

AA203: Optimal and Learning-based Control

Combined Course Notes

James Harrison*

April 23, 2019

Introduction

These notes accompany the newly revised (Spring 2019) version of AA203 at Stanford. The goal of this new course is to present a unified treatment of optimal control and reinforcement learning (RL), with an emphasis on model-based reinforcement learning. The goal of the instructors is to unify the subjects as much as possible, and to concretize connections between these research communities.

How is this course different from a standard class on Optimal Control? First, we will emphasize practical computational tools for real world optimal control problems, such as model predictive control and sequential convex programming. Beyond this, the last third of the course focuses on the case in which an exact model of the system is not available. We will discuss this setting both in the online context (typically referred to as adaptive optimal control) and in the episodic context (the typical setting for reinforcement learning).

How is this course different from a standard class on Reinforcement Learning? Many classes on reinforcement learning focus primarily on the setting of discrete Markov Decision Processes (MDPs), whereas we will focus primarily on continuous MDPs. More importantly, the focus on discrete MDPs leads planning with a known model (which is typically referred to as “planning” or “control” in RL) to be relatively simple. In this course, we will spend considerably more time focusing on planning with a known model in both continuous and discrete time. Finally, the focus of this course will primarily be on model-based methods. We will touch briefly on model-free methods at the end, and combinations of model-free and model-based approaches.

*Contact: jharrison@stanford.edu

A Note on Notation

The notation and language used in the control theory and reinforcement learning communities vary substantially, as so we will state all of the notational choices we make in this section. First, optimal control problems are typically stated in terms of minimizing a cost function, whereas reinforcement learning problems aim to maximize a reward. These are mathematically identical statements, where one is simply the negation of the other. Herein, we will use the control theoretic approach of cost minimization. We write c for the cost function, f for the system dynamics, and denote the state and action at time t as \mathbf{x}_t and \mathbf{u}_t respectively. We write scalars as lower case letters, vectors as bold lower case letters, and matrices as upper case letters. We write a deterministic policy as $\pi(\mathbf{x})$, and a stochastic policy as $\pi(\mathbf{u} \mid \mathbf{x})$. We write the cost-to-go (negation of the value function) associated with policy π at time t and state \mathbf{x} as $J_t^\pi(\mathbf{x})$. We will also sometimes refer to the cost-to-go as the value, but in these notes we are always referring to the expected sum of future costs. For an in-depth discussion of the notational and language differences between the artificial intelligence and control theory communities, we refer the reader to [Pow12].

For notational convenience, we will write the Hessian of a function $f(x)$, evaluated at x^* , as $\nabla^2 f(x^*)$.

Prerequisites

While these notes aim to be almost entirely self contained, familiarity with undergraduate level calculus, differential equations, and linear algebra (equivalent to CME 102 and EE 263 at Stanford) are assumed. We will briefly review nonlinear optimization in the first section of these notes, but previous experience with optimization (e.g. EE 364A) will be helpful. Finally, previous experience with machine learning (at the level of CS 229) is beneficial.

Contents

1	Nonlinear Optimization	3
1.1	Unconstrained Nonlinear Optimization	3
1.1.1	Necessary Conditions for Optimality	4
1.1.2	Sufficient Conditions for Optimality	5
1.1.3	Special case: Convex Optimization	5
1.1.4	Computational Methods	6
1.2	Constrained Nonlinear Optimization	8
1.2.1	Equality Constrained Optimization	8
1.2.2	Inequality Constrained Optimization	10
1.3	Further Reading	11
2	Dynamic Programming and the Linear Quadratic Regulator	12
2.1	The Optimal Control Problem	12
2.2	Dynamic Programming and the Principle of Optimality	13

2.3	Discrete LQR	14
2.3.1	LQR with (Bi)linear Cost and Affine Dynamics	16
2.3.2	LQR Tracking around a Linear Trajectory	17
2.3.3	LQR Tracking around a Nonlinear Trajectory	18
2.4	Iterative LQR and Differential Dynamic Programming	19
2.4.1	Iterative LQR	19
2.4.2	Differential Dynamic Programming	20
2.4.3	Algorithmic Details for iLQR and DDP	22
2.5	Stochastic Optimal Control	22
2.5.1	The Stochastic Optimal Control Problem	23
2.5.2	The Principle of Optimality: Stochastic Case	23
2.6	Stochastic LQR and LQG	24
2.6.1	LQR with Additive Noise	24
2.6.2	Problems with Imperfect State Information	24
2.6.3	LQG and the Separation Principle	25
2.7	Further Reading	26
3	The HJB and HJI Equations	26
3.1	The Principle of Optimality in Continuous Time	26
3.1.1	Hamilton-Jacobi-Bellman	26
3.1.2	Continuous-Time LQR	28
3.2	Differential Games	29
3.2.1	Differential Games and Information Patterns	29
3.2.2	Hamilton-Jacobi-Isaacs	30
3.2.3	Reachability	31
3.3	Further Reading	32

1 Nonlinear Optimization

In this section we discuss the generic nonlinear optimization problem that forms the basis for the rest of the material presented in this class. We write the minimization problem as

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$$

where f is the cost function, usually assumed twice continuously differentiable, $x \in \mathbb{R}^n$ is the optimization variable, and $\mathcal{X} \subset \mathbb{R}^n$ is the constraint set. The special case in which the cost function is linear and the constraint set is specified by linear equations and/or inequalities is *linear optimization*, which we will not discuss.

1.1 Unconstrained Nonlinear Optimization

We will first address the unconstrained case, in which $\mathcal{X} = \mathbb{R}^n$. A vector \mathbf{x}^* is said to be an unconstrained *local minimum* if there exists $\epsilon > 0$ such that $f(\mathbf{x}^*) \leq f(\mathbf{x})$ for all $\mathbf{x} \in \{\mathbf{x} \mid$

$\|\mathbf{x} - \mathbf{x}^*\| \leq \epsilon\}$, and \mathbf{x}^* is said to be an unconstrained *global minimum* if $f(\mathbf{x}^*) \leq f(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^n$.

1.1.1 Necessary Conditions for Optimality

For a differentiable cost function, we can compare the cost of a point to its neighbors by considering a small variation $\Delta\mathbf{x}$ from \mathbf{x}^* . By using Taylor expansions, this yields a first order cost variation

$$f(\mathbf{x}^* + \Delta\mathbf{x}) - f(\mathbf{x}^*) \approx \nabla f(\mathbf{x}^*)^T \Delta\mathbf{x} \quad (1)$$

and a second order cost variation

$$f(\mathbf{x}^* + \Delta\mathbf{x}) - f(\mathbf{x}^*) \approx \nabla f(\mathbf{x}^*)^T \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^T \nabla^2 f(\mathbf{x}^*) \Delta\mathbf{x}. \quad (2)$$

Setting $\Delta\mathbf{x}$ to be equal to positive and negative multiples of the unit coordinate vector, we have

$$\frac{\partial f(\mathbf{x}^*)}{\partial x_i} \geq 0 \quad (3)$$

where x_i denotes the i 'th coordinate of \mathbf{x} , and

$$\frac{\partial f(\mathbf{x}^*)}{\partial x_i} \leq 0 \quad (4)$$

for all i , which is only satisfied by $\nabla f(\mathbf{x}^*) = 0$. This is referred to as the *first order necessary condition for optimality*. Looking at the second order variation, and noting that $\nabla f(\mathbf{x}^*) = 0$, we expect

$$f(\mathbf{x}^* + \Delta\mathbf{x}) - f(\mathbf{x}^*) \geq 0 \quad (5)$$

and thus

$$\Delta\mathbf{x}^T \nabla^2 f(\mathbf{x}^*) \Delta\mathbf{x} \geq 0 \quad (6)$$

which implies $\nabla^2 f(\mathbf{x}^*)$ is positive semidefinite. This is referred to as the *second order necessary condition for optimality*. Stating these conditions formally,

Theorem 1.1 (Necessary Conditions for Optimality (NOC)). *Let \mathbf{x}^* be an unconstrained local minimum of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $f \in C^1$ in an open set S containing \mathbf{x}^* . Then,*

$$\nabla f(\mathbf{x}^*) = 0. \quad (7)$$

If $f \in C^2$ within S , $\nabla^2 f(\mathbf{x}^)$ is positive semidefinite.*

Proof. See section 1.1 of [Ber16].

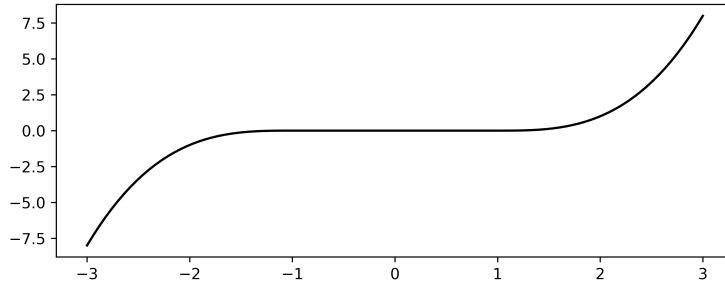


Figure 1: An example of a function for which the necessary conditions of optimality are satisfied but the sufficient conditions are not.

1.1.2 Sufficient Conditions for Optimality

If we strengthen the second order condition to $\nabla^2 f(\mathbf{x}^*)$ being positive definite, we have the sufficient conditions for \mathbf{x}^* being a local minimum. Why is the second order necessary conditions not sufficient? An example function is given in figure 1. Formally,

Theorem 1.2 (Sufficient Conditions for Optimality (SOC)). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be C^2 in an open set S . Suppose a vector \mathbf{x}^* satisfies the conditions $\nabla f(\mathbf{x}^*) = 0$ and $\nabla^2 f(\mathbf{x}^*)$ is positive definite. Then \mathbf{x}^* is a strict unconstrained local minimum of f .*

Proof is again given in Section 1.1 of [Ber16]. There are several reasons why the optimality conditions are important. In a general nonlinear optimization setting, they can be used to filter candidates for global minima. They can be used for sensitivity analysis, in which the sensitivity of \mathbf{x}^* to model parameters can be quantified [Ber16]. This is common in e.g. microeconomics. Finally, these conditions often provide the basis for the design and analysis of optimization algorithms.

1.1.3 Special case: Convex Optimization

A special case within nonlinear optimization is the set of *convex optimization* problems. A set $S \subset \mathbb{R}^n$ is called *convex* if

$$\alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \in S, \quad \forall \mathbf{x}, \mathbf{y} \in S, \forall \alpha \in [0, 1]. \quad (8)$$

For S convex, a function $f : S \rightarrow \mathbb{R}$ is called convex if

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}). \quad (9)$$

This class of problems has several important characteristics. If f is convex, then

- A local minimum of f over S is also a global minimum over S . If in addition f is strictly convex (the inequality in (9) is strict), there exists at most one global minimum of f .

- If $f \in C^1$ and convex, and the set S is open, $\nabla f(\mathbf{x}^*) = 0$ is a necessary and sufficient condition for a vector $\mathbf{x}^* \in S$ to be a global minimum over S .

Convex optimization problems have several nice properties that make them (usually) computationally efficient to solve, and the first property above gives a certificate of having obtained global optimality that is difficult or impossible to obtain in the general nonlinear optimization setting. For a thorough treatment of convex optimization theory and algorithms, see [BV04].

1.1.4 Computational Methods

In this subsection we will discuss the class of algorithms known as *gradient methods* for finding local minima in nonlinear optimization problems. These approaches, rely (roughly) on following the gradient of the function “downhill”, toward the minima. More concretely, these algorithms rely on taking steps of the form

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k \quad (10)$$

where if $\nabla f(\mathbf{x}) \neq 0$, \mathbf{d}^k is chosen so that

$$\nabla f(\mathbf{x})^T \mathbf{d}^k < 0 \quad (11)$$

and $\alpha > 0$. Typically, the step size α^k is chosen such that

$$f(\mathbf{x}^k + \alpha^k \mathbf{d}^k) < f(\mathbf{x}^k), \quad (12)$$

but generally, the step size and the direction of descent (\mathbf{d}^k) are tuning parameters.

