now look at this QP subproblem in the detailed formulation. It is remarkably similar to the original OCP. To reduce notational overhead, let us define a few abbreviations: first, the diagonal blocks of the Hessian of the Lagrangian

$$Q_k = \nabla_{w_k}^2 \mathcal{L}(w, \lambda, \mu), \quad k = 0, \dots, N,$$

second, the objective gradients

$$g_k = \nabla_{(x,u)} L(x_k, u_k), \quad k = 0, \dots, N-1, \quad \text{and} \quad g_N = \nabla_x E(x_N),$$

third the system discontinuities (that can be non-zero in the simultaneous approach)

$$a_k = f_k(x_k, u_k) - x_{k+1}, \quad k = 0, \dots, N-1,$$

and fourth the transition matrices

$$A_k = \frac{\partial f_k}{\partial x_k}(x_k, u_k), \quad B_k = \frac{\partial f_k}{\partial u_k}(x_k, u_k), \quad k = 0, \dots, N-1,$$

fifth the residual of the coupled constraints

$$r = \sum_{k=0}^{N-1} r_k(x_k, u_k) + r_N(x_N),$$

as well as its derivatives

$$R_k = \frac{\partial r_k}{\partial (x_k, u_k)}(x_k, u_k), \quad k = 0, \dots, N-1, \quad \text{and} \quad R_N = \frac{\partial r_N}{\partial x}(x_N),$$

and last the inequality constraint residuals and their derivatives

$$h_k = h_k(x_k, u_k), \quad H_k = \frac{\partial h_k}{\partial (x_k, u_k)}(x_k, u_k) \quad \text{and} \quad h_N = h_N(x_N), \quad H_N = \frac{\partial h_N}{\partial x}(x_N).$$

With all the above abbreviations, the detailed form of the QP subproblem is finally given as follows.

$$\begin{aligned} & \underset{\Delta x_0, \Delta u_0, \dots, \Delta x_N}{\text{minimize}} && \frac{1}{2} \sum_{k=0}^{N-1} \begin{bmatrix} \Delta x_k \\ \Delta u_k \end{bmatrix}^T Q_k \begin{bmatrix} \Delta x_k \\ \Delta u_k \end{bmatrix} + \frac{1}{2} \Delta x_N^T Q_N \Delta x_N + \sum_{k=0}^N \begin{bmatrix} \Delta x_N \\ \Delta u_N \end{bmatrix}^T g_k + \Delta x_N^T g_N \\ & \text{subject to} && a_k + A_k \Delta x_k + B_k \Delta u_k - \Delta x_{k+1} = 0, \quad \text{for } k = 0, \dots, N-1, \end{aligned} \tag{6.4}$$

$$\begin{aligned} & r + \sum_{k=0}^{N-1} R_k \begin{bmatrix} \Delta x_k \\ \Delta u_k \end{bmatrix} + R_N \Delta x_N = 0, \\ & h_k + H_k \begin{bmatrix} \Delta x_k \\ \Delta u_k \end{bmatrix} \leq 0, \quad \text{for } k = 0, \dots, N-1, \\ & h_N + H_N \Delta x_N \leq 0. \end{aligned} \tag{6.5}$$

This is again an optimal control problem, but a linear-quadratic one. It is a convex QP if the Hessian blocks Q_k are positive definite, and can be solved by a variety of sparsity exploiting QP solvers.

6.3 Sparsity Exploitation in QP Solvers

When regarding the QP (6.4) one way would be to apply a sparse interior point QP solver like OOQP to it, or a sparse active set method. This can be very efficient. Another way would be to first reduce, or *condense*, the variable space of the QP, and then apply a standard dense QP solver to the reduced problem. Let us treat this way first.

6.3.1 Condensing

When we regard the linearized dynamic system equations (6.5) they correspond to an affine time variant system in the steps Δx_k , namely

$$\Delta x_{k+1} = a_k + A_k \Delta x_k + B_k \Delta u_k. \quad (6.7)$$

If the values for Δx_0 as well as for all $\{\Delta u_k\}_{k=0}^{N-1}$ would be known, then also the values for $\{\Delta x_k\}_{k=1}^N$ can be obtained by a forward simulation of this linear system. Due to its linearity, the resulting map will be linear, i.e. we can write

$$\begin{bmatrix} \Delta x_1 \\ \vdots \\ \Delta x_N \end{bmatrix} = v + M \begin{bmatrix} \Delta x_0 \\ \Delta u_0 \\ \vdots \\ \Delta u_{N-1} \end{bmatrix}, \quad (6.8)$$

$$\Delta w_{\text{dep}} = v + M \Delta w_{\text{ind}} \quad (6.9)$$

with a vector $v \in \mathbb{R}^{N \cdot n_x}$ and a matrix $M \in \mathbb{R}^{(N \cdot n_x) \times (n_x + N \cdot n_u)}$, and dividing the variables into a dependent and an independent part, $\Delta w = (\Delta w_{\text{dep}}, \Delta w_{\text{ind}})$.

The vector v can be generated recursively by simulating the affine dynamic system (6.7) with all inputs set to zero, i.e. $\Delta w_{\text{ind}} = 0$. This yields the forward recursion

$$v_1 = a_0, \quad v_{k+1} = a_k + A_k v_k, \quad k = 1, \dots, N-1$$

for the components of the vector $v = (v_1, \dots, v_N)$. The subblocks of the matrix M can be obtained recursively as well in a straightforward way. Note that the matrix is lower triangular because the states Δx_j do not depend on Δu_k if $k \geq j$. On the other hand, if $k < j$, the corresponding matrix blocks are $A_{j-1} \cdots A_{k+1} B_k$. Finally, the dependence of Δx_j on Δx_0 is $A_{j-1} \cdots A_0$. In this way, all blocks of the matrix M are defined.

To get a notationally different, but equivalent view on condensing, note that the linear dynamic system equations (6.5) are nothing else than the linear system

$$\begin{bmatrix} A_0 & B_0 & -\mathbb{I} & & & \\ & A_1 & B_1 & -\mathbb{I} & & \\ & & \ddots & & \ddots & \\ & & & A_{N-1} & B_{N-1} & -\mathbb{I} \end{bmatrix} \begin{bmatrix} \Delta x_0 \\ \Delta u_0 \\ \Delta x_1 \\ \Delta u_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_{N-1} \\ \Delta u_{N-1} \\ \Delta x_N \end{bmatrix} = - \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_N \end{bmatrix} \quad (6.10)$$

After reordering the variables into dependent and independent ones, this system can be written as

$$\left[\begin{array}{ccc|ccc} A_0 & B_0 & & -\mathbb{I} & & \\ & A_1 & B_1 & -\mathbb{I} & & \\ & & \ddots & & \ddots & \\ & & & A_{N-1} & B_{N-1} & -\mathbb{I} \end{array} \right] \begin{bmatrix} \Delta x_0 \\ \Delta u_0 \\ \vdots \\ \Delta u_{N-1} \\ \Delta x_1 \\ \vdots \\ \Delta x_N \end{bmatrix} = - \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_N \end{bmatrix} \quad (6.11)$$

which we can summarize as

$$[X|Y] \begin{bmatrix} \Delta w_{\text{ind}} \\ \Delta w_{\text{dep}} \end{bmatrix} = -a$$

so that we get the explicit solution

$$\Delta w_{\text{dep}} = \underbrace{(-Y^{-1}a)}_{=v} + \underbrace{(-Y^{-1}X)}_{=M} \Delta w_{\text{ind}}.$$

Note that the submatrix Y is always invertible due the fact that it is lower triangular and has (negative) unit matrices on its diagonal.

Once the vector v and matrix M are computed, we can formulate a *condensed QP* which has only the independent variables Δw_{ind} as degrees of freedom. This condensed QP can be solved by a dense QP solver, and the resulting solution Δw_{ind}^* can be expanded again to yield also the QP solution for $w_{\text{dep}}^* = v + M\Delta w_{\text{ind}}^*$. The QP multipliers $\lambda_{\text{dep}} = (\lambda_1, \dots, \lambda_N)$ for the constraints (6.5) can be obtained from the dense QP solution in a slightly more complex way. The trick is to regard the Lagrangian of the original QP (6.5), $\mathcal{L}^{\text{QP}}(\Delta w_{\text{ind}}, \Delta w_{\text{dep}}, \lambda_{\text{dep}}, \lambda_r, \mu)$ and note that the condensed QP yields also the multipliers λ_r^*, μ^* , which turn out to be the correct multipliers also for the uncondensed QP. Thus, the only missing quantity is λ_{dep}^* . It can be obtained by using the following two observations: first, for the true QP solution must hold that the Lagrange gradient is zero, also with respect to Δw_{dep} . Second, this Lagrange gradient depends linearly on the unknown multipliers λ_{dep} which contribute to it via the term $Y^T \lambda_{\text{dep}}$, i.e. we have

$$0 = \nabla_{\Delta w_{\text{dep}}} \mathcal{L}^{\text{QP}}(\Delta w_{\text{ind}}^*, \Delta w_{\text{dep}}^*, \lambda_{\text{dep}}^*, \lambda_r^*, \mu^*) = \nabla_{\Delta w_{\text{dep}}} \mathcal{L}^{\text{QP}}(\Delta w_{\text{ind}}^*, \Delta w_{\text{dep}}^*, 0, \lambda_r^*, \mu^*) + Y^T \lambda_{\text{dep}}^*.$$

It is a favourable fact that the Lagrange gradient depends on the missing multipliers via the matrix Y^T , because this matrix is invertible. Thus, we obtain an explicit equation for obtaining the missing multipliers, namely

$$\lambda_{\text{dep}}^* = -Y^{-T} \nabla_{\Delta w_{\text{dep}}} \mathcal{L}^{\text{QP}}(\Delta w_{\text{ind}}^*, \Delta w_{\text{dep}}^*, 0, \lambda_r^*, \mu^*).$$

Note that the multipliers would not be needed within a Gauss-Newton method.

Summarizing, condensing reduces the original QP to a QP that has the size of the QP in the sequential approach. Nearly all sparsity is lost, but the dimension of the QP is much reduced. Condensing is favourable if the horizon length N and the control dimension n_u are relatively small compared to the state dimension n_x . If the initial value is fixed, then also Δx_0 can be eliminated from the condensed QP before passing it to a dense QP solver, further reducing the dimension.

On the other hand, if the state dimension n_x is very small compared to $N \cdot n_u$, condensing is not favourable due to the fact that it destroys sparsity. This is most easily seen in the Hessian. In the original sparse QP, the block sparse Hessian has $N(n_x + n_u)^2 + n_x^2$ nonzero elements. This is linear in N . In contrast to this, the condensed Hessian is dense and has $(n_x + Nn_u)^2$ elements, which is quadratic in N . Thus, if N is large, not only might the condensed Hessian need more (!) storage than the original one, also the solution time of the QP becomes cubic in N (factorization costs of the Hessian).

6.3.2 Sparse KKT System

A different way to exploit the sparsity present in the QP (6.4) is to keep all variables in the problem and use within the QP solver linear algebra routines that exploit sparsity of matrices. This can be realized within both, interior point (IP) methods as well as in active set methods, but is much easier to illustrate at the example of IP methods. For illustration, let us assume a problem without coupled constraints (6.6) and assume that all inequalities have been transformed into primal barrier terms that are added to the objective. Then, in each interior point iteration, an equality constrained QP of the following simple form needs to be solved.

$$\begin{aligned} \underset{\Delta x_0, \Delta u_0, \dots, \Delta x_N}{\text{minimize}} \quad & \frac{1}{2} \sum_{k=0}^{N-1} \begin{bmatrix} \Delta x_k \\ \Delta u_k \end{bmatrix}^T \begin{bmatrix} Q_k^x & Q_k^{xu} \\ (Q_k^{xu})^T & Q_k^u \end{bmatrix} \begin{bmatrix} \Delta x_k \\ \Delta u_k \end{bmatrix} + \frac{1}{2} \Delta x_N^T Q_N \Delta x_N \\ & + \sum_{k=0}^N \begin{bmatrix} \Delta x_N \\ \Delta u_N \end{bmatrix}^T g_k + \Delta x_N^T g_N \end{aligned} \quad (6.12)$$

$$\text{subject to} \quad a_k + A_k \Delta x_k + B_k \Delta u_k - \Delta x_{k+1} = 0, \quad \text{for } k = 0, \dots, N-1. \quad (6.13)$$

Formulating the Lagrangian of this QP and differentiating it with respect to all its primal and dual variables $y = (\Delta x_0, \Delta u_0, \lambda_1, \Delta x_1, \Delta u_1, \dots, \lambda_N, \Delta x_N)$ in this order we obtain a linear system of the following block tridiagonal

form

$$\begin{bmatrix}
 Q_0^x & Q_0^{xu} & A_0^T & & & & & \\
 (Q_0^{xu})^T & Q_0^u & B_0^T & & & & & \\
 A_0 & B_0 & 0 & -\mathbb{I} & & & & \\
 & & -\mathbb{I} & Q_1^x & Q_1^{xu} & A_1^T & & \\
 & & & (Q_1^{xu})^T & Q_1^u & B_1^T & & \\
 & & & A_1 & B_1 & 0 & -\mathbb{I} & \\
 & & & & & -\mathbb{I} & & \\
 & & & & & & \ddots & \ddots \\
 & & & & & & A_{N-1} & B_{N-1} & 0 & -\mathbb{I} \\
 & & & & & & & -\mathbb{I} & Q_N &
 \end{bmatrix}
 \begin{bmatrix}
 \Delta x_0 \\
 \Delta u_0 \\
 \lambda_1 \\
 \Delta x_1 \\
 \Delta u_1 \\
 \lambda_2 \\
 \vdots \\
 \lambda_N \\
 \Delta x_N
 \end{bmatrix}
 =
 \begin{bmatrix}
 * \\
 * \\
 * \\
 * \\
 * \\
 * \\
 * \\
 * \\
 * \\
 *
 \end{bmatrix} \quad (6.14)$$

This linear system can be solved with a banded direct factorization routine, whose runtime is proportional to $N(n_x + n_u)^3$. We will see in the next chapter that a particularly efficient way to solve the above linear system can be obtained by applying the principle of dynamic programming to the equality constrained quadratic subproblem (6.12).

Summarizing, the approach to directly solve the sparse QP without condensing is advantageous if Nn_u is large compared to n_x . It needs, however, sparse linear algebra routines within the QP solver. This is easier to implement in the case of IP methods than for active set methods.

Chapter 7

Dynamic Programming

”In view of all that we have said in the foregoing sections, the many obstacles we appear to have surmounted. What casts the pall over our victory celebration? It is the curse of dimensionality, a malediction that has plagued the scientist from earliest days.”

Richard E. Bellman

Dynamic programming (DP) is a very different approach to solve optimal control problems than the ones presented previously. The methodology was developed in the fifties and sixties of the 19th century, most prominently by Richard Bellman [Bel57] who also coined the term dynamic programming. Interestingly, dynamic programming is easiest to apply to systems with discrete state and control spaces, so that we will introduce this case first. When DP is applied to discrete time systems with continuous state spaces, some approximations have to be made, usually by discretization. Generally, this discretization leads to exponential growth of computational cost with respect to the dimension n_x of the state space, what Bellman called the “curse of dimensionality”. It is the only but major drawback of DP and limits its practical applicability to systems with $n_x \approx 6$. In the continuous time case, DP is formulated as a partial differential equation in the state space, the Hamilton-Jacobi-Bellman (HJB) equation, suffering from the same limitation. On the positive side, DP can easily deal with all kinds of hybrid systems or non-differentiable dynamics, and it even allows us to treat stochastic optimal control with recourse, or minimax games, without much additional effort. An excellent textbook on discrete time optimal control and dynamic programming is [Ber05]. Let us now start with discrete control and state spaces.

