

Appendix C

Optimization Problems and Optimality Conditions

C.1 Unconstrained Optimization

The mathematical model of an (unconstrained) optimization problem can be generally described by a domain or constraint set \mathcal{D} in \mathbb{R}^n and an objective function $f : \mathcal{D} \rightarrow \mathbb{R}$ that maps an element of \mathcal{D} to a real value. The optimization problem seeks an optimal solution $\mathbf{x}_\star \in \mathcal{D}$ such that the value of f is minimized:

$$f(\mathbf{x}_\star) \leq f(\mathbf{x}), \quad \text{for all } \mathbf{x} \in \mathcal{D}.$$

In particular, if $\mathcal{D} = \mathbb{R}^n$, it is called an unconstrained optimization problem.

Definition C.1.1 (Local and Global Minima). *A variable \mathbf{x}_\star is a local minimum of f if there exists a neighborhood $\mathcal{B}_\epsilon(\mathbf{x}_\star) \doteq \{\mathbf{x} \in \mathcal{D} \mid \|\mathbf{x} - \mathbf{x}_\star\|_2 < \epsilon\}$ for some $\epsilon > 0$ such that*

$$f(\mathbf{x}_\star) \leq f(\mathbf{x}), \quad \text{for all } \mathbf{x} \in \mathcal{B}_\epsilon(\mathbf{x}_\star).$$

The variable \mathbf{x}_\star is a global minimum of f if $\mathcal{B}_\epsilon(\mathbf{x}_\star) = \mathcal{D}$. The above local and global minima are said to be strict if the corresponding inequalities are also strict for $\mathbf{x} \neq \mathbf{x}_\star$.

If the objective function f is differentiable, then conditions for the optimality can be expressed in terms of its derivatives. In particular, if \mathbf{x}_\star is a local minimum, then within a small neighborhood $\mathcal{B}_\epsilon(\mathbf{x}_\star)$, for any given vector $\mathbf{v} \in \mathbb{R}^n$, we have

$$f(\mathbf{x}_\star + t \cdot \mathbf{v}) \geq f(\mathbf{x}_\star)$$

for sufficiently small $t > 0$ such that $t \cdot \mathbf{v} \in \mathcal{B}_\epsilon(\mathbf{0})$. Hence we have

$$\lim_{t \rightarrow 0} \frac{f(\mathbf{x}_\star + t \cdot \mathbf{v}) - f(\mathbf{x}_\star)}{t} = \nabla f(\mathbf{x}_\star)^T \mathbf{v} \geq 0.$$

Notice that this must be true for both \mathbf{v} and $-\mathbf{v}$. Then for the inequality to hold for all $\mathbf{v} \in \mathbb{R}^n$, we must have

$$\nabla f(\mathbf{x}_\star) = \mathbf{0}. \quad (\text{C.1.1})$$

Definition C.1.2 (Stationary Point). *A point \mathbf{x}_\star that satisfies the condition $\nabla f(\mathbf{x}_\star) = \mathbf{0}$ is referred to as a stationary point of $f(\mathbf{x})$.*

If f is twice continuously differentiable and \mathbf{x}_\star is a stationary point with $\nabla f(\mathbf{x}_\star) = \mathbf{0}$, we have:

$$f(\mathbf{x}_\star + t \cdot \mathbf{v}) \approx f(\mathbf{x}_\star) + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{x}_\star) \mathbf{v} t^2 + o(t^2).$$

If \mathbf{x}_\star is a local minimum, we have

$$f(\mathbf{x}_\star + t \cdot \mathbf{v}) - f(\mathbf{x}_\star) \geq 0 \quad \Rightarrow \quad \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{x}_\star) \mathbf{v} t^2 \geq 0$$

for all $\mathbf{v} \in \mathbb{R}^n$. This implies the matrix $\nabla^2 f(\mathbf{x}_\star)$ is necessarily positive semi-definite, namely,

$$\nabla^2 f(\mathbf{x}_\star) \succeq 0. \quad (\text{C.1.2})$$

It is then not difficult to show the following sufficient condition for local minima:

Proposition C.1.3 (Second-Order Sufficient Optimality Condition). *Let $f : \mathcal{D} \rightarrow \mathbb{R}$ be twice continuously differentiable. If \mathbf{x}_\star satisfies the conditions*

$$\nabla f(\mathbf{x}_\star) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\mathbf{x}_\star) \succ 0,$$

Then \mathbf{x}_\star is a strict local minimum of $f(\mathbf{x})$.

In general, a local minimum is not necessarily a global minimum in the domain of $f(\mathbf{x})$. Therefore, the global minimum can be found by exhaustively comparing the values of f at all local minima. However, when the objective function f is convex, the following proposition shows that any local minimum is also a global minimum.

Proposition C.1.4 (Global Optimality of Convex Functions). *Let $f : \mathcal{D} \rightarrow \mathbb{R}$ be a convex function over convex set \mathcal{D} . Then*

1. *A local minimum of f is also a global minimum. Furthermore, if f is strictly convex, then the global minimum, if exists, is unique.*
2. *A point $\mathbf{x}_\star \in \mathcal{D}$ is a global minimum of f if $\mathbf{0} \in \partial f(\mathbf{x}_\star)$. In the case that f is differentiable, $\nabla f(\mathbf{x}_\star) = \mathbf{0}$ implies that \mathbf{x}_\star is a global minimum.*

Finally, we note that given an objective function f , a local minimum need not exist. For example, the simple scalar function $f(x) = x$ does not have a minimal value in the domain of real numbers as $\inf_{x \in \mathbb{R}} f(x) = -\infty$. Therefore, a sufficient condition for f to have at least one local minimum is that the set $\{f(\mathbf{x}) | \mathbf{x} \in \mathcal{D}\}$ is bounded below. Alternatively, according to the Weierstrass theorem, if f is continuous and the domain set $\mathcal{D} \subseteq \mathbb{R}^n$ is compact (i.e. closed and bounded), then f has at least one local minimum.