We will look at the general class of descent directions of the form

$$\mathbf{d}^k = -D^k \nabla f(\mathbf{x}^k) \quad (13)$$

where $D^k > 0$ (note that this guarantees $\nabla f(\mathbf{x}^k)^T \mathbf{d}^k < 0$).

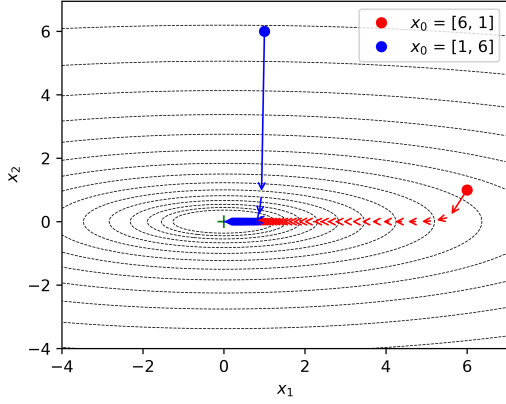
Steepest descent, $D^k = I$. The simplest choice of descent direction is directly following the gradient, and ignoring second order function information. In practice, this often leads to slow convergence (figure 2a) and possible oscillation (figure 2b).

Newton’s Method, $D^k = (\nabla^2 f(\mathbf{x}^k))^{-1}$. The underlying idea of this approach is to at each iteration, minimize the quadratic approximation of f around \mathbf{x}^k ,

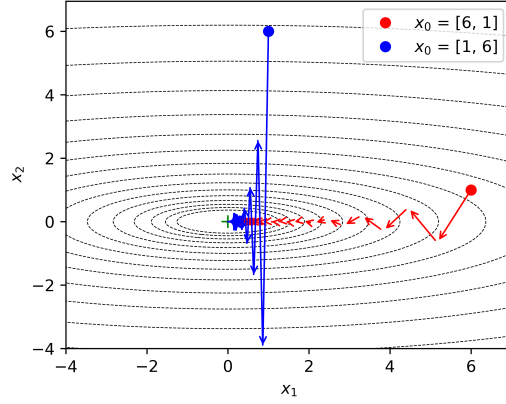
$$f^k(\mathbf{x}) = f(\mathbf{x}^k) + \nabla f(\mathbf{x}^k)^T (\mathbf{x} - \mathbf{x}^k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^k)^T \nabla^2 f(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k). \quad (14)$$

Setting the derivative of this to zero, we obtain

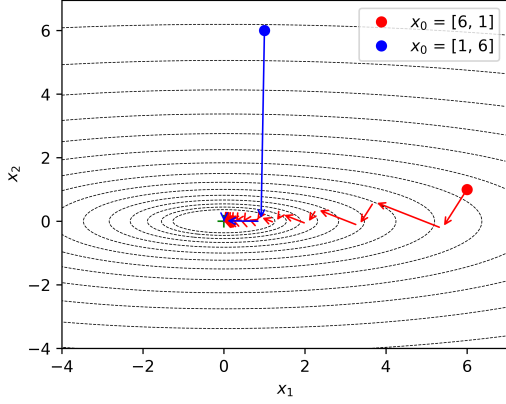
$$\nabla f(\mathbf{x}^k) + \nabla^2 f(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k) = 0 \quad (15)$$



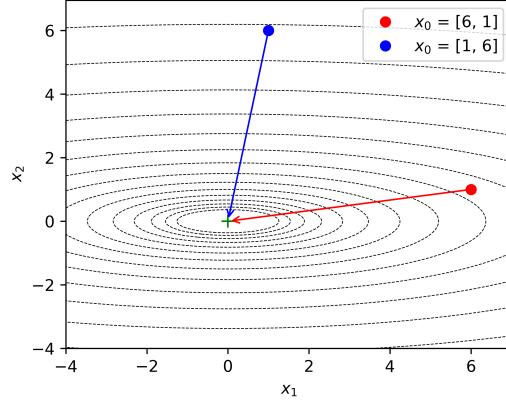
(a) Steepest descent, small fixed step size.



(b) Steepest descent, large fixed step size.



(c) Steepest descent, step size chosen via line search.



(d) Newton's method. Note that the method converges in one step.

Figure 2: Comparison of steepest descent methods with various step sizes, and Newton's method, on the same quadratic cost function.

and thus, by setting \mathbf{x}^{k+1} to be the \mathbf{x} that satisfies the above, we get the

$$\mathbf{x}^{k+1} = \mathbf{x}^k - (\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k) \quad (16)$$

or more generally,

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha (\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k). \quad (17)$$

Note that this update is only valid for $\nabla^2 f(\mathbf{x}^k) > 0$. When this condition doesn't hold, \mathbf{x}^{k+1} is not a minimizer of the second order approximation (as a result of the SOCs). See figure 2d for an example where Newton's method converges in one step, as a result of the cost function being quadratic.

Diagonally scaled steepest descent, $D^k = \text{diag}(d_1^k, \dots, d_n^k)$. Have $d_i^k > 0 \forall i$. A popular choice is

$$d_i^k = \left(\frac{\partial^2 f(\mathbf{x}^k)}{\partial x_i^2} \right)^{-1} \quad (18)$$

which is a diagonal approximation of the Hessian.

Modified Newton's method, $D^k = (\nabla^2 f(\mathbf{x}^0))^{-1}$. Requires $\nabla^2 f(\mathbf{x}^0) > 0$. For cases in which one expects $\nabla^2 f(\mathbf{x}^0) \approx \nabla^2 f(\mathbf{x}^k)$, this removes having to compute the Hessian at each step.

In addition to choosing the descent direction, there also exist a variety of methods to choose the step size α . A computationally intensive but efficient (in terms of the number of steps taken) is using a minimization rule of the form

$$\alpha^k = \text{argmin}_{\alpha \geq 0} f(\mathbf{x}^k + \alpha \mathbf{d}^k) \quad (19)$$

which is usually solved via line search (figure 2c). Alternative approaches include a limited minimization rule, in which you constrain $\alpha^k \in [0, s]$ during the line search, or simpler approach such as a constant step size (which may not guarantee convergence), or a diminishing scheduled step size. In this last case, schedules are typically chosen such that $\alpha^k \rightarrow 0$ as $k \rightarrow \infty$, while $\sum_{k=0}^{\infty} \alpha^k = +\infty$.

1.2 Constrained Nonlinear Optimization

In this section we will address the general constrained nonlinear optimization problem,

$$\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$$

which may equivalently be written

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & \mathbf{x} \in \mathcal{X} \end{aligned}$$

where the set \mathcal{X} is usually specified in terms of equality and inequality constraints. To operate within this problem structure, we will develop a set of optimality conditions involving auxiliary variables called *Lagrange multipliers*.

1.2.1 Equality Constrained Optimization

We will first look at optimization with equality constraints of the form

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h_i(\mathbf{x}) = 0, \quad i = 1, \dots, m \end{aligned}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $h_i : \mathbb{R}^n \rightarrow \mathbb{R}$ are C^1 . We will write $\mathbf{h} = [h_1, \dots, h_m]^T$. For a given local minimum \mathbf{x}^* , there exist scalars $\lambda_1, \dots, \lambda_m$ called Lagrange multipliers such that

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla h_i(\mathbf{x}^*) = 0. \quad (20)$$

There are several possible interpretations for Lagrange multipliers. First, note that the cost gradient $\nabla f(\mathbf{x}^*)$ is in the subspace spanned by the constraint gradients at \mathbf{x}^* . Equivalently, $\nabla f(\mathbf{x}^*)$ is orthogonal to the subspace of first order feasible variations

$$V(\mathbf{x}^*) = \{\Delta \mathbf{x} \mid \nabla h_i(\mathbf{x}^*)^T \Delta \mathbf{x} = 0, i = 1, \dots, m\}. \quad (21)$$

This subspace is the space of variations $\Delta \mathbf{x}$ for which $\mathbf{x} = \mathbf{x}^* + \Delta \mathbf{x}$ satisfies the constraint $\mathbf{h}(\mathbf{x}) = 0$ up to first order. Therefore, at a local minimum, the first order cost variation $\nabla f(\mathbf{x}^*)^T \Delta \mathbf{x}$ is zero for all variations $\Delta \mathbf{x}$ in this space.

Given this informal understanding, we may now precisely state the necessary conditions for optimality in constrained optimization.

Theorem 1.3 (NOC for equality constrained optimization). *Let \mathbf{x}^* be a local minimum of f subject to $\mathbf{h}(\mathbf{x}) = 0$, and assume that the constraint gradients $\nabla h_1(\mathbf{x}^*), \dots, \nabla h_m(\mathbf{x}^*)$ are linearly independent. Then there exists a unique vector $\boldsymbol{\lambda}^* = [\lambda_1^*, \dots, \lambda_m^*]^T$ called a Lagrange multiplier vector, such that*

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla h_i(\mathbf{x}^*) = 0. \quad (22)$$

If in addition f and \mathbf{h} are C^2 , we have

$$\mathbf{y}^T (\nabla^2 f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla^2 h_i(\mathbf{x}^*)) \mathbf{y} \geq 0, \quad \forall \mathbf{y} \in V(\mathbf{x}^*) \quad (23)$$

where

$$V(\mathbf{x}^*) = \{\mathbf{y} \mid \nabla h_i(\mathbf{x}^*)^T \mathbf{y} = 0, i = 1, \dots, m\}. \quad (24)$$

Proof. See [Ber16] Section 3.1.1 and 3.1.2. □

We will sketch two possible proofs for the NOC for equality constrained optimization.

Penalty approach. This approach relies on adding to the cost function a large penalty term for constraint violation. This is the same approach that will be used in proving the necessary conditions for inequality constrained optimization, and is the basis of a variety of practical numerical algorithms.

Elimination approach. This approach views the constraints as a system of m equations with n unknowns, for which m variables can be expressed in terms of the remaining $m - n$ variables. This reduces the problem to an unconstrained optimization problem.

Note that in theorem 1.3, we assumed the gradients of the constraint functions were linearly independent. A feasible vector for which this holds is called *regular*. If this condition is violated, a Lagrange multiplier for a local minimum may not exist.

For convenience, we will write the necessary conditions in terms of the Lagrangian function $L : \mathbb{R}^{m+n} \rightarrow \mathbb{R}$,

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i h_i(\mathbf{x}). \quad (25)$$

This function allows the NOC conditions to be succinctly stated as

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*) = 0 \quad (26)$$

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*) = 0 \quad (27)$$

$$\mathbf{y}^T \nabla_{\mathbf{x}\mathbf{x}}^2 L(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{y} \geq 0, \quad \forall \mathbf{y} \in V(\mathbf{x}^*). \quad (28)$$

which form a system of $n + m$ equations with $n + m$ unknowns. Given this notation, we can state the sufficient conditions.

Theorem 1.4 (SOC for equality constrained optimization). *Assume that f and \mathbf{h} are C^2 and let $\mathbf{x}^* \in \mathbb{R}^n$ and $\boldsymbol{\lambda}^* \in \mathbb{R}^m$ satisfy*

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*) = 0 \quad (29)$$

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*) = 0 \quad (30)$$

$$\mathbf{y}^T \nabla_{\mathbf{x}\mathbf{x}}^2 L(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{y} > 0, \quad \forall \mathbf{y} \neq 0, \mathbf{y} \in V(\mathbf{x}^*). \quad (31)$$

Proof. See [Ber16] Section 3.2. □

Note that the SOC does not include regularity of \mathbf{x}^* .

