

# Lectures on Optimization

Sinho Chewi

April 14, 2025

## Contents

<b>1</b>	<b>[1/14] Introduction and basics of convex functions</b>	<b>3</b>
1.1	Overview of the course . . . . .	3
1.2	Preliminaries on convexity and smoothness . . . . .	8
<b>2</b>	<b>[1/16] Gradient flow</b>	<b>12</b>
<b>3</b>	<b>[1/21] Gradient descent: smooth case</b>	<b>17</b>
<b>4</b>	<b>[1/23] Lower bounds for smooth optimization</b>	<b>23</b>
4.1	Reductions between the convex and strongly convex settings . . . . .	23
4.2	Lower bounds . . . . .	25
<b>5</b>	<b>[1/28–1/30] Acceleration</b>	<b>28</b>
5.1	Quadratic case: the conjugate gradient method . . . . .	28
5.2	General case: continuous time . . . . .	33
5.3	General case: discrete time . . . . .	35
<b>6</b>	<b>[2/4–2/13] Non-smooth convex optimization</b>	<b>38</b>
6.1	Convex analysis . . . . .	39
6.2	Projected subgradient methods . . . . .	44
6.3	Cutting plane methods . . . . .	48
6.4	Lower bounds . . . . .	51
<b>7</b>	<b>[2/18] Frank–Wolfe</b>	<b>55</b>

<b>8</b>	<b>[2/20] Proximal methods</b>	<b>58</b>
8.1	Algorithms and examples . . . . .	59
8.2	Convergence analysis . . . . .	62
<b>9</b>	<b>[2/25–2/27] Fenchel duality</b>	<b>64</b>
9.1	(Optional) Connection with classical mechanics . . . . .	65
9.2	Duality correspondences . . . . .	69
<b>10</b>	<b>[3/4–3/6] Mirror methods</b>	<b>74</b>
10.1	Bregman divergences and relative convexity/smoothness . . . . .	75
10.2	Algorithms and convergence analysis . . . . .	77
10.3	Online algorithms and multiplicative weights . . . . .	82
<b>11</b>	<b>[3/25–4/3] Alternating minimization</b>	<b>87</b>
11.1	Alternating projections . . . . .	88
11.2	Convergence analysis for alternating minimization . . . . .	90
11.3	Case study: entropic optimal transport . . . . .	95
<b>12</b>	<b>[4/8–4/17] Stochastic optimization</b>	<b>100</b>
12.1	Stochastic mirror proximal gradient descent . . . . .	100
12.2	Implications for statistical generalization . . . . .	104
12.3	Central limit theorem for Polyak–Ruppert averaging . . . . .	107
12.4	Variance reduction . . . . .	117

# 1 [1/14] Introduction and basics of convex functions

These lecture notes supplement S&DS 432/632 (Advanced Optimization Techniques), taught in Spring 2025. They are not meant to be comprehensive.

The notes are primarily based on the books [Bub15; Nes18], as well as my personal understanding of the subject formed through discussions with many people over the years. Please send me corrections and feedback via email. I thank Linghai Liu, Leda Wang, Ruixiao Wang, Ilias Zadik, and Matthew S. Zhang for correcting my mistakes.

**Logistics.** The problem sets, syllabus, and all other information can be found at the [course website](#) and the Canvas page. Grading is based on six problem sets and one take-home final exam, each of which counts for 1/7 of the total grade. All questions related to logistics should be directed to my email: [sinho.chewi@yale.edu](mailto:sinho.chewi@yale.edu).

**Audience.** This course focuses on the theory of optimization. In particular, the course is **mathematical** in nature and taught in a theorem–proof format. The course assumes familiarity with basic proofs and logical reasoning, as well as linear algebra, multivariate calculus, and probability theory.

The reader should also be familiar with asymptotic notions (big- $O$  notation). We use the shorthand notation  $a \lesssim b$  (resp.  $a \gtrsim b$ ) to mean that  $a \leq Cb$  (resp.  $a \geq b/C$ ) for an absolute constant  $C > 0$  (i.e., a constant that does not depend on other parameters of the problem), and  $a \asymp b$  to mean that both  $a \lesssim b$  and  $a \gtrsim b$  hold. We use  $a = O(b)$  and  $a \lesssim b$  interchangeably.

## 1.1 Overview of the course

The basic problem of optimization is to compute an approximate minimizer of a given function  $f : \mathcal{X} \rightarrow \mathbb{R}$ . In this course,  $\mathcal{X}$  is always taken to be a subset of  $\mathbb{R}^d$ , although generalizations are possible (e.g., to manifolds).

**Black-box optimization and the oracle model.** What does it mean to “compute”? The answer depends on the representation of  $f$  and our model of computation. We start by studying *black-box optimization*. In this model, we presume that we can *evaluate*  $f$ , and possibly its derivatives, at any chosen point  $x \in \mathcal{X}$ .

The advantage of the black-box model is that it applies very *generally*: it is difficult to find situations in which we need to optimize a function but we cannot even evaluate

it! Consequently, algorithms developed in this model can be applied to the majority of problems encountered in practice<sup>1</sup>—witness the ubiquity of gradient descent.

The disadvantage is that by its very generality, it cannot take advantage of additional structural information about  $f$  which can bring computational savings. That is why, later in the course, we turn toward the study of *structured* optimization problems.

It is easy, at least at an intuitive level, to describe algorithms which are valid in the black-box model. Namely, they are algorithms which only “interact” with  $f$  through evaluations of  $f$  and its derivatives. The existence of an algorithm, together with a corresponding mathematical analysis of the number of iterations to reach an approximate minimizer contingent upon assumptions on  $f$ , provide an *upper bound* on the complexity of the optimization task. In this course, we are also interested in *lower bounds*, which delineate fundamental limitations encountered by *any* algorithm. In order to prove such a lower bound, we need to formalize the notion of “interaction” alluded to above, and this leads to the important concept of an *oracle*.

First, observe that it does not make sense to discuss the complexity of optimizing a *single* function  $f$ . For if  $x_\star$  is the minimizer of  $f$ , we can consider the algorithm “output  $x_\star$ ”, which yields the correct answer in one iteration. But this algorithm is silly, since it utterly fails at optimizing any other function whose minimizer does not happen to be  $x_\star$ . Reflecting upon this situation, we do not consider an optimization algorithm to be sensible when it happens to succeed for one particular problem; rather, we expect it to succeed on many similar problems. Hence, we talk about a *class* of functions  $\mathcal{F}$  of interest, and we require our algorithms to succeed on *every*  $f \in \mathcal{F}$ .