C.2 Constrained Optimization

In the previous section, the constraint set of the optimization problems is assumed to be rather general. However, in most optimization problems considered in this book, the constraints are formulated as equality or inequality conditions. For example, the domain $\mathcal{D} \subseteq \mathbb{R}^n$ of a polyhedron can be specified by a set of equality and inequality conditions. *Lagrange multipliers* are a set of supportive variables to facilitate the derivation of optimality conditions for such constrained optimization problems. Arguably, Lagrange multiplier theory is the most influential theory in constrained optimization. In duality theory that we will discuss in the next section, the same Lagrange multiplier variables are also called *dual variables*, which will play a central role as the optimization variables of the *dual problems*.

First, we consider the optimization problem with equality constraints:

$$\min f(\mathbf{x}) \quad \text{subj. to} \quad h_i(\mathbf{x}) = 0, \quad i = 1, \dots, m, \quad (\text{C.2.1})$$

where f and each h_i are assumed to be continuously differentiable.¹ Conveniently, we further assume the gradients of the equality conditions at a feasible solution \mathbf{x}

$$\nabla h_1(\mathbf{x}), \nabla h_2(\mathbf{x}), \dots, \nabla h_m(\mathbf{x})$$

are linearly independent. Such a solution \mathbf{x} is also called *regular*.

The optimality conditions for (C.2.1) can be conveniently derived in terms of the Lagrangian function $L : \mathbb{R}^{n+m} \rightarrow \mathbb{R}$ as

$$L(\mathbf{x}, \boldsymbol{\lambda}) \doteq f(\mathbf{x}) + \sum_{i=1}^m \lambda_i h_i(\mathbf{x}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{h}(\mathbf{x}) \rangle, \quad (\text{C.2.2})$$

where λ_i are the Lagrange multipliers for the equality conditions, and $\boldsymbol{\lambda} = [\lambda_1, \lambda_2, \dots, \lambda_m]^T \in \mathbb{R}^m$ is the corresponding Lagrange multiplier vector; and for brevity, we denote $\mathbf{h} = [h_1, h_2, \dots, h_m]^T$ as a map from \mathbb{R}^n to \mathbb{R}^m .

The basic Lagrange multiplier theory states the following necessary condition for the optimality of a regular solution.

¹In the main text, we need to generalize to cases when f is not differentiable.

Proposition C.2.1 (Necessary Conditions). *Let \mathbf{x}_\star be a local minimum of function $f(\mathbf{x})$ subject to $h_i(\mathbf{x}) = 0$, $i = 1, \dots, m$. Further assume \mathbf{x}_\star is regular. Then there exists a Lagrange multiplier vector $\boldsymbol{\lambda}_\star = (\lambda_{\star,1}, \lambda_{\star,2}, \dots, \lambda_{\star,m}) \in \mathbb{R}^m$, such that*

$$\begin{aligned}\nabla_{\mathbf{x}} L(\mathbf{x}_\star, \boldsymbol{\lambda}_\star) &= \nabla f(\mathbf{x}_\star) + \sum_{i=1}^m \lambda_{\star,i} \nabla h_i(\mathbf{x}_\star) = 0, \\ \nabla_{\boldsymbol{\lambda}} L(\mathbf{x}_\star, \boldsymbol{\lambda}_\star) &= \mathbf{h}(\mathbf{x}_\star) = 0.\end{aligned}\tag{C.2.3}$$

Furthermore, if f and \mathbf{h} are twice continuously differentiable, we have

$$\begin{aligned}\mathbf{v}^T \nabla_{\mathbf{x}\mathbf{x}}^2 L(\mathbf{x}_\star, \boldsymbol{\lambda}_\star) \mathbf{v}^T &= \mathbf{v}^T (\nabla^2 f(\mathbf{x}_\star) + \sum_{i=1}^m \lambda_{\star,i} \nabla^2 h_i(\mathbf{x}_\star)) \mathbf{v}^T \\ &\geq 0, \quad \forall \mathbf{v} : \mathbf{v}^T \nabla h_i(\mathbf{x}_\star) = 0, \quad i = 1, \dots, m.\end{aligned}\tag{C.2.4}$$

In (C.2.4), the conditions for vector $\mathbf{v} \in \mathbb{R}^n$ that satisfies $\mathbf{v}^T \nabla h_i(\mathbf{x}_\star) = 0$ can be understood as follows. If we consider a new point $\mathbf{x}' = \mathbf{x}_\star + \mathbf{v}$, due to the fact that $\mathbf{v}^T \nabla h_i(\mathbf{x}_\star) = 0$, the small variation \mathbf{v} will not change the value of $\mathbf{h}(\mathbf{x}') = 0$. Therefore, we can define

$$V(\mathbf{x}_\star) = \{\mathbf{v} \mid \mathbf{v}^T \nabla h_i(\mathbf{x}_\star) = 0, \quad i = 1, \dots, m\}.\tag{C.2.5}$$

as the *subspace of first-order feasible variations*.

In summary, the first-order condition (C.2.3) implies the gradient $\nabla f(\mathbf{x}_\star)$ is orthogonal to $V(\mathbf{x}_\star)$, which resembles the first-order condition $\nabla f(\mathbf{x}_\star) = 0$ in unconstrained optimization. The second-order condition (C.2.4) implies the Hessian of the Lagrangian function $L(\mathbf{x}_\star, \boldsymbol{\lambda}_\star)$ is positive semidefinite when constrained in $V(\mathbf{x}_\star)$.