1.2.2 Inequality Constrained Optimization

We will now address the general case, including inequality constraints,

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{s.t.} \quad & h_i(\mathbf{x}) = 0, \quad i = 0, \dots, m \\ & g_j(\mathbf{x}) \leq 0, \quad j = 0, \dots, r \end{aligned}$$

where f, h_i, g_i are C^1 . The key intuition for the case of inequality constraints is based on realizing that for any feasible point, some subset of the constraints will be active (for which $g_j(\mathbf{x}) = 0$), while the complement of this set will be inactive. We define the active set of inequality constraints, which we denote

$$A(\mathbf{x}) = \{j \mid g_j(\mathbf{x}) = 0\}. \quad (32)$$

A constraint is active at \mathbf{x} if it is in $A(\mathbf{x})$, otherwise it is inactive. Note that if \mathbf{x}^* is a local minimum of the inequality constrained problem, then \mathbf{x}^* is a local minimum of the identical problem with the inactive constraints removed. Moreover, at this local minimum, the constraints may be treated as equality constraints. Thus, if \mathbf{x}^* is regular, there exists Lagrange multipliers $\lambda_1^*, \dots, \lambda_m^*$ and $\mu_j^*, j \in A(\mathbf{x}^*)$ such that

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i \nabla h_i(\mathbf{x}^*) + \sum_{j \in A(\mathbf{x}^*)} \mu_j^* \nabla g_j(\mathbf{x}^*) = 0. \quad (33)$$

We will define the Lagrangian

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i h_i(\mathbf{x}) + \sum_{j=1}^r \mu_j^* g_j(\mathbf{x}), \quad (34)$$

which we will use to state the necessary and sufficient conditions.

Theorem 1.5 (Karush-Kuhn-Tucker NOC). *Let \mathbf{x}^* be a local minimum for the inequality constrained problem where f, h_i, g_j are C^1 and assume \mathbf{x}^* is regular (equality and active inequality constraint gradients are linearly independent). Then, there exists unique Lagrange multiplier vectors $\boldsymbol{\lambda}^*$ and $\boldsymbol{\mu}^*$ such that*

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = 0 \quad (35)$$

$$\boldsymbol{\mu} \geq 0 \quad (36)$$

$$\mu_j^* = 0, \quad \forall j \notin A(\mathbf{x}^*) \quad (37)$$

If in addition, $f, \mathbf{h}, \mathbf{g}$ are C^2 , we have

$$\mathbf{y}^T \nabla_{\mathbf{xx}}^2 L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \mathbf{y} \geq 0 \quad (38)$$

for all \mathbf{y} such that

$$\nabla h_i(\mathbf{x}^*)^T \mathbf{y} = 0, \quad i = 1, \dots, m \quad (39)$$

$$\nabla g_j(\mathbf{x}^*)^T \mathbf{y} = 0, \quad j \in A(\mathbf{x}^*) \quad (40)$$

Proof. See [Ber16] Section 3.3.1. □

The SOC are obtained similarly to the equality constrained case.

1.3 Further Reading

In this section we have addressed the necessary and sufficient conditions for constrained and unconstrained nonlinear optimization. This section is based heavily on [Ber16], and we refer the reader to this book for further details. We have avoided discussing linear programming, which is itself a large topic of study, about which many books have been written (we refer the reader to [BT97] as a good reference on the subject).

Convex optimization has become a powerful and widespread tool in modern optimal control. While we have only addressed it briefly here, [BV04] offers a fairly comprehensive treatment of the theory and practice of convex optimization. For a succinct overview with a focus on machine learning, we refer the reader to [Kol08].

2 Dynamic Programming and the Linear Quadratic Regulator

2.1 The Optimal Control Problem

In this section, we will outline the deterministic continuous-time optimal control problem that we will aim to solve. We will denote the state at time t as $\mathbf{x}(t) \in \mathbb{R}^n$, and the control as $\mathbf{u}(t) \in \mathbb{R}^m$. We will also occasionally write these as \mathbf{x}_t and \mathbf{u}_t , respectively. We will write the continuous-time systems dynamics as

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t), t). \quad (41)$$

We will refer to a history of control input values during an interval $[t_0, t_f]$ as a control history, and we will refer to a history of state values over this interval as a state trajectory.

Different control problems may call for various constraints. For example, we may constrain a quadrotor to only fly in space not occupied by obstacles. Examples of constraints we will see are

- Initial and final conditions, $\mathbf{x}(t_0) = \mathbf{x}_0$, $\mathbf{x}(t_f) = \mathbf{x}_f$
- Trajectory constraints, $\mathbf{x} \leq \mathbf{x}(t) \leq \bar{\mathbf{x}}$
- Control limits, $\underline{\mathbf{u}} \leq \mathbf{u}(t) \leq \bar{\mathbf{u}}$.

A state trajectory and control history that satisfy the constraints during the entire time interval $[t_0, t_f]$ are called admissible trajectories and admissible controls, respectively.

Finally, we will define the performance measure,

$$J = c_f(\mathbf{x}(t_f), t_f) + \int_{t_0}^{t_f} c(\mathbf{x}(t), \mathbf{u}(t), t) dt \quad (42)$$

where c is the instantaneous cost function, and c_f is the terminal state cost. We are now able to state the continuous-time optimal control problem. We aim to find an admissible control, \mathbf{u}^* , which causes the system (41) to follow an admissible trajectory, \mathbf{x}^* , that minimizes the performance measure given by (42). The minimizer $(\mathbf{x}^*, \mathbf{u}^*)$ is called an optimal trajectory-control pair.

Note, first of all, that this is an extremely general problem formulation. We have not fixed our system dynamics, cost function, or specific constraints. We can't, in general, guarantee the existence or uniqueness of the optimal solution.

There are two possible solution forms for the optimal control. The first, $\mathbf{u}^* = e(\mathbf{x}(t_0), t)$ is referred to as an open-loop solution. This is an input function that is applied to the system, without using feedback. Practically, such solutions usually require augmentation with a feedback controller, as small model mismatch may lead to compounding errors. The second possible solution form is a feedback policy, $\mathbf{u}^* = \pi(\mathbf{x}(t), t)$. This feedback law maps all state-time pairs to an action and thus is usually more robust to possible model mismatch. However, depending on the particular problem formulation, open-loop solutions may be easier to compute.

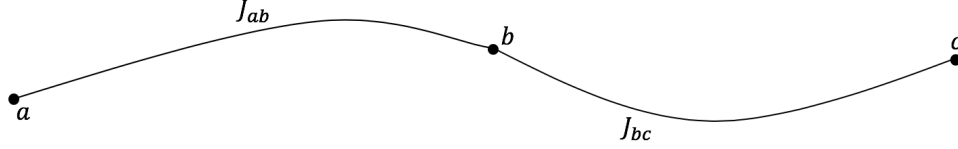


Figure 3: An optimal trajectory connecting point a to point c . There are no better (lower cost) trajectories than the sub-trajectory connecting b and c , by the principle of optimality.

2.2 Dynamic Programming and the Principle of Optimality

In this chapter we will outline the principle of optimality, and the method of dynamic programming (DP), one of two main approaches to solving the optimal control problem. The second, so-called variational approaches based on Pontryagin's Maximum Principle (PMP) will be discussed in future chapters. Dynamic programming has the strong advantage of yielding a feedback policy, however, exactly solving the dynamic programming problem is infeasible for many systems. We will address special cases in which the DP problem can be solved exactly, and approximate methods that work for a wide variety of systems.

Despite having just introduced the optimal control problem in continuous time, we will be operating in discrete time here, in which we aim to minimize

$$J_f(\mathbf{x}_0) = c_f(\mathbf{x}_N) + \sum_{k=0}^{N-1} c(\mathbf{x}_k, \mathbf{u}_k, k). \quad (43)$$

We will extend the methods we develop in this chapter to continuous time in the next chapter.

The principle of optimality is as follows. Figure 3 shows a trajectory from point a to c . If the cost of the trajectory, $J_{ac} = J_{ab} + J_{bc}$, is minimal, then J_{bc} is also a minimum cost trajectory connecting b and c . The proof of this principle, stated informally, is simple. Assume there exists an alternative trajectory connecting b and c , for which we will write the cost as \tilde{J}_{bc} , that achieves $\tilde{J}_{bc} < J_{bc}$. Then, we have

$$\tilde{J}_{ac} = J_{ab} + \tilde{J}_{bc} \quad (44)$$

$$< J_{ab} + J_{bc} \quad (45)$$

$$= J_{ac}, \quad (46)$$

and thus J_{ac} isn't minimal. More formally,

Theorem 2.1 (Discrete-time Principle of Optimality: Deterministic Case). *Let $\pi^* = (\pi_0^*, \dots, \pi_{N-1}^*)$ be an optimal policy. Assume state \mathbf{x}_k is reachable. Consider the subproblem whereby we are at \mathbf{x}_k at time k and we wish to minimize the cost-to-go from time k to time N . Then the truncated policy $(\pi_k^*, \dots, \pi_{N-1}^*)$ is optimal for the subproblem.*

Dynamic programming, intuitively, proceeds backwards in time, first solving simpler shorter horizon problems. If we have found the optimal policy for times $k+1$ to $N-1$, along with the associated cost-to-go for each state, choosing the optimal policy for time k

is a one step optimization problem. More concretely, we will assume we have dynamics of the form $\mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k, k)$ with $\mathbf{u}_k \in \mathcal{U}(\mathbf{x}_k)$, and the cost given by (43). Then, dynamic programming iterates backward in time, from $N - 1$ to 0, with

$$J_N(\mathbf{x}_N) = c_T(\mathbf{x}_N) \quad (47)$$

$$J_k(\mathbf{x}_k) = \min_{\mathbf{u}_k \in \mathcal{U}(\mathbf{x}_k)} \{c_k(\mathbf{x}_k, \mathbf{u}_k, k) + J_{k+1}(f(\mathbf{x}_k, \mathbf{u}_k, k))\}. \quad (48)$$

Note that here we have considered only deterministic dynamical systems (there is no stochastic disturbance). Equation (48) is one form of the *Bellman equation*, one of the most important relations in optimal control.

Dynamic programming raises many practical issues if one were to attempt to apply it directly. To perform the recursion, J_{k+1} must be known for all \mathbf{x}_{k+1} (or more precisely, all \mathbf{x}_{k+1} that are reachable from \mathbf{x}_k). If the state space is discrete (and relatively small), this is tractable as the cost-to-go may just be maintained in tabular form. In the next subsection, we will discuss an extremely important case in continuous space in which the cost-to-go can be computed exactly for all states. However, for general systems, we can not expect to be able to compute the cost-to-go for all states. Possible approaches to make the DP approach tractable are discretizing the state space, approximating the cost-to-go (i.e. restricting the family of functions that J_{k+1} may be in), or interpolating between cost-to-go computed for a finite set of states.

2.3 Discrete LQR

An important instance in which dynamic programming can be solved analytically for continuous state-action systems is the *linear quadratic regulator* problem. We will fix the dynamics of the system to be (possibly time-varying) linear,

$$\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k \quad (49)$$

and the cost function as quadratic

$$c(\mathbf{x}_k, \mathbf{u}_k) = \frac{1}{2}(\mathbf{x}_k^T Q_k \mathbf{x}_k + \mathbf{u}_k^T R_k \mathbf{u}_k) \quad (50)$$

$$c_N(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T Q_N \mathbf{x}_k \quad (51)$$

where $Q_k \in \mathbb{R}^{n \times n}$ is positive semi-definite and $R_k \in \mathbb{R}^{m \times m}$ is positive definite for all $k = 0, \dots, N$. Importantly, we assume \mathbf{x}_k and \mathbf{u}_k are unconstrained for all k . To perform DP recursion, we initialize

$$J_N^*(\mathbf{x}_N) = \frac{1}{2} \mathbf{x}_N^T Q_N \mathbf{x}_N := \frac{1}{2} \mathbf{x}_N^T V_N \mathbf{x}_N. \quad (52)$$

Then, applying (48), we have

$$J_{N-1}^*(\mathbf{x}_{N-1}) = \frac{1}{2} \min_{\mathbf{u}_{N-1} \in \mathbb{R}^m} \{ \mathbf{x}_{N-1}^T Q_{N-1} \mathbf{x}_{N-1} + \mathbf{u}_{N-1}^T R_{N-1} \mathbf{u}_{N-1} + \mathbf{x}_N^T V_N \mathbf{x}_N \} \quad (53)$$

which, applying the dynamics,

$$J_{N-1}^*(\mathbf{x}_{N-1}) = \frac{1}{2} \min_{\mathbf{u}_{N-1} \in \mathbb{R}^m} \{ \mathbf{x}_{N-1}^T Q_{N-1} \mathbf{x}_{N-1} + \mathbf{u}_{N-1}^T R_{N-1} \mathbf{u}_{N-1} + (A_{N-1} \mathbf{x}_{N-1} + B_{N-1} \mathbf{u}_{N-1})^T V_N (A_{N-1} \mathbf{x}_{N-1} + B_{N-1} \mathbf{u}_{N-1}) \}. \quad (54)$$