7.1 Dynamic Programming in Discrete State Space

Let us regard a dynamic system

$$x_{k+1} = f(x_k, u_k)$$

with $f : \mathbb{X} \times \mathbb{U} \rightarrow \mathbb{X}$, i.e. $x_k \in \mathbb{X}$ and $u_k \in \mathbb{U}$, where we do not have to specify the sets \mathbb{X} and \mathbb{U} yet. We note, however, that we need to assume they are finite for a practical implementation of DP. Thus, let us in this section assume they are finite with $n_{\mathbb{X}}$ and $n_{\mathbb{U}}$ elements, respectively. Let us also define a stage cost $L(x, u)$ and terminal cost $E(x)$ that take values from $\mathbb{R}_{\infty} = \mathbb{R} \cup \{\infty\}$, where infinity denotes infeasible pairs (x, u) or x . The optimal control problem that we first address can be stated as

$$\begin{aligned} & \underset{x_0, u_0, x_1, \dots, u_{N-1}, x_N}{\text{minimize}} && \sum_{k=0}^{N-1} L(x_k, u_k) + E(x_N) \end{aligned} \quad (7.1a)$$

$$\text{subject to} \quad f(x_k, u_k) - x_{k+1} = 0, \quad \text{for } k = 0, \dots, N-1, \quad (7.1b)$$

$$\bar{x}_0 - x_0 = 0. \quad (7.1c)$$

Given the fact that the initial value is fixed and the controls $\{u_k\}_{k=0}^{N-1}$ are the only true degrees of freedom, and given that each $u_k \in \mathbb{U}$ takes one of the $n_{\mathbb{U}}$ elements of \mathbb{U} , there exist exactly $n_{\mathbb{U}}^N$ different trajectories, each with a specific value of the objective function, where infinity denotes an infeasible trajectory. Assuming that the evaluation of f and of L takes one computational unit, and noting that each trajectory needs N such evaluations,

the overall complexity of simple enumeration is $O(Nn_{\mathbb{U}}^N)$. Simple enumeration of all possible trajectories thus has a complexity that grows exponentially with the horizon length N .

Dynamic programming is just a more intelligent way to enumerate all possible trajectories. It starts from the *principle of optimality*, i.e. the observation that each subtrajectory of an optimal trajectory is an optimal trajectory as well. More specifically, in DP we define the *value function* or *cost-to-go function* as the optimal cost that would be obtained if at time $k \in \{0, \dots, N\}$ and at state \bar{x}_k we solve the optimal control problem on a shortened horizon:

$$J_k(\bar{x}_k) = \min_{x_k, u_k, \dots, u_{N-1}, x_N} \sum_{i=k}^{N-1} L(x_i, u_i) + E(x_N) \quad (7.2a)$$

$$\text{subject to} \quad f(x_i, u_i) - x_{i+1} = 0, \quad \text{for } i = k, \dots, N-1, \quad (7.2b)$$

$$\bar{x}_k - x_k = 0. \quad (7.2c)$$

Thus, each function $J_k : \mathbb{X} \rightarrow \mathbb{R}_{\infty}$ summarizes the cost-to-go to the end when starting at a given state. For the case $k = N$ we trivially have $J_N(x) = E(x)$. The principle of optimality states now that for any $k \in \{0, \dots, N-1\}$ holds

$$J_k(\bar{x}_k) = \min_u L(\bar{x}_k, u) + J_{k+1}(f(\bar{x}_k, u)). \quad (7.3)$$

This immediately allows us to perform a recursion to compute all functions J_k one after the other, starting with $k = N-1$ and then reducing k in each recursion step by one, until we have obtained J_0 . This recursion is called the *dynamic programming recursion*. Once all the value functions J_k are computed, the *optimal feedback control* for a given state x_k at time k is given by

$$u_k^*(x_k) = \arg \min_u L(x_k, u) + J_{k+1}(f(x_k, u))$$

This allows us to reconstruct the optimal trajectory by a forward simulation that starts at $x_0 = \bar{x}_0$ and then proceeds as follows:

$$x_{k+1} = f(x_k, u_k^*(x_k)), \quad k = 0, \dots, N-1.$$

In this way, DP allows us to solve the optimal control problem up to global optimality, but with a different complexity than simple enumeration. To assess its complexity, let us remark that the most cost intensive step is the generation of the N cost-to-go functions J_k . Each recursion step (7.3) needs to go through all $n_{\mathbb{X}}$ states x . For each state it needs to test $n_{\mathbb{U}}$ controls u by evaluating once the system $f(x, u)$ and stage cost $L(x, u)$, which by definition costs one computational unit. Thus, the overall computational complexity is $O(Nn_{\mathbb{X}}n_{\mathbb{U}})$. Compared with simple enumeration, where we had $O(Nn_{\mathbb{U}}^N)$, DP is often much better even for moderately sized horizons N . Let us for example assume an optimal control problem with $n_{\mathbb{U}} = 10$, $n_{\mathbb{X}} = 1000$, $N = 100$. Then simple enumeration has a cost of 10^{102} while DP has a cost of 10^6 .

One of the main advantages of dynamic programming, that can likewise be defined for continuous state spaces, is that we do not need to make any assumptions (such as differentiability or convexity) on the functions f, L, E defining the problem, and still it solves the problem up to global optimality. On the other hand, if it shall be applied to a continuous state space, we have to represent the functions J_k on the computer, e.g. by tabulation on a grid in state space. If the continuous state space \mathbb{X}_{cont} is a box in dimension n_x , and if we use a rectangular grid with m intervals in each dimension, then the total number of grid points is m^{n_x} . If we perform DP on this grid, then the above complexity estimate is still valid, but with $n_{\mathbb{X}} = m^{n_x}$. Thus, when DP is applied to systems with continuous state spaces, it has exponential complexity in the dimension of the state space; it suffers from what Bellman called the *curse of dimensionality*. There exist many ways to approximate the value function, e.g. by neural networks or other functional representations [BT96], but the global optimality guarantee of dynamic programming is lost in these cases. On the other hand, there exists one special case where DP can be performed exactly in continuous state spaces, that we treat next.

7.2 Linear Quadratic Problems

Let us regard now linear quadratic optimal control problems of the form

$$\underset{x, u}{\text{minimize}} \quad \sum_{i=0}^{N-1} \begin{bmatrix} x_i \\ u_i \end{bmatrix}^T \begin{bmatrix} Q_i & S_i^T \\ S_i & R_i \end{bmatrix} \begin{bmatrix} x_i \\ u_i \end{bmatrix} + x_N^T P_N x_N \quad (7.4)$$

subject to

$$\begin{aligned} x_0 - \bar{x}_0 &= 0, \\ x_{i+1} - A_i x_i - B_i u_i &= 0, \quad i = 0, \dots, N-1. \end{aligned}$$

Let us apply dynamic programming to this case. In each recursion step, we have to solve, for a time varying stage

cost $L_k(x, u) = \begin{bmatrix} x_k \\ u_k \end{bmatrix}^T \begin{bmatrix} Q_k & S_k^T \\ S_k & R_k \end{bmatrix} \begin{bmatrix} x_k \\ u_k \end{bmatrix}$ and a dynamic system $f_k(x, u) = A_k x + B_k u$ the recursion step

$$J_k(x) = \min_u L_k(x, u) + J_{k+1}(f_k(x, u)),$$

where we start with $J_N(x) = x^T P_N x$. Fortunately, it can be shown that under these circumstances, each J_k is quadratic, i.e. it again has the form $J_k(x) = x^T P_k x$. More specifically, the following theorem holds, where we drop the index k for simplicity.

Theorem 10 (Quadratic Representation of Value Function) *If $R + B^T P B$ is positive definite, then the minimum $J_{\text{new}}(x)$ of one step of the DP recursion*

$$J_{\text{new}}(x) = \min_u \begin{bmatrix} x \\ u \end{bmatrix}^T \left(\begin{bmatrix} Q & S^T \\ S & R \end{bmatrix} + [A \mid B]^T P [A \mid B] \right) \begin{bmatrix} x \\ u \end{bmatrix}$$

is a quadratic function explicitly given by $J_{\text{new}}(x) = x^T P_{\text{new}} x$ with

$$P_{\text{new}} = Q + A^T P A - (S^T + A^T P B)(R + B^T P B)^{-1}(S + B^T P A). \quad (7.5)$$

The proof starts by noting that the optimization problem for a specific x is given by

$$J_{\text{new}}(x) = \min_u \begin{bmatrix} x \\ u \end{bmatrix}^T \begin{bmatrix} Q + A^T P A & S^T + A^T P B \\ S + B^T P A & R + B^T P B \end{bmatrix} \begin{bmatrix} x \\ u \end{bmatrix}.$$

Then it uses the fact that for invertible $\bar{R} = R + B^T P B$ this problem can be solved explicitly, yielding the formula (7.5), by a direct application of the *Schur complement lemma*, that can easily be verified by direct calculation.

Lemma 1 (Schur Complement Lemma) *If \bar{R} is positive definite then*

$$\min_u \begin{bmatrix} x \\ u \end{bmatrix}^T \begin{bmatrix} \bar{Q} & \bar{S}^T \\ \bar{S} & \bar{R} \end{bmatrix} \begin{bmatrix} x \\ u \end{bmatrix} = x^T (\bar{Q} - \bar{S}^T \bar{R}^{-1} \bar{S}) x \quad (7.6)$$

and the minimizer $u^*(x)$ is given by $u^*(x) = -\bar{R}^{-1} \bar{S} x$.

The above theorem allows us to solve the optimal control problem by first computing explicitly all matrices P_k , and then performing the forward closed loop simulation. More explicitly, starting with P_N , we iterate for $k = N-1, \dots, 0$ backwards

$$P_k = Q_k + A_k^T P_{k+1} A_k - (S_k^T + A_k^T P_{k+1} B_k)(R_k + B_k^T P_{k+1} B_k)^{-1}(S_k + B_k^T P_{k+1} A_k). \quad (7.7)$$

This is sometimes called the *Difference Riccati Equation*. Then, we obtain the optimal feedback $u_k^*(x_k)$ by

$$u_k^*(x_k) = -(R_k + B_k^T P_{k+1} B_k)^{-1}(S_k + B_k^T P_{k+1} A_k)x_k,$$

and finally, starting with $x_0 = \bar{x}_0$ we perform the forward recursion

$$x_{k+1} = A_k x_k + B_k u_k^*(x_k),$$

which delivers the complete optimal trajectory of the linear quadratic optimal control problem.

An important and more general case are problems with linear quadratic costs and affine linear systems, i.e. problems of the form

$$\underset{x, u}{\text{minimize}} \quad \sum_{i=0}^{N-1} \begin{bmatrix} 1 \\ x_i \\ u_i \end{bmatrix}^T \begin{bmatrix} * & q_i^T & s_i^T \\ q_i & Q_i & S_i^T \\ s_i & S_i & R_i \end{bmatrix} \begin{bmatrix} 1 \\ x_i \\ u_i \end{bmatrix} + \begin{bmatrix} 1 \\ x_N \end{bmatrix}^T \begin{bmatrix} * & p_N^T \\ p_N & P_N \end{bmatrix} \begin{bmatrix} 1 \\ x_N \end{bmatrix} \quad (7.8)$$

$$\begin{aligned} \text{subject to} \quad & x_0 - x_0^{\text{fix}} = 0, \\ & x_{i+1} - A_i x_i - B_i u_i - c_i = 0, \quad i = 0, \dots, N-1. \end{aligned}$$

These optimization problems appear at many occasions, for example as linearizations of nonlinear optimal control problems, as in Chapter 6, in reference tracking problems with $L_i(x_i, u_i) = \|x_i - x_i^{\text{ref}}\|_Q^2 + \|u_i\|_R^2$, or in *moving horizon estimation* (MHE) with cost $L_i(x_i, u_i) = \|C x_i - y_i^{\text{meas}}\|_Q^2 + \|u_i\|_R^2$. They can be treated by exactly the same recursion formulae as above, by augmenting the system states x_k to

$$\tilde{x}_k = \begin{bmatrix} 1 \\ x_k \end{bmatrix}$$

and replacing the dynamics by

$$\tilde{x}_{k+1} = \begin{bmatrix} 1 & 0 \\ c_k & A_k \end{bmatrix} \tilde{x}_k + \begin{bmatrix} 0 \\ B_k \end{bmatrix} u_k$$

with initial value

$$\tilde{x}_0^{\text{fix}} = \begin{bmatrix} 1 \\ x_0^{\text{fix}} \end{bmatrix}$$

Then the problem (7.8) can be reformulated in the form of problem (7.4) and can be solved using exactly the same difference Riccati equation formula as before!

7.3 Infinite Horizon Problems

Dynamic programming can easily be generalized to infinite horizon problems of the form

$$\begin{aligned} & \text{minimize}_{x, u} \quad \sum_{i=0}^{\infty} L(x_i, u_i) \\ & \text{subject to} \end{aligned}$$

$$\begin{aligned} x_0 - \bar{x}_0 &= 0, \\ x_{i+1} - f(x_i, u_i) &= 0, \quad i = 0, \dots, \infty. \end{aligned}$$

Interestingly, the cost-to-go function $J_k(x_k)$ defined in Equation (7.2) becomes independent of the index k , i.e it holds that $J_k = J_{k+1}$ for all k . This directly leads to the *Bellman Equation*:

$$J(x) = \min_u \underbrace{L(x, u) + J(f(x, u))}_{=\tilde{J}(x, u)}$$

The optimal controls are obtained by the function

$$u^*(x) = \arg \min_u \tilde{J}(x, u).$$

This feedback is called the *stationary optimal feedback control*. It is a static state feedback law.

7.4 The Linear Quadratic Regulator

An important special case is again the case of a linear system with quadratic cost. It is the solution to an infinite horizon problem with a linear system $f(x, u) = Ax + Bu$ and quadratic cost

$$L(x, u) = \begin{bmatrix} x \\ u \end{bmatrix}^T \begin{bmatrix} Q & S^T \\ S & R \end{bmatrix} \begin{bmatrix} x \\ u \end{bmatrix}.$$

For its solution, we just require a stationary solution of the Riccati recursion (7.7), setting $P_k = P_{k+1}$, which yields the so called *algebraic Riccati equation in discrete time*

$$P = Q + A^T P A - (S^T + A^T P B)(R + B^T P B)^{-1}(S + B^T P A).$$

This is a nonlinear matrix equation in the symmetric matrix P , i.e. with $n_x(n_x + 1)/2$ unknowns. It can either be solved by an iterative application of the difference Riccati recursion (7.7) starting with e.g. a zero matrix $P = 0$,

or by faster converging procedures such as Newton-type methods, where, however, care has to be taken to avoid possible shadow solutions that are not positive definite. Once the solution matrix P is found, the optimal feedback control $u^*(x)$ is given by

$$u^*(x) = - \underbrace{(R + B^T P B)^{-1} (S + B^T P A)}_{=K} x$$

This feedback is called the *Linear Quadratic Regulator (LQR)*, and K is the LQR gain.

7.5 Robust and Stochastic Dynamic Programming

One of its most interesting characteristics is that DP can easily be applied to games like chess, or to *robust optimal control problems*. Here, an adverse player choses counter-actions, or disturbances, w_k against us. They influence both the stage costs L_k as well as the system dynamics f_k and while we want to minimize, our adversary wants to maximize. The robust DP recursion for such a minimax game is simply:

$$J_k(x) = \min_u \max_w \underbrace{L_k(x, u, w) + J_{k+1}(f_k(x, u, w))}_{=\tilde{J}_k(x, u)}$$

starting with

$$J_N(x) = E(x).$$

The solution obtained by DP takes into account that we can react to the actions by the adversary, i.e. that we can apply feedback, and in the model predictive control (MPC) literature such a feedback law is sometimes called Closed-Loop Robust Optimal Control [BBM03].

Alternatively, we might have a stochastic system and the aim is to find the feedback law that gives us the best expected value. Here, instead of the maximum, we take an *expectation* over the disturbances w_k . The stochastic DP recursion is simply given by

$$J_k(x) = \min_u \mathbb{E}_w \{ \underbrace{L_k(x, u, w) + J_{k+1}(f_k(x, u, w))}_{=\tilde{J}_k(x, u)} \}$$

where $\mathbb{E}_w\{\cdot\}$ is the expectation operator, i.e. the integral over w weighted with the probability density function $\rho(w|x, u)$ of w given x and u :

$$\mathbb{E}_w\{\phi(x, u, w)\} = \int \phi(x, u, w) \rho(w|x, u) dw.$$

In case of finitely many disturbances, this is just a weighted sum. Note that DP avoids the combinatorial explosion of scenario trees that are often used in stochastic programming, but of course suffers from the curse of dimensionality. It is the preferred option for long horizon problems with small state spaces.