The algorithm is designed to succeed on  $\mathcal{F}$  and thus, in an anthropomorphic sense, it “knows”  $\mathcal{F}$ . However, it does not know which particular  $f \in \mathcal{F}$  it is trying to optimize. (If it possessed knowledge of  $f$ , then we run into the issue from before, namely it could simply output the minimizer.) The role of the oracle is to act as an intermediary between the algorithm and the function. Namely, we assume that the algorithm is allowed to ask certain questions (“queries”) to the oracle for  $f$ , and this is the only means by which the algorithm can gather more information about  $f$ . The allowable queries and responses determine the nature of the oracle, e.g.:

- a **zeroth-order oracle** accepts a query point  $x \in \mathbb{R}^d$  and outputs  $f(x)$ ;
- a **first-order oracle** accepts a query point  $x \in \mathbb{R}^d$  and outputs  $(f(x), \nabla f(x))$ .

Most of the course focuses on optimization with a first-order oracle, but other oracles are possible (e.g., *linear optimization oracles* and *proximal oracles*). The zeroth-order and

---

<sup>1</sup>There is a caveat: in this course, we solely consider continuous optimization problems. Combinatorial optimization is an entirely different beast.

first-order oracles are easy to justify, as they correspond to the black-box model described above. As the oracles become more exotic, it becomes necessary to show that they are reasonable, by describing important applications in which such access to  $f$  is feasible.

The *query complexity* of  $\mathcal{F}$  for a particular choice of oracle, as a function of the prescribed tolerance  $\varepsilon$ , is then (informally) defined to be the minimum number  $N$  such that there exists an algorithm which, for any  $f \in \mathcal{F}$ , makes  $N$  queries to the oracle for  $f$  and outputs a point  $x$  with  $f(x) - \min f \leq \varepsilon$ .

It is worth noting that query complexity is not the same as computational complexity. Indeed, query complexity only counts the number of interactions with the oracle, and the algorithm is allowed to perform unlimited computations between interactions. In principle, this could lead to a situation in which query complexity is wholly unrepresentative of the true computational cost of optimization—this would be the case if optimal algorithms in the oracle model were contrived and impractical. Thankfully, this is not the case. The oracle model is widely adopted as the standard model for optimization because it is the setting in which we can make precise claims about complexity, and because it generally aligns with optimization in practice.

This summarizes the conceptual framework for optimization theory—the “identity cards of the field” [Nes18], although a careful treatment of the framework only becomes necessary when discussing lower bounds (and hence we elaborate on the details then). As a branch of mathematics, the theory of optimization could be defined as the quest to characterize the query complexity of various classes  $\mathcal{F}$ , under various oracle models, and thereby identify optimal algorithms. This indeed remains a core element of the field, but as query complexity reaches maturity, research has shifted toward different types of questions, often inspired by practical developments.

**The role of convexity.** In order to optimize efficiently, we need to place assumptions on  $f$ , ideally minimal ones. For example, we can assume that  $f$  is continuous. In this course, however, we are interested in *quantitative* rates of convergence for algorithms, and for this purpose, a *qualitative* assumption such as continuity is not enough. A quantitative form of continuity is to assume that  $f$  is *L-Lipschitz in the  $\ell_\infty$  norm*:

$$|f(x) - f(y)| \leq L \max_{i \in [d]} |x[i] - y[i]| \quad \text{for all } x, y \in \mathcal{X}. \quad (1.1)$$

Also, for concreteness, let us take  $\mathcal{X}$  to be the cube,  $\mathcal{X} = [0, 1]^d$ . In the language of the framework above, we consider the class

$$\mathcal{F} = \{f : [0, 1]^d \rightarrow \mathbb{R} \mid f \text{ satisfies (1.1)}\}. \quad (1.2)$$

One can then prove the following negative result.

**Theorem 1.1.** For any  $0 < \varepsilon < L/2$  and any deterministic algorithm, the complexity of  $\varepsilon$ -approximately minimizing functions in the class defined in (1.2) to within  $\varepsilon$  using a zeroth-order oracle is at least  $\lfloor \frac{L}{2\varepsilon} \rfloor^d$ .

Thus, for  $\varepsilon < L/4$ , the complexity grows *exponentially* with the dimension. The proof is not difficult; see, e.g., [Nes18, Theorem 1.1.2]. It is also robust: variants of the result can be proven when the notion of Lipschitzness is w.r.t. the  $\ell_2$  norm; when the oracle is taken to be a first-order oracle; when the algorithm is allowed to be randomized; etc. The message is clear: in order for optimization to be tractable in the worst case, we must impose some structural assumptions.

The black-box oracles we have been considering are *local* in nature: given a query point  $x \in \mathbb{R}^d$ , the oracle reveals some information about the behavior of  $f$  in a local neighborhood of  $x$ . Assumptions such as Lipschitzness effectively govern how large this local neighborhood is. But ultimately, to render optimization tractable, we must ensure that local information yields global consequences. As justified in the next subsection, a key assumption that makes this possible is *convexity*.

Of course, not every problem is convex, and non-convex optimization often still succeeds. But for the purpose of understanding the core principles underlying optimization, there is no better starting place. It is important to remember that convex problems abound in every application domain; here, we give two classical examples from statistics.

**Example 1.2 (logistic regression).** The data consists of  $n$  pairs  $(X_i, Y_i) \in \mathbb{R}^d \times \{0, 1\}$ , where  $X_i$  is a vector of covariates and  $Y_i$  is a binary response. The statistical model assumes that the pairs are independently drawn, the covariates are deterministic, and  $Y_i$  has a Bernoulli distribution with parameter  $\exp(\langle \theta, X_i \rangle) / \{1 + \exp(\langle \theta, X_i \rangle)\}$ . The goal is to infer the parameter  $\theta$ .