Rearranging, we have

$$J_{N-1}^*(\mathbf{x}_{N-1}) = \frac{1}{2} \min_{\mathbf{u}_{N-1} \in \mathbb{R}^m} \{ \mathbf{x}_{N-1}^T (Q_{N-1} + A_{N-1}^T V_N A_{N-1}) \mathbf{x}_{N-1} + \mathbf{u}_{N-1}^T (R_{N-1} + B_{N-1}^T V_N B_{N-1}) \mathbf{u}_{N-1} + 2 \mathbf{u}_{N-1}^T (B_{N-1}^T V_N A_{N-1}) \mathbf{x}_{N-1} \}. \quad (55)$$

Note that this optimization problem is convex in \mathbf{u}_{N-1} as $R_{N-1} + B_{N-1}^T V_N B_{N-1} > 0$. Therefore, any local minima is a global minima, and therefore we can simply apply the first order optimality conditions. Differentiating,

$$\frac{\partial J_{N-1}^*}{\partial \mathbf{u}_{N-1}}(\mathbf{x}_{N-1}) = (R_{N-1} + B_{N-1}^T V_N B_{N-1}) \mathbf{u}_{N-1} + (B_{N-1}^T V_N A_{N-1}) \mathbf{x}_{N-1} \quad (56)$$

and setting this to zero yields

$$\mathbf{u}_{N-1}^* = -(R_{N-1} + B_{N-1}^T V_N B_{N-1})^{-1} (B_{N-1}^T V_N A_{N-1}) \mathbf{x}_{N-1} \quad (57)$$

which we write

$$\mathbf{u}_{N-1}^* = L_{N-1} \mathbf{x}_{N-1} \quad (58)$$

which is a time-varying linear feedback policy. Plugging this feedback policy into (54),

$$J_{N-1}^*(\mathbf{x}_{N-1}) = \mathbf{x}_{N-1}^T (Q_{N-1} + L_{N-1}^T R_{N-1} L_{N-1} + (A_{N-1} + B_{N-1} L_{N-1})^T V_N (A_{N-1} + B_{N-1} L_{N-1})) \mathbf{x}_{N-1}. \quad (59)$$

Critically, this implies that the cost-to-go is always a positive semi-definite quadratic function of the state. Because the optimal policy is always linear, and the optimal cost-to-go is always quadratic, the DP recursion may be recursively performed backward in time and the minimization may be performed analytically.

Following the same procedure, we can write the DP recursion for the discrete-time LQR controller:

1. $V_N = Q_N$
2. $L_k = -(R_k + B_k^T V_{k+1} B_k)^{-1} (B_k^T V_{k+1} A_k)$
3. $V_k = Q_k + L_k^T R_k L_k + (A_k + B_k L_k)^T V_{k+1} (A_k + B_k L_k)$
4. $\mathbf{u}_k^* = L_k \mathbf{x}_k$

$$5. J_k^*(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T V_k \mathbf{x}_k$$

There are several implications of this recurrence relation. First, even if A, B, Q, R are all constant (not time-varying), the policy is still time-varying. Why is this the case? Control effort invested early in the problem will yield dividends over the remaining length of the horizon, in terms of lower state cost for all future time steps. However, as the remaining length of the episode becomes shorter, this tradeoff is increasingly imbalanced, and the control effort will decrease. However, for a linear time-invariant system, if (A, B) is controllable, the feedback gain L_k approach a constant as the episode length approaches infinity. This time-invariant policy is practical for long horizon control problems, and may be approximately computed by running the DP recurrence relation until approximate convergence.

2.3.1 LQR with (Bi)linear Cost and Affine Dynamics

In the previous subsection, we have derived the common formulation of the LQR controller. In this subsection, we will derive the discrete time LQR controller for a more general system with bilinear/linear terms in the cost and affine terms in the dynamics. This derivation will be the basis of algorithms we will build up in the following subsections. More concretely, we consider systems with stage-wise cost

$$c(\mathbf{x}_k, \mathbf{u}_k) = \frac{1}{2} \mathbf{x}_k^T Q_k \mathbf{x}_k + \frac{1}{2} \mathbf{u}_k^T R_k \mathbf{u}_k + \mathbf{u}_k^T H_k \mathbf{x}_k + \mathbf{q}_k^T \mathbf{x}_k + \mathbf{r}_k \mathbf{u}_k + q_k, \quad (60)$$

terminal cost

$$c_N(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T Q_N \mathbf{x}_k + \mathbf{q}_N^T \mathbf{x}_k + q_N, \quad (61)$$

and dynamics

$$\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k + \mathbf{d}_k. \quad (62)$$

The cost-to-go will take the form

$$J_k(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T V_k \mathbf{x}_k + \mathbf{v}_k^T \mathbf{x}_k + v_k. \quad (63)$$

Repeating our approach from the last subsection, we have

$$\begin{aligned} J_k^*(\mathbf{x}_k) = \min_{\mathbf{u}_k \in \mathbb{R}^m} \{ & \frac{1}{2} \mathbf{x}_k^T Q_k \mathbf{x}_k + \frac{1}{2} \mathbf{u}_k^T R_k \mathbf{u}_k + \mathbf{u}_k^T H_k \mathbf{x}_k + \mathbf{q}_k^T \mathbf{x}_k + \mathbf{r}_k \mathbf{u}_k + q_k \\ & + \frac{1}{2} (A_k \mathbf{x}_k + B_k \mathbf{u}_k + \mathbf{d}_k)^T V_{k+1} (A_k \mathbf{x}_k + B_k \mathbf{u}_k + \mathbf{d}_k) \\ & + \mathbf{v}_{k+1}^T (A_k \mathbf{x}_k + B_k \mathbf{u}_k + \mathbf{d}_k) + v_{k+1} \}. \end{aligned} \quad (64)$$

Rearranging, have

$$\begin{aligned} J_k^*(\mathbf{x}_k) = \min_{\mathbf{u}_k \in \mathbb{R}^m} \{ & \frac{1}{2} \mathbf{x}_k^T (Q_k + A_k^T V_{k+1} A_k) \mathbf{x}_k + \frac{1}{2} \mathbf{u}_k^T (R_k + B_k^T V_{k+1} B_k) \mathbf{u}_k \\ & + \mathbf{u}_k^T (H_k + B_k^T V_{k+1} A_k)^T \mathbf{x}_k + (\mathbf{q}_k + A_k^T V_{k+1} \mathbf{d}_k + A_k^T \mathbf{v}_{k+1})^T \mathbf{x}_k \\ & + (\mathbf{r}_k + B_k^T V_{k+1} \mathbf{d}_k + B_k^T \mathbf{v}_{k+1}) \mathbf{u}_k + (v_{k+1} + \frac{1}{2} \mathbf{d}_k^T V_{k+1} \mathbf{d}_k + \mathbf{v}_{k+1}^T \mathbf{d}_k) \}. \end{aligned} \quad (65)$$

Solving this minimization problem, we see that our optimal controller takes the form

$$\mathbf{u}_k^* = \mathbf{l}_k + L_k \mathbf{x}_k. \quad (66)$$

We will define the following useful terms which will be used throughout the remainder of this section

$$S_{\mathbf{u},k} = \mathbf{r}_k + B_k^T \mathbf{v}_{k+1} + B_k^T V_{k+1} \mathbf{d}_k \quad (67)$$

$$S_{\mathbf{uu},k} = R_k + B_k^T V_{k+1} B_k \quad (68)$$

$$S_{\mathbf{ux},k} = H_k + B_k^T V_{k+1} A_k. \quad (69)$$

Given this notation, all necessary terms can be computed via the following relations

$$1. \ V_N = Q_N; \ \mathbf{v}_N = \mathbf{q}_N; \ v_N = q_N$$

2.

$$L_k = -S_{\mathbf{uu},k}^{-1} S_{\mathbf{ux},k} \quad (70)$$

$$\mathbf{l}_k = -S_{\mathbf{uu},k}^{-1} S_{\mathbf{u},k} \quad (71)$$

3.

$$V_k = Q_k + A_k^T V_{k+1} A_k - L_k^T S_{\mathbf{uu},k} L_k \quad (72)$$

$$\mathbf{v}_k = \mathbf{q}_k + A_k^T (\mathbf{v}_{k+1} + V_{k+1} \mathbf{d}_k) + S_{\mathbf{ux},k}^T \mathbf{l}_k \quad (73)$$

$$v_k = v_{k+1} + q_k + \mathbf{d}_k^T \mathbf{v}_{k+1} + \frac{1}{2} \mathbf{d}_k^T V_{k+1} \mathbf{d}_k + \frac{1}{2} \mathbf{l}_k^T S_{\mathbf{u},k} \quad (74)$$

$$4. \ \mathbf{u}_k^* = \mathbf{l}_k + L_k \mathbf{x}_k$$

$$5. \ J_k(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^T V_k \mathbf{x}_k + \mathbf{v}_k^T \mathbf{x}_k + v_k.$$

Note that in the following subsections (specifically in our discussion of differential dynamic programming) we will introduce more convenient (and compact) notation.

2.3.2 LQR Tracking around a Linear Trajectory

In the previous subsections we have considered the generic linear quadratic control problem, in which we want to regulate to a fixed point, and deviations from this point are penalized. In this section, we will address the case in which we want to track a pre-specified trajectory. Let us assume (for now) that we have been given a nominal trajectory of the form $(\bar{\mathbf{x}}_0, \dots, \bar{\mathbf{x}}_N)$ and $(\bar{\mathbf{u}}_0, \dots, \bar{\mathbf{u}}_{N-1})$. We will also assume that this trajectory satisfies our given dynamics, such that

$$\bar{\mathbf{x}}_{k+1} = A_k \bar{\mathbf{x}}_k + B_k \bar{\mathbf{u}}_k + \mathbf{d}_k, \ \forall k = 0, \dots, N-1. \quad (75)$$

Then, we can rewrite our dynamics in terms of deviations from the nominal trajectory,

$$\delta \mathbf{x}_k = \mathbf{x}_k - \bar{\mathbf{x}}_k \quad (76)$$

$$\delta \mathbf{u}_k = \mathbf{u}_k - \bar{\mathbf{u}}_k. \quad (77)$$

Rewriting, we have

$$\delta \mathbf{x}_{k+1} = A_k \delta \mathbf{x}_k + B_k \delta \mathbf{u}_k. \quad (78)$$

Thus, tracking the nominal trajectory reduces to driving the state deviation, $\delta \mathbf{x}_k$, to zero. Note that solving this problem requires rewriting the original cost function in terms of the deviations $\delta \mathbf{x}_k, \delta \mathbf{u}_k$.

2.3.3 LQR Tracking around a Nonlinear Trajectory

Despite LQR being an extremely powerful approach to optimal control, it suffers from a handful of limitations. First and foremost, it assumes the dynamics are (possibly time-varying) linear, and the cost function is quadratic. While most systems are in fact nonlinear, a typical approach to designing feedback controllers is to linearize around some operating point. This is an effective method for designing regulators, which aim to control the system to some particular state. If, in contrast, we wish to track a trajectory, we must instead linearize around this trajectory. We will assume we are given a nominal trajectory which satisfies the nonlinear dynamics, such that

$$\bar{\mathbf{x}}_{k+1} = f(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k), \quad \forall k = 0, \dots, N-1. \quad (79)$$

Given this, we can linearize our system at each timestep by Taylor expanding,

$$\mathbf{x}_{k+1} \approx f(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) + \underbrace{\frac{\partial f}{\partial \mathbf{x}}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)}_{A_k} (\mathbf{x}_k - \bar{\mathbf{x}}_k) + \underbrace{\frac{\partial f}{\partial \mathbf{u}}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)}_{B_k} (\mathbf{u}_k - \bar{\mathbf{u}}_k) \quad (80)$$

which allows us to again rewrite the system in terms of deviations, to get

$$\delta \mathbf{x}_{k+1} = A_k \delta \mathbf{x}_k + B_k \delta \mathbf{u}_k \quad (81)$$

which is linear in $\delta \mathbf{x}_k, \delta \mathbf{u}_k$. Note that design of systems of this type often require careful design and analysis, as deviating from the nominal trajectory results in the loss of accuracy of the local model linearization.