7.6 Interesting Properties of the DP Operator

Let us define the *dynamic programming operator* T_k acting on one value function, J_{k+1} , and giving another one, J_k , by

$$T_k[J](x) = \min_u L_k(x, u) + J(f_k(x, u)).$$

Note that the operator T_k maps from the space of functions $\mathbb{X} \rightarrow \mathbb{R}_\infty$ into itself. With this operator, the dynamic programming recursion is compactly written as $J_k = T_k[J_{k+1}]$, and the stationary Bellman equation would just be $J = T[J]$. Let us for notational simplicity drop the index k in the following. An interesting property of the DP operator T is its *monotonicity*, as follows.

Theorem 11 (Monotonicity of DP) *Regard two value functions J and J' . If $J \geq J'$ in the sense that for all $x \in \mathbb{X}$ holds that $J(x) \geq J'(x)$ then also*

$$T[J] \geq T[J'].$$

The proof is

$$T[J](x) = \min_u L(x, u) + \underbrace{J(f(x, u))}_{\geq J'(f(x, u))} \geq \min_u L(x, u) + J'(f(x, u)) = T[J'](x)$$

This monotonicity property holds also for robust or stochastic dynamic programming, and is for example used in existence proofs for solutions of the stationary Bellman equation, or in stability proofs of model predictive control (MPC) schemes [MRRS00].

Another interesting observation is that certain DP operators T preserve convexity of the value function J .

Theorem 12 (Convex dynamic programming) *If the system is affine in (x, u) , i.e. $f(x, u, w) = A(w)x + B(w)u + c(w)$, and if the stage cost $L(x, u, w)$ is convex in (x, u) , then the DP, the robust DP, and the stochastic DP operators T preserve convexity of J , i.e. if J is a convex function, then $T[J]$ is again a convex function.*

Proof: It is interesting to note that no restrictions are given on how the functions depend on w . The proof of the convexity preservation starts by noting that for fixed w , $L(x, u, w) + J(f(x, u, w))$ is a convex function in (x, u) . Because also the maximum over all w , or the positively weighted sum of an expectation value computation, preserve convexity, the function $\tilde{J}(x, u)$ is in all three cases convex in both x and u . Finally, the minimization of a convex function over one of its arguments preserves convexity, i.e. the resulting value function $T[J]$ defined by

$$T[J](x) = \min_u \tilde{J}(x, u)$$

is convex.

But why would convexity be important in the context of DP? First, convexity of $\tilde{J}(x, u)$ implies that the computation of the feedback law $\arg \min_u \tilde{J}(x, u)$ is a convex parametric program and could reliably be solved by local optimization methods. Second, it might be possible to represent the value function $J(x)$ more efficiently than by tabulation on a grid, for example as the pointwise maximum of affine functions

$$J(x) = \max_i a_i^T \begin{bmatrix} 1 \\ x \end{bmatrix}.$$

It is an interesting fact that that for piecewise linear convex costs and constraints and polyhedral uncertainty this representation is exact and leads to an exact robust DP algorithm that might be called *polyhedral DP* [BBM03, DB04]. The polyhedral convex representability of the cost-to-go for linear systems with piecewise linear cost is indirectly exploited in some explicit MPC approaches [PDB⁺00, BBM02]. Polyhedral representations with a limited number of facets can also be used to approximate a convex cost-to-go and still yield some guarantees on the closed-loop system [BD06, BD08, JM10]. Finally, note that also the linear quadratic regulator is a special case of convex dynamic programming.

7.7 The Gradient of the Value Function

The meaning of the cost-to-go, or the value function, J_k is that it is the cost incurred on the remainder of the horizon for the best possible strategy. In order to make an interesting connection between the value function and the multipliers λ_k that we encountered in derivative based optimization methods, let us now regard a discrete time optimal control problem as in the previous chapters, but without coupled constraints, as these cannot directly be treated with dynamic programming. We assume further that the initial value is fixed and that all inequality and terminal constraints are subsumed in the stage cost $L(x, u)$ and terminal cost $E(x_N)$ by barrier functions that take infinite values outside the feasible domain but are differentiable inside. For terminal equality constraints, e.g. a fixed terminal state, assume for the moment that these are approximated by a terminal region of non-zero volume on which again a barrier can be defined. Thus, we regard the following problem.

$$\begin{aligned} & \underset{x_0, u_0, x_1, \dots, u_{N-1}, x_N}{\text{minimize}} && \sum_{k=0}^{N-1} L(x_k, u_k) + E(x_N) \end{aligned} \tag{7.9a}$$

$$\text{subject to} \quad f(x_k, u_k) - x_{k+1} = 0, \quad \text{for } k = 0, \dots, N-1, \tag{7.9b}$$

$$\bar{x}_0 - x_0 = 0. \tag{7.9c}$$

The dynamic programming recursion for this problem is given by:

$$J_N(x) = E(x), \quad J_k(x) = \min_u L(x, u) + J_{k+1}(f(x, u)), \quad k = N-1, \dots, 0. \quad (7.10)$$

We remember that we obtained the optimal solution by the forward recursion

$$x_0 = \bar{x}_0, \quad x_{k+1} = f(x_k, u_k), \quad k = 0, \dots, N-1,$$

where u_k is defined by

$$u_k = \arg \min_u L(x_k, u) + J_{k+1}(f(x_k, u)). \quad (7.11)$$

The solution of this optimization problem in u necessarily satisfies the first order necessary optimality condition

$$\nabla_u L(x_k, u_k) + \frac{\partial f}{\partial u}(x_k, u_k)^T \nabla J_{k+1}(f(x_k, u_k)) = 0 \quad (7.12)$$

which defines u_k locally if the problem is locally strictly convex, i.e., if objective has a positive definite Hessian at (x_k, u_k) . We now formulate simple conditions on x_k and u_k that follow necessarily from the DP recursion. For this aim we first note that on the optimal trajectory holds $x_{k+1} = f(x_k, u_k)$ and that we trivially obtain along the optimal trajectory

$$J_N(x_N) = E(x_N), \quad J_k(x_k) = L(x_k, u_k) + J_{k+1}(x_{k+1}), \quad k = N-1, \dots, 0.$$

This implies for example that the value function remains constant on the whole trajectory for problems with zero stage costs. However, it is even more interesting to regard the gradient $\nabla J_k(x_k)$ along the optimal state trajectory. If we differentiate (7.10) at the point x_k with respect to x we obtain

$$\nabla J_N(x_k) = \nabla E(x_k), \quad \nabla J_k(x_k)^T = \frac{d}{dx} \underbrace{L(x_k, u_k) + J_{k+1}(f(x_k, u_k))}_{=: \tilde{J}_k(x_k, u_k)} \quad k = N-1, \dots, 0. \quad (7.13)$$

In the evaluation of the total derivative it is needed to observe that the optimal u_k is via (7.12) an implicit function of x_k . However, it turns out that the derivative does not depend on $\frac{du_k}{dx_k}$ because of

$$\frac{d}{dx} \tilde{J}_k(x_k, u_k) = \frac{\partial \tilde{J}_k}{\partial x}(x_k, u_k) + \underbrace{\frac{\partial \tilde{J}_k}{\partial u}(x_k, u_k)}_{=0} \frac{du_k}{dx_k}, \quad (7.14)$$

where the partial derivative with respect to u is zero because of (7.12). Thus, the gradients of the value function at the optimal trajectory have to satisfy the recursion

$$\nabla J_k(x_k) = \nabla_x L(x_k, u_k) + \frac{\partial f}{\partial x}(x_k, u_k)^T \nabla J_{k+1}(x_{k+1}) \quad k = N-1, \dots, 0. \quad (7.15)$$

This recursive condition on the gradients $\nabla J_k(x_k)$ is equivalent to the first order necessary condition (FONC) for optimality that we obtained previously for differentiable optimal control problems, if we identify the gradients with the multipliers, i.e. set

$$\lambda_k = \nabla J_k(x_k). \quad (7.16)$$

This is a very important interpretation of the multipliers λ_k : they are nothing else than the gradients of the value function along the optimal trajectory!

7.8 A Discrete Time Minimum Principle

Collecting all necessary conditions of optimality that we just derived, but substituting $\nabla J_k(x_k)$ by λ_k we arrive indeed exactly to the same conditions (5.16) that we derived in Chapter 5 in a completely different way.

$$x_0 = \bar{x}_0 \quad (7.17a)$$

$$x_{k+1} = f(x_k, u_k), \quad k = 0, \dots, N-1, \quad (7.17b)$$

$$\lambda_N = \nabla_{x_N} E(x_N) \quad (7.17c)$$

$$\lambda_k = \nabla_x L(x_k, u_k) + \frac{\partial f}{\partial x}(x_k, u_k)^T \lambda_{k+1}, \quad k = N-1, \dots, 1, \quad (7.17d)$$

$$0 = \nabla_u L(x_k, u_k) + \frac{\partial f}{\partial u}(x_k, u_k)^T \lambda_{k+1}, \quad k = 0, \dots, N-1. \quad (7.17e)$$

In the context of continuous time problems, we will arrive at a very similar formulation, which has the interesting features that the recursion for λ becomes a differential equation that can be integrated forward in time if desired, and that the optimization problem in (7.11) does only depend on the gradient of J . This will facilitate the formulation and numerical solution of the necessary optimality conditions as a boundary value problem.

Chapter 8

Continuous Time Optimal Control Problems

When we are confronted with a problem whose dynamic system lives in continuous time and whose control inputs are a function, we speak of a *continuous time optimal control problem*. This type of problem is the focus of this and the following chapter. We will encounter variations of the same concepts as in the discrete time setting, such as Lagrange multipliers λ , the value function J , or the difference between sequential or simultaneous methods. Some numerical methods and details, however, are only relevant in the continuous time setting, such as Pontryagin's Maximum Principle (not covered in this lecture) or the ODE solvers with sensitivity generation described in Section 9.4.

8.1 Formulation of Continuous Time Optimal Control Problems

In an ODE setting, a continuous time optimal control problem can be stated as follows.

$$\begin{aligned}
 & \underset{x(\cdot), u(\cdot)}{\text{minimize}} && \int_0^T L(x(t), u(t)) dt + E(x(T)) && (8.1) \\
 & \text{subject to} && && \\
 & x(0) - x_0 = 0, && && \text{(fixed initial value)} \\
 & \dot{x}(t) - f(x(t), u(t)) = 0, && t \in [0, T], && \text{(ODE model)} \\
 & h(x(t), u(t)) \leq 0, && t \in [0, T], && \text{(path constraints)} \\
 & r(x(T)) \leq 0, && && \text{(terminal constraints).}
 \end{aligned}$$

The problem and its variables are visualized in Figure 8.1.

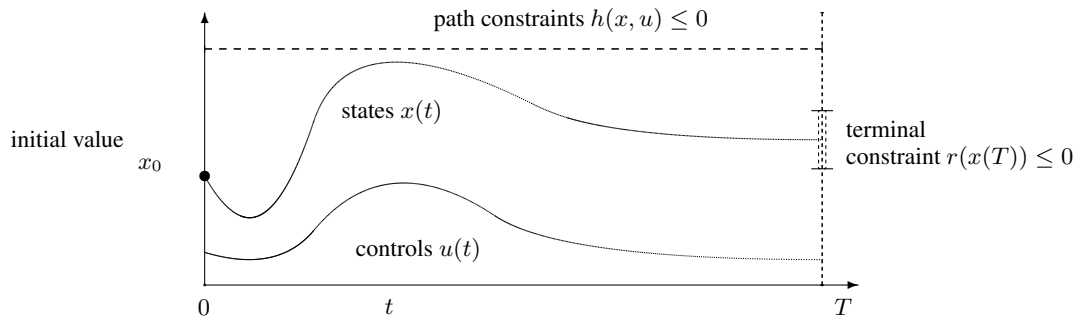


Figure 8.1: The variables and constraints of a continuous time optimal control problem.

The integral cost contribution $L(x, u)$ is sometimes called the *Lagrange term* (which should not be confused with the Lagrange function) and the terminal cost $E(x(T))$ is sometimes called a *Mayer term*. The combination of both, like here, is called a *Bolza objective*.

Note that any Lagrange objective term can be reformulated as a Mayer term, if we add an additional “cost state” c that has to satisfy the differential equation $\dot{c} = L(x, u)$, and then simply take $c(T)$ as the terminal Mayer cost term. Conversely, every differentiable Mayer term can be replaced by a Lagrange term, namely by $L(x, u) = \nabla E(x)^T f(x, u)$, as the cost integral then satisfies the equality $\int_0^T L(x, u) dt = \int_0^T \frac{dE}{dt} dt = E(x(T)) - E(x_0)$. These two equivalences mean that it would be no restriction of generality to take only one of the two cost contributions, Lagrange or Mayer term, in the above formulation; however, in this script we choose to use the full Bolza objective.

So far, we wrote all functions L, E, f, h independent of time t or of parameters p , and we will leave both of these generalizations away in the remainder of this script. However, all the methods presented in the following chapters can easily be adapted to these two cases, using again state augmentation, as follows. If time dependence occurs, we just introduce a “clock state” t with differential equation $\dot{t} = 1$, and work with the augmented system $\dot{\tilde{x}} = \tilde{f}(\tilde{x}, u)$:

$$\tilde{x} = \begin{bmatrix} x \\ t \end{bmatrix}, \tilde{f}(\tilde{x}, u) = \begin{bmatrix} f(x, u, t) \\ 1 \end{bmatrix}$$

Likewise, in the case that time constant, but free optimization parameters p occur, they can be incorporated as “parameter state” p with differential equation $\dot{p} = 0$ and free initial value.

Another interesting case that is specific to continuous time problems is when the duration T of the problem is free. As an example, we might think of a robot arm that should move an object in minimal time from its current state to some desired terminal position. In this case, we might rescale the time horizon to the interval $[0, 1]$ by a time constant but free variable T that is treated like an optimization parameter. Then we regard a scaled problem $\dot{\tilde{x}} = \tilde{f}(\tilde{x}, u)$

$$\tilde{x} = \begin{bmatrix} x \\ T \end{bmatrix}, \tilde{f}(\tilde{x}, u) = \begin{bmatrix} T \cdot f(x, u) \\ 0 \end{bmatrix}$$

with pseudo time $\tau \in [0, 1]$, where the initial condition $T(0)$ for the “state” T is free and T satisfies again $\dot{T} = 0$.

We note that although all the above reformulations make it possible to transfer the methods in this script to the respective special cases, an efficient numerical implementation should exploit the structures inherent in these special cases.

8.2 Overview of Numerical Approaches

Generally speaking, there are three basic families of approaches to address continuous time optimal control problems, (a) state-space, (b) indirect, and (c) direct approaches, cf. the top row of Fig. 8.2. We follow here the outline given in [DBDW06].

State-space approaches use the principle of optimality that states that each subarc of an optimal trajectory must be optimal. While this was the basis of dynamic programming in discrete time, in the continuous time case this leads to the so-called *Hamilton-Jacobi-Bellman (HJB) equation*, a partial differential equation (PDE) in the state space. Methods to numerically compute solution approximations exist, but the approach, like dynamic programming, suffers from Bellmans “curse of dimensionality” and is restricted to small state dimensions. This approach is not treated in this script.