Proposition C.2.2 (Sufficient Conditions). *Assume f and h are twice continuously differentiable. Let $(\mathbf{x}_\star, \boldsymbol{\lambda}_\star) \in \mathbb{R}^{n+m}$ satisfy*

$$\begin{aligned}\nabla_{\mathbf{x}} L(\mathbf{x}_\star, \boldsymbol{\lambda}_\star) &= 0, \\ \nabla_{\boldsymbol{\lambda}} L(\mathbf{x}_\star, \boldsymbol{\lambda}_\star) &= 0, \\ \mathbf{v}^T \nabla_{\mathbf{x}\mathbf{x}}^2 L(\mathbf{x}_\star, \boldsymbol{\lambda}_\star) \mathbf{v}^T &> 0 \quad \text{for all } \mathbf{v} \in V(\mathbf{x}_\star), \mathbf{v} \neq \mathbf{0}.\end{aligned}\tag{C.2.6}$$

Then \mathbf{x}_\star is a strict local minimum of $f(\mathbf{x})$ subject to $\mathbf{h}(\mathbf{x}) = 0$.

C.3 Basic Duality Theory

Recall the Lagrangian function for the above equality-constrained optimization problem:

$$L(\mathbf{x}, \boldsymbol{\lambda}) \doteq f(\mathbf{x}) + \sum_{i=1}^m \lambda_i h_i(\mathbf{x}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{h}(\mathbf{x}) \rangle,\tag{C.3.1}$$

where $\boldsymbol{\lambda} = [\lambda_1, \lambda_2, \dots, \lambda_m]^T \in \mathbb{R}^m$ are the Lagrangian multipliers.

In duality theory, the vector $\boldsymbol{\lambda}$ is also called the *dual variables* for the so-called *dual function*:

$$q(\boldsymbol{\lambda}) \doteq \inf_{\mathbf{x} \in \mathcal{D}} L(\mathbf{x}, \boldsymbol{\lambda}).\tag{C.3.2}$$

Correspondingly, $f(\mathbf{x})$ is referred to as the *primal function* and \mathbf{x} the *primal variables*.

A simple property of the dual function q is that it is a concave function regardless whether the primal problem is convex or not, since q is the point-wise infimum of a family of affine functions with respect to $(\boldsymbol{\lambda})$.

Another important property of the dual function is that $q(\boldsymbol{\lambda})$ is a lower bound of $f(\mathbf{x}')$ for any feasible solution \mathbf{x}' . In particular, $q(\boldsymbol{\lambda})$ is a lower bound of the optimal value $f(\mathbf{x}_*)$. This can be easily verified since for a feasible \mathbf{x}' satisfying $\mathbf{h}(\mathbf{x}') = 0$, we have

$$q(\boldsymbol{\lambda}) = \inf_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{h}(\mathbf{x}) \rangle \leq \inf_{\mathbf{x} \in \mathcal{D}, \mathbf{h}(\mathbf{x})=0} f(\mathbf{x}) \leq f(\mathbf{x}').$$

For the dual function $q(\boldsymbol{\lambda})$ to provide a meaningful lower bound for $f(\mathbf{x}_*)$, it is natural to avoid trivial cases when $q(\boldsymbol{\lambda}) = -\infty$. So we normally restrict the domain of the dual function q to:

$$\mathcal{C} \doteq \{\boldsymbol{\lambda} \mid q(\boldsymbol{\lambda}) > -\infty\}. \quad (\text{C.3.3})$$

More specifically, the dual variables $(\boldsymbol{\lambda})$ that satisfy above conditions are called *dual feasible solutions*.

A very useful concept in duality theory is the so-called *duality gap* between the primal and dual functions

$$f(\mathbf{x}) - q(\boldsymbol{\lambda}). \quad (\text{C.3.4})$$

Since q is a lower bound of f , the duality gap is always nonnegative (over the set of feasible solutions). More importantly, when the duality gap is zero, namely, there exists a feasible solution \mathbf{x}_* and $\boldsymbol{\lambda}_*$ such that $f(\mathbf{x}_*) = q(\boldsymbol{\lambda}_*)$, then \mathbf{x}_* is the optimal primal solution and $\boldsymbol{\lambda}_*$ is the optimal dual solution.

Since the dual function $q(\boldsymbol{\lambda})$ is a lower bound of the primal function $f(\mathbf{x})$, in particular of its minimal value $f(\mathbf{x}_*)$. Naturally, when we want to achieve the best lower-bound estimation of the minimal value, we can consider the following optimization problem in the dual space:

$$\max_{\boldsymbol{\lambda}} q(\boldsymbol{\lambda}). \quad (\text{C.3.5})$$

The problem (C.3.5) is called the *Lagrange dual problem* associated with the original *primal problem* (C.2.1).

Since the optimal solution $q(\boldsymbol{\lambda}_*)$ is the best lower-bound approximation of the global minimum $f(\mathbf{x}_*)$, the following inequality condition holds trivially:

$$q(\boldsymbol{\lambda}_*) \leq f(\mathbf{x}_*). \quad (\text{C.3.6})$$

The condition is known as the *weak duality condition*. Furthermore, when the equality can be obtained in (C.3.6), the duality gap between f and q becomes zero, and we say the primal and dual function pair satisfy the *strong duality condition*.

The strong duality condition can be achieved for convex objective functions subject to linear constraints.