In designing this tracking system, a second question now occurs: how do we choose our cost function? One possible option is arbitrary choice of Q and R by the system designer. This has the advantage of being easily customizable to change system behavior, and we can guarantee the necessary conditions on these matrices. A second option, if we are given some arbitrary (possibly non-quadratic) cost function c , is to locally quadratize the cost function.

Writing

$$c_k := c(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) \quad (82)$$

$$c_{i,k} := \frac{\partial c}{\partial i}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) \quad (83)$$

$$c_{ij,k} := \frac{\partial^2 c}{\partial i \partial j}(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k) \quad (84)$$

we can second order Taylor expand our cost function around our nominal trajectory

$$c(\delta \mathbf{x}_k, \delta \mathbf{u}_k) \approx \frac{1}{2} \begin{bmatrix} 1 \\ \delta \mathbf{x}_k \\ \delta \mathbf{u}_k \end{bmatrix}^T \begin{bmatrix} 2c_k & c_{\mathbf{x},k}^T & c_{\mathbf{u},k}^T \\ c_{\mathbf{x},k} & c_{\mathbf{x}\mathbf{x},k} & c_{\mathbf{u}\mathbf{x},k}^T \\ c_{\mathbf{u},k} & c_{\mathbf{u}\mathbf{x},k} & c_{\mathbf{u}\mathbf{u},k} \end{bmatrix} \begin{bmatrix} 1 \\ \delta \mathbf{x}_k \\ \delta \mathbf{u}_k \end{bmatrix}. \quad (85)$$

Here $c_{\mathbf{x}\mathbf{x},k}$ and $c_{\mathbf{u}\mathbf{u},k}$ replace Q_k and R_k from the previous section, respectively. There are two primary concerns with this approach to choosing the cost function. First, we require the quadratic form in (85) to be positive semi-definite and $c_{\mathbf{u}\mathbf{u},k}$ to be positive definite, for all k . Second, we have an implicit cost that we would like to stay close to the nominal trajectory to ensure our linearized model does not become inaccurate. As a result of this implicit cost, we may wish to tune the cost terms to yield tracking that is better suited to the nonlinear model that we are tracking.

2.4 Iterative LQR and Differential Dynamic Programming

2.4.1 Iterative LQR

We have addressed the case in which we wish to track a given trajectory with LQR. A natural question, now, is whether we can use LQR to improve on this nominal trajectory? Iterative LQR augments tracking LQR with a forward pass in which the nominal trajectory is updated. As a consequence, it can be used to improve trajectories and in most cases, can be used as a practical trajectory generation and control algorithm for nonlinear systems. We will define the following useful terms

$$Q_k = c_k + v_{k+1} \quad (86)$$

$$Q_{\mathbf{x},k} = c_{\mathbf{x},k} + f_{\mathbf{x},k}^T v_{k+1} \quad (87)$$

$$Q_{\mathbf{u},k} = c_{\mathbf{u},k} + f_{\mathbf{u},k}^T v_{k+1} \quad (88)$$

$$Q_{\mathbf{x}\mathbf{x},k} = c_{\mathbf{x}\mathbf{x},k} + f_{\mathbf{x},k}^T V_{k+1} f_{\mathbf{x},k} \quad (89)$$

$$Q_{\mathbf{u}\mathbf{u},k} = c_{\mathbf{u}\mathbf{u},k} + f_{\mathbf{u},k}^T V_{k+1} f_{\mathbf{u},k} \quad (90)$$

$$Q_{\mathbf{u}\mathbf{x},k} = c_{\mathbf{u}\mathbf{x},k} + f_{\mathbf{u},k}^T V_{k+1} f_{\mathbf{x},k} \quad (91)$$

where $f_{\mathbf{x},k} = A_k$ and $f_{\mathbf{u},k} = B_k$. In this form, the optimal control perturbation is

$$\delta \mathbf{u}_k^* = \mathbf{l}_k + L_k \delta \mathbf{x}_k \quad (92)$$

Algorithm 1 iLQR

Require: Nominal control sequence, $(\bar{\mathbf{u}}_0, \dots, \bar{\mathbf{u}}_{N-1})$

- 1: $\delta \mathbf{u}_k = 0$ for all k
 - 2: **while** not converged **do**
 Forward pass:
 - 3: Compute nominal trajectory $\bar{\mathbf{x}}_{k+1} = f(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k + \delta \mathbf{u}_k)$ and set $\bar{\mathbf{u}}_k \leftarrow \bar{\mathbf{u}}_k + \delta \mathbf{u}_k$
 - Backward pass:
 - 4: Compute Q terms around $(\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k)$ for all k via (86 – 91)
 - 5: Update feedback law via (93 – 94)
 - 6: Update value approximation via (95 – 97)
 - 7: **end while**
 - 8: Compute control law $\pi_k(\mathbf{x}_k) = \bar{\mathbf{u}}_k + \mathbf{l}_k + L_k(\mathbf{x}_k - \bar{\mathbf{x}}_k)$
 - 9: **return** $\{\pi_k\}_{k=0}^{N-1}$
-

where

$$\mathbf{l}_k = -Q_{\mathbf{u}\mathbf{u},k}^{-1} Q_{\mathbf{u},k} \quad (93)$$

$$L_k = -Q_{\mathbf{u}\mathbf{u},k}^{-1} Q_{\mathbf{u}\mathbf{x},k}. \quad (94)$$

Finally, the local backward recursion can be completed by updating the value function terms via

$$v_k = Q_k - \frac{1}{2} \mathbf{l}_k^T Q_{\mathbf{u}\mathbf{u},k} \mathbf{l}_k \quad (95)$$

$$\mathbf{v}_k = Q_{\mathbf{x},k} - L_k^T Q_{\mathbf{u}\mathbf{u},k} \mathbf{l}_k \quad (96)$$

$$V_k = Q_{\mathbf{x}\mathbf{x},k} - L_k^T Q_{\mathbf{u}\mathbf{u},k} L_k. \quad (97)$$

So far, we have simply derived an alternative method for performing a quadratic approximation of the DP recursion around some nominal trajectory. The iterative LQR (iLQR) algorithm differs by introducing a forward pass that updates the trajectory that is being tracked. The algorithm alternates between forward passes, in which the control policy is applied to the nonlinear dynamics, and backward passes in which the cost function and dynamics are linearized around the new nominal trajectory, and the quadratic approximation of the value, as well as the new control law, is computed. The iterative LQR algorithm is outlined in Algorithm 1. Critically, note that this algorithm returns both a nominal trajectory, in terms of the $\bar{\mathbf{x}}_k, \bar{\mathbf{u}}_k$, as well as a feedback policy that stabilizes around this trajectory.

2.4.2 Differential Dynamic Programming

Iterative LQR performs trajectory optimization by first linearizing the dynamics and quadratizing the cost function, and then performing the dynamic programming recursion to compute optimal controls. While this linearization/quadratization approach is sufficient for approximating the Bellman equation such that it may be solved analytically, an alternative approach

is to directly approximate the Bellman equation. *Differential dynamic programming* (DDP) directly builds a quadratic approximation of the right hand side of the Bellman equation (as opposed to first approximating the dynamics and the cost function), which may then be solved analytically. We will first define the change in the value of J_k under a perturbation $\delta \mathbf{x}_k, \delta \mathbf{u}_k$,

$$Q(\delta \mathbf{x}_k, \delta \mathbf{u}_k) := c(\bar{\mathbf{x}}_k + \delta \mathbf{x}_k, \bar{\mathbf{u}}_k + \delta \mathbf{u}_k) + J_{k+1}(f(\bar{\mathbf{x}}_k + \delta \mathbf{x}_k, \bar{\mathbf{u}}_k + \delta \mathbf{u}_k)). \quad (98)$$

Note that Q here is different from the Q matrix in Section 2.3. Using the same notation as in (82), we can write the quadratic expansion of (98) as

$$Q(\delta \mathbf{x}_k, \delta \mathbf{u}_k) \approx \frac{1}{2} \begin{bmatrix} 1 \\ \delta \mathbf{x}_k \\ \delta \mathbf{u}_k \end{bmatrix}^T \begin{bmatrix} 2Q_k & Q_{\mathbf{x},k}^T & Q_{\mathbf{u},k}^T \\ Q_{\mathbf{x},k} & Q_{\mathbf{x}\mathbf{x},k} & Q_{\mathbf{u}\mathbf{x},k}^T \\ Q_{\mathbf{u},k} & Q_{\mathbf{u}\mathbf{x},k} & Q_{\mathbf{u}\mathbf{u},k} \end{bmatrix} \begin{bmatrix} 1 \\ \delta \mathbf{x}_k \\ \delta \mathbf{u}_k \end{bmatrix} \quad (99)$$

where

$$Q_k = c_k + v_{k+1} \quad (100)$$

$$Q_{\mathbf{x},k} = c_{\mathbf{x},k} + f_{\mathbf{x},k}^T \mathbf{v}_{k+1} \quad (101)$$

$$Q_{\mathbf{u},k} = c_{\mathbf{u},k} + f_{\mathbf{u},k}^T \mathbf{v}_{k+1} \quad (102)$$

$$Q_{\mathbf{x}\mathbf{x},k} = c_{\mathbf{x}\mathbf{x},k} + f_{\mathbf{x},k}^T V_{k+1} f_{\mathbf{x},k} + \mathbf{v}_{k+1} \cdot f_{\mathbf{x}\mathbf{x},k} \quad (103)$$

$$Q_{\mathbf{u}\mathbf{u},k} = c_{\mathbf{u}\mathbf{u},k} + f_{\mathbf{u},k}^T V_{k+1} f_{\mathbf{u},k} + \mathbf{v}_{k+1} \cdot f_{\mathbf{u}\mathbf{u},k} \quad (104)$$

$$Q_{\mathbf{u}\mathbf{x},k} = c_{\mathbf{u}\mathbf{x},k} + f_{\mathbf{u},k}^T V_{k+1} f_{\mathbf{x},k} + \mathbf{v}_{k+1} \cdot f_{\mathbf{u}\mathbf{x},k}. \quad (105)$$

Note that these terms differ only from iLQR via the last term in (103 – 105), which are second order approximation of the dynamics. Note that the dot notation denotes tensor contraction.

Given this, we can partially minimize this quadratic form over the control deviation,

$$\delta \mathbf{u}_k^* = \operatorname{argmin}_{\delta \mathbf{u}} Q(\delta \mathbf{x}_k, \delta \mathbf{u}) = \mathbf{l}_k + L_k \delta \mathbf{x}_k \quad (106)$$

where

$$\mathbf{l}_k = -Q_{\mathbf{u}\mathbf{u},k}^{-1} Q_{\mathbf{u},k} \quad (107)$$

$$L_k = -Q_{\mathbf{u}\mathbf{u},k}^{-1} Q_{\mathbf{u}\mathbf{x},k}. \quad (108)$$

The DDP algorithm is identical to Algorithm 1, just with the alternative definitions for $Q_{\mathbf{x}\mathbf{x},k}$, $Q_{\mathbf{u}\mathbf{u},k}$ and $Q_{\mathbf{u}\mathbf{x},k}$. The main philosophical difference between iLQR and DDP is that iLQR first approximates the dynamics and cost, and then solves the Bellman equation directly, whereas DDP directly approximates the Bellman equation. While DDP yields a more accurate approximation, computing the second order dynamics terms is expensive in practice. Practically, iLQR is sufficient for most applications.