Indirect Methods use the necessary conditions of optimality of the infinite problem to derive a boundary value problem (BVP) in ordinary differential equations (ODE). This BVP must numerically be solved, and the approach is often sketched as “first optimize, then discretize”. The class of indirect methods encompasses also the well known calculus of variations and the Euler-Lagrange differential equations, and the so-called *Pontryagin Maximum Principle*. The numerical solution of the BVP is performed by shooting techniques or by collocation. The two major drawbacks are that the underlying differential equations are often difficult to solve due to strong nonlinearity and instability, and that changes in the control structure, i.e. the sequence of arcs where different constraints are active, are difficult to handle: they usually require a completely new problem setup. Moreover, on so called singular arcs, higher index differential-algebraic equations (DAE) arise which necessitate specialized solution techniques. This approach is not explained in detail in this script.

Direct methods transform the original infinite optimal control problem into a finite dimensional nonlinear programming problem (NLP) which is then solved by structure exploiting numerical optimization methods. Roughly speaking, direct methods transform the continuous time dynamic system into a discrete time system and then proceed as described in the preceding chapters of this script. The approach is therefore often sketched as “first discretize, then optimize”. One of the most important advantages of direct compared to indirect methods is that they

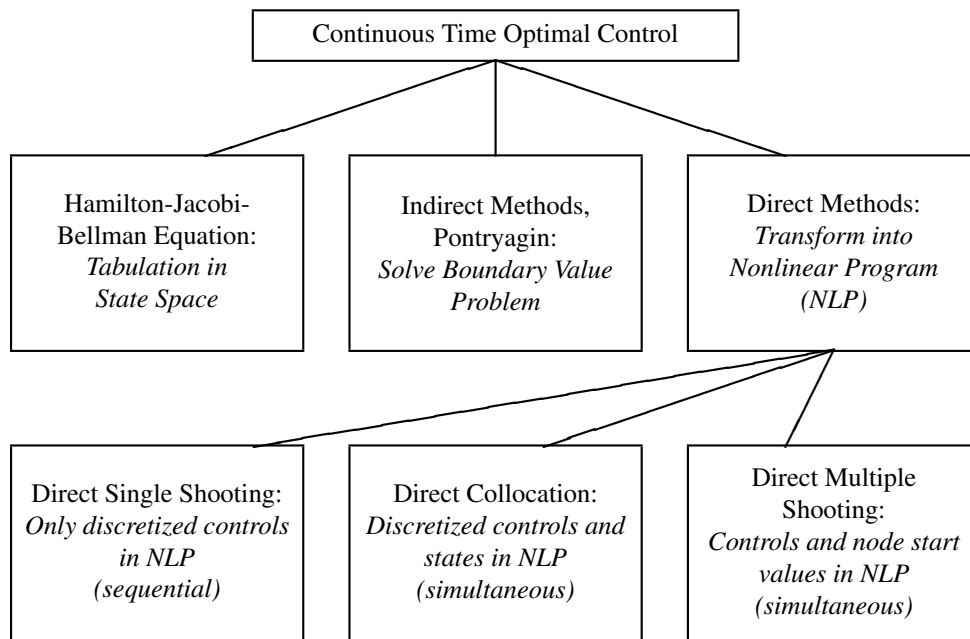


Figure 8.2: The optimal control family tree.

can easily treat inequality constraints, like the inequality path constraints in the formulation above. This is because structural changes in the active constraints during the optimization procedure are treated by well developed NLP methods that can deal with inequality constraints and active set changes. All direct methods are based on a finite dimensional parameterization of the control trajectory, but differ in the way the state trajectory is handled, cf. the bottom row of Fig. 8.2. For solution of constrained optimal control problems in real world applications, direct methods are nowadays by far the most widespread and successfully used techniques, and are therefore the focus of this script. Brief descriptions of three of the direct methods – single shooting, multiple shooting, and collocation – and some algorithmic details are given in Chapter 9, while we point out that the first chapters of the script covering finite dimensional optimization and discrete time dynamic systems have already covered most of the algorithmic ideas relevant for direct approaches to optimal control.

Chapter 9

Direct Approaches to Continuous Optimal Control

Direct methods to continuous optimal control finitely parameterize the infinite dimensional decision variables, notably the controls $u(t)$, such that the original problem is approximated by a finite dimensional nonlinear program (NLP). This NLP can then be addressed by structure exploiting numerical NLP solution methods. For this reason, the approach is often characterized as “First discretize, then optimize.” The direct approach connects easily to all optimization methods developed in the continuous optimization community, such as the methods described in Chapters 2 and 3. Most successful direct methods even parameterize the problem such that the resulting NLP has the structure of a discrete time optimal control problem, such that all the techniques and structures described in Chapters 5 and 6 are applicable. For this reason, the current chapter is kept relatively short; its major aim is to outline the major concepts and vocabulary in the field.

We start by describing *direct single shooting*, *direct multiple shooting*, and *direct collocation* and a variant *pseudospectral methods*. We also discuss how sensitivities are computed in the context of shooting methods. The optimization problem formulation we address in this chapter is the same as (8.1) in Chapter 8. The direct methods differ in how they transcribe this problem into a finite NLP. The problem (8.1) has a fixed initial value which simplifies in particular the single shooting method, but all concepts can in a straightforward way be generalized to other OCP formulations with free initial values.

9.1 Direct Single Shooting

All shooting methods use an embedded ODE or DAE solver in order to eliminate the continuous time dynamic system. They do so by first parameterizing the control function $u(t)$, e.g. by polynomials, by piecewise constant functions, or, more generally, by piecewise polynomials. We denote the finite control parameters by the vector q , and the resulting control function by $u(t; q)$. The most widespread parameterization are piecewise constant controls, for which we choose a fixed grid $0 = t_0 < t_1 < \dots < t_N = T$, and N parameters $q_i \in \mathbb{R}^{n_u}$, $i = 0, \dots, N-1$, and then we set

$$u(t; q) = q_i \quad \text{if } t \in [t_i, t_{i+1}].$$

Thus, the dimension of the vector $q = (q_0, \dots, q_{N-1})$ is Nn_u . In single shooting, which is a *sequential approach* earliest presented in [HR71, SS78], we then regard the states $x(t)$ on $[0, T]$ as dependent variables that are obtained by a forward integration of the dynamic system, starting at x_0 and using the controls $u(t; q)$. We denote the resulting trajectory as $x(t; q)$. In order to discretize inequality path constraints, we choose a grid, typically the same as for the control discretization, at which we check the inequalities. Thus, in single shooting, we transcribe the OCP (8.1) into the following NLP, that is visualized in Figure 9.1.

$$\begin{aligned} & \underset{q \in \mathbb{R}^{Nn_u}}{\text{minimize}} && \int_0^T L(x(t; q), u(t; q)) dt + E(x(T; q)) \end{aligned} \tag{9.1}$$

subject to

$$\begin{aligned} h(x(t_i; q), u(t_i; q)) & \leq 0, & i = 0, \dots, N-1, & \text{(discretized path constraints)} \\ r(x(T; q)) & \leq 0. & & \text{(terminal constraints)} \end{aligned}$$

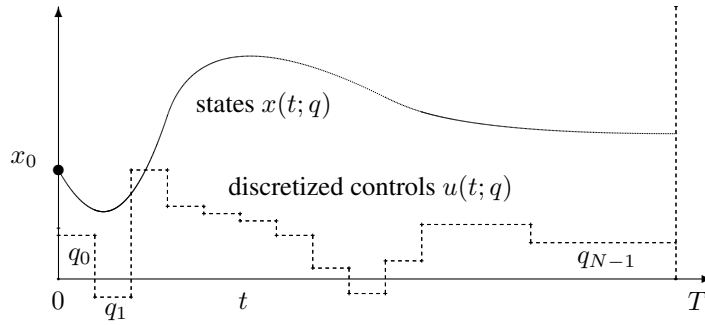


Figure 9.1: The NLP variables in the direct single shooting method.

As the only variable of this NLP is the vector $q \in \mathbb{R}^{Nn_u}$ that influences nearly all problem functions, the above problem can usually be solved by a dense NLP solver in a black-box fashion. As the problem functions and their derivatives are expensive to compute, while a small QP is cheap to solve, often Sequential Quadratic Programming (SQP) is used, e.g. the codes NPSOL or SNOPT. Let us first assume the Hessian needs not be computed but can be obtained e.g. by BFGS updates.

The computation of the derivatives can be done in different ways with a different complexity: first, we can use forward derivatives, using finite differences or algorithmic differentiation. Taking the computational cost of integrating one time interval as one computational unit, this means that one complete forward integration costs N units. Given that the vector q has Nn_u components, this means that the computation of all derivatives costs $(Nn_u + 1)N$ units when implemented in the most straightforward way. This number can still be reduced by one half if we take into account that controls at the end of the horizon do not influence the first part of the trajectory. We might call this way the *reduced derivative computation* as it computes directly only the reduced quantities needed in each reduced QP.

Second, if the number of output quantities such as objective and inequality constraints is not big, we can use the principle of reverse automatic differentiation in order to generate the derivatives. In the extreme case that no inequality constraints are present and we only need the gradient of the objective, this gradient can cheaply be computed by reverse AD, as done in the so called *gradient methods*. Note that in this case the same adjoint differential equations of the indirect approach can be used for reverse computation of the gradient, but that in contrast to the indirect method we do not eliminate the controls, and we integrate the adjoint equations backwards in time. The complexity for one gradient computation is only $4N$ computational units. However, each additional state constraint necessitates a further backward sweep.

Third, in the case that we have chosen piecewise controls, as here, we might use the fact that after the piecewise control discretization we have basically transformed the continuous time OCP into a discrete time OCP (see next section). Then we can compute the derivatives with respect to both s_i and q_i on each interval separately, which costs $(n_x + n_u + 1)$ units. This means a total derivative computation cost of $N(n_x + n_u + 1)$ units. In contrast to the second (adjoint) approach, this approach can handle an arbitrary number of path inequality constraints, like the first one. Note that it has the same complexity that we obtain in the standard implementation of the multiple shooting approach, as explained next. We remark here already that both shooting methods can each implement all the above ways of derivative generation, but differ in one respect only, namely that single shooting is a sequential and multiple shooting a simultaneous approach.

9.2 Direct Multiple Shooting

The direct multiple shooting method that was originally developed by Bock and Plitt [BP84] performs first a piecewise control discretization on a grid, exactly as we did in single shooting, i.e. we set

$$u(t) = q_i \quad \text{for } t \in [t_i, t_{i+1}].$$

But then, it solves the ODE separately on each interval $[t_i, t_{i+1}]$, starting with artificial initial values s_i :

$$\begin{aligned}\dot{x}_i(t; s_i, q_i) &= f(x_i(t; s_i, q_i), q_i), \quad t \in [t_i, t_{i+1}], \\ x_i(t_i; s_i, q_i) &= s_i.\end{aligned}$$

Thus, we obtain trajectory pieces $x_i(t; s_i, q_i)$. Likewise, we numerically compute the integrals

$$l_i(s_i, q_i) := \int_{t_i}^{t_{i+1}} L(x_i(t; s_i, q_i), q_i) dt.$$

Finally, we choose a grid at which we check the inequality path constraints; here we choose the same as for the controls and states, but note that a much finer sampling would be possible as well, which, however, requires continuous output from the integrator. Thus, the NLP that is solved in multiple shooting and that is visualized in Figure 9.2 is given by

$$\begin{aligned}& \underset{s, q}{\text{minimize}} && \sum_{i=0}^{N-1} l_i(s_i, q_i) + E(s_N) && (9.2) \\& \text{subject to} && && \\& x_0 - s_0 &= 0, && \text{(initial value)} \\& x_i(t_{i+1}; s_i, q_i) - s_{i+1} &= 0, \quad i = 0, \dots, N-1, && \text{(continuity)} \\& h(s_i, q_i) &\leq 0, \quad i = 0, \dots, N, && \text{(discretized path constraints)} \\& r(s_N) &\leq 0. && \text{(terminal constraints)}\end{aligned}$$

Note that by setting $f_i(s_i, q_i) := x_i(t_{i+1}; s_i, q_i)$ the continuity conditions can be interpreted a discrete time dynamic system $s_{i+1} = f_i(s_i, q_i)$ and the above optimal control problem has exactly the same structure as the discrete time optimal control problem (6.1) discussed in detail in Chapter 6. Most important, we can and should employ a sparsity exploiting NLP solver. Regarding the derivative computation, nearly all cost resides in the derivatives of the discrete time dynamic system, i.e. the matrices A_i and B_i in (6.5). If again the simulation on one interval, i.e. one evaluation of f_i , costs one unit, then the computation of these matrices by finite differences costs $(n_x + n_u + 1)$, and as we need N of them, we have a total derivative computation cost of $N(n_x + n_u + 1)$ per Newton-type iteration.

Remark on Schlöder’s Reduction Trick: We point out here that the derivatives of the condensed QP could also directly be computed, using the reduced way, as explained as first variant in the context of single shooting. It exploits the fact that the initial value x_0 is fixed in the NMPC problem, changing the complexity of the derivative computations. It is only advantageous for large state but small control dimensions as it has a complexity of $N^2 n_u$. It was originally developed by Schlöder [Sch88] in the context of Gauss-Newton methods and generalized to general SQP shooting methods by [Sch05]. A further generalization of this approach to solve a “lifted” (larger, but equivalent) system with the same computational cost per iteration is the so called *lifted Newton method* [AD10] where also an analysis of the benefits of lifting is made.

The main advantages of lifted Newton approaches such as multiple shooting compared with single shooting are the facts that (a) we can also initialize the state trajectory, and (b), that they show superior local convergence properties in particular for unstable systems. An interesting remark is that if the original system is linear, continuity is perfectly satisfied in all SQP iterations, and single and multiple shooting would be identical. Also, it is interesting to recall that the Lagrange multipliers λ_i for the continuity conditions are an approximation of the adjoint variables, and that they indicate the costs of continuity.

Finally, it is interesting to note that a direct multiple shooting algorithm can be made a single shooting algorithm easily: we only have to overwrite, before the derivative computation, the states s by the result of a forward simulation using the controls q obtained in the last Newton-type iteration. From this perspective, we can regard single shooting as a variant of multiple shooting where we perturb the result of each iteration by a “feasibility improvement” that makes all continuity conditions feasible by the forward simulation, implicitly giving priority to the control guess over the state guess [TWR04].

9.3 Direct Collocation

A third important class of direct methods are the direct transcription methods, most notably *direct collocation*. Here we discretize the infinite OCP in both controls and states on a fixed and relatively fine grid $\{t_k\}_{k=0}^N$; recall that

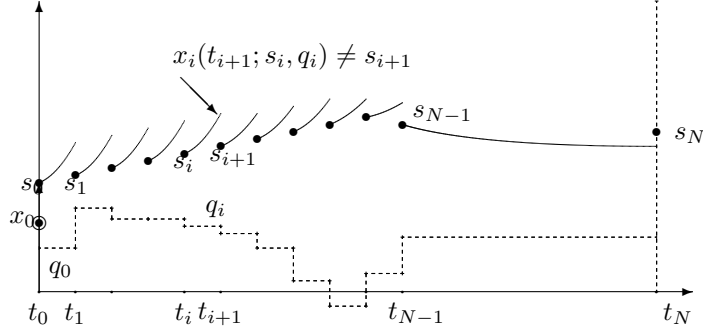


Figure 9.2: The NLP variables in the direct multiple shooting method.

each collocation interval corresponds to an integrator step. We denote the states on the grid points by $s_k \approx x(t_k)$. We choose a parameterization of the controls on the same grid, e.g. piecewise constant or piecewise polynomials, with control parameters q_k that yield on each interval a function $u_k(t; q)$

On each collocation interval $[t_k, t_{k+1}]$ a set of m collocation points $t_k^{(1)}, \dots, t_k^{(m)}$ is chosen and the trajectory is approximated by a polynomial $p_k(t; v_k)$ with coefficient vector v_k . As equalities of the optimization problem we now require that the collocation conditions (9.3) are met at the collocation points.

$$s_k = p_k(t_k; v_k) \quad (9.3a)$$

$$f(p_k(t_k^{(1)}; v_k), u_k(t_k^{(1)}; q_k)) = p'_k(t_k^{(1)}; v) \quad (9.3b)$$

$$\vdots \quad (9.3c)$$

$$f(p_k(t_k^{(m)}; v_k), u_k(t_k^{(m)}; q_k)) = p'_k(t_k^{(m)}; v) \quad (9.3d)$$

We summarize this system by the vector equation $c_k(s_k, v_k, q_k) = 0$ that has as many components as the vector v_i . Additionally, we also require continuity across interval boundaries, i.e. we add the constraints $p_k(t_{k+1}; v_k) - s_{k+1} = 0$. We also approximate the integrals $\int_{t_k}^{t_{k+1}} L(x, u) dt$ on the collocation intervals by a quadrature formula using the same collocation points, which we denote by the term $l_k(s_k, v_k, q_k)$. Path constraints are enforced on a grid, e.g. the interval boundaries, which we do here. We point out, that much finer sampling is possible as well, e.g. on the collocation nodes or even more often. Thus, we obtain a large scale, but sparse NLP:

$$\underset{s, v, q}{\text{minimize}} \quad \sum_{k=0}^{N-1} l_k(s_k, v_k, q_k) + E(s_N) \quad (9.4)$$

subject to

$$\begin{aligned} s_0 - x_0 &= 0, & & \text{(fixed initial value)} \\ c_k(s_k, v_k, q_k) &= 0, & k = 0, \dots, N-1, & \text{(collocation conditions)} \\ p_k(t_{k+1}; v_k) - s_{k+1} &= 0, & k = 0, \dots, N-1, & \text{(continuity conditions)} \\ h(s_k, q_k) &\leq 0, & k = 0, \dots, N-1, & \text{(discretized path constraints)} \\ r(s_N) &\leq 0. & & \text{(terminal constraints)} \end{aligned}$$

This large sparse NLP needs to be solved by structure exploiting solvers, and due to the fact that the problem functions are typically relatively cheap to evaluate compared with the cost of the linear algebra, nonlinear interior point methods are often the most efficient approach here. A widespread combination is to use collocation with IPOPT using the AMPL interface. It is interesting to note that, like in direct multiple shooting, the multipliers of the continuity conditions are again an approximation of the adjoint variables.