Theorem C.3.1 (Strong Duality Theorem). *Let the objective function $f(\mathbf{x})$ in (C.2.1) be convex and $\mathbf{h}(\mathbf{x})$ be linear. If the optimal value f_* is finite, then the optimal solution for its dual problem exists and there is no duality gap.*

Under the strong duality condition, the minimal value of f can be found by optimizing the dual problem $q(\boldsymbol{\lambda})$, and the optimal primal solution can also be obtained by minimizing the Lagrangian function $L(\mathbf{x}, \boldsymbol{\lambda}_*)$ over \mathbf{x} . In other words, the optimal $(\mathbf{x}_*, \boldsymbol{\lambda}_*)$ is the saddle point of the Lagrangian function:

$$\max_{\boldsymbol{\lambda}} \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}). \quad (\text{C.3.7})$$

In the above, we have assumed all functions are differentiable. In this book, we often need to optimize a convex function that is not differentiable and the type of constraints are of the form $\mathbf{Ax} = \mathbf{y}$.

Lemma C.3.2. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ convex, $\mathbf{y} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, and let \mathbf{x}_* be some point satisfying $\mathbf{Ax}_* = \mathbf{y}$. If there exists $\boldsymbol{\nu}$ such that*

$$\mathbf{A}^* \boldsymbol{\nu} \in \partial f(\mathbf{x}_*), \quad (\text{C.3.8})$$

then \mathbf{x}_ is a solution to the optimization problem*

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{Ax} = \mathbf{y} \end{aligned} \quad (\text{C.3.9})$$

Proof. Consider any \mathbf{x}' satisfying $\mathbf{Ax}' = \mathbf{y}$. By the subgradient inequality B.3.3,

$$\begin{aligned} f(\mathbf{x}') & \geq f(\mathbf{x}_*) + \langle \mathbf{A}^* \boldsymbol{\nu}, \mathbf{x}' - \mathbf{x}_* \rangle \\ & = f(\mathbf{x}_*) + \langle \boldsymbol{\nu}, \mathbf{A}(\mathbf{x}' - \mathbf{x}_*) \rangle \\ & = f(\mathbf{x}_*), \end{aligned} \quad (\text{C.3.10})$$

since $\mathbf{Ax}' = \mathbf{Ax}_*$. Thus, \mathbf{x}_* is optimal. \square

Appendix D

Methods for Optimization

In this chapter, we review classical approaches to solving optimization problems of the form

$$\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}), \quad (\text{D.0.1})$$

in which we seek to minimize an objective function f over some domain \mathcal{D} . All of the algorithms we describe are *iterative methods* of optimization, which produce a sequence of points

$$\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k, \dots \quad (\text{D.0.2})$$

starting from some initialization \mathbf{x}_0 . The goal is to generate a sequence $\{\mathbf{x}_k\}$ which quickly converges to a minimizer \mathbf{x}_\star of f over \mathcal{D} . The total time an iterative method requires to produce an acceptable answer depends chiefly on two quantities:

- (i) **Per iteration cost:** how long it takes to generate the next point \mathbf{x}_{k+1} given the previous points $\mathbf{x}_0, \dots, \mathbf{x}_k$.
- (ii) **Convergence rate:** how quickly the iterates \mathbf{x}_k improve in quality. This dictates how many iterations are required to produce a sufficiently good solution. This may be measured either in terms of the distance of the iterate \mathbf{x}_k to a minimizer,

$$\|\mathbf{x}_k - \mathbf{x}_\star\|, \quad (\text{D.0.3})$$

or in terms of the sub-optimality in objective value,

$$|f(\mathbf{x}_k) - f(\mathbf{x}_\star)|. \quad (\text{D.0.4})$$

Criteria (i) and (ii) are usually in tension: we can fast convergence rate at the price of very expensive iterations, or we can have very cheap iterations at the price of a relatively slow convergence.

In the era of big data or large models, many practical problems involve optimizing over very large number of model parameters or training over large-scale datasets. Due to computation limitations, we typically can only afford to do fairly simple calculations in each iteration. Hence we are mainly interested in methods that achieve the fastest possible convergence rate out of methods that only work with *first order* information (values of $f(\mathbf{x})$ and $\nabla f(\mathbf{x})$). Sometimes due to memory limitation and time requirement, we need to store the data and conduct the calculation over many *parallel* processes or a *distributed* network of machines. To reduce communication cost and delay, we often prefer algorithms that are amenable to parallel or distributed implementation and require minimal exchange of data and information across different processes or machines. In this appendix, we sketch basic ideas of some of the most popular and effective techniques that enhance the performance of first order methods, especially those are suitable for solving large-scale problems. We also provide references where the reader can find more complete exposition and analysis of these techniques.

D.1 Gradient Descent

Perhaps the simplest iterative method of optimization is *gradient descent*, also known as the gradient method, which applies to *differentiable* functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$. The method comes from the simplest idea that from the current state \mathbf{x}_k , one would like to take a small step $t \geq 0$ in the direction $\mathbf{v} \in \mathbb{R}^n$ to $\mathbf{x}_{k+1} = \mathbf{x}_k + t \cdot \mathbf{v}$ such that the value of f decreases:

$$f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k).$$

Since f is differentiable, we know that up to first-order approximation:

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) = f(\mathbf{x}_k + t \cdot \mathbf{v}) - f(\mathbf{x}_k) \approx t \cdot \nabla f(\mathbf{x}_k)^T \mathbf{v}.$$

The gradient $\nabla f(\mathbf{x}_k)$ points in the direction of steepest increase of the objective f ; the negative gradient is the direction of steepest decent. So in order for $f(\mathbf{x}_{k+1})$ to be smaller than $f(\mathbf{x}_k)$, it is natural to take the direction in which the value of f drops the fastest: $\mathbf{v} \propto -\nabla f(\mathbf{x}_k)$.