2.4.3 Algorithmic Details for iLQR and DDP

Algorithm 1 leaves out several details that would be critical for implementing the algorithm. First, what convergence criteria should we use? In [TL05], the authors stop when the update to the nominal control action sequence is sufficiently small. In [LK14], the authors iterate until the cost of the trajectory (with some additional penalty terms) increases. Finally, a variety of convergence criteria are based on expected trajectory improvement, computed via line search [JM70, TET12]. In the forward pass, standard iLQR computes an updated nominal control sequence via $\bar{\mathbf{u}}_k \leftarrow \bar{\mathbf{u}}_k + \mathbf{l}_k + L_k \delta \mathbf{x}_k$. Instead we can weight \mathbf{l}_k with a scalar $\alpha \in [0, 1]$ for which we perform line search. This results in increased stability (as with standard line search for step size determination in nonlinear optimization) and possibly faster convergence. When α is close to zero, or alternative conditions (such as expected improvement being small) are met, we terminate. For a further discussion of this approach, we refer the reader to [TET12], which also features a discussion of step size determination in the DDP literature.

Iterative LQR and DDP rely on minimizing a second order approximation of the cost-to-go perturbation. However, we do not have any guarantees on the convexity of $Q(\delta \mathbf{x}_k, \delta \mathbf{u}_k)$ for arbitrary cost functions. Note that DDP is performing a Newton step [LS92] (iLQR is performing a Newton step with an approximation of the Hessian) via decomposing the optimization problem over controls into N smaller optimization problems. As such, standard approaches from Newton methods for regularization have been applied, such as replacing $Q_{\mathbf{u}\mathbf{u},k}$ with $Q_{\mathbf{u}\mathbf{u},k} + \mu I$, which is convex for sufficiently large μ . Alternative approaches have been explored in [TET12, TMT14], based on regularizing the quadratic term in the approximate cost-to-go.

Both iLQR and DDP are local methods. Full dynamic programming approaches yield globally optimal feedback policies. In contrast, iLQR and DDP yield nominal trajectories and local stabilizing controllers. However, these local controllers are often sufficient for tracking the trajectory. As they are local method, choice of initial control sequence is important, and poor choice may result in poor convergence. Additionally, we have not considered constraints on either state or action in the derivation of iLQR or DDP. This is currently an active area of research [XLH17, TMT14, GB17].

2.5 Stochastic Optimal Control

We began this chapter by discussing the generic optimal control problem that we consider in this course, and took a small detour to discuss the special case of discrete-time LQR. We will now generalize the previously stated optimal control problem to consider the stochastic case, in which random noise is added to the state transitions and/or observations. In the remainder of this chapter we will discuss the discrete time case.

2.5.1 The Stochastic Optimal Control Problem

We consider systems of the form

$$\mathbf{x}_{k+1} = f_k(\mathbf{x}_k, \mathbf{u}_k, \boldsymbol{\omega}_k), \quad k = 0, \dots, N-1 \quad (109)$$

where $\boldsymbol{\omega}_k \sim p(\cdot \mid \mathbf{x}_k, \mathbf{u}_k)$ is the disturbance or noise. We write the expected cost under policy $\pi = \{\pi_0, \dots, \pi_{N-1}\}$ as

$$J_\pi(\mathbf{x}_0) = \mathbb{E}_{\boldsymbol{\omega}_{0:N-1}} \left[c_N(\mathbf{x}_N) + \sum_{k=0}^{N-1} c_k(\mathbf{x}_k, \pi_k(\mathbf{x}_k), \boldsymbol{\omega}_k) \right]. \quad (110)$$

Then, the stochastic control problem we wish to solve is to find

$$J^*(\mathbf{x}_0) = \min_{\pi} J_\pi(\mathbf{x}_0). \quad (111)$$

In contrast to the deterministic optimal control problem, we are specifically interested in finding the optimal closed-loop *policy* in the stochastic case. Closed-loop policies can achieve lower cost than open-loop action sequences, as they take advantage of current state information in their action selection. We have also, importantly, assumed the cost is strictly additive. Finally, we have assumed we are operating in a *risk-neutral* setting. This risk neutrality corresponds to the expectation in (110), which could be replaced by any *risk metric*, which maps a distribution over outcomes to a scalar.

2.5.2 The Principle of Optimality: Stochastic Case

We can now state the principle of optimality for the stochastic optimal control problem. Note that this is a strict generalization of the deterministic case.

Theorem 2.2 (Discrete-time Principle of Optimality: Stochastic Case). *Let $\pi^* = (\pi_0^*, \dots, \pi_{N-1}^*)$ be an optimal policy. Assume state \mathbf{x}_k is reachable. Consider the tail subproblem*

$$\mathbb{E}_{\boldsymbol{\omega}_{i:N-1}} \left[c_T(\mathbf{x}_N) + \sum_{k=i}^{N-1} c_k(\mathbf{x}_k, \pi(\mathbf{x}_k), \boldsymbol{\omega}_k) \right]. \quad (112)$$

Then the truncated policy $(\pi_i^, \dots, \pi_{N-1}^*)$ is optimal for the subproblem.*

The intuition behind the stochastic principle of optimality is effectively the same as for the deterministic, and the proof is also based on decomposition of the total cost into two cost terms. This is possible due to the linearity of expectation. Stated simply, if a better policy existed for the tail problem, this would imply π^* is suboptimal.

The stochastic version of the principle of optimality leads to a concomitant dynamic programming algorithm, which takes the form

$$J_N(\mathbf{x}_N) = c_T(\mathbf{x}_N) \quad (113)$$

$$J_k(\mathbf{x}_k) = \min_{\mathbf{u}_k \in \mathcal{U}(\mathbf{x}_k)} \mathbb{E}_{\boldsymbol{\omega}_k} [c_k(\mathbf{x}_k, \mathbf{u}_k, \boldsymbol{\omega}_k) + J_{k+1}(f(\mathbf{x}_k, \mathbf{u}_k, \boldsymbol{\omega}_k))] \quad (114)$$

and the optimal policy is

$$\pi_k^*(\mathbf{x}_k) = \operatorname{argmin}_{\mathbf{u}_k \in \mathcal{U}(\mathbf{x}_k)} \mathbb{E}_{\boldsymbol{\omega}_k} [c_k(\mathbf{x}_k, \mathbf{u}_k, \boldsymbol{\omega}_k) + J_{k+1}(f(\mathbf{x}_k, \mathbf{u}_k, \boldsymbol{\omega}_k))]. \quad (115)$$

2.6 Stochastic LQR and LQG

2.6.1 LQR with Additive Noise

We will first address the stochastic LQR problem. The system dynamics are

$$\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k + \boldsymbol{\omega}_k \quad (116)$$

where $\boldsymbol{\omega}_k \sim \mathcal{N}(0, \Sigma_\omega)$, and the stage-wise cost is

$$c_k(\mathbf{x}_k, \mathbf{u}_k) = \frac{1}{2}(\mathbf{x}_k^T Q_k \mathbf{x}_k + \mathbf{u}_k^T R_k \mathbf{u}_k). \quad (117)$$

with terminal cost $\frac{1}{2} \mathbf{x}_N^T Q_N \mathbf{x}_N$. We wish to minimize the expected cost. The cost-to-go, as in the deterministic case, will be quadratic. Thus, plugging into the Bellman equation, we have

$$\begin{aligned} J_k^*(\mathbf{x}_k) &= \min_{\mathbf{u}_k \in \mathbb{R}^m} \mathbb{E} \left[\frac{1}{2} \mathbf{x}_k^T Q_k \mathbf{x}_k + \frac{1}{2} \mathbf{u}_k^T R_k \mathbf{u}_k \right. \\ &\quad \left. + \frac{1}{2} (A_k \mathbf{x}_k + B_k \mathbf{u}_k + \boldsymbol{\omega}_k)^T V_{k+1} (A_k \mathbf{x}_k + B_k \mathbf{u}_k + \boldsymbol{\omega}_k) \right] \end{aligned} \quad (118)$$

$$\begin{aligned} &= \min_{\mathbf{u}_k \in \mathbb{R}^m} \left\{ \frac{1}{2} \mathbf{x}_k^T Q_k \mathbf{x}_k + \frac{1}{2} \mathbf{u}_k^T R_k \mathbf{u}_k \right. \\ &\quad \left. + \mathbb{E} \left[\frac{1}{2} (A_k \mathbf{x}_k + B_k \mathbf{u}_k + \boldsymbol{\omega}_k)^T V_{k+1} (A_k \mathbf{x}_k + B_k \mathbf{u}_k + \boldsymbol{\omega}_k) \right] \right\}. \end{aligned} \quad (119)$$

Following the same minimization procedure as for LQR, we see that the policy is identical to that in Section 2.3. The Riccati equation, however, is

$$V_k = Q_k + L_k^T R_k L_k + \mathbb{E}[(A_k + B_k L_k + \boldsymbol{\omega}_k)^T V_{k+1} (A_k + B_k L_k + \boldsymbol{\omega}_k)] \quad (120)$$

$$= Q_k + L_k^T R_k L_k + (A_k + B_k L_k)^T V_{k+1} (A_k + B_k L_k) + \text{tr}(\Sigma_\omega V_{k+1}) \quad (121)$$

where $\text{tr}(\cdot)$ denotes the trace. The equality between (120) and (121) holds as

$$\mathbb{E}[(A_k + B_k L_k)^T V_{k+1} \boldsymbol{\omega}_k] = 0 \quad (122)$$

for zero-mean $\boldsymbol{\omega}_k$, and $\mathbb{E}[\boldsymbol{\omega}_k^T V_{k+1} \boldsymbol{\omega}_k] = \text{tr}(\Sigma_\omega V_{k+1})$. Note that this is identical to the deterministic case, other than the additive trace term at the end.

2.6.2 Problems with Imperfect State Information

We now consider the case in which direct, perfect state information isn't available. Instead, we have a noise-corrupted measurements

$$\mathbf{z}_0 = h_0(\mathbf{x}_0, \boldsymbol{\nu}_0) \quad (123)$$

$$\mathbf{z}_k = h(\mathbf{x}_k, \boldsymbol{\nu}_k), \quad k = 0, \dots, N-1. \quad (124)$$

The observation disturbance is characterized by distribution

$$p(\cdot \mid \mathbf{x}_k, \dots, \mathbf{x}_0, \mathbf{u}_{k-1}, \dots, \mathbf{u}_0, \boldsymbol{\omega}_{k-1}, \dots, \boldsymbol{\omega}_0, \boldsymbol{\nu}_{k-1}, \dots, \boldsymbol{\nu}_0) \quad (125)$$

and the initial state \mathbf{x}_0 is distributed according to $p(\mathbf{x}_0)$. We will define the information vector as

$$\mathbf{i}_k = [\mathbf{z}_0^T, \dots, \mathbf{z}_k^T, \mathbf{u}_0^T, \dots, \mathbf{u}_{k-1}^T]^T. \quad (126)$$

Armed with this, we will consider *admissible* policies $\pi(\mathbf{i}_k) \in \mathcal{U}_k$, which implies they are *causal* — they do not rely on information only available in the future. The goal of the control problem, then is to minimize

$$\mathbb{E}_{\mathbf{x}_0, \boldsymbol{\omega}_{0:N-1}, \boldsymbol{\nu}_{0:N-1}} \left[c_T(\mathbf{x}_N) + \sum_{k=0}^{N-1} c(\mathbf{x}_k, \pi(\mathbf{i}_k), \boldsymbol{\omega}_k) \right]. \quad (127)$$

Treating this problem as a perfect state information problem and attempting to solve directly results in several problems. In addition to the standard difficulties associated with applying DP to generic problems, \mathbf{i}_k has expanding dimension over the length of the problem. Alternatively, we may reason in terms of sufficient statistics: quantities that summarize all of the informational content of \mathbf{i}_k . For example, if we can construct a conditional distribution over state, $p(\mathbf{x}_k \mid \mathbf{i}_k)$, we can design a policy of the form $\pi_k(p(\mathbf{x}_k \mid \mathbf{i}_k))$.