An interesting variant of orthogonal collocation methods that is often called the *pseudospectral optimal control method* uses only one collocation interval but on this interval it uses an extremely high order polynomial. State constraints are then typically enforced at all collocation points.

9.4 Sensitivity Computation in Shooting Methods

In all shooting methods we need to compute derivatives of the result of an ODE integration routine, or, in the more general case, of a DAE solver, on a given time interval. Let us for notational simplicity regard just the autonomous ODE case $\dot{x} = f(x)$ on a time interval $[0, T]$. The case of control or other parameters on which this ODE depends as well as time dependence can conceptually be covered by state augmentation. Thus, we regard a starting point s and the evolution of the ODE

$$\dot{x} = f(x), \quad t \in [0, T], \quad x(0) = s. \quad (9.5)$$

This gives a solution $x(t; s)$, $t \in [0, T]$, and we are most interested in the terminal value $x(T; s)$ and in the sensitivity matrix

$$G(t) = \frac{\partial x(t; s)}{\partial s}, \quad t \in [0, T],$$

and in particular its terminal value. This matrix $G(T) \in \mathbb{R}^{n_x \times n_x}$ can be computed in many different ways, five of which we briefly sketch here.

1. External Numerical Differentiation (END)
2. Solution of the Variational Differential Equations
3. Algorithmic Differentiation (AD) of the Integrator
4. Internal Algorithmic Differentiation within the Integrator
5. Internal Numerical Differentiation (IND)

In all five methods we assume that the integrator to be differentiated is a state-of-the-art integrator with inbuilt error control and adaptive step-size selection.

The first approach, *External Numerical Differentiation (END)*, just treats the integrator as a black box function and uses finite differences. We perturb s by some quantity $\epsilon > 0$ in the direction of the unit vectors e_i and call the integrator several times in order to compute directional derivatives by finite differences:

$$G(T)e_i \approx \frac{x(T; s + \epsilon e_i) - x(T; s)}{\epsilon} \quad (9.6)$$

The cost of this approach to compute $G(T)$ is $(n_x + 1)$ times the cost of a forward simulation. The approach is very easy to implement, but suffers from one serious problem: due to integrator adaptivity, each call might have a different discretization grid. This error control of each trajectory does not only create an overhead, but worse, it might result in discontinuous perturbations even for small ϵ . It is important to note that due to adaptivity, the output $x(T; s)$ is not a differentiable function in s , but only guaranteed to be close to the true solution within the integrator accuracy TOL, e.g. TOL = 10^{-4} . Thus, we need to use, as a rule of thumb, $\epsilon = \sqrt{\text{TOL}}$ in order to make large enough perturbations. As finite differences always mean that we lose half the digits of accuracy, we might easily end e.g. with a derivative that has only two valid digits.

A completely different approach is to formulate and solve the *variational differential equations* along with the nominal trajectory. This means that we solve, together with $\dot{x} = f(x)$, the additional matrix differential equation

$$\dot{G}(t) = \frac{\partial f}{\partial x}(x(t))G(t), \quad t \in [0, T], \quad G(0) = \mathbb{I}.$$

This is much more accurate than the first approach at a similar computational cost, but we have to get analytic expressions for $\frac{\partial f}{\partial x}(x(t))$. Also, it is interesting to note that the computed sensitivity $G(T)$ might not be 100% identical with the derivative of the (discretized) integrator result $x(T; s)$.

This last disadvantage is avoided in the third approach, *Algorithmic Differentiation (AD) of the Integrator*, where we first freeze the discretization scheme at the current nominal trajectory and then apply an AD tool to the whole integrator. This is up to machine precision 100% identical with the derivative of the numerical solution $x(T; s)$ for a given fixed discretization grid. In a practical implementation, the integrator and right hand side function $f(x)$ need to be in the same or in compatible computer languages that are treated by the corresponding AD tool (e.g. C++ when using ADOL-C). Also, if an implicit integrator is used, it should be noted that the underlying Newton iterations will be differentiated, which might create considerable and avoidable overhead compared to the variational differential equation approach.

A fourth approach, *Internal Algorithmic Differentiation (AD) of the Integrator* can be seen as a combination of the variational differential equation and AD. Here, AD is applied to each step of the integrator in a custom implementation of the integrator, but care is taken that no components of the algorithm are differentiated that need not be differentiated, such as Newton matrices. The approach is illustrated for an Euler scheme (where it is identical to both the variational differential equation and external AD). If the grid is given by $\{t_k\}_{k=0}^N$ and the Euler iterates

$$x_{k+1} = x_k + (t_{k+1} - t_k)f(x_k), \quad k = 0, \dots, N-1, \quad x_0 = s.$$

then this approach generates matrices

$$G_{k+1} = G_k + (t_{k+1} - t_k) \frac{\partial f}{\partial x}(x_k) G_k, \quad k = 0, \dots, N-1, \quad G_0 = \mathbb{I}.$$

This approach is usually the most computationally efficient of the exact differentiation approaches but requires a custom implementation of an ODE/DAE solver that is explicitly designed for the generation of sensitivities. Note that as in the previous two approaches, we cannot deal with black-box right hand side functions $f(x)$ as we need to compute their derivatives symbolically or algorithmically, though the matrix $\frac{\partial f}{\partial x}(x_k)$ could of course also be computed by finite differences.

This last idea can be generalized to the concept of *Internal Numerical Differentiation (IND)* [Boc87]. At first sight it is similar to END, but needs a custom implementation and differs in several respects. First, all trajectories are computed simultaneously; only the nominal trajectory is adaptive, while the perturbed trajectories use the nominal, frozen grid. In implicit methods, also matrix factorizations etc. will be frozen. At the end of the interval, we use the finite difference formula (9.6) but with a much smaller perturbation, namely $\epsilon = \sqrt{\text{PREC}}$ where PREC is the machine precision, typically 10^{-16} . The derivatives will have the accuracy $\sqrt{\text{PREC}}$, i.e. usually 10^{-8} , which is much higher than for END.

Again, we illustrate IND at hand of the explicit Euler integration scheme, where each perturbed trajectory with index $i = 1, \dots, n_x$ just satisfies

$$x_{k+1}^i = x_k^i + (t_{k+1} - t_k)f(x_k^i), \quad k = 0, \dots, N-1, \quad x_0^i = s + \epsilon e_i.$$

Note that due to the fact that adaptivity and possible matrix factorizations are switched off for the perturbed trajectories, IND is not only more accurate, but also cheaper than END.

9.5 A Classification of Direct Optimal Control Methods

It is an interesting exercise to try to classify Newton type optimal control algorithms, where we follow the presentation given in [DFH09]. Let us have a look at how nonlinear optimal control algorithms perform their major algorithmic components, each of which comes in several variants:

- (a) Treatment of Inequalities: Nonlinear IP vs. SQP
- (b) Nonlinear Iterations: Simultaneous vs. Sequential
- (c) Derivative Computations: Full vs. Reduced
- (d) Linear Algebra: Banded vs. Condensing

In the last two of these categories, we observe that the first variants each exploit the specific structures of the simultaneous approach, while the second variant reduces the variable space to the one of the sequential approach. Note that reduced derivatives imply condensed linear algebra, so the combination [Reduced,Banded] is excluded. In the first category, we might sometimes distinguish two variants of SQP methods, depending on how they solve their underlying QP problems, via active set QP solvers (SQP-AS) or via interior point methods (SQP-IP).

Based on these four categories, each with two alternatives, and one combination excluded, we obtain 12 possible combinations. In these categories, the classical single shooting method [SS78] could be classified as [SQP,Sequential,Reduced] or as [SQP,Sequential,Full,Condensing] because some variants compute directly the reduced derivatives, while others compute first the stagewise derivative matrices A_i and B_i and condense then. Tenny's feasibility perturbed SQP method [TWR04] could be classified as [SQP,Sequential,Full,Banded], and Bock's multiple shooting [BP84] as well as the classical reduced SQP collocation methods [THE75, Bie84, Bet01] as [SQP,Simultaneous,Full,Condensing]. The band structure exploiting SQP variants from Steinbach [Ste94] and Franke [Fra98] are classified as [SQP-IP,Simultaneous,Full,Banded], while the widely used interior point direct collocation method in conjunction with IPOPT by Biegler and Wächter [Wae02] as [IP,Simultaneous,Full,Banded]. The reduced Gauss-Newton method of Schlöder [Sch88] would here be classified as [SQP,Simultaneous,Reduced].

Chapter 10

Nonlinear Model Predictive Control

So far, we have regarded one single optimal control problem and focussed on ways to numerically solve this problem. Once we have computed such a solution, we might try to control the corresponding real process with the obtained control trajectory. This approach to use a precomputed control trajectory is called *open-loop control*. Unfortunately, the result will most probably be very dissatisfying, as the real process will typically not coincide completely with the model that we have used for optimization. If we wanted for example move a robot arm to a terminal point, the robot arm might end at a very different location than the model predicted. This is due to the difference of the model with the reality, sometimes called *model-plant-mismatch*. This mismatch might be due to modelling errors or external, unforeseen disturbances.

On the other hand, we might be able to observe the real process during its time development, and notice, for example, that the robot arm moves differently than predicted. This will allow us to correct the control inputs online in order to get a better performance; this procedure is called *feedback control* or *closed-loop control*. Feedback allows us to improve the practical performance of optimal control enormously. In its most basic form, we could use ad-hoc implementations of feedback that react to deviations from the planned state trajectory by basic control schemes such as a *proportional-integral (PI)* controller. On the other hand, we might use again optimal control techniques in order to react to disturbances of the state, by using *optimal feedback control*, which we had outlined in the Chapter 7 on dynamic programming (DP). In the case of the moving robot arm this would result in the following behaviour: if during its motion the robot arm is strongly pushed by an external disturbance, it will not try to come back to its planned trajectory but instead adapt to the new situation and follow the new optimal trajectory. This is straightforward in the case of DP or HJB, where we have the optimal feedback control precomputed for all possible states. But as said, these approaches are impossible to use for nontrivial state dimensions, i.e. systems with more than, say, 3-8 states. Thus, typically we cannot precompute the optimal feedback control in advance.

A possible remedy is to compute the optimal feedback control in *real-time*, or *online*, during the runtime of the process. In the case of the robot arm this means that after the disturbance, we would call our optimization solver again in order to quickly compute the new optimal trajectory. If we could solve this problem exactly and infinitely fast, we would get exactly the same feedback as in optimal feedback control. In reality, we have to work with approximations: first, we might simplify the optimal control problem in order to allow faster computation, e.g. by predicting only a limited amount of time into the future, and second, we might adapt our algorithms to the new task, namely that we have to solve optimization problems again and again. This task is called *real-time optimization* or *embedded optimization*, due to the fact that in many cases, the numerical optimization will be carried out on *embedded hardware*, i.e. processors that reside not in a desktop computer but e.g. in a feedback control system.

While this idea of *optimal feedback control via real-time optimization* sounds challenging or even impossible for the fast motion of robot arms, it is since decades industrial practice in the process control industry under the name of *Model Predictive Control (MPC)*. There, time scales are often in the range of minutes and allow ample time for each optimization. The main stream implementation of MPC can in discrete time roughly be formulated as follows: (1) observe the current state of the system \bar{x}_0 , (2) predict and optimize the future behaviour of the process on a limited time window of N steps by solving an open-loop optimization problem starting at the state \bar{x}_0 , (3) implement the first control action u_0^* at the real process, (4) move the optimization horizon one time step forward and repeat the procedure. MPC is sometimes also called *receding horizon control* due to this movement of the *prediction horizon*. The name *nonlinear MPC*, short *NMPC*, is reserved for the special case of MPC with underlying nonlinear dynamic systems, while linear MPC refers to MPC with linear system models. Note that

NMPC leads typically to non-convex optimization problems while nearly all linear MPC formulations use convex cost and constraints.

Note that in the case of a time-invariant system and cost, the subsequent optimization problems differ only by the initial value \bar{x}_0 and nothing else, and thus, the MPC feedback is time-invariant as well. If we would be able to solve the problem with an infinite prediction horizon, we would obtain the stationary optimal feedback control. The limitation of the horizon to a finite length N allows us to solve the problem numerically. If we choose N large enough, it will be a good approximation to the infinite horizon problem.

In this script, we do not focus on the different ways to formulate the MPC problem, but on its numerical solution by suitable real-time optimization methods. This and the next chapter follows the presentation given in [DFH09] and [Die02] and focusses on the MPC optimal control problem.

10.1 NMPC Optimization Problem

Let us in this chapter regard the following simplified optimal control problem in discrete time augmented with algebraic equations.

$$\begin{aligned} \underset{\bar{x}, z, u}{\text{minimize}} \quad & \sum_{i=0}^{N-1} L(x_i, z_i, u_i) \quad + \quad E(x_N) \end{aligned} \quad (10.1a)$$

$$\text{subject to} \quad x_0 - \bar{x}_0 = 0, \quad (10.1b)$$

$$x_{i+1} - f(x_i, z_i, u_i) = 0, \quad i = 0, \dots, N-1, \quad (10.1c)$$

$$g(x_i, z_i, u_i) = 0, \quad i = 0, \dots, N-1, \quad (10.1d)$$

$$h(x_i, z_i, u_i) \leq 0, \quad i = 0, \dots, N-1, \quad (10.1e)$$

$$r(x_N) \leq 0. \quad (10.1f)$$

Here, $x_i \in \mathbb{R}^{n_x}$ is the differential state, $z_i \in \mathbb{R}^{n_z}$ the algebraic state, and $u_i \in \mathbb{R}^{n_u}$ is the control. Functions f and g are assumed twice differentiable and map into \mathbb{R}^{n_x} and \mathbb{R}^{n_z} , respectively. The algebraic state z_i is uniquely determined by (10.1d) when x_i and u_i are fixed, as we assume that $\frac{\partial g}{\partial z}$ is invertible everywhere.

We choose to regard this difference-algebraic system form because it covers several parametrization schemes for continuous time dynamic systems in differential algebraic equation (DAE) form, in particular direct multiple shooting with DAE relaxation [LBS⁺03] and direct collocation [THE75, Bie84]. Note that in the case of collocation, all collocation equations on a collocation interval would be collected within the function g and the collocation node values in the variables z_i , see the formulation in formula (9.4).