The maximum likelihood estimator (MLE) for this model is the solution to the convex optimization problem

$$\hat{\theta}_{\text{MLE}} \in \arg \min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (\log(1 + \exp \langle \theta, X_i \rangle) - Y_i \langle \theta, X_i \rangle).$$

**Example 1.3 (LASSO).** The data consists of  $n$  pairs  $(X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$ . The statistical model assumes that the pairs are independently drawn, and that  $Y_i = \langle \theta, X_i \rangle + \xi_i$ , where the  $\xi_i$ 's are i.i.d. noise variables independent of the  $X_i$ 's. When the parameter  $\theta$  is assumed to be sparse, it is standard to use the LASSO estimator, which is the solution to the convex optimization problem

$$\hat{\theta}_{\text{LASSO}} = \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{2n} \sum_{i=1}^n (Y_i - \langle \theta, X_i \rangle)^2 + \lambda \|\theta\|_1 \right\}.$$

Here,  $\lambda > 0$  is the regularization parameter and  $\|\cdot\|_1$  denotes the  $\ell_1$  norm, defined via  $\|\theta\|_1 := \sum_{i=1}^d |\theta[i]|$ .

In these examples, the estimator is defined as the solution to a convex problem which is not solvable in closed form, necessitating the use of numerical optimization. Actually, it is not that most problems in the “wild” are convex and hence there was a need to develop convex optimization. In fact, it often goes the other way around: convex optimization is such a powerful tool that problems are intentionally formulated to be convex. This is the case for the LASSO estimator, which can be motivated as a convex relaxation of the (statistically superior)  $\ell_0$ -constrained least-squares estimator.

**First-order methods.** This course largely focuses on first-order methods, namely, gradient descent and its variants. This class of methods is natural from the perspective of the theory. Equally importantly, first-order methods are lightweight and therefore scalable to large problem sizes, making them the method of choice even for highly non-convex settings which fall squarely outside of the theory.

**Beyond the black-box model.** After developing results for the black-box model, we study structured problems which admit more efficient solutions. The LASSO estimator of [Example 1.3](#) can be treated as a “composite” optimization problem (a sum of a smooth and a non-smooth function), and the estimators in both [Example 1.2](#) and [Example 1.3](#) (and empirical risk minimization more generally) are “finite sum” problems whose computation can be sped up via the use of stochastic gradients. Other examples include the use of alternative geometries (mirror descent) and the use of coordinate-wise structure (alternating maximization/coordinate descent).

We also study interior-point methods, which are a practically effective suite of algorithms which solve linear programs (LPs) and semidefinite programs (SDPs) with polynomial iteration complexities.

Further topics are considered as time permits.

## 1.2 Preliminaries on convexity and smoothness

We assume familiarity with the basic notion of convexity, and we briefly review it here.

**Definition 1.4.** A subset  $\mathcal{C} \subseteq \mathbb{R}^d$  is **convex** if for all  $x, y \in \mathcal{C}$  and all  $t \in [0, 1]$ , the point  $(1 - t)x + ty$  also lies in  $\mathcal{C}$ .

**Definition 1.5.** Let  $\mathcal{C}$  be convex and let  $\alpha \geq 0$ . A function  $f : \mathcal{C} \rightarrow \mathbb{R}$  is  **$\alpha$ -convex** if for all  $x, y \in \mathcal{C}$  and all  $t \in [0, 1]$ ,

$$f((1 - t)x + ty) \leq (1 - t)f(x) + tf(y) - \frac{\alpha}{2} t(1 - t) \|y - x\|^2. \quad (1.3)$$

When  $\alpha = 0$ , this is just the usual definition of a convex function. When  $\alpha > 0$ , we say that the function is *strongly* convex.

The definition above has the advantage that it does not require  $f$  to be differentiable. However, for the purposes of checking and utilizing convexity, it is convenient to have the following equivalent reformulations, which should be committed to memory. For simplicity, we focus on the case  $\mathcal{C} = \mathbb{R}^d$ .

**Proposition 1.6 (convexity equivalences).** Let  $\mathcal{C} = \mathbb{R}^d$  and  $\alpha \geq 0$ .

1. If  $f$  is continuously differentiable, (1.3) is equivalent to each of the following:

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\alpha}{2} \|y - x\|^2 \quad \text{for all } x, y \in \mathbb{R}^d. \quad (1.4)$$

$$\langle \nabla f(y) - \nabla f(x), y - x \rangle \geq \alpha \|y - x\|^2 \quad \text{for all } x, y \in \mathbb{R}^d. \quad (1.5)$$

2. If  $f$  is twice continuously differentiable, (1.3) is equivalent to

$$\langle v, \nabla^2 f(x) v \rangle \geq \alpha \|v\|^2 \quad \text{for all } v, x \in \mathbb{R}^d. \quad (1.6)$$

*Proof.* Assume that  $f$  is continuously differentiable.

(1.3)  $\Rightarrow$  (1.4): Rearranging (1.3) yields, for  $t > 0$ ,

$$f(y) \geq f(x) + \frac{f((1 - t)x + ty) - f(x)}{t} + \frac{\alpha(1 - t)}{2} \|y - x\|^2.$$

Sending  $t \searrow 0$  yields (1.4).



(1.4)  $\Rightarrow$  (1.5): Swap  $x$  and  $y$  in (1.4) and add the resulting inequality back to (1.4).

(1.5)  $\Rightarrow$  (1.3): By the fundamental theorem of calculus, for  $v := y - x$ ,

$$f(y) = f(x) + \int_0^1 \langle \nabla f(x + sv), v \rangle ds,$$

$$f((1-t)x + ty) = f(x) + \int_0^1 \langle \nabla f(x + stv), tv \rangle ds.$$

Hence, (1.5) yields

$$\begin{aligned} f((1-t)x + ty) - (1-t)f(x) - tf(y) &= -t \int_0^1 \langle \nabla f(x + sv) - \nabla f(x + stv), v \rangle ds \\ &\leq -t \int_0^1 \alpha s(1-t) \|v\|^2 ds = -\frac{\alpha}{2} t(1-t) \|v\|^2. \end{aligned}$$

Finally, assume that  $f$  is twice continuously differentiable. Letting  $y = x + \varepsilon v$  in (1.5) and sending  $\varepsilon \searrow 0$  establishes (1.6). Conversely, the fundamental theorem of calculus shows that

$$\langle \nabla f(y) - \nabla f(x), y - x \rangle = \int_0^1 \langle \nabla^2 f(x + t(y-x)) (y-x), y-x \rangle dt,$$

and hence (1.6) implies (1.5).  $\square$

The equivalent statements each have their own interpretation: for  $\alpha = 0$ , (1.3) states that  $f$  lies below each of its secant lines between the intersection points; (1.4) states that  $f$  globally lies above each of its tangent lines; (1.6) states that  $\nabla f$  is a monotone vector field; and (1.6) is a statement about curvature.