2.6.3 LQG and the Separation Principle

Given the generic imperfect state information control problem, we will address an extremely important special case. We will again consider quadratic cost of the form

$$\frac{1}{2} \mathbb{E} \left[\mathbf{x}_N^T Q_N \mathbf{x}_N + \sum_{k=0}^{N-1} \mathbf{x}_k^T Q_k \mathbf{x}_k + \mathbf{u}_k^T R_k \mathbf{u}_k \right] \quad (128)$$

subject to dynamics

$$\mathbf{x}_{k+1} = A_k \mathbf{x}_k + B_k \mathbf{u}_k + \boldsymbol{\omega}_k \quad (129)$$

and measurements

$$\mathbf{z}_k = C_k \mathbf{x}_k + \boldsymbol{\nu}_k. \quad (130)$$

The initial state \mathbf{x}_0 , and process and measurements noise $\boldsymbol{\omega}_{0:N-1}, \boldsymbol{\nu}_{0:N-1}$ independent, zero-mean Gaussians. We will write the covariance of $\boldsymbol{\omega}_k$ and $\boldsymbol{\nu}_k$ as $\Sigma_{\boldsymbol{\omega},k}$ and $\Sigma_{\boldsymbol{\nu},k}$, respectively. We will write $\Sigma_{\mathbf{x},0}$ for the covariance of \mathbf{x}_0 .

We will skip the lengthy derivation, but the optimal control policy takes the form

$$\mathbf{u}_k^* = L_k \hat{\mathbf{x}}_k \quad (131)$$

where $\hat{\mathbf{x}}_k$ is the state estimate from the Kalman filter, and L_k is the standard gain from LQR. Note that the Kalman filter maintains a Gaussian estimate of the state, and so the sufficient statistics are the mean and variance. However, the policy depends only on the mean. The policy can be designed as if access to perfect state information is available, while the estimator can be designed without considering the controller. This is known as the *separation principle*.

2.7 Further Reading

A comprehensive coverage of linear quadratic methods for optimal control is Anderson and Moore [AM07]. LQG is covered in discrete time in [Ber12]. The original, comprehensive reference on DDP is [JM70], but a large body of literature on the method has been produced since then. The original papers on iLQR are [TL05, LT04].

3 The HJB and HJI Equations

In this section, we will extend the ideas of dynamic programming to the continuous time setting. Restating the continuous time optimal control problem, we assume dynamics

$$\dot{\mathbf{x}}(t) = f(\mathbf{x}(t), \mathbf{u}(t), t) \quad (132)$$

and cost

$$J(\mathbf{x}(0)) = c_f(\mathbf{x}(t_f), t_f) + \int_0^{t_f} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau. \quad (133)$$

where t_f is fixed.

3.1 The Principle of Optimality in Continuous Time

3.1.1 Hamilton-Jacobi-Bellman

As in the discrete time principle of optimality, consider the tail problem

$$J(\mathbf{x}(t), \{\mathbf{u}(\tau)\}_{\tau=t}^{t_f}, t) = c_f(\mathbf{x}(t_f), t_f) + \int_t^{t_f} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau \quad (134)$$

where $t \leq t_f$ and $\mathbf{x}(t)$ is an admissible state value. The optimal solution to this tail problem comes from the functional minimization

$$J^*(\mathbf{x}(t), t) = \min_{\{\mathbf{u}(\tau)\}_{\tau=t}^{t_f}} \left\{ c_f(\mathbf{x}(t_f), t_f) + \int_t^{t_f} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau \right\}. \quad (135)$$

Note, then, that due to the additivity of cost we can split the problem up over time,

$$J^*(\mathbf{x}(t), t) = \min_{\{\mathbf{u}(\tau)\}_{\tau=t}^{t_f}} \left\{ \int_t^{t+\Delta t} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau + c_f(\mathbf{x}(t_f), t_f) + \int_{t+\Delta t}^{t_f} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau \right\} \quad (136)$$

which by applying the principle of optimality to the tail cost,

$$J^*(\mathbf{x}(t), t) = \min_{\{\mathbf{u}(\tau)\}_{\tau=t}^{t+\Delta t}} \left\{ \int_t^{t+\Delta t} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \tau) d\tau + J^*(\mathbf{x}(t+\Delta t), t+\Delta t) \right\}. \quad (137)$$

Let $J_t^*(\mathbf{x}(t), t) = \nabla_t J^*(\mathbf{x}(t), t)$ and $J_x^*(\mathbf{x}(t), t) = \nabla_x J^*(\mathbf{x}(t), t)$. Taylor expanding, we have

$$J^*(\mathbf{x}(t), t) = \min_{\{\mathbf{u}(\tau)\}_{\tau=t}^{t+\Delta t}} \{c(\mathbf{x}(t), \mathbf{u}(t), t)\Delta t + J^*(\mathbf{x}(t), t) + (J_t^*(\mathbf{x}(t), t))\Delta t + (J_x^*(\mathbf{x}(t), t))^T(\mathbf{x}(t + \Delta t) - \mathbf{x}(t)) + o(\Delta t)\} \quad (138)$$

for small Δt . The first term is a result of Taylor expanding the integral and applying the fundamental theorem of calculus. Note that we can pull $J^*(\mathbf{x}(t), t)$ out of the minimization over cost, as this quantity will not vary under different choices of future actions. Dividing through by Δt and taking the limit $\Delta t \rightarrow 0$, we obtain the *Hamilton-Jacobi-Bellman* equation

$$0 = J_t^*(\mathbf{x}(t), t) + \min_{\mathbf{u}(t)} \{c(\mathbf{x}(t), \mathbf{u}(t), t) + (J_x^*(\mathbf{x}(t), t))^T f(\mathbf{x}(t), \mathbf{u}(t), t)\} \quad (139)$$

with terminal condition

$$J^*(\mathbf{x}(t_f), t_f) = c_f(\mathbf{x}(t_f), t_f). \quad (140)$$

For convenience, we will define the Hamiltonian

$$\mathcal{H}(\mathbf{x}(t), \mathbf{u}(t), J_x^*, t) := c(\mathbf{x}(t), \mathbf{u}(t), t) + (J_x^*(\mathbf{x}(t), t))^T f(\mathbf{x}(t), \mathbf{u}(t), t) \quad (141)$$

which allow us to compactly write the HJB equation as

$$0 = J_t^*(\mathbf{x}(t), t) + \min_{\mathbf{u}(t)} \{\mathcal{H}(\mathbf{x}(t), \mathbf{u}(t), J_x^*, t)\}. \quad (142)$$

The HJB equation is a partial differential equation that, for cost-to-go $J^*(\mathbf{x}(t), t)$, will satisfy all time-state pairs $(\mathbf{x}(t), t)$. The previous informal derivation assumed differentiability of $J^*(\mathbf{x}(t), t)$, which we do not know a priori. This assumption is rectified by the following theorem on solutions to the HJB equation.

Theorem 3.1 (Sufficiency Theorem). *Suppose $V(\mathbf{x}, t)$ is a solution to the HJB equation, that $V(\mathbf{x}, t)$ is C^1 in \mathbf{x} and t , and that*

$$0 = V_t(\mathbf{x}, t) + \min_{\mathbf{u} \in \mathcal{U}} \{c(\mathbf{x}, \mathbf{u}, t) + (V_x(\mathbf{x}, t))^T f(\mathbf{x}, \mathbf{u}, t)\} \\ V(\mathbf{x}, t_f) = c_f(\mathbf{x}, t_f) \quad \forall \mathbf{x}$$

Suppose also that $\pi^(\mathbf{x}, t)$ attains the minimum in this equation for all t and \mathbf{x} . Let $\{\mathbf{x}^*(t) \mid t \in [t_0, t_f]\}$ be the state trajectory obtained from the given initial condition $\mathbf{x}(0)$ when the control trajectory $\mathbf{u}^*(t) = \pi^*(\mathbf{x}^*(t), t), t \in [t_0, t_f]$ is used. Then V is equal to the optimal cost-to-go function, i.e.,*

$$V(\mathbf{x}, t) = J^*(\mathbf{x}, t) \quad \forall \mathbf{x}, t. \quad (143)$$

Furthermore, the control trajectory $\{\mathbf{u}^(t) \mid t \in [t_0, t_f]\}$ is optimal..*

Proof. [Ber12], Volume 1, Section 7.2. □

3.1.2 Continuous-Time LQR

As a useful result of the HJB equations, we will derive LQR in continuous time. We aim to minimize

$$J(\mathbf{x}(0)) = \frac{1}{2} \mathbf{x}^T(t_f) Q_f \mathbf{x}(t_f) + \frac{1}{2} \int_0^{t_f} \mathbf{x}^T(t) Q(t) \mathbf{x}(t) + \mathbf{u}^T(t) R(t) \mathbf{u}(t) dt \quad (144)$$

subject to dynamics

$$\dot{\mathbf{x}}(t) = A(t) \mathbf{x}(t) + B(t) \mathbf{u}(t). \quad (145)$$

As in discrete LQR, we will assume $Q_f, Q(t)$ are positive semidefinite, and $R(t)$ is positive definite. We will also assume t_f is fixed, and the state and action are unconstrained.

We will write the Hamiltonian,

$$\mathcal{H} = \frac{1}{2} \mathbf{x}^T(t) Q(t) \mathbf{x}(t) + \frac{1}{2} \mathbf{u}^T(t) R(t) \mathbf{u}(t) + J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T (A(t) \mathbf{x}(t) + B(t) \mathbf{u}(t)) \quad (146)$$

which yields necessary optimality conditions

$$0 = \nabla_{\mathbf{u}} \mathcal{H} = R(t) \mathbf{u}(t) + B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t). \quad (147)$$

Since $\nabla_{\mathbf{u}\mathbf{u}}^2 \mathcal{H} = R(t) > 0$, the control that satisfies the necessary conditions is the global minimizer. Rearranging, we have

$$\mathbf{u}^*(t) = -R^{-1}(t) B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t) \quad (148)$$

which we can plug back into the Hamiltonian to yield

$$\begin{aligned} \mathcal{H} &= \frac{1}{2} \mathbf{x}^T(t) Q(t) \mathbf{x}(t) + \frac{1}{2} J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T B(t) R^{-1}(t) B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t) \\ &\quad + J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T A(t) \mathbf{x}(t) - J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T B(t) R^{-1}(t) B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t) \\ &= \frac{1}{2} \mathbf{x}^T(t) Q(t) \mathbf{x}(t) - \frac{1}{2} J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T B(t) R^{-1}(t) B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t) + J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T A(t) \mathbf{x}(t). \end{aligned} \quad (149)$$

This gives the HJB equation

$$\begin{aligned} 0 &= J_t^*(\mathbf{x}(t), t) + \frac{1}{2} \mathbf{x}^T(t) Q(t) \mathbf{x}(t) - \frac{1}{2} J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T B(t) R^{-1}(t) B^T(t) J_{\mathbf{x}}^*(\mathbf{x}(t), t) \\ &\quad + J_{\mathbf{x}}^*(\mathbf{x}(t), t)^T A(t) \mathbf{x}(t) \end{aligned} \quad (151)$$

with boundary condition

$$J^*(\mathbf{x}(t_f), t_f) = \frac{1}{2} \mathbf{x}^T(t_f) Q_f \mathbf{x}(t_f). \quad (152)$$

It may appear as if we are stuck here, as this form of the HJB doesn't immediately yield $J^*(\mathbf{x}(t), t)$. Armed with the knowledge that the discrete time LQR problem has a quadratic cost-to-go, we will cross our fingers and guess a solution of the form

$$J^*(\mathbf{x}(t), t) = \frac{1}{2} \mathbf{x}^T(t) V(t) \mathbf{x}(t). \quad (153)$$

Substituting, we have

$$0 = \frac{1}{2}\mathbf{x}^T(t)\dot{V}(t)\mathbf{x}(t) + \frac{1}{2}\mathbf{x}^T(t)Q(t)\mathbf{x}(t) - \frac{1}{2}\mathbf{x}^T(t)V(t)B(t)R^{-1}(t)B^T(t)V(t)\mathbf{x}(t) + \mathbf{x}^T(t)V(t)A(t)\mathbf{x}(t) \quad (154)$$

Note that we will decompose

$$\mathbf{x}^T(t)V(t)A(t)\mathbf{x}(t) = \frac{1}{2}\mathbf{x}^T(t)V(t)A(t)\mathbf{x}(t) + \frac{1}{2}\mathbf{x}^T(t)A^T(t)V(t)\mathbf{x}(t) \quad (155)$$

which yields

$$0 = \frac{1}{2}\mathbf{x}^T(t) \left(\dot{V}(t) + Q(t) - V(t)B(t)R^{-1}(t)B^T(t)V(t) + V(t)A(t) + A^T(t)V(t) \right) \mathbf{x}(t). \quad (156)$$

This equation must hold for all $\mathbf{x}(t)$, so

$$-\dot{V}(t) = Q(t) - V(t)B(t)R^{-1}(t)B^T(t)V(t) + V(t)A(t) + A^T(t)V(t) \quad (157)$$

with boundary condition $V(t_f) = Q_f$.