Here, the free variables are the differential state vector $x = (x_0^T, x_1^T, \dots, x_{N-1}^T, x_N^T)^T$ at all considered time points and the algebraic and control vector on all but the last time points: $z = (z_0^T, z_1^T, \dots, z_{N-1}^T)^T$ and $u = (u_0^T, u_1^T, \dots, u_{N-1}^T)^T$.

The task in real-time optimization for NMPC is now the following: for a given value of \bar{x}_0 , we need to approximately solve the above optimization problem as fast as possible, and of the obtained solution, it is the optimal value u_0 that we need fastest in order to provide the NMPC feedback. We might call the exact solution $u_0^*(\bar{x}_0)$ in order to express its dependence on the initial value \bar{x}_0 . The only reason why we formulate and optimize the large optimization problem is because it delivers us this map $u_0^* : \mathbb{R}^{n_x} \rightarrow \mathbb{R}^{n_u}$, which is an approximation to the optimal feedback control.

Remark on fixed and free parameters: In most NMPC applications there are some *constant* parameters \bar{p} that are assumed constant for the NMPC optimization, but that change for different problems, like \bar{x}_0 . We do not regard them here for notational convenience, but note that they can be treated by state augmentation, i.e. regarded as constant system states with fixed initial value \bar{p} .

10.2 Nominal Stability of NMPC

Very often, one is interested in stabilizing the nonlinear dynamic system at a given set point for states and controls, which we might without loss of generality set to zero here. This steady state, that satisfies $f(0, 0, 0) = 0$, $g(0, 0, 0) = 0$ must be assumed to be feasible, i.e. $h(0, 0, 0) \leq 0$. One then often uses as stage cost the quadratic deviation from this set point, i.e., $L(x, u) = x^T Q x + u^T R u$ with positive definite matrices Q, R . It is important to

note that this function is positive definite, i.e., $L(0, 0) = 0$ and $L(x, u) > 0$ otherwise. In this case, one would ideally like to solve the infinite horizon problem with $N = \infty$ in order to obtain the true stationary optimal feedback control; this would automatically ensure stability, as the value function $J(x)$ can be shown to decrease along the trajectory of the nominal system in each time step by $-L(x_0, u^*(x_0))$ and can thus serve as a Lyapunov function. But as we have in practice to choose a finite N , the question arises how we can ensure nominal stability of NMPC nevertheless. One way due to [KG88, MM90] is to impose a *zero terminal constraint* i.e. to require $x_N = 0$ as terminal boundary condition (10.1f) in the NMPC problem and to employ no terminal cost, i.e. $E(x_N) = 0$.

In this case of a zero terminal constraint, it can be shown that the value function J_0 of the finite horizon problem is a Lyapunov function that decreases by at least $-L(\bar{x}_0, u^*(\bar{x}_0))$ in each time step. To prove this, let us assume that $(x_0^*, z_0^*, u_0^*, x_1^*, z_1^*, u_1^*, \dots, x_N^*)$ is the solution of the NMPC problem (10.1a)-(10.1f) starting with initial value \bar{x}_0 . After application of this feedback to the nominal system, i.e. without model-plant-mismatch, the system will evolve exactly as predicted, and for the next NMPC problem the initial value \bar{x}'_0 will be given by $\bar{x}'_0 = x_1^*$. For this problem, the *shifted* version of the previous solution $(x_1^*, z_1^*, u_1^*, \dots, x_N^*, 0, 0, 0)$ is a feasible point, and due to the zero values at the end, no additional cost arises at the end of the horizon. However, because the first stage cost term moved out of the horizon, we have that the cost of this feasible point of the next NMPC problem is reduced by exactly $-L(\bar{x}_0, u^*(\bar{x}_0))$. After further optimization, the cost can only be further reduced. Thus, we have proven that the value function J_0 is reduced along the trajectory, i.e. $J_0(\bar{x}'_0) \leq J_0(\bar{x}_0) - L(\bar{x}_0, u^*(\bar{x}_0))$. More generally, one can relax the zero terminal constraint and construct combinations of terminal cost $E(x_N)$ and terminal inequalities $r(x_N) \leq 0$ that have the same property but are less restrictive, cf. e.g. [CA98, DNMS00, MRRS00].

10.3 Online Initialization via Shift

For exploiting the fact that NMPC requires the solution of a whole sequence of neighboring NLPs and not just a number of stand-alone problems, we have first the possibility to *initialize* subsequent problems efficiently based on previous information.

A first and obvious way to transfer solution information from one solved NMPC problem to the initialization of the next one is employing the shift that we used already in the proof of nominal stability above. It is motivated by the principle of optimality of subarcs, which, in our context, states the following: Let us assume we have computed an optimal solution $(x_0^*, z_0^*, u_0^*, x_1^*, z_1^*, u_1^*, \dots, x_N^*)$ of the NMPC problem (10.1a)-(10.1f) starting with initial value \bar{x}_0 . If we regard a shortened NMPC problem without the first interval, which starts with the initial value \bar{x}_1 chosen to be x_1^* , then for this shortened problem the vector $(x_1^*, z_1^*, u_1^*, \dots, x_N^*)$ is the optimal solution.

Based on the expectation that the measured or observed true initial value for the shortened NMPC problem differs not much from x_1^* – i.e. we believe our prediction model and expect no disturbances – this “shrinking” horizon initialization is canonical, and it is used in MPC of batch or finite time processes, see e.g. [HAM98, DBS05].

However, in the case of moving horizon problems, the horizon is not only shortened by removing the first interval, but also prolonged at the end by appending a new terminal interval – i.e. the horizon is moved forward in time. In the moving horizon case, the principle of optimality is thus not strictly applicable, and we have to think about how to initialize the appended new variables z_N, u_N, x_{N+1} . Often, they are obtained by setting $u_N := u_{N-1}$ or setting u_N as the steady state control. The states z_N and x_{N+1} are then obtained by forward simulation. In the case that zero is the steady state and we had a zero terminal constraint, this would just result in zero values to be appended, as in the proof in the previous section. In any case, this transformation of the variables from one problem to the next is called “shift initialization”. It is not as canonical as the “shrinking horizon” case, because the shifted solution is not optimal for the new undisturbed problem. However, in the case of long horizon lengths N we can expect the shifted solution to be a good initial guess for the new solution. Moreover, for most NMPC schemes with stability guarantee (for an overview see e.g. [MRRS00]) there exists a canonical choice of u_N that implies feasibility (but not optimality) of the shifted solution for the new, undisturbed problem. The shift initialization is very often used e.g. in [LB90, BR91, MHAM96, Die].

A comparison of shifted vs. non-shifted initializations was performed in [BDLS99] with the result that for autonomous NMPC problems that shall regulate a system to steady state, there is usually no advantage of a shift initialization compared to the “primitive” warm start initialization that leaves the variables at the previous solution. In the extreme case of short horizon lengths, it turns out to be even advantageous NOT to shift the previous solution, as subsequent solutions are less dominated by the initial values than by the terminal conditions. On the other hand, shift initialization are a crucial prerequisite in periodic tracking applications [Die] and whenever the system or cost

function are not autonomous.

10.4 Outline of Real-Time Optimization Strategies

In NMPC we would dream to have the solution to a new optimal control problem instantly, which is impossible due to computational delays. Several ideas help us to deal with this issue.

Offline precomputations: As consecutive NMPC problems are similar, some computations can be done once and for all before the controller starts. In the extreme case, this leads to an explicit precomputation of the NMPC control law that has raised much interest in the linear MPC community [BBM02], or a solution of the Hamilton-Jacobi-Bellman Equation, both of which are prohibitive for state and parameter dimensions above ten. But also when online optimization is used, code optimization for the model routines is often essential, and it is in some cases even possible to precompute and factorize Hessians or even Jacobians in Newton type Optimization routines, in particular in the case of neighboring feedback control along reference trajectories [KEB87, BM00]. Also, pre-optimized compilable computer code can be auto-generated that is specific to the family of optimization problems, which is e.g. in convex optimization pursued in [MWB10].

Delay compensation by prediction: When we know how long our computations for solving an NMPC problem will take, it is a good idea *not* to address a problem starting at the current state but to simulate at which state the system will be when we will have solved the problem. This can be done using the NMPC system model and the open-loop control inputs that we will apply in the meantime [FA03]. This feature is used in many practical NMPC schemes with non-negligible computation time.

Division into preparation and feedback phase: A third ingredient of several NMPC algorithms is to divide the computations in each sampling time into a preparation phase and a feedback phase [DBS⁺02]. The more CPU intensive preparation phase (a) is performed with an old predicted state \bar{x}_0 before the new state estimate, say \bar{x}'_0 , is available, while the feedback phase (b) then delivers quickly an *approximate* solution to the optimization problem for \bar{x}'_0 . Often, this approximation is based on one of the tangential predictors discussed in the next chapter.

Iterating while the problem changes: A fourth important ingredient of some NMPC algorithms is the idea to work on the optimization problem while it changes, i.e., to never iterate the Newton type procedure to convergence for an NMPC problem getting older and older during the iterations, but to rather work with the most current information in each new iteration. This idea is used in [LB90, DBS⁺02, Oht04].

As a historical note, one of the first true online algorithms for nonlinear MPC was the *Newton-Type Controller of Li and Biegler* [LB89]. It is based on a sequential optimal control formulation, thus it iterates in the space of controls $u = (u_0, u_1, \dots, u_{N-1})$ only. It uses an SQP type procedure with Gauss-Newton Hessian and line search, and in each sampling time, only one SQP iteration is performed. The transition from one problem to the next uses a shift of the controls $u^{\text{new}} = (u_1, \dots, u_{N-1}, u_N^{\text{new}})$. The result of each SQP iterate is used to give an approximate feedback to the plant. As a sequential scheme without tangential predictor, it seems to have had sometimes problems with nonlinear convergence, though closed-loop stability was proven for open-loop stable processes [LB90].

In the next chapter, we will discuss several other real-time optimization algorithms in more detail that are all based on ideas from the field of parametric nonlinear optimization.

Chapter 11

Moving Horizon Estimation

In order to predict and optimize the future behaviour of a dynamic system, one needs to know the state and possibly some unknown parameters of the system. Aim of this chapter is to present methods that estimate the current state and system parameters from a series of measurements in the past. It turns out that many estimation formulations naturally lead to optimization problems that have nearly the same structure as the optimal control problems treated earlier in this course. One powerful method for online state and parameter estimation uses the measurements on a moving time window in the past, and is called moving horizon estimation. It is the main topic of this chapter, and a technology often combined with nonlinear model predictive control (NMPC), with which its optimization problems share many characteristics. However, one can also use all available information to estimate the state and parameters of a system. In the first part of this chapter we will look at an estimation problem on a fixed time horizon.

11.1 State and Parameter Estimation Problem Formulation

Throughout this chapter we regard a dynamic system of the following form

$$x_{k+1} = f_k(x_k, w_k), \quad (11.1a)$$

$$y_k = g_k(x_k, w_k) + v_k, \quad k = 0, \dots, N-1. \quad (11.1b)$$

Here, f_k describes the time varying system dynamics, g_k models the measurement process, x_k are the unknown system states, and w_k are unknown disturbances. The measurement noise is also unknown and given by v_k , while the only quantities that we know are the measurements y_k . We assume that we have some important other piece of information, namely some knowledge - or an educated guess - on the probability density functions (PDF) for the noises v_k and disturbances w_k for $k = 0, \dots, N-1$, as well as for the initial state x_0 .

For ease of notation, we sloppily denote by $P(x)$ the PDF of a random variable X at the point x , i.e. we have $P(x) \geq 0$, $\int P(x)dx = 1$, and the expectation of variable X is computed as $E\{X\} = \int xP(x)dx$. Without loss of generality, we assume the following form of PDFs:

$$P(v_k) = \exp(-\Phi_k(v_k)) \cdot \text{const}, \quad k = 0, \dots, N-1, \quad (11.2a)$$

$$P(w_k) = \exp(-\beta_k(w_k)) \cdot \text{const}, \quad k = 0, \dots, N-1, \quad \text{and} \quad (11.2b)$$

$$P(x_0) = \exp(-\alpha_0(x_0)) \cdot \text{const}, \quad (11.2c)$$

where the constants are just for normalization and will later not be of further interest. Note that any PDF can be brought into this form by taking the negative logarithm, and that a zero value of the PDF corresponds to a value $+\infty$ for the negative logarithm.

Remark: Note that if $(x_0, w_0, w_1, \dots, w_{N-1})$ would be known, they would uniquely determine all states (x_1, \dots, x_N) . The reason why we like to give a-priori PDFs for all variables $(x_0, w_0, w_1, \dots, w_{N-1})$ is that this helps us to ensure that a unique optimal solution exists for the resulting estimation problems, independent of the observability properties of the system. If additional a-priori information would be known, e.g. for some of the states (x_1, \dots, x_N) , it could be added easily to the estimation problem formulations that follow.

11.1.1 Generality of the Considered Dynamic System Class

Though the dynamic system setting in Eqs. (11.1) is a rather compact formulation, it comprises many estimation settings of practical interest. We discuss a few of them.

Systems with known inputs

If we would have a system described by $x_{k+1} = \tilde{f}(x_k, u_k, w_k)$ with known inputs u_k , we can bring it into the form (11.1) by defining

$$f_k(x_k, w_k) := \tilde{f}(x_k, u_k, w_k),$$

i.e. the dependence of the system on the known controls makes the system time variant.

Systems with measured inputs

How could we deal with a system described by $x_{k+1} = \tilde{f}(x_k, u_k, \tilde{w}_k)$ with inputs u_k that we do not know exactly, but for which we have measurements \tilde{u}_k ? If the measurement noise on the input measurements is denoted by \tilde{v}_k , we define a disturbance vector $w_k = (\tilde{w}_k, \tilde{v}_k)$ and bring the system into the form (11.1) by setting

$$f_k(x_k, w_k) := \tilde{f}(x_k, \tilde{u}_k + \tilde{v}_k, \tilde{w}_k).$$

Systems with unknown parameters

Very often we do not only want to know the system states but also some parameters that are unknown, but constant in time. If the original system state would be given by \tilde{x}_k and the original dynamics by $\tilde{x}_{k+1} = \tilde{f}(\tilde{x}_k, p)$, we can proceed as follows to bring the system into the form (11.1). First, we introduce an individual parameter value p_k for each time interval. Second, we define the augmented system state $x_k = (\tilde{x}_k, p_k)$. Third, we define the augmented dynamical system (11.1) as

$$f_k(x_k, w_k) := \begin{bmatrix} \tilde{f}(\tilde{x}_k, p_k) \\ p_k \end{bmatrix},$$

such that the second part of the system dynamics equation, $p_{k+1} = p_k$, ensures that the “parameter state” p_k remains constant over time.