Therefore, gradient descent generates its next iterate by stepping in the direction of the negative gradient

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \nabla f(\mathbf{x}_k). \quad (\text{D.1.1})$$

Here, $t_k \geq 0$ is a scalar, often called the *step size*.¹ The step size t_k can either be determined analytically from the properties of the function f , or numerically by performing a *line search*, which produces an approximate solution² to the one dimensional problem:

$$\min_{t \geq 0} f(\mathbf{x}_k + t \nabla f(\mathbf{x}_k)). \quad (\text{D.1.2})$$

Convergence of gradient descent.

A principal virtue of gradient descent is that for many problems, ∇f can be computed efficiently. To understand the overall properties of the method, we need to know how many iterations it requires to obtain a solution of a given desired quality. This depends in turn on the properties of the objective function f .

We begin by assuming that f is a convex, differentiable function, and that the gradient $\nabla f(\mathbf{x})$ is L -Lipschitz:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}')\|_2 \leq L \|\mathbf{x} - \mathbf{x}'\|_2, \quad \forall \mathbf{x}, \mathbf{x}'. \quad (\text{D.1.3})$$

This condition states that the gradient does not change too rapidly as we move from point to point. Intuitively, this means that a first-order model for the objective function generated by taking a Taylor expansion at point \mathbf{x} will be valid over a relatively large portion of the space. Indeed, it turns out that under these hypotheses, we can take t_k to a uniform $\frac{1}{L}$; smaller L allows larger steps. Moreover, it can be shown that with this choice,

$$\begin{aligned} f(\mathbf{x}_{k+1}) &\leq f(\mathbf{x}_k) - \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|^2 \\ &\leq f(\mathbf{x}_k). \end{aligned} \quad (\text{D.1.4})$$

Thus, with this choice, the gradient method is a *descent method*: it strictly decreases the objective at each iteration, until \mathbf{x}_k reaches a minimizer. The following theorem gives an overall control on the rate of convergence, measured in function values:

Theorem D.1.1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable function with $\nabla f(\mathbf{x})$ L -Lipschitz. Let $\mathbf{X}_\star \neq \emptyset$ denote the set of minimizers of f , and f_\star the minimum value of f over \mathbb{R}^n . Consider the gradient method with constant step size $t_k = \frac{1}{L}$. Then*

$$f(\mathbf{x}_k) - f_\star \leq \frac{L}{2} \frac{\|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2}{k}. \quad (\text{D.1.5})$$

Moreover, as $k \rightarrow \infty$, $\mathbf{x}_k \rightarrow \mathbf{X}_\star$.

A proof of this theorem (actually a more generalized version) can be found in Chapter 8, Section 8.2.

¹or the *learning rate* in learning algorithms.

²Typically, this is done by *backtracking*: starting from some nominal value of t , we reduce t until ... backtracking, Armijo and Wolfe conditions

Several aspects of this result are worth noting. First, the suboptimality in function values decreases as $1/k$. In particular, as $k \rightarrow \infty$, $f(\mathbf{x}_k) \rightarrow f_\star$. Second, the rate of convergence depends on the Lipschitz constant L – the smaller L is, faster f approaches f_\star . Finally, the rate of convergence depends on the distance of the initialization to \mathbf{x}_\star . A strength of this result is that it is nonasymptotic (the bound works for all k , not just k large) and does not depend on dimension n . For applications, we care not just about function values, but about the quality of the iterates \mathbf{x}_k . Here, we are guaranteed that \mathbf{x}_k approaches \mathbf{x}_\star . However, no general, dimension-independent bound on the rate of convergence is known.

D.2 Rates of Convergence and Acceleration

How good the gradient method is? More generally, if we restrict ourselves to relatively simple methods that only use gradient and function value information, what rate can we obtain? This fundamental question motivates the study of lower bounds for the computational efficiency of methods. This requires a model of computation. One simply model for first order methods assumes that at each iteration, the next point \mathbf{x}_{k+1} is generated based only on the previous points $\mathbf{x}_0, \dots, \mathbf{x}_k$, their function values $f(\mathbf{x}_0), \dots, f(\mathbf{x}_k)$, and gradients $\nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_k)$:

$$f(\mathbf{x}_{k+1}) = \mathcal{F}_{k+1}(\mathbf{x}_0, \dots, \mathbf{x}_k, f(\mathbf{x}_0), \dots, f(\mathbf{x}_k), \nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_k)). \quad (\text{D.2.1})$$

This is sometimes referred to as a *black box model*, since the method only accesses the function f through its value and gradient.

It has been shown that:

Theorem D.2.1. *For every L and R , there exists a convex differentiable function f with ∇f L -Lipschitz, and an initial point \mathbf{x}_0 satisfying $\|\mathbf{x}_0 - \mathbf{x}_\star\|_2 \leq R$ such that*

$$f(\mathbf{x}_k) - f_\star \geq c \frac{LR^2}{k^2}, \quad (\text{D.2.2})$$

where $c > 0$ is a numerical constant.

This result can be read as saying that for the class of functions with Lipschitz gradients, the best generic rate of convergence that any gradient-like method can achieve is $O(1/k^2)$. Notice that Theorem D.1.1 implies that the gradient method converges at a rate of $O(1/k)$. For large k , this is *much* worse!