As noted above, the key feature of convexity is that local information yields global conclusions. Before describing this, let us first recall some basic facts about optimization. For simplicity, we consider unconstrained optimization throughout.

**Lemma 1.7 (existence of minimizer).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be continuous and its level sets be bounded. Then, there exists a global minimizer of  $f$ .

*Proof.* The proof uses some analysis. Let  $x_0 \in \mathbb{R}^d$  and let  $\mathcal{K} := \{f \leq f(x_0)\}$  denote the level set. By the continuity assumption,  $\mathcal{K}$  is closed and bounded, thus compact. Let  $\{x_n\}_{n \in \mathbb{N}}$  be a minimizing sequence,  $f(x_n) \rightarrow \inf f$ . By compactness, it admits a subsequence, still denoted  $\{x_n\}_{n \in \mathbb{N}}$ , which converges to some  $x_\star \in \mathbb{R}^d$ . By continuity,  $f(x_\star) = \lim_{n \rightarrow \infty} f(x_n) = \inf f$ .  $\square$

**Lemma 1.8** (necessary conditions for optimality). Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be minimized at  $x_\star$ .

1. If  $f$  is continuously differentiable, then  $\nabla f(x_\star) = 0$ .
2. If  $f$  is twice continuously differentiable, then  $\nabla^2 f(x_\star) \geq 0$ .

*Proof.* Let  $v \in \mathbb{R}^d$  and  $\varepsilon > 0$ ; then,  $f(x_\star + \varepsilon v) - f(x_\star) \geq 0$ . If  $f$  is continuously differentiable, this yields  $\int_0^1 \langle \nabla f(x_\star + \varepsilon t v), v \rangle dt \geq 0$ . By continuity of  $\nabla f$ , sending  $\varepsilon \searrow 0$  proves that  $\langle \nabla f(x_\star), v \rangle \geq 0$  for all  $v \in \mathbb{R}^d$ , which entails  $\nabla f(x_\star) = 0$ .

If  $f$  is twice continuously differentiable, we can expand once more to obtain  $0 \leq \int_0^1 \int_0^1 \langle \nabla^2 f(x_\star + \varepsilon s t v), v \rangle ds dt$ . By continuity of  $\nabla^2 f$ , sending  $\varepsilon \searrow 0$  then proves that  $\langle \nabla^2 f(x_\star) v, v \rangle \geq 0$  for all  $v \in \mathbb{R}^d$ .  $\square$

The conditions  $\nabla f(x_\star) = 0$ ,  $\nabla^2 f(x_\star) \geq 0$  are necessary for optimality, but not sufficient in general. The issue is that the proof of Lemma 1.8 is entirely local, so the same conclusion holds even if  $x_\star$  is only assumed to be a *local* minimizer. On the other hand, under the assumption of convexity, the first-order necessary condition becomes sufficient.

**Lemma 1.9** (sufficient condition for optimality). Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be *convex* and continuously differentiable, and let  $\nabla f(x_\star) = 0$ . Then,  $x_\star$  is a global minimizer of  $f$ .

In particular, every local minimizer of  $f$  is a global minimizer.

*Proof.* This easily follows from (1.4) with  $x = x_\star$ .  $\square$

Next, we note that the minimizer is unique if  $f$  is strictly convex.

**Lemma 1.10** (uniqueness of minimizer). Assume that  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is strictly convex, i.e., for all distinct  $x, y \in \mathbb{R}^d$  and  $t \in (0, 1)$ ,  $f((1-t)x + ty) < (1-t)f(x) + tf(y)$ . Then, if  $f$  admits a minimizer  $x_\star$ , it is unique.

*Proof.* If we had two distinct minimizers  $x_\star, \tilde{x}_\star$ , so that  $f(x_\star) = f(\tilde{x}_\star)$ , then strict convexity would imply  $f(\frac{1}{2}x_\star + \frac{1}{2}\tilde{x}_\star) < f(x_\star)$ , which is a contradiction.  $\square$

If  $f$  is strongly convex, then it is strictly convex. Also, from, e.g., (1.4), we see that  $f$  grows at least quadratically at  $\infty$ , which implies that it has bounded level sets. We can therefore conclude:

**Corollary 1.11.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be strongly convex and continuously differentiable. Then, it admits a unique minimizer  $x_\star$ , which is characterized by  $\nabla f(x_\star) = 0$ .

Finally, when discussing algorithms, we also need a dual condition—an *upper* bound on the Hessian—which in this context is called *smoothness*.<sup>2</sup>

**Definition 1.12.** Let  $\beta \geq 0$ . We say that  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\beta$ -**smooth** if it is continuously differentiable and

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|y - x\|^2 \quad \text{for all } x, y \in \mathbb{R}^d. \quad (1.7)$$

The following proposition is established in the same way as [Proposition 1.6](#), so we omit the proof.

**Proposition 1.13 (smoothness equivalences).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be continuously differentiable and  $\beta \geq 0$ . Then,  $f$  is  $\beta$ -smooth if and only if

$$\langle \nabla f(y) - \nabla f(x), y - x \rangle \leq \beta \|y - x\|^2 \quad \text{for all } x, y \in \mathbb{R}^d.$$

If  $f$  is twice continuously differentiable, this is also equivalent to

$$\langle v, \nabla^2 f(x) v \rangle \leq \beta \|v\|^2 \quad \text{for all } v, x \in \mathbb{R}^d.$$

If  $f$  is convex,  $\beta$ -smooth, and twice continuously differentiable, then  $0 \leq \nabla^2 f \leq \beta I$ , which implies that the gradient  $\nabla f$  is  $\beta$ -Lipschitz:

$$\|\nabla f(y) - \nabla f(x)\| \leq \beta \|y - x\| \quad \text{for all } x, y \in \mathbb{R}^d. \quad (1.8)$$

This remains true even without assuming twice differentiability ([Exercise 3.1](#)).

## Bibliographical notes

For further discussion on the oracle model, see [\[NY83, §1\]](#).

## Exercises

**Exercise 1.1.** Let  $f = \frac{\alpha}{2} \|\cdot\|^2$ , where  $\alpha \geq 0$ . Show via direct computation that (1.3) holds with equality.