11.2 The Trajectory Estimation Problem

A first question one might want to answer is the following: given the measurements $y = (y_0, \dots, y_{N-1})$, what are the most probable state and disturbance trajectories $x = (x_0, \dots, x_N)$ and $w = (w_0, \dots, w_{N-1})$? We decide to work in a Bayesian estimation framework, and our aim is to find the maximum a-posteriori (MAP) estimate that maximizes the conditional PDF $P(x, w|y)$ of the trajectory, given the measurements. Using Bayes’ formula, this PDF is given by

$$P(x, w|y) = \frac{P(x, w, y)}{P(y)} \quad (11.3a)$$

$$= \frac{P(y|x, w) \cdot P(x, w)}{P(y)} \quad (11.3b)$$

$$= P(y|x, w) \cdot P(x, w) \cdot \text{const.} \quad (11.3c)$$

Instead of maximizing the conditional PDF, one can equivalently minimize the negative logarithm of the PDF. Thus, the MAP estimate is given by

$$\arg \min_{x, w} -\log P(x, w) - \log P(y|x, w).$$

Fortunately, we can find explicit expressions for both terms. First, we note that

$$\begin{aligned} P(x, w) &= P(x_0, \dots, x_N, w_0, \dots, w_{N-1}) \\ &= \begin{cases} 0, & \text{if not } x_{k+1} = f_k(x_k, w_k) \text{ for all } k = 0, \dots, N-1, \\ P(x_0)P(w_0) \cdots P(w_{N-1}), & \text{else.} \end{cases} \end{aligned}$$

This means that

$$-\log P(x, w) = \begin{cases} \infty, & \text{if not } x_{k+1} = f_k(x_k, w_k) \text{ for all } k = 0, \dots, N-1, \\ \alpha_0(x_0) + \sum_{k=0}^{N-1} \beta_k(w_k) + \text{const}, & \text{else.} \end{cases}$$

For the other term, we use the fact that the conditional probability $P(y_k|x, w)$ to obtain a measurement y_k only depends on the state x_k and disturbance w_k at the same time point. Because of $y_k = g_k(x_k, w_k) + v_k$, it is given by $P(y_k|x_k, w_k) = P(v_k)$, with $v_k = y_k - g_k(x_k, w_k)$. Thus, the following identities hold:

$$\begin{aligned} P(y|x, w) &= P(y_0, \dots, y_{N-1}|x_0, \dots, x_N, w_0, \dots, w_{N-1}) \\ &= \prod_{k=0}^{N-1} P(y_k|x_k, w_k) \\ &= \prod_{k=0}^{N-1} P(v_k), \quad \text{with } v_k = y_k - g_k(x_k, w_k) \text{ for } k = 0, \dots, N-1, \\ &= \prod_{k=0}^{N-1} \exp(-\Phi_k(y_k - g_k(x_k, w_k))) \cdot \text{const.} \end{aligned}$$

Therefore, we obtain the compact expression

$$-\log P(y|x, w) = \sum_{k=0}^{N-1} \Phi_k(y_k - g_k(x_k, w_k)).$$

Taking both expressions together, we obtain the MAP estimate as solution of the following minimization problem, where we exclude the infinite objective values by the corresponding constraints:

$$\begin{aligned} \underset{x, w}{\text{minimize}} \quad & \alpha_0(x_0) + \sum_{k=0}^{N-1} [\Phi_k(y_k - g_k(x_k, w_k)) + \beta_k(w_k)] \end{aligned} \quad (11.4a)$$

$$\text{subject to} \quad x_{k+1} - f_k(x_k, w_k) = 0, \quad \text{for } k = 0, \dots, N-1. \quad (11.4b)$$

We will often call the term $\alpha_0(x_0)$ the “arrival cost”, as it measures the “cost” for arriving at x_0 . For notational convenience, we also define the shorthand

$$\varphi_k(x_k, w_k) := \Phi_k(y_k - g_k(x_k, w_k)) + \beta_k(w_k)$$

and call this term, as in the previous chapters, the “stage cost”. Note that the optimization problem (11.4) is of exactly the same form as the optimal control problems discussed previously in this lecture.

11.2.1 Examples for the stage and arrival costs

Very often the cost terms $\alpha_0(x_0)$, $\beta_k(w_k)$ and $\Phi_k(v_k)$ are chosen as quadratic penalties. For notational convenience we define $\|x\|_P^2 := x^\top P x$ for positive definite matrices $P \succ 0$. Note that quadratic penalties correspond to weighted ℓ_2 -norms, as $\|x\|_P^2 = \|P^{\frac{1}{2}} x\|_2^2$, where $P^{\frac{1}{2}}$ is the unique symmetric matrix square root such that $P^{\frac{1}{2}} \cdot P^{\frac{1}{2}} = P$. A typical choice for the arrival cost is $\alpha_0(x_0) = \|x_0 - \bar{x}_0\|_P^2$, where \bar{x}_0 is an a-priori guess for the initial state, and P an inverse covariance matrix expressing the confidence we have for this guess. For the disturbances, a penalty $\beta_k(w_k) = \|w_k\|_R^2$ expresses how unlikely we expect them to be. For the measurement errors, the quadratic penalty $\Phi_k(v_k) = \|v_k\|_Q^2$ is often used, where Q^{-1} is the covariance matrix we expect for the measurement errors.

Instead of quadratic penalties, that correspond to the assumption of Gaussian distributions, other choices are possible as well. Mostly, one uses convex functions, because of their beneficial properties for optimization. Two other popular convex penalty functions are the (possibly weighted) ℓ_1 -norm $\|v\|_1 = \sum_{i=1}^{n_v} |v_i|$, which corresponds to a Laplace distribution, and the Huber penalty, that is for a scalar input $v \in \mathbb{R}$ defined as

$$\Phi_{\text{Huber}, \sigma}(v) = \begin{cases} v^2 & \text{if } |v| < \sigma, \\ 2\sigma|v| - \sigma^2 & \text{else.} \end{cases}$$

The Huber penalty corresponds to a distribution that looks like a Gaussian in the neighborhood of zero, but which has “fatter tails” than a Gaussian. These fat tails can express our expectation that outliers might appear, i.e.

that we expect that large residuals have a higher probability than a normal distribution would suggest. From the penalty function perspective, both the ℓ_1 - and the Huber-penalty have the property that they penalize large error residuals less than a quadratic penalty would do. Thus, using ℓ_1 - or Huber-penalties for the measurement error functions $\Phi_k(v_k)$ allows one to design estimators that are more robust against outliers than the usual ℓ_2 -norm based estimators.

Remark on parameter jumps: An interesting other application of the ℓ_1 -norm arises in the case when we want to detect jumps in some parameter p , but we expect these jumps to occur only rarely. In addition to the usual system dynamics and measurement equation, one can then model the parameter dynamics by $p_{k+1} = p_k + w_k$ and penalize the parameter jumps with an ℓ_1 -norm, i.e. choose $\beta_k(w_k) := \|w_k\|_1$. This discourages changes in p_k , and nonzero values for w_k , i.e. changes in p_k , will only occur in the optimal solution if there is a significant benefit in terms of the other optimization objective terms.

11.3 Dynamic Programming for the Trajectory Estimation Problem

Because the trajectory estimation problem is an optimal control problem, it can also be solved by dynamic programming. In this context, it is interesting to observe that dynamic programming can in principle be performed in forward as well as in backward direction. In estimation problems, in contrast to standard optimal control problems, one usually choses to go in forward direction. The reason is that dynamic programming then allows us to “forget the past” and to just summarize the contribution of the past in one function, which we call the “arrival cost”. The arrival cost is the equivalent to the “cost-to-go” in the usual backwards dynamic programming recursion. We define the arrival cost $\alpha_n(x_n)$ for any $n \leq N$ as the cost to arrive after n steps at state x_n :

$$\alpha_n(x_n) := \min_{x_0, w_0, \dots, x_{n-1}, w_{n-1}} \alpha_0(x_0) + \sum_{k=0}^{n-1} \varphi_k(x_k, w_k) \quad \text{s.t. } x_{k+1} = f_k(x_k, w_k), \text{ for } k = 0, \dots, n-1. \quad (11.5)$$

Note that x_n is not a variable, but a fixed parameter for the optimization problem. By the dynamic programming principle, one can compute the arrival cost recursively, using the fact that the only connection between time n and $n+1$ is via the state x_{n+1} . The dynamic programming recursion proceeds as follows, for $n = 0, \dots, N-1$:

$$\alpha_{n+1}(x_{n+1}) = \min_{x_n, w_n} \alpha_n(x_n) + \varphi_n(x_n, w_n) \quad \text{s.t. } x_{n+1} = f_n(x_n, w_n). \quad (11.6)$$

Again, note that x_{n+1} is a fixed parameter to the optimization problem. To use dynamic programming to solve the trajectory estimation problem, one proceeds as follows:

1. Start with the given arrival cost $\alpha_0(\cdot)$.
2. Compute $\alpha_1(\cdot)$ up to $\alpha_N(\cdot)$, using the dynamic programming recursion (11.6)
3. Compute $x_N^* = \arg \min_{x_N} \alpha_N(x_N)$.
4. For $n = N-1, \dots, 0$, compute $(x_n^*, w_n^*) = \arg \min_{x_n, w_n} \alpha_n(x_n) + \varphi_n(x_n, w_n) \quad \text{s.t. } x_{n+1}^* = f_n(x_n, w_n)$.

Note that very often one is only interested in the estimate for the last state, x_N^* , which is already obtained after Step 3. Thus, Step 4 is optional, and only needed if one wants to know an estimate of the complete trajectory. However, if one is really only interested in the last state x_N , why should one first try to maximize the MAP $P(x, w|y)$ of the complete trajectory? In this case, one should rather maximize directly the PDF $P(x_N|y)$ of the last state, as we will do in Section 11.5. It will later turn out that both estimation formulations, the trajectory estimation and the direct estimation of the last state, lead to the same results for linear quadratic estimation problems.

11.4 Linear Quadratic Trajectory Estimation

Let us specialize the trajectory estimation problem to the special case of linear dynamic systems with quadratic costs, i.e. with underlying Gaussian PDFs for disturbances and measurement errors. In this case we deal with the

following quantities.

$$f_k(x_k, w_k) = A_k x_k + b_k + w_k, \quad (11.7a)$$

$$g_k(x_k, w_k) = C_k x_k, \quad (11.7b)$$

$$\beta_k(w_k) = \frac{1}{2} \|w_k\|_R^2, \quad (11.7c)$$

$$\Phi_k(v_k) = \frac{1}{2} \|v_k\|_Q^2, \quad \text{for } k = 0, \dots, N-1, \quad \text{and} \quad (11.7d)$$

$$\alpha_0(x_0) = \frac{1}{2} \|x_0 - \bar{x}_0\|_{P_0}^2. \quad (11.7e)$$

Note that we have chosen a formulation for the system dynamics in which the disturbances affect every state directly. This will allow us to simplify some later expressions. The optimal control problem resulting from this linear quadratic estimation setup is the following.

$$\begin{aligned} \underset{x, w}{\text{minimize}} \quad & \frac{1}{2} \|x_0 - \bar{x}_0\|_{P_0}^2 + \sum_{k=0}^{N-1} \left[\frac{1}{2} \|y_k - C_k x_k\|_Q^2 + \frac{1}{2} \|w_k\|_R^2 \right] \end{aligned} \quad (11.8a)$$

$$\text{subject to} \quad x_{k+1} - A_k x_k - b_k - w_k = 0, \quad \text{for } k = 0, \dots, N-1. \quad (11.8b)$$

One can easily eliminate all w_k using the equality constraints, and then one obtains the following unconstrained quadratic optimization problem.

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \|x_0 - \bar{x}_0\|_{P_0}^2 + \sum_{k=0}^{N-1} \left[\frac{1}{2} \|y_k - C_k x_k\|_Q^2 + \frac{1}{2} \|x_{k+1} - A_k x_k - b_k\|_R^2 \right] \quad (11.9a)$$

To solve it, one might just differentiate the objective function with respect to x and set the gradient to zero, which results in a sparse linear equation system for the optimal state trajectory x^* . On the other hand, one could also use dynamic programming to solve it. To formulate the dynamic programming recursion, we first state a useful lemma and corollary.

Lemma 2 (Schur Complement Lemma) *If $R \succ 0$, the following identity holds*

$$\begin{bmatrix} x \\ u \end{bmatrix}^\top \begin{bmatrix} Q & S^\top \\ S & R \end{bmatrix} \begin{bmatrix} x \\ u \end{bmatrix} = x^\top (Q - S^\top R^{-1} S) x + \|R^{-1} S x + u\|_R^2. \quad (11.10)$$

In particular,

$$\min_u \begin{bmatrix} x \\ u \end{bmatrix}^\top \begin{bmatrix} Q & S^\top \\ S & R \end{bmatrix} \begin{bmatrix} x \\ u \end{bmatrix} = x^\top (Q - S^\top R^{-1} S) x.$$

If in addition $\begin{bmatrix} Q & S^\top \\ S & R \end{bmatrix} \succ 0$, then also $Q - S^\top R^{-1} S \succ 0$.

The proof of the lemma uses the matrix decomposition

$$\begin{bmatrix} Q & S^\top \\ S & R \end{bmatrix} = \begin{bmatrix} Q - S^\top R^{-1} S & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} S^\top R^{-1} S & S^\top \\ S & R \end{bmatrix}$$

and the fact that the second term can be expressed as

$$\begin{bmatrix} x \\ u \end{bmatrix}^\top \begin{bmatrix} S^\top R^{-1} S & S^\top \\ S & R \end{bmatrix} \begin{bmatrix} x \\ u \end{bmatrix} = x^\top S^\top R^{-1} S x + 2u^\top S x + u^\top R u = \|R^{-1} S x + u\|_R^2.$$

From this we also obtain the following corollary.

Corollary 3 (Summarizing Linear Quadratic Costs) *If $R \succ 0$ then*

$$\min_u \begin{bmatrix} 1 \\ x \\ u \end{bmatrix}^\top \begin{bmatrix} c & q^\top & s^\top \\ q & Q & S^\top \\ s & S & R \end{bmatrix} \begin{bmatrix} 1 \\ x \\ u \end{bmatrix} = c - s^\top R^{-1} s + 2x^\top (q - S^\top R^{-1} s) + x^\top (Q - S^\top R^{-1} S) x.$$

The proof of the corollary uses the previous lemma with $\tilde{S} = [s|S]$, $\tilde{Q} = \begin{bmatrix} c & q^\top \\ q & Q \end{bmatrix}$ and $\tilde{x} = \begin{bmatrix} 1 \\ x \end{bmatrix}$, and the fact that

$$\tilde{Q} - \tilde{S}^\top R^{-1} \tilde{S} = \begin{bmatrix} c & q^\top \\ q & Q \end{bmatrix} - \begin{bmatrix} s^\top R^{-1} s & s^\top R^{-1} S \\ S^\top R^{-1} s & S^\top R^{-1} S \end{bmatrix}.$$

To formulate the dynamic programming recursion, we assume that $\alpha_k(x_k) = \frac{1}{2} \|x_k - \bar{x}_k\|_{P_k}^2$ and eliminate w_k , which results in the following formula.

$$\alpha_{k+1}(x_{k+1}) = \min_{x_k} \frac{1}{2} \|x_k - \bar{x}_k\|_{P_k}^2 + \frac{1}{2} \|y_k - C_k x_k\|_Q^2 + \frac{1}{2} \|x_{k+1} - A_k x_k - b_k\|_R^2$$

Using the above corollary, we know that the solution is a quadratic function. We use the identity

$$\begin{aligned} & \|x_k - \bar{x}_k\|_{P_k}^2 + \|y_k - C_k x_k\|_Q^2 + \|x_{k+1} - A_k x_k - b_k\|_R^2 \\ &= \text{const} + \begin{bmatrix} 1 \\ x_{k+1} \\ x_k \end{bmatrix}^\top \begin{bmatrix} \text{const} & (-Rb_k)^\top & (P_k \bar{x}_k - C_k^\top Q y_k + A_k^\top R b_k)^\top \\ (-Rb_k) & R & (-A_k^\top R)^\top \\ (P_k \bar{x}_k - C_k^\top Q y_k + A_k^\top R b_k) & (-A_k^\top R) & (P_k + C_k^\top Q C_k + A_k^\top R A_k) \end{bmatrix} \begin{bmatrix} 1 \\ x_{k+1} \\ x_k \end{bmatrix} \end{aligned}$$

Based on the corollary, with $\tilde{u} = x_k$, $\tilde{R} := P_k + C_k^\top Q C_k + A_k^\top R A_k$ and $\tilde{s} := (P_k \bar{x}_k - C_k^\top Q y_k + A_k^\top R b_k)$, the quadratic function is explicitly given by

$$\alpha_{k+1}(x_{k+1}) = \text{const} + \frac{1}{2} x_{k+1}^\top \left(R - (A_k^\top R)^\top \tilde{R}^{-1} A_k^\top R \right) x_{k+1} + x_{k+1}^\top \left(-Rb_k + (A_k^\top R)^\top \tilde{R}^{-1} \tilde{s} \right).$$

We define the matrix $P_{k+1} := \left(R - (A_k^\top R)^\top \tilde{R}^{-1} A_k^\top R \right)$, which is positive definite due to the fact that the original quadratic function was positive definite in (x_k, x_{k+1}) . To bring $\alpha_{k+1}(x_{k+1})$ into a more compact form, we define $\bar{x}_{k+1} = -P_{k+1}^{-1} \left(-Rb_k + (A_k^\top R)^\top \tilde{R}^{-1} \tilde{s} \right)$. We can then show that

$$\alpha_{k+1}(x_{k+1}) = \frac{1}{2} \|x_{k+1} - \bar{x}_{k+1}\|_{P_{k+1}}^2 + \text{const}$$

as an immediate consequence of the following basic lemma.