Could the gradient method be suboptimal? Figure D.1 shows the behavior of gradient descent on two different problems. The figure plots the level sets $S_\beta = \{\mathbf{x} \mid f(\mathbf{x}) = \beta\}$ of the objective f as well as the iterates \mathbf{x}_k . Because the gradient $\nabla f(\mathbf{x})$ is orthogonal to the level set containing \mathbf{x} ,

[[FIGURE HERE]]

Figure D.1. **Gradient descent on ill-conditioned problems** John: show the behavior of the gradient method on two problems – a nearly spherical quadratic and a more ellipsoidal quadratic

[[FIGURE HERE]]

Figure D.2. **The gradient method and heavy ball method.**

the gradient method moves orthogonal to the level sets. At left, we show a function $f(\mathbf{x})$ whose level sets are nearly circular. The gradient method makes rapid progress. At right is a function $f(\mathbf{x})$ whose level sets are more elongated. The iterates “chatter” repeatedly changing direction and making slow progress towards \mathbf{x}_\star .

The heavy ball method.

The bad behavior in Figure D.1 can be mitigated by preventing the steps $\mathbf{x}_{k+1} - \mathbf{x}_k$ from changing direction too rapidly. An intuitive way to accomplish this is to treat the iterate \mathbf{x}_k as a particle with some amount of momentum, which causes it to continue moving in the same direction. This suggests an update of the form

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}). \quad (\text{D.2.3})$$

Because it treats \mathbf{x}_k as a particle with nonzero mass, this method is aptly called the *heavy ball method*. Figure D.2 compares the heavy ball method to the gradient method on an ill-conditioned quadratic. Notice that the heavy ball method takes far fewer iterations to reach the vicinity of \mathbf{x}_\star .

Nesterov’s accelerated method.

Although the heavy ball method improves over the gradient method, its worst case rate of convergence is still $O(1/k)$. However, by using momentum in a clever way, it is possible to achieve a better rate of convergence of $O(1/k^2)$, which matches the lower bound in Theorem D.2.1. This means, perhaps surprisingly, that there is a gradient-like method that is fundamentally better than gradient descent!

The method that achieves this optimal rate is known as *Nesterov’s accelerated gradient method*. Strictly speaking, it is not a momentum method. Rather, it uses two sequences of iterates \mathbf{x}_k and \mathbf{p}_k . The auxiliary point \mathbf{p}_k is extrapolated from \mathbf{x}_k in a form similar to that in the heavy ball method:

$$\mathbf{p}_{k+1} \doteq \mathbf{x}_k + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}).$$

At each iteration, we move to this new point, compute the gradient at this point, and descend from it (instead of \mathbf{x}_k):

$$\mathbf{x}_{k+1} = \mathbf{p}_{k+1} - \alpha \nabla f(\mathbf{p}_{k+1}). \quad (\text{D.2.4})$$

As we will show in Section 8.3 of Chapter 8, with properly chosen weights β_k and α , the gradient method is indeed accelerated and can achieve the optimal convergence rate of $O(1/k^2)$, for the class of functions with Lipschitz gradients.

Theorem D.2.2. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable function with $\nabla f(\mathbf{x})$ being L -Lipschitz. Let $\mathbf{X}_\star \neq \emptyset$ denote the set of minimizers of f and f_\star the minimum value of f over \mathbb{R}^n . The iterates \mathbf{x}_k produced by the accelerated gradient method satisfy*

$$f(\mathbf{x}_k) - f_\star \leq \frac{L}{2} \frac{\|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2}{(k+1)^2}. \quad (\text{D.2.5})$$

Moreover, as $k \rightarrow \infty$, $\mathbf{x}_k \rightarrow \mathbf{x}_\star$.

Strongly convex functions.

Notice that Theorem D.2.1 characterizes the best possible rate of convergence for gradient-like methods for the class of functions with Lipschitz gradients; and Theorem D.2.2 states that this rate can be achieved with the accelerated gradient methods. Nevertheless, this does not mean this is the best one can do for more restricted classes of functions. If, in addition to Lipschitz gradients, the functions satisfy additional nice properties (such as strongly convex defined in Appendix B), one can show that gradient-like methods converge at a linear rate [Boyd and Vandenberghe, 2004].

Theorem D.2.3. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable strongly convex function with constant μ and $\nabla f(\mathbf{x})$ being L -Lipschitz. Let f_\star be the minimum value of f over \mathbb{R}^n . Then the iterates \mathbf{x}_k produced by the gradient-descent $\mathbf{x}_{k+1} = \mathbf{x}_k - t \nabla f(\mathbf{x}_k)$ with $t = \frac{2}{L+\mu}$ satisfy*

$$f(\mathbf{x}_k) - f_\star \leq \frac{L}{2} e^{-\alpha k} \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2 \quad (\text{D.2.6})$$

for some constant $\alpha > 0$.

That is, $f(\mathbf{x}_k) - f_\star$ converges to zero exponentially in the order of $O(e^{-\alpha k})$, much faster than $O(1/k^2)$. In this book, the class of functions that we often encounter are not necessarily strongly convex. Nevertheless, they may satisfy certain weaker notion of strong convexity, known as *restricted strong convexity*. We will see that under such conditions, one may also expect gradient-like methods to achieve linear rate of convergence.

Nondifferentiable functions.

The main assumption of gradient descent methods is that the objective function $f(\mathbf{x})$ is differentiable in \mathbf{x} . In this book, we often need to minimize

functions that are not everywhere differentiable, such as functions involving the ℓ^1 norm $\|\mathbf{x}\|_1$. In such cases, we need to generalize the notion of gradient to “subgradients” (see Definition 2.3.4 in Chapter 2). Essentially, subgradients at a point \mathbf{x} is the set of vectors $\mathbf{u} \in \mathbb{R}^n$ such that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{u}, \mathbf{y} - \mathbf{x} \rangle, \quad \forall \mathbf{y}.$$

We often denote the set of subgradients as $\partial f(\mathbf{x})$. To minimize such a function $f(\mathbf{x})$, we may generalize the gradient descent method by replacing the gradient $\nabla f(\mathbf{x})$ with any subgradient:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \mathbf{g}_k, \quad \mathbf{g}_k \in \partial f(\mathbf{x}_k).$$

A main disadvantage of such subgradient descent methods is their relatively poor convergence rate. In general, the convergence rate of subgradient descent for non-smooth objective functions is

$$f(\mathbf{x}_k) - f_\star = O\left(1/\sqrt{k}\right).$$

The reader can refer to [Nemirovski, 1995, Nemirovski, 2007, Nesterov, 2003b] for more detailed analysis of subgradient descent algorithms.