---

<sup>2</sup>This is not to be confused with the mathematical usage of “smoothness” as “infinitely differentiable”.

## 2 [1/16] Gradient flow

Before we turn toward our main first-order algorithm of interest, namely gradient descent, we first study the situation in continuous time via the gradient flow. Throughout this section, we let  $(x_t)_{t \geq 0}$  denote the gradient flow for  $f$ :

$$\dot{x}_t = -\nabla f(x_t). \quad (\text{GF})$$

This is an ordinary differential equation (ODE), and since the main purpose of this section is to develop intuition, we assume that  $f$  is twice continuously differentiable and do not worry about showing that (GF) is well-posed. We use the following notation throughout these notes:

$$x_\star \in \arg \min f, \quad f_\star := \min f = f(x_\star).$$

Generally, we always assume that  $f$  admits a minimizer.

The most basic property of GF is that it always decreases the function value.

**Lemma 2.1** (descent property of GF). For any  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , the gradient flow  $(x_t)_{t \geq 0}$  of  $f$  satisfies

$$\partial_t f(x_t) = -\|\nabla f(x_t)\|^2 \leq 0.$$

*Proof.* By the chain rule,  $\partial_t f(x_t) = \langle \nabla f(x_t), \dot{x}_t \rangle = -\|\nabla f(x_t)\|^2$ .  $\square$

To obtain quantitative convergence results, we now use the assumption of convexity. Our first result shows that under strong convexity, the gradient flow *contracts*.

**Theorem 2.2** (contraction of GF). Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be  $\alpha$ -convex. Let  $(y_t)_{t \geq 0}$  be another gradient flow for  $f$ , i.e.,  $\dot{y}_t = -\nabla f(y_t)$ . Then, for all  $t \geq 0$ ,

$$\|y_t - x_t\| \leq \exp(-\alpha t) \|y_0 - x_0\|.$$

*Proof.* We differentiate the squared distance between the two flows:

$$\partial_t (\|y_t - x_t\|^2) = 2 \langle y_t - x_t, \dot{y}_t - \dot{x}_t \rangle = -2 \langle y_t - x_t, \nabla f(y_t) - \nabla f(x_t) \rangle \leq -2\alpha \|y_t - x_t\|^2,$$

where the last inequality is (1.5). The proof is concluded by applying Grönwall's lemma (see Lemma 2.3) below.  $\square$

The proof above arrives at what is called a *differential inequality*, that is, an inequality which holds between a quantity and its derivative(s). This is a common strategy for analyzing ODEs/PDEs, and it can be loosely viewed as the continuous-time analogue of induction. The following standard lemma is useful for handling such inequalities.

**Lemma 2.3 (Grönwall).** Suppose that  $u : [0, T] \rightarrow \mathbb{R}$  is a continuously differentiable curve that satisfies the differential inequality

$$\dot{u}(t) \leq Au(t) + B(t), \quad t \in [0, T].$$

Then, it holds that

$$u(t) \leq u(0) \exp(At) + \int_0^t B(s) \exp(A(t-s)) ds, \quad t \in [0, T].$$

*Proof.* The idea is to differentiate  $t \mapsto \exp(-At) u(t)$ :

$$\partial_t [\exp(-At) u(t)] = \exp(-At) \{-Au(t) + \dot{u}(t)\} \leq B(t) \exp(-At).$$

By the fundamental theorem of calculus,

$$\exp(-At) u(t) - u(0) \leq \int_0^t B(s) \exp(-As) ds.$$

Rearranging yields the result.  $\square$

There are many variants of Grönwall's lemma that can be proven in similar ways, e.g., we can allow time-varying  $A$  as well.

Returning to [Theorem 2.2](#), we can apply [Lemma 2.3](#) with  $A = -2\alpha$  and  $B = 0$  to conclude that  $\|y_t - x_t\|^2 \leq \exp(-2\alpha t) \|y_0 - x_0\|^2$ , which proves the theorem. Note in particular that we can take  $y_t = x_\star$  for all  $t \geq 0$ , so it yields the following statement about convergence to the minimizer:  $\|x_t - x_\star\| \leq \exp(-\alpha t) \|x_0 - x_\star\|$ .

The next result is about convergence in function value, and unlike [Theorem 2.2](#), it yields convergence for the case  $\alpha = 0$  as well.

**Theorem 2.4 (convergence of GF in function value).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be  $\alpha$ -convex,  $\alpha \geq 0$ . Then, for all  $t \geq 0$ ,

$$f(x_t) - f_\star \leq \frac{\alpha}{2(\exp(\alpha t) - 1)} \|x_0 - x_\star\|^2.$$

When  $\alpha = 0$ , the right-hand side should be interpreted as its limiting value as  $\alpha \rightarrow 0$ , namely,  $\frac{1}{2t} \|x_0 - x_\star\|^2$ .

*Proof.* We differentiate  $t \mapsto \|x_t - x_\star\|^2$ , but this time we apply (1.4):

$$\partial_t(\|x_t - x_\star\|^2) = -2 \langle \nabla f(x_t), x_t - x_\star \rangle \leq -\alpha \|x_t - x_\star\|^2 - 2(f(x_t) - f_\star).$$

Applying Grönwall's lemma (Lemma 2.3) with  $A = -\alpha$ ,  $B(t) = -2(f(x_t) - f_\star)$ ,

$$0 \leq \|x_t - x_\star\|^2 \leq \exp(-\alpha t) \|x_0 - x_\star\|^2 - 2 \int_0^t \exp(-\alpha(t-s)) (f(x_s) - f_\star) ds.$$

By the descent property (Lemma 2.1),  $f(x_s) \geq f(x_t)$ , so that

$$\begin{aligned} \int_0^t \exp(-\alpha(t-s)) (f(x_s) - f_\star) ds &\geq (f(x_t) - f_\star) \int_0^t \exp(-\alpha(t-s)) ds \\ &= (f(x_t) - f_\star) \frac{1 - \exp(-\alpha t)}{\alpha}. \end{aligned}$$

Rearranging yields the result.  $\square$

When  $\alpha > 0$ , Theorem 2.4 shows that  $f(x_t) - f_\star = O(\exp(-\alpha t))$ . When  $\alpha = 0$ , the rate becomes  $f(x_t) - f_\star = O(1/t)$ . Actually, the rate in Theorem 2.4 is not sharp (see Exercise 2.1 and Exercise 2.2). However, the statement and proof are chosen because they form the basis of our approach in discrete time.