Lemma 3 If $P \succ 0$ and $\bar{x} = -P^{-1}g$ then $\frac{1}{2}x^\top P x + g^\top x = \|x - \bar{x}\|_P^2 + \text{const}$.

11.5 Recursive Bayesian Estimation of the Last State

Very often, one is only interested in estimating the last state x_N , not in the whole trajectory. For this aim, a technique that is very similar to dynamic programming can be used that is called *Recursive Bayesian Estimation*. The idea is to recursively compute the conditional PDF of the state x_{n+1} given all measurements y_0, \dots, y_n . We note that the only memory of the system is the state x_n , and that the latest measurement y_n helps us to learn more about the PDF of x_n . For these reason, one can derive the following identity.

$$\begin{aligned} P(x_{n+1}|y_0, \dots, y_n) &= \int P(x_{n+1}|x_n) P(x_n|y_0, \dots, y_n) dx_n \\ &= \int P(x_{n+1}|x_n, w_n) P(x_n, w_n|y_0, \dots, y_n) dx_n dw_n \\ &= \int_{f_n(x_n, w_n)=x_{n+1}} P(x_n, w_n|y_0, \dots, y_n) dx_n dw_n \\ &= \int_{f_n(x_n, w_n)=x_{n+1}} \frac{P(x_n, w_n|y_0, \dots, y_{n-1}) P(y_n|x_n, w_n)}{P(y_n|y_0, \dots, y_{n-1})} dx_n dw_n \\ &= \text{const} \cdot \int_{f_n(x_n, w_n)=x_{n+1}} P(x_n, w_n|y_0, \dots, y_{n-1}) P(y_n|x_n, w_n) dx_n dw_n \\ &= \text{const} \cdot \int_{f_n(x_n, w_n)=x_{n+1}} P(w_n) P(x_n|y_0, \dots, y_{n-1}) P(y_n|x_n, w_n) dx_n dw_n \\ &= \text{const} \cdot \int_{f_n(x_n, w_n)=x_{n+1}} e^{-\beta_n(w_n)} P(x_n|y_0, \dots, y_{n-1}) e^{-\Phi(y_n - g_n(x_n, w_n))} dx_n dw_n \end{aligned}$$

The result is a recursive formula to compute $P(x_{n+1}|y_0, \dots, y_n)$ from the last measurement y_n and from $P(x_n|y_0, \dots, y_{n-1})$.

Bibliography

- [AD10] J. Albersmeyer and M. Diehl. The Lifted Newton Method and its Application in Optimization. *SIAM Journal on Optimization*, 20(3):1655–1684, 2010.
- [BBM02] A. Bemporad, F. Borrelli, and M. Morari. Model Predictive Control Based on Linear Programming - The Explicit Solution. *IEEE Transactions on Automatic Control*, 47(12):1974–1985, 2002.
- [BBM03] A. Bemporad, F. Borrelli, and M. Morari. Min-max Control of Constrained Uncertain Discrete-Time Linear Systems. *IEEE Transactions on Automatic Control*, 2003. in press.
- [BD06] J. Björnberg and M. Diehl. Approximate robust dynamic programming and robustly stable MPC. *Automatica*, 42(5):777–782, May 2006.
- [BD08] J. Björnberg and M. Diehl. Approximate Dynamic Programming for Generation of Robustly Stable Feedback Controllers. In *Modeling, Simulation and Optimization of Complex Processes. International Conference on High Performance Scientific Computing*, pages 69–86, Heidelberg, 2008. Springer.
- [BDLS99] H.G. Bock, M. Diehl, D.B. Leineweber, and J.P. Schlöder. Efficient direct multiple shooting in nonlinear model predictive control. In F. Keil, W. Mackens, H. Voß, and J. Werther, editors, *Scientific Computing in Chemical Engineering II*, volume 2, pages 218–227, Berlin, 1999. Springer.
- [Bel57] R. Bellman. *Dynamic programming*. Princeton University Press, 1957.
- [Ber05] D. Bertsekas. *Dynamic Programming and Optimal Control*, volume 1. Athena Scientific, 3rd edition, 2005.
- [Bet01] J.T. Betts. *Practical Methods for Optimal Control Using Nonlinear Programming*. SIAM, Philadelphia, 2001.
- [Bie84] L.T. Biegler. Solution of dynamic optimization problems by successive quadratic programming and orthogonal collocation. *Computers and Chemical Engineering*, 8:243–248, 1984.
- [BM00] C. Büskens and H. Maurer. SQP-methods for solving optimal control problems with control and state constraints: adjoint variables, sensitivity analysis and real-time control. *Journal of Computational and Applied Mathematics*, 120:85–108, 2000.
- [Boc87] H.G. Bock. *Randwertproblemmethoden zur Parameteridentifizierung in Systemen nichtlinearer Differentialgleichungen*, volume 183 of *Bonner Mathematische Schriften*. Universität Bonn, Bonn, 1987.
- [BP84] H.G. Bock and K.J. Plitt. A multiple shooting algorithm for direct solution of optimal control problems. In *Proceedings 9th IFAC World Congress Budapest*, pages 242–247. Pergamon Press, 1984.
- [BR91] L.T. Biegler and J.B. Rawlings. Optimization approaches to nonlinear model predictive control. In W.H. Ray and Y. Arkun, editors, *Proc. 4th International Conference on Chemical Process Control - CPC IV*, pages 543–571. AIChE, CACHE, 1991.
- [BT96] D.P. Bertsekas and J.N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, Belmont, MA, 1996.

- [BTN01] A. Ben-Tal and A. Nemirovski. *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*, volume 3 of *MPS/SIAM Series on Optimization*. SIAM, 2001.
- [BV04] S. Boyd and L. Vandenberghe. *Convex Optimization*. University Press, Cambridge, 2004.
- [CA98] H. Chen and F. Allgöwer. A quasi-infinite horizon nonlinear model predictive control scheme with guaranteed stability. *Automatica*, 34(10):1205–1218, 1998.
- [Dan63] G. B. Dantzig. *Linear Programming and Extensions*. Princeton University Press, 1963.
- [DB04] M. Diehl and J. Björnberg. Robust Dynamic Programming for Min-Max Model Predictive Control of Constrained Uncertain Systems. *IEEE Transactions on Automatic Control*, 49(12):2253–2257, December 2004.
- [DBDW06] Moritz Diehl, Hans Georg Bock, Holger Diedam, and Pierre-Brice Wieber. *Fast Motions in Biomechanics and Robotics*, volume 340, chapter Fast Direct Multiple Shooting Algorithms for Optimal Robot Control, pages 65–93. Springer, 2006.
- [DBS⁺02] M. Diehl, H.G. Bock, J.P. Schlöder, R. Findeisen, Z. Nagy, and F. Allgöwer. Real-time optimization and Nonlinear Model Predictive Control of Processes governed by differential-algebraic equations. *Journal of Process Control*, 12(4):577–585, 2002.
- [DBS05] M. Diehl, H.G. Bock, and J.P. Schlöder. A real-time iteration scheme for nonlinear optimization in optimal feedback control. *SIAM Journal on Control and Optimization*, 43(5):1714–1736, 2005.
- [DFH09] M. Diehl, H. J. Ferreau, and N. Haverbeke. *Nonlinear model predictive control*, volume 384 of *Lecture Notes in Control and Information Sciences*, chapter Efficient Numerical Methods for Nonlinear MPC and Moving Horizon Estimation, pages 391–417. Springer, 2009.
- [Die] NOW NAMED TO Diehl2004f.
- [Die02] M. Diehl. *Real-Time Optimization for Large Scale Nonlinear Processes*, volume 920 of *Fortschr.-Ber. VDI Reihe 8, Meß-, Steuerungs- und Regelungstechnik*. VDI Verlag, Düsseldorf, 2002. Download also at: <http://www.ub.uni-heidelberg.de/archiv/1659/>.
- [DNMS00] G. De Nicolao, L. Magni, and R. Scattolini. Stability and Robustness of Nonlinear Receding Horizon Control. In F. Allgöwer and A. Zheng, editors, *Nonlinear Predictive Control*, volume 26 of *Progress in Systems Theory*, pages 3–23, Basel Boston Berlin, 2000. Birkhäuser.
- [FA03] R. Findeisen and F. Allgöwer. Computational Delay in Nonlinear Model Predictive Control. Proc. Int. Symp. Adv. Control of Chemical Processes, ADCHEM, 2003.
- [FBD08] H. J. Ferreau, H. G. Bock, and M. Diehl. An online active set strategy to overcome the limitations of explicit MPC. *International Journal of Robust and Nonlinear Control*, 18(8):816–830, 2008.
- [Fra98] R. Franke. *Integrierte dynamische Modellierung und Optimierung von Systemen mit saisonaler Wärmespeicherung*. PhD thesis, Technische Universität Ilmenau, Germany, 1998.
- [GMS97] P.E. Gill, W. Murray, and M.A. Saunders. SNOPT: An SQP algorithm for large-scale constrained optimization. Technical report, Numerical Analysis Report 97-2, Department of Mathematics, University of California, San Diego, La Jolla, CA, 1997.
- [GT82] A. Griewank and Ph.L. Toint. Partitioned variable metric updates for large structured optimization problems. *Numerische Mathematik*, 39:119–137, 1982.
- [GW08] A. Griewank and A. Walther. *Evaluating Derivatives*. SIAM, 2 edition, 2008.
- [HAM98] A. Helbig, O. Abel, and W. Marquardt. Model Predictive Control for On-line Optimization of Semi-batch Reactors. In *Proc. Amer. Contr. Conf.*, pages 1695–1699, Philadelphia, 1998.
- [HFD11] B. Houska, H.J. Ferreau, and M. Diehl. ACADO Toolkit – An Open Source Framework for Automatic Control and Dynamic Optimization. *Optimal Control Applications and Methods*, 32(3):298–312, 2011.

- [HR71] G.A. Hicks and W.H. Ray. Approximation methods for optimal control systems. *Can. J. Chem. Engng.*, 49:522–528, 1971.
- [JM10] C.N. Jones and M. Morari. Polytopic approximation of explicit model predictive controllers. *IEEE Transactions on Automatic Control*, 55(11):2542–2553, 2010.
- [Kar39] W. Karush. Minima of Functions of Several Variables with Inequalities as Side Conditions. Master’s thesis, Department of Mathematics, University of Chicago, 1939.
- [KEB87] P. Krämer-Eis and H.G. Bock. Numerical Treatment of State and Control Constraints in the Computation of Feedback Laws for Nonlinear Control Problems. In P. Deufhard et al., editor, *Large Scale Scientific Computing*, pages 287–306. Birkhäuser, Basel Boston Berlin, 1987.
- [KG88] S.S. Keerthi and E.G. Gilbert. Optimal infinite-horizon feedback laws for a general class of constrained discrete-time systems: Stability and moving-horizon approximations. *Journal of Optimization Theory and Applications*, 57(2):265–293, 1988.
- [KT51] H.W. Kuhn and A.W. Tucker. Nonlinear programming. In J. Neyman, editor, *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*, Berkeley, 1951. University of California Press.
- [LB89] W.C. Li and L.T. Biegler. Multistep, Newton-Type Control Strategies for Constrained Nonlinear Processes. *Chem. Eng. Res. Des.*, 67:562–577, 1989.
- [LB90] W.C. Li and L.T. Biegler. Newton-Type Controllers for Constrained Nonlinear Processes with Uncertainty. *Industrial and Engineering Chemistry Research*, 29:1647–1657, 1990.
- [LBS⁺03] D.B. Leineweber, I. Bauer, A.A.S. Schäfer, H.G. Bock, and J.P. Schlöder. An Efficient Multiple Shooting Based Reduced SQP Strategy for Large-Scale Dynamic Process Optimization (Parts I and II). *Computers and Chemical Engineering*, 27:157–174, 2003.
- [LSBS03] D.B. Leineweber, A.A.S. Schäfer, H.G. Bock, and J.P. Schlöder. An Efficient Multiple Shooting Based Reduced SQP Strategy for Large-Scale Dynamic Process Optimization. Part II: Software Aspects and Applications. *Computers and Chemical Engineering*, 27:167–174, 2003.
- [MHAM96] A. M’hamdi, A. Helbig, O. Abel, and W. Marquardt. Newton-type Receding Horizon Control and State Estimation. In *Proc. 13rd IFAC World Congress*, pages 121–126, San Francisco, 1996.
- [MM90] D. Q. Mayne and H. Michalska. Receding horizon control of nonlinear systems. *IEEE Transactions on Automatic Control*, 35(7):814–824, 1990.
- [MRRS00] D.Q. Mayne, J.B. Rawlings, C.V. Rao, and P.O.M. Scokaert. Constrained model predictive control: stability and optimality. *Automatica*, 26(6):789–814, 2000.
- [MWB10] J. Mattingley, Y. Wang, and Stephen Boyd. Code generation for receding horizon control. In *Proceedings of the IEEE International Symposium on Computer-Aided Control System Design*, Yokohama, Japan, 2010.
- [NW06] J. Nocedal and S.J. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer, 2 edition, 2006.
- [Oht04] T. Ohtsuka. A Continuation/GMRES Method for Fast Computation of Nonlinear Receding Horizon Control. *Automatica*, 40(4):563–574, 2004.
- [PDB⁺00] E. N. Pistikopoulos, V. Dua, N. A. Bozinis, A. Bemporad, and M. Morari. On-line optimization via off-line parametric optimization tools. *Computers and Chemical Engineering*, 24:183–188, 2000.
- [Sch88] J.P. Schlöder. *Numerische Methoden zur Behandlung hochdimensionaler Aufgaben der Parameteridentifizierung*, volume 187 of *Bonner Mathematische Schriften*. Universität Bonn, Bonn, 1988.
- [Sch05] A.A.S. Schäfer. *Efficient reduced Newton-type methods for solution of large-scale structured optimization problems with application to biological and chemical processes*. PhD thesis, Universität Heidelberg, 2005.

- [SS78] R.W.H. Sargent and G.R. Sullivan. The development of an efficient optimal control package. In J. Stoer, editor, *Proceedings of the 8th IFIP Conference on Optimization Techniques (1977), Part 2*, Heidelberg, 1978. Springer.
- [Ste94] M.C. Steinbach. A structured interior point SQP method for nonlinear optimal control problems. In R. Bulirsch and D. Kraft, editors, *Computation Optimal Control*, pages 213–222, Basel Boston Berlin, 1994. Birkhäuser.
- [THE75] T.H. Tsang, D.M. Himmelblau, and T.F. Edgar. Optimal control via collocation and non-linear programming. *International Journal on Control*, 21:763–768, 1975.
- [TWR04] M.J. Tenny, S.J. Wright, and J.B. Rawlings. Nonlinear model predictive control via feasibility-perturbed sequential quadratic programming. *Computational Optimization and Applications*, 28(1):87–121, April 2004.
- [Wae02] A. Waechter. *An Interior Point Algorithm for Large-Scale Nonlinear Optimization with Applications in Process Engineering*. PhD thesis, Carnegie Mellon University, 2002.
- [WB06] A. Wächter and L.T. Biegler. On the Implementation of a Primal-Dual Interior Point Filter Line Search Algorithm for Large-Scale Nonlinear Programming. *Mathematical Programming*, 106(1):25–57, 2006.
- [WB09] A. Wächter and L. Biegler. IPOPT - an Interior Point OPTimizer. <https://projects.coin-or.org/Ipopt>, 2009.