Nevertheless, as we will see in Chapter 8, in many of our problems, the objective function $f(\mathbf{x})$ is of the form $f_1(\mathbf{x}) + f_2(\mathbf{x})$ with f_1 being smooth and f_2 nonsmooth. If for the nonsmooth part f_2 , the so called *proximal operator*:

$$\min_{\mathbf{x}} f_2(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{w}\|_2^2 \quad (\text{D.2.7})$$

has a closed-form solution or can be solved efficiently, then the subgradient descent method can be properly modified so that it would enjoy the same convergence rate as the smooth case. See the *proximal gradient* method in Section 8.2 of Chapter 8.

Constrained optimization.

It is very common in practice that we want to minimize a function $f(\mathbf{x})$ while the desired solution \mathbf{x}_\star is constrained to some subset $C \subset \mathbb{R}^n$:

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in C. \end{aligned} \quad (\text{D.2.8})$$

Solutions in the subset C are called *feasible* solutions. Notice that if we still apply the gradient descent method to minimize $f(\mathbf{x})$. Then, at each descent iteration $\mathbf{p}_{k+1} = \mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)$, even if \mathbf{x}_k is feasible, the new state \mathbf{p}_{k+1} may step outside of the constrained set: $\mathbf{p}_{k+1} \notin C$. A natural and simple fix to this issue is to “project” \mathbf{p}_{k+1} back to the set C :

$$\mathbf{x}_{k+1} = \mathcal{P}_C[\mathbf{p}_{k+1}] = \arg \min_{\mathbf{x} \in C} \frac{1}{2} \|\mathbf{x} - \mathbf{p}_{k+1}\|_2^2 \quad (\text{D.2.9})$$

where \mathbf{x}_{k+1} is the point in C closest to \mathbf{p}_{k+1} . This will ensure the new iterate \mathbf{x}_{k+1} is always feasible. This method is called *projected gradient descent*,

and we use it to provide a simplest algorithm for minimizing the ℓ^1 norm in Chapter 2. This simple method is also the inspiration for other first-order constrained optimization methods such as the *Frank-Wolfe* method introduced in Section 8.6 of Chapter 8.

One disadvantage of such projected gradient descent methods is their relatively poor convergence rate. In the case the constraints are equality constraints: $C = \{\mathbf{x} \mid \mathbf{h}(\mathbf{x}) = 0\}$, one could try to convert the constrained optimization

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \mathbf{h}(\mathbf{x}) = 0 \quad (\text{D.2.10})$$

to an unconstrained one by penalizing any deviation of $\mathbf{h}(\mathbf{x})$ from 0:

$$\min f(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{h}(\mathbf{x})\|_2^2. \quad (\text{D.2.11})$$

This is known as the *penalty method*. One can show that as $\mu \rightarrow +\infty$, the solution to the unconstrained optimization approaches that of the constrained one. However, in practice, as μ becomes large, the unconstrained problem becomes increasingly harder to solve as its gradient Lipschitz becomes increasingly large. See Section 8.4 of Chapter 8 for an example.

As we have discussed in Appendix C, another way to convert the constrained optimization problem is through the Lagrangian formulation. The optimal (feasible) solution \mathbf{x}_\star to the above constrained optimization is also the optimal solution $(\mathbf{x}_\star, \boldsymbol{\lambda}_\star)$ to the unconstrained optimization:

$$\max_{\boldsymbol{\lambda}} \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}) \quad (\text{D.2.12})$$

where the Lagrangian function is defined as

$$L(\mathbf{x}, \boldsymbol{\lambda}) \doteq f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{h}(\mathbf{x}) \rangle.$$

It is natural to consider solving the above min-max problem through the following alternating optimization scheme:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k), \quad (\text{D.2.13})$$

$$\boldsymbol{\lambda}_{k+1} = \arg \max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}_{k+1}, \boldsymbol{\lambda}). \quad (\text{D.2.14})$$

Although the saddle point of the Lagrangian is the desired optimal solution, there is no guarantee that each step of the above iteration would produce feasible iterates. As we see in Section 8.4 of Chapter 8, even for some simple problems, the above subproblems might fail to have a solution (the value of the objective function can be unbounded).

To remedy this problem, one could augment the Lagrangian $L(\mathbf{x}, \boldsymbol{\lambda})$ with an extra quadratic penalty term for the constraint:

$$L_\mu(\mathbf{x}, \boldsymbol{\lambda}) \doteq f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{h}(\mathbf{x}) \rangle + \frac{\mu}{2} \|\mathbf{h}(\mathbf{x})\|_2^2,$$

which is known as the *Augmented Lagrangian*. As we will see in Section 8.4 of Chapter 8, the augmented Lagrangian leads to much better conditioned

subproblems for the alternating scheme:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \mathcal{L}_{\mu}(\mathbf{x}, \boldsymbol{\lambda}_k), \quad (\text{D.2.15})$$

$$\boldsymbol{\lambda}_{k+1} = \arg \max_{\boldsymbol{\lambda}} \mathcal{L}_{\mu}(\mathbf{x}_{k+1}, \boldsymbol{\lambda}), \quad (\text{D.2.16})$$

and the sequence of iterates $\{(\mathbf{x}_k, \boldsymbol{\lambda}_k)\}$ typically converge to the desired optimal solution $(\mathbf{x}_{\star}, \boldsymbol{\lambda}_{\star})$ for a properly chosen μ or a sequence $\{\mu_k\}$.