Therefore, the HJB PDE has been reduced to a set of matrix ordinary differential equations (the Riccati equation). This is integrated backwards in time to find the full control policy as a function of time. Once we have found $V(t)$, the control policy is

$$\mathbf{u}^*(t) = -R^{-1}(t)B^T(t)V(t)\mathbf{x}(t). \quad (158)$$

Similarly to the discrete case, the feedback gains tend toward constant in the limit of the infinite horizon problem, under some technical assumptions.

3.2 Differential Games

We have so far addressed the case in which we aim to solve the optimal control problem for a single agent. We will now consider an adversarial game setting, in which there exists another player that aims to maximally harm the first agent. In particular, we will consider zero sum games in which the second agent aims to maximize the cost of the first agent. While the differential game setting is not restricted to this case — agents may have separate cost functions that partially interfere or aid each other — the zero-sum case lends itself to useful analytical tools.

3.2.1 Differential Games and Information Patterns

We consider the two player differential game with dynamics

$$\dot{\mathbf{x}}(t) = f(\mathbf{x}(t), \mathbf{u}(t), \mathbf{d}(t)) \quad (159)$$

where the first player takes action $\mathbf{u}(t)$ at time t , and the second player takes action $\mathbf{d}(t)$. The state $\mathbf{x}(t)$ is the joint state of both players. We write the cost as

$$J(\mathbf{x}(t)) = c_f(\mathbf{x}(0)) + \int_t^0 c(\mathbf{x}(\tau), \mathbf{u}(\tau), \mathbf{d}(\tau)) d\tau \quad (160)$$

which the first agent aims to maximize, and the second agent aims to minimize.

To fully specify the differential game, we must specify what each agent knows, and when. This is referred to as the *information pattern* of the game. In addition to capturing the knowledge of the state available to each agent, the information pattern also captures the knowledge of each other agents' strategies available to each agent.

3.2.2 Hamilton-Jacobi-Isaacs

The key idea in building the multi-agent equivalent of the HJB equation will again be to apply the principle of optimality. We consider the information pattern in which the adversary has access to the instantaneous control action of the first agent, so the cost takes the form

$$J(\mathbf{x}(t), t) = \min_{\Gamma(\mathbf{u})(\cdot)} \max_{\mathbf{u}(\cdot)} \left\{ \int_t^0 c(\mathbf{x}(\tau), \mathbf{u}(\tau), \mathbf{d}(\tau)) d\tau + c_f(\mathbf{x}(0)) \right\}. \quad (161)$$

Applying the dynamic programming principle, we have

$$J(\mathbf{x}(t), t) = \min_{\Gamma(\mathbf{u})(\cdot)} \max_{\mathbf{u}(\cdot)} \left\{ \int_t^{t+\Delta t} c(\mathbf{x}(\tau), \mathbf{u}(\tau), \mathbf{d}(\tau)) d\tau + J(\mathbf{x}(t+\Delta t), t+\Delta t) \right\}. \quad (162)$$

We can take the same strategy as with the informal derivation of the HJB equation, and Taylor expand both terms to yield

$$\begin{aligned} J(\mathbf{x}(t), t) = \min_{\Gamma(\mathbf{u})(\cdot)} \max_{\mathbf{u}(\cdot)} \{ & c(\mathbf{x}(\tau), \mathbf{u}(\tau), \mathbf{d}(\tau)) \Delta t + J(\mathbf{x}(t), t) \\ & + (J_{\mathbf{x}}(\mathbf{x}(t), t))^T f(\mathbf{x}(t), \mathbf{u}(t), \mathbf{d}(t)) \Delta t + J_t(\mathbf{x}(t), t) \Delta t \}. \end{aligned} \quad (163)$$

Note that we are optimizing over instantaneous actions, and so we optimizing over finite dimensional quantities as opposed to functions. Dividing through by Δt and removing redundant terms, we get the *Hamilton-Jacobi-Isaacs* (HJI) equation

$$0 = J_t(\mathbf{x}, t) + \max_{\mathbf{u}} \min_{\mathbf{d}} \{ c(\mathbf{x}, \mathbf{u}, \mathbf{d}) + (J_{\mathbf{x}}(\mathbf{x}, \mathbf{u}, \mathbf{d}))^T f(\mathbf{x}, \mathbf{d}, \mathbf{u}) \} \quad (164)$$

with boundary condition

$$J(\mathbf{x}, 0) = c_f(\mathbf{x}). \quad (165)$$

Note that we have switched the order of the min/max.

3.2.3 Reachability

Differential games have applications in multi-agent modeling (both in the context of autonomous systems engineering and, e.g., economics and operations research). One concrete application in engineering is reachability analysis. In this setting, an agent aims to compute the set of states in which there exists a policy that either avoids a target set or enters a target set, subject to adversarial disturbances. The former case, in which we would like to avoid a target set, is useful for safety verification. If we are able to, even in the worst case, guarantee e.g. collision avoidance, we have guarantees on safety (subject of course to our system assumptions). The latter case is useful for task satisfaction. For example, we would like a quadrotor to reach a set of safe hovering poses, even under adversarial disturbances. Finding the backward reachable set in this case would find all states such that there exists a policy that succeeds in reaching the target set.

More concretely, the first case aims to find a set

$$\mathcal{A}(t) = \{\bar{\mathbf{x}} : \exists \Gamma(\mathbf{u})(\cdot), \forall \mathbf{u}(\cdot), \dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u}, \mathbf{d}), \mathbf{x}(t) = \bar{\mathbf{x}}, \mathbf{x}(0) \in \mathcal{T}\} \quad (166)$$

where \mathcal{T} is the unsafe set which we aim to avoid. Breaking this down, $\mathcal{A}(t)$ is the set of states at time t such that there exists $\Gamma(\mathbf{u})$ that maps action \mathbf{u} to a disturbance such that, following the dynamics induced by the disturbance and the action sequence, the state is in \mathcal{T} at time 0 (note that we are considering $t \leq 0$).

The second case aims to find a set

$$\mathcal{R}(t) = \{\bar{\mathbf{x}} : \forall \Gamma(\mathbf{u})(\cdot), \exists \mathbf{u}(\cdot), \dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u}, \mathbf{d}), \mathbf{x}(t) = \bar{\mathbf{x}}, \mathbf{x}(0) \in \mathcal{T}\}, \quad (167)$$

where in this case \mathcal{T} is the set that we wish to reach. In this setting, we wish to find all states that, no matter what strategy the disturbance takes, there exist control actions that can steer the system to the goal state. Because the disturbance is adversarial (we reason over all adversary strategies), this is an extremely conservative form of safety analysis.

Computation of the backward reachable set results from solving a differential *game of kind* in which the outcome is Boolean (i.e. whether or not $\mathbf{x}(0) \in \mathcal{T}$). This boolean outcome can be encoded by removing the running cost and choosing a particular form for the final cost. In particular, we can choose a final cost where

$$\mathbf{x} \in \mathcal{T} \iff c_f(\mathbf{x}) \leq 0. \quad (168)$$

As a result, the agent should aim to maximize c_f to avoid \mathcal{T} , whereas the disturbance should aim to minimize it. The two settings then take the following forms:

- Set avoidance: $J(\mathbf{x}, t) = \min_{\Gamma(\mathbf{u})} \max_{\mathbf{u}} c_f(\mathbf{x}(0))$
- Set reaching: $J(\mathbf{x}, t) = \max_{\Gamma(\mathbf{u})} \min_{\mathbf{u}} c_f(\mathbf{x}(0))$

Sets vs. Tubes. We have so far considered avoidance or reachability problems for which we care about set membership at time $t = 0$. However, for something like collision avoidance, we would like to stay collision free at every time as opposed to a particular time. *Backward reachable sets* capture the case in which only the final time set membership matters, and states for times $t < 0$ do not matter. *Backward reachable tubes* capture the entire time duration of the problem. Any state that passes through the target at any time in the problem duration is included. This yields a modified value function of the form

$$J(\mathbf{x}, t) = \min_{\Gamma(\mathbf{u})} \max_{\mathbf{u}} \min_{\tau \in [t, 0]} c_f(\mathbf{x}(\tau)). \quad (169)$$

If the target set membership holds at any time τ' , then $\min_{\tau \in [t, 0]} c_f(\mathbf{x}(\tau)) \leq c_f(\mathbf{x}(\tau')) \leq 0$.

3.3 Further Reading

Our coverage of reachability analysis is based on the [MBT05], which is an important early work in the field, in addition to being a relatively comprehensive coverage of the method. For a review of differential games with a (slight) emphasis on economics and management science, we refer the reader to [Bre10]. For a review of HJB and continuous time LQR, we refer the reader to [Ber12] and [Kir12].

References

- [AM07] Brian DO Anderson and John B Moore. *Optimal control: linear quadratic methods*. Courier Corporation, 2007.
- [Ber12] Dimitri P Bertsekas. *Dynamic programming and optimal control*. Number 1. 4 edition, 2012.
- [Ber16] Dimitri P Bertsekas. *Nonlinear programming*. Athena Scientific, 2016.
- [Bre10] Alberto Bressan. Noncooperative differential games. a tutorial. 2010.
- [BT97] Dimitris Bertsimas and John N Tsitsiklis. *Introduction to linear optimization*. Athena Scientific, 1997.
- [BV04] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [GB17] Markus Giffthaler and Jonas Buchli. A projection approach to equality constrained iterative linear quadratic optimal control. In *IEEE-RAS International Conference on Humanoid Robotics (Humanoids)*, 2017.
- [JM70] David Jacobson and David Mayne. *Differential Dynamic Programming*. Elsevier, 1970.

- [Kir12] Donald E Kirk. *Optimal control theory: an introduction*. Courier Corporation, 2012.
- [Kol08] Zico Kolter. Convex optimization overview. *CS 229 Lecture Notes*, 2008.
- [LK14] Sergey Levine and Vladlen Koltun. Learning complex neural network policies with trajectory optimization. In *International Conference on Machine Learning (ICML)*, 2014.
- [LS92] Li-zhi Liao and Christine A Shoemaker. Advantages of differential dynamic programming over newton’s method for discrete-time optimal control problems. Technical report, Cornell University, 1992.
- [LT04] Weiwei Li and Emanuel Todorov. Iterative linear quadratic regulator design for nonlinear biological movement systems. In *International Conference on Informatics in Control, Automation, and Robotics*, 2004.
- [MBT05] Ian M Mitchell, Alexandre M Bayen, and Claire J Tomlin. A time-dependent hamilton-jacobi formulation of reachable sets for continuous dynamic games. *IEEE Transactions on Automatic Control*, 2005.
- [Pow12] Warren B Powell. AI, OR and control theory: A rosetta stone for stochastic optimization. 2012.
- [TET12] Yuval Tassa, Tom Erez, and Emanuel Todorov. Synthesis and stabilization of complex behaviors through online trajectory optimization. In *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2012.
- [TL05] Emanuel Todorov and Weiwei Li. A generalized iterative LQG method for locally-optimal feedback control of constrained nonlinear stochastic systems. In *American Control Conference (ACC)*, 2005.
- [TMT14] Yuval Tassa, Nicolas Mansard, and Emo Todorov. Control-limited differential dynamic programming. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2014.
- [XLH17] Zhaoming Xie, C Karen Liu, and Kris Hauser. Differential dynamic programming with nonlinear constraints. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2017.