

AA203: Optimal and Learning-based Control

Course Notes

James Harrison

April 8, 2019

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 $f(\mathbf{x}^*) = 0$, 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].

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.

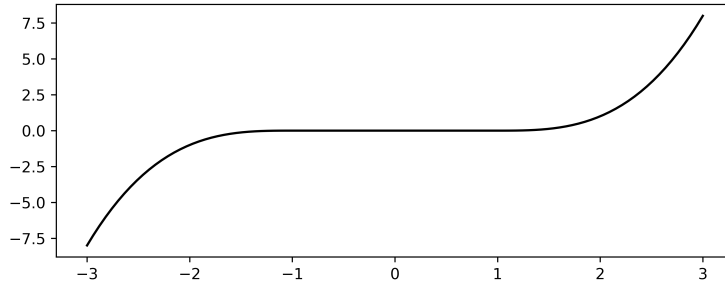


Figure 1: An example of a function for which the necessary conditions of optimality are satisfied but the sufficient conditions are not.

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

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

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,

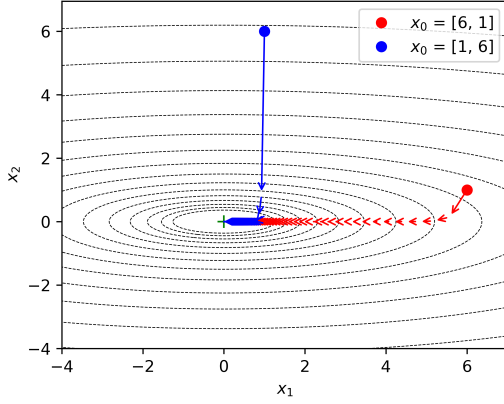
$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha (\nabla^2 f(\mathbf{x}^k))^{-1} \nabla f(\mathbf{x}^k). \quad (17)$$

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.

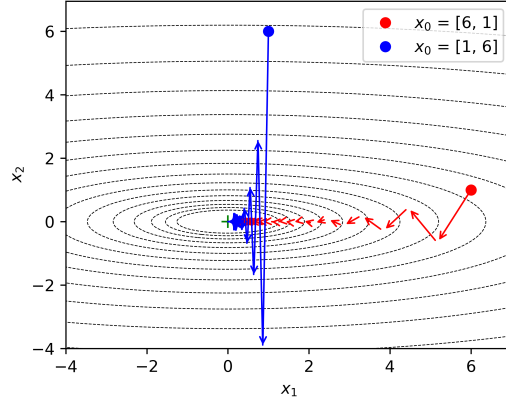
Diagonally scaled steepest descent, $D^k = \text{diag}(d_1^k, \dots, d_n^k)$. Have $d_i^k > 0 \forall i$. A popular choice is

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

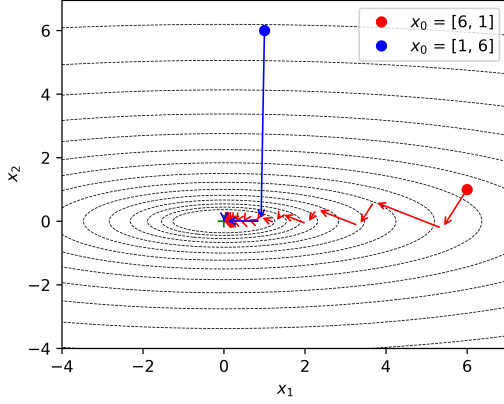
which is a diagonal approximation of the Hessian.



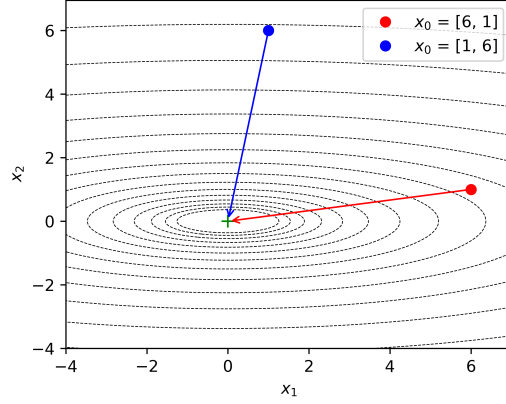
(a) Steepest descent, small fixed step size.



(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.

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

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

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

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

Proof. See [Ber16] Section 3.1.1 and 3.1.2. \square

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

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

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].

References

- [Ber16] Dimitri P Bertsekas. *Nonlinear programming*. Athena Scientific, 2016.
- [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.
- [Kol08] Zico Kolter. Convex optimization overview. *CS 229 Lecture Notes*, 2008.