D.3 Block Coordinate Descent and Alternating Direction Method of Multipliers

In many optimization problems we may encounter in practice, the dimension of \mathbf{x} could be so high that we might not even afford to conduct gradient descent to minimize $f(\mathbf{x})$ for all the variables together. Very often the objective function $f(\mathbf{x})$ has certain decomposable structures such as:

$$\min f(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x}^i). \quad (\text{D.3.1})$$

For example, the ℓ^1 -norm function $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$ is such a decomposable function. In such cases, we may conduct the so-called *block coordinate descent* to take advantage of such decomposable structures by iteratively minimizing the objective function with respect to one block of variables at a time.

More specifically, assume the domain \mathcal{D} can be written as a Cartesian product

$$\mathcal{D} = \mathcal{D}_1 \times \mathcal{D}_2 \times \cdots \times \mathcal{D}_m,$$

where each $\mathcal{D}_i \subseteq \mathbb{R}^{n_i}$, $n_1 + n_2 + \cdots + n_m = n$. The variables can be also partitioned into m blocks as $\mathbf{x} = (\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m) \in \mathbb{R}^n$ with each $\mathbf{x}^i \in \mathcal{D}_i$. The block coordinate descent scheme proceeds as follows:

1. Initialize $\mathbf{x}_0 = (\mathbf{x}_0^1, \mathbf{x}_0^2, \dots, \mathbf{x}_0^m)$.
2. In the k -th iteration, for every $i = 1, \dots, m$,

$$\mathbf{x}_k^i = \arg \min_{\mathbf{x} \in \mathcal{D}_i} f(\mathbf{x}_k^1, \dots, \mathbf{x}_k^{i-1}, \bar{\mathbf{x}}, \mathbf{x}_k^{i+1}, \dots, \mathbf{x}_k^m).$$

3. Repeat Step 2 until the solution converges.

In the literature, the convergence of block coordinate descent methods can be proven under different conditions. A most natural condition is when the objective function $f(\mathbf{x}^1, \dots, \mathbf{x}^{i-1}, \mathbf{x}^i, \mathbf{x}^{i+1}, \dots, \mathbf{x}^m)$ is *strictly convex* with respect to each block \mathbf{x}^i . This guarantees the minimal solution \mathbf{x}_{\star}^i is also unique. For a more detailed discussion about conditions under which such methods converge, the reader is referred to [Bertsekas, 2003].

In compressive sensing or statistical learning,³ very often we need to deal with an objective function $f(\mathbf{x})$ that is composed of multiple terms:

$$f(\mathbf{x}) = f_1(\mathbf{x}) + f_2(\mathbf{x}) + \cdots + f_m(\mathbf{x}). \quad (\text{D.3.2})$$

To obtain more scalable algorithms such that we can optimize each term in a parallel or distributed fashion, we could rewrite this problem in terms of a set of local variables $\mathbf{x}^i \in \mathbb{R}^n$ and one global variable \mathbf{z} :

$$\text{minimize } \sum_{i=1}^m f_i(\mathbf{x}^i) \quad \text{subject to } \mathbf{x}^i = \mathbf{z}, \quad i = 1, \dots, m. \quad (\text{D.3.3})$$

In the literature, this is also known as the *consensus optimization*. To solve such a constrained optimization problem, we could apply the above block descent method to its augment Lagrangian:

$$L(\mathbf{x}^1, \dots, \mathbf{x}^m, \mathbf{z}, \boldsymbol{\lambda}) = \sum_{i=1}^m f_i(\mathbf{x}^i) + \langle \boldsymbol{\lambda}^i, \mathbf{x}^i - \mathbf{z} \rangle + \frac{\mu}{2} \|\mathbf{x}^i - \mathbf{z}\|_2^2. \quad (\text{D.3.4})$$

This leads to the following iterative process:

$$\begin{aligned} \mathbf{x}_{k+1}^i &= \arg \min_{\mathbf{x}^i} f_i(\mathbf{x}^i) + \langle \boldsymbol{\lambda}^i, \mathbf{x}^i - \mathbf{z} \rangle + \frac{\mu}{2} \|\mathbf{x}^i - \mathbf{z}\|_2^2, \\ \mathbf{z}_{k+1} &= \frac{1}{m} \sum_{i=1}^m \left(\mathbf{x}_{k+1}^i + \frac{1}{\mu} \boldsymbol{\lambda}_k^i \right), \\ \boldsymbol{\lambda}_{k+1}^i &= \boldsymbol{\lambda}_k^i + \mu (\mathbf{x}_{k+1}^i - \mathbf{z}_{k+1}). \end{aligned}$$

This is known as the *Alternating Direction Method of Multipliers* (ADMM). Notice that the above scheme is rather amenable to distributed implementation as each local process can solve in parallel a subproblem for \mathbf{x}^i and then share the information through the common variable \mathbf{z} . In Chapter 8, we will study the ADMM scheme for the case with $m = 2$ in great detail. For a more detailed exposition of ADMM and more general variants, the reader may refer to [Boyd et al., 2011b].

D.4 Nonconvex Problems

nonconvex problems – convergence to a stationary point? avoid saddle points? trust region methods?

³Say training a deep neural networks over a very large set of training samples, where \mathbf{x} are the network parameters.