Next, we observe that convexity is not needed for convergence in function value. Due to the descent property (Lemma 2.1), it suffices to have a lower bound on the norm of the gradient to ensure that we make sufficient progress. For example, we can impose the following condition.

**Definition 2.5.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be continuously differentiable and  $\alpha > 0$ . We say that  $f$  satisfies a **Polyak–Łojasiewicz (PL) inequality** with constant  $\alpha$  if

$$\|\nabla f(x)\|^2 \geq 2\alpha (f(x) - f_\star) \quad \text{for all } x \in \mathbb{R}^d. \quad (\text{PL})$$

The next statement is an immediate corollary of Lemma 2.1, (PL), and Grönwall's lemma (Lemma 2.3).

**Corollary 2.6 (convergence of GF under PL).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfy (PL) with constant  $\alpha > 0$ . Then, for all  $t \geq 0$ ,

$$f(x_t) - f_\star \leq (f(x_0) - f_\star) \exp(-2\alpha t).$$

We present a few key properties of the PL inequality.

**Proposition 2.7** (strong convexity  $\Rightarrow$  PL  $\Rightarrow$  quadratic growth). Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\alpha > 0$ . The following implications hold.

1. If  $f$  is  $\alpha$ -convex, then  $f$  satisfies (PL) with constant  $\alpha$ .
2. If  $f$  satisfies (PL) with constant  $\alpha$ , then it satisfies the following **quadratic growth** property:

$$f(x) - f_\star \geq \frac{\alpha}{2} \inf_{x_\star \in \mathcal{X}_\star} \|x - x_\star\|^2, \quad \text{for all } x \in \mathbb{R}^d, \quad (\text{QG})$$

where  $\mathcal{X}_\star$  denotes the set of minimizers of  $f$ .

*Proof.*

1. Setting  $y = x_\star$  in (1.4), we obtain

$$\begin{aligned} -(f(x) - f_\star) &\geq \langle \nabla f(x), x_\star - x \rangle + \frac{\alpha}{2} \|x - x_\star\|^2 \\ &\geq -\|\nabla f(x)\| \|x_\star - x\| + \frac{\alpha}{2} \|x - x_\star\|^2 \geq -\frac{1}{2\alpha} \|\nabla f(x)\|^2, \end{aligned}$$

where the last inequality uses  $ab \leq \frac{\lambda}{2} a^2 + \frac{1}{2\lambda} b^2$  for all  $\lambda > 0$ .

2. Let  $(x_t)_{t \geq 0}$  denote the gradient flow for  $f$  started at  $x_0 = x$ . For simplicity, we present a proof *assuming* that the gradient flow converges to a point  $x_\star$ , although this assumption can be avoided (cf. [KNS16]). By Corollary 2.6, we see that  $x_\star \in \mathcal{X}_\star$ .

We start by observing that

$$\partial_t (\|x_t - x_0\|^2) = -2 \langle \nabla f(x_t), x_t - x_0 \rangle \leq 2 \|\nabla f(x_t)\| \|x_t - x_0\|$$

and hence

$$\partial_t \|x_t - x_0\| \leq \|\nabla f(x_t)\|.$$

We differentiate the following quantity:  $\mathcal{L}_t := \sqrt{\frac{\alpha}{2}} \|x_t - x_0\| + \sqrt{f(x_t) - f_\star}$ .

$$\dot{\mathcal{L}}_t \leq \sqrt{\frac{\alpha}{2}} \|\nabla f(x_t)\| - \frac{\|\nabla f(x_t)\|^2}{2\sqrt{f(x_t) - f_\star}} \leq 0,$$

where we applied (PL). Since  $\mathcal{L}_0 = \sqrt{f(x) - f_\star}$  and  $\mathcal{L}_\infty = \sqrt{\frac{\alpha}{2}} \|x - x_\star\|$ , we deduce the result from  $\mathcal{L}_0 \geq \mathcal{L}_\infty$ .

□

Hence, strong convexity implies (PL), but is (PL) truly weaker than convexity? Indeed, there are examples. In particular, the PL condition has been of interest in recent years because it holds for certain overparametrized models (Exercise 2.3).

We conclude this section by studying the implication of Lemma 2.1 alone. The fundamental theorem of calculus shows that

$$\frac{1}{t} \int_0^t \|\nabla f(x_s)\|^2 ds \leq \frac{f(x_0) - f(x_t)}{t} \leq \frac{f(x_0) - f_\star}{t}.$$

We therefore arrive at the following simple consequence.

**Corollary 2.8 (convergence of GF in gradient norm).** For any  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ ,

$$\min_{s \in [0, t]} \|\nabla f(x_s)\| \leq \sqrt{\frac{f(x_0) - f_\star}{t}}.$$

(In contrast, note that if we additionally assume convexity, then Exercise 2.1 shows that  $\|\nabla f(x_t)\| = O(1/t)$ .)

This implies there exists a sequence of times  $\{t_n\}_{n \in \mathbb{N}} \nearrow \infty$  such that  $\|\nabla f(x_{t_n})\| \rightarrow 0$ . (Indeed,  $\min_{s \in [n, 2n]} \|\nabla f(x_s)\| = O(1/n^{1/2})$ , so we can choose  $t_n \in [n, 2n]$ .) However, the gradient flow may not converge. Famously, it is a result of [Loj63] that for *real analytic*  $f$ , if the gradient flow remains bounded, then it does converge, and hence necessarily to a stationary point. Of course, such a stationary point may not be a global minimizer.

The idea of subsequent sections is to replicate the preceding analysis in discrete time.

## Bibliographical notes

My understanding of Theorem 2.4, Exercise 2.1, and Exercise 2.2 is based on extensive discussions with Jason M. Altschuler, Adil Salim, Andre Wibisono, and Ashia Wilson. The proof in Exercise 2.1 is taken from [OV01], and the extension in Exercise 2.2 to  $\alpha > 0$  is recorded in [LMW24, §F]. Both of these references pertain to the Langevin diffusion, but underneath the hood they make use of principles from optimization; see [Che25] for an introduction to this perspective.

The PL inequality is attributed to [Loj63; Pol63] and it was popularized in [KNS16]. The proof that (PL) implies the quadratic growth inequality goes back at least to the celebrated work of [OV00].



## Exercises

**Exercise 2.1.** Let  $f$  be convex. Show that the following quantity is decreasing,  $\mathcal{L}_t \leq 0$ :

$$\mathcal{L}_t := t^2 \|\nabla f(x_t)\|^2 + 2t(f(x_t) - f_\star) + \|x_t - x_\star\|^2.$$

Deduce the following gradient bound:

$$\|\nabla f(x_t)\|^2 \leq \frac{1}{t^2} \|x_0 - x_\star\|^2.$$

Moreover, use (1.4) to argue that  $2t(f(x_t) - f_\star) \leq t^2 \|\nabla f(x_t)\|^2 + \|x_t - x_\star\|^2$ , hence

$$f(x_t) - f_\star \leq \frac{1}{4t} \|x_0 - x_\star\|^2. \quad (2.1)$$

Note that this improves upon Theorem 2.4 by a factor of 2. Furthermore, show that (2.1) is sharp, as follows: for any  $R, t > 0$ , let  $f : x \mapsto \frac{R}{2t} \max\{0, x\}$ ,  $x_0 = R$ , and show that (2.1) holds with equality.

**Exercise 2.2.** Extend Exercise 2.1 to the case  $\alpha > 0$ . Toward this end, consider

$$\mathcal{L}_t := A_t \|\nabla f(x_t)\|^2 + 2B_t(f(x_t) - f_\star) + \|x_t - x_\star\|^2.$$

Choose  $A_t, B_t$  carefully to ensure that  $\mathcal{L}_t \leq -\alpha \mathcal{L}_t$ , and thereby deduce the following sharp bounds:

$$\|\nabla f(x_t)\|^2 \leq \frac{\alpha^2 \|x_0 - x_\star\|^2}{\exp(2\alpha t) (1 - \exp(-\alpha t))^2}, \quad f(x_t) - f_\star \leq \frac{\alpha \|x_0 - x_\star\|^2}{2(\exp(2\alpha t) - 1)}.$$

**Exercise 2.3.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be  $\alpha$ -convex with  $\alpha > 0$ , and let  $g : \mathbb{R}^d \rightarrow \mathbb{R}^n$  with  $d \geq n$ . Assume that  $g$  is surjective and that for all  $x \in \mathbb{R}^d$ , if  $\nabla g(x)$  denotes the Jacobian at  $x$  (interpreted as a  $d \times n$  matrix), then  $\nabla g(x)^\top \nabla g(x) \geq \sigma I_n$ . Show that the composition  $f \circ g$  satisfies (PL) with constant  $\alpha\sigma$ . Note that for  $d > n$ , there are typically multiple minimizers of  $f \circ g$ .

## 3 [1/21] Gradient descent: smooth case

In this section, we study the **gradient descent** algorithm:

$$x_{n+1} := x_n - h \nabla f(x_n). \quad (\text{GD})$$

From the perspective of numerical analysis, this is the *Euler* or *forward* discretization of (GF). Our aim is to show that if  $f$  is smooth, and the step size is sufficiently small (as a function of the smoothness), then the conclusions for (GF) transfer to (GD). Throughout this section, we assume that  $f$  is twice continuously differentiable and  $\beta$ -smooth.

Some of the results in this section pertain to a single step of (GD), so we use the following notation:

$$x^+ := x - h \nabla f(x) .$$

The first step is to establish the descent property.

**Lemma 3.1** (descent lemma). For any  $\beta$ -smooth  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , if  $h \leq 1/\beta$ , then

$$f(x^+) - f(x) \leq -\frac{h}{2} \|\nabla f(x)\|^2 .$$

*Proof.* By the smoothness inequality (1.7),

$$f(x^+) \leq f(x) + \langle \nabla f(x), x^+ - x \rangle + \frac{\beta}{2} \|x^+ - x\|^2 = f(x) - h \|\nabla f(x)\|^2 + \frac{\beta h^2}{2} \|\nabla f(x)\|^2 .$$

If  $h \leq 1/\beta$ , then  $-h(1 - \beta h/2) \leq -h/2$ . □

It is natural to state the subsequent results in terms of the following parameter.

**Definition 3.2.** Let  $f$  be  $\alpha$ -convex and  $\beta$ -smooth. Then, the **condition number** of  $f$  is defined to be the ratio  $\kappa := \beta/\alpha \geq 1$ .

When  $f$  is quadratic,  $f(x) = \frac{1}{2} \langle x, Ax \rangle$  with  $A$  symmetric, then  $\alpha, \beta$  correspond to the minimum and maximum eigenvalues of  $A$  respectively, and the ratio  $\beta/\alpha$  is known in numerical linear algebra as the condition number of the matrix  $A$ . Thus, Definition 3.2 provides a natural generalization of this notion. With this definition in hand, we now arrive at our first convergence result for (GD).

**Theorem 3.3** (contraction of GD). Let  $f$  be  $\alpha$ -convex and  $\beta$ -smooth. For all  $x, y \in \mathbb{R}^d$  and step size  $h \leq 1/\beta$ ,

$$\|y^+ - x^+\| \leq (1 - \alpha h)^{1/2} \|y - x\| .$$

*Proof.* Expanding the square,

$$\|y^+ - x^+\|^2 = \|y - x\|^2 - 2h \langle y - x, \nabla f(y) - \nabla f(x) \rangle + h^2 \|\nabla f(y) - \nabla f(x)\|^2.$$

By (3.5) in [Exercise 3.1](#) below, for  $h \leq 1/\beta$  and from (1.5) we have

$$\|y^+ - x^+\|^2 \leq \|y - x\|^2 - h \langle \nabla f(y) - \nabla f(x), y - x \rangle \leq (1 - \alpha h) \|y - x\|^2. \quad \square$$

In particular, if we take  $y = x_\star$ ,  $h = 1/\beta$ , and iterate, it yields

$$\|x_N - x_\star\| \leq \left(1 - \frac{1}{\kappa}\right)^{N/2} \|x_0 - x_\star\| \leq \exp\left(-\frac{N}{2\kappa}\right) \|x_0 - x_\star\|.$$

Thus, to obtain  $\|x_N - x_\star\| \leq \varepsilon$ , it suffices to take  $N \geq 2\kappa \log(\|x_0 - x_\star\|/\varepsilon)$ .

The essence of these proofs is that the first-order term (scaling as  $h$ ) replicates the continuous-time calculation, and we must apply smoothness in an appropriate way to control the second-order term (scaling as  $h^2$ ). In the above proof, note that if we naïvely use Lipschitzness of the gradient (1.8) to control the second-order term, it leads to the suboptimal choice of step size  $h = 1/(\beta\kappa)$ , and a contraction factor of  $(1 - 1/\kappa^2)^{1/2}$ . To obtain  $\|x_N - x_\star\| \leq \varepsilon$ , we would then have the estimate  $N \geq 2\kappa^2 \log(\|x_0 - x_\star\|/\varepsilon)$ , which is substantially worse. In conclusion, a bit of finesse is necessary. (In fact, [Theorem 3.3](#) can also be improved, and the sharp rate is derived in [Exercise 3.2](#).)

Next, we turn toward the analogue of [Theorem 2.4](#).

**Theorem 3.4 (convergence of GD in function value).** Let  $f$  be  $\alpha$ -convex and  $\beta$ -smooth. For any step size  $h \leq 1/\beta$ ,

$$\|x^+ - x_\star\|^2 \leq (1 - \alpha h) \|x - x_\star\|^2 - 2h (f(x^+) - f_\star). \quad (3.1)$$

Therefore,

$$f(x_N) - f_\star \leq \frac{\alpha}{2 \{(1 - \alpha h)^{-N} - 1\}} \|x_0 - x_\star\|^2. \quad (3.2)$$

When  $\alpha = 0$ , the right-hand side should be interpreted as its limiting value as  $\alpha \rightarrow 0$ , namely,  $\frac{1}{2Nh} \|x_0 - x_\star\|^2$ .

*Proof.* Expanding the square and applying convexity via (1.4),

$$\begin{aligned} \|x^+ - x_\star\|^2 &= \|x - x_\star\|^2 - 2h \langle \nabla f(x), x - x_\star \rangle + h^2 \|\nabla f(x)\|^2 \\ &\leq (1 - \alpha h) \|x - x_\star\|^2 - 2h (f(x) - f_\star) + h^2 \|\nabla f(x)\|^2. \end{aligned}$$

For  $h \leq 1/\beta$ , the descent lemma (Lemma 3.1) now implies (3.1).

The proof of (3.2), based on iterating the recursive inequality (3.1), is justified after Lemma 3.5 below.  $\square$

We remark for later use that the proof of (3.1) goes through even if we replace  $x_\star$  with any other point  $z \in \mathbb{R}^d$ , i.e.,

$$\|x^+ - z\|^2 \leq (1 - \alpha h) \|x - z\|^2 - 2h (f(x^+) - f(z)), \quad \text{for all } z \in \mathbb{R}^d. \quad (3.3)$$

Iterating (3.1) is a matter of unrolling the recursion, but in order to maintain the analogy with continuous time, we refer to the lemma below as “discrete Grönwall”.

**Lemma 3.5 (discrete Grönwall).** Suppose that for some  $A > 0$ ,

$$u_{n+1} \leq Au_n + B_n \quad \text{for } n = 0, 1, \dots, N-1.$$

Then,

$$u_N \leq A^N u_0 + \sum_{n=1}^N A^{N-n} B_{n-1}.$$

*Proof.* We multiply the given inequality by  $A^{-(n+1)}$  to form a telescoping sum:

$$A^{-N} u_N - u_0 = \sum_{n=0}^{N-1} A^{-(n+1)} (u_{n+1} - Au_n) \leq \sum_{n=0}^{N-1} A^{-(n+1)} B_n.$$

Rearrange to obtain the result.  $\square$

To complete the proof of Theorem 3.4, we apply Lemma 3.5 with  $u_n = \|x_n - x_\star\|^2$ ,  $A = 1 - \alpha h$ , and  $B_n = -2h (f(x_{n+1}) - f_\star)$ , yielding

$$2h \sum_{n=1}^N (1 - \alpha h)^{N-n} (f(x_n) - f_\star) \leq (1 - \alpha h)^N \|x_0 - x_\star\|^2.$$

For  $h \leq 1/\beta$ , the descent lemma (Lemma 3.1) implies  $f(x_n) - f_\star \geq f(x_N) - f_\star$ , so

$$f(x_N) - f_\star \leq \frac{\|x_0 - x_\star\|^2}{2h \sum_{n=1}^N (1 - \alpha h)^{-n}} = \frac{\alpha \|x_0 - x_\star\|^2}{2 \{(1 - \alpha h)^{-N} - 1\}}.$$

In particular, let us set  $h = 1/\beta$ . For  $\alpha > 0$  it yields

$$f(x_N) - f_\star \leq \frac{\alpha \|x_0 - x_\star\|^2}{2 \{(1 - 1/\kappa)^{-N} - 1\}}$$

and for  $\alpha = 0$ , it yields

$$f(x_N) - f_\star \leq \frac{\beta \|x_0 - x_\star\|^2}{2N}.$$

The proof of convergence under (PŁ) is strikingly easy.

**Theorem 3.6** (convergence of GD under PŁ). Let  $f$  be  $\beta$ -smooth and satisfy (PŁ) with constant  $\alpha > 0$ . Then, for all  $h \leq 1/\beta$ ,

$$f(x_N) - f_\star \leq (1 - \alpha h)^N (f(x_0) - f_\star).$$

*Proof.* By the descent lemma (Lemma 3.1) and (PŁ),

$$\begin{aligned} f(x^+) - f_\star &= f(x) - f_\star + f(x^+) - f(x) \leq f(x) - f_\star - \frac{h}{2} \|\nabla f(x)\|^2 \\ &\leq (1 - \alpha h) (f(x) - f_\star). \end{aligned} \quad \square$$

Finally, we present the result for obtaining a stationary point.

**Theorem 3.7.** Let  $f$  be  $\beta$ -smooth and  $h \leq 1/\beta$ . Then,

$$\min_{n=0,1,\dots,N-1} \|\nabla f(x_n)\| \leq \sqrt{\frac{2(f(x_0) - f_\star)}{Nh}}.$$

*Proof.* Telescope the descent lemma (Lemma 3.1):

$$\frac{h}{2N} \sum_{n=0}^{N-1} \|\nabla f(x_n)\|^2 \leq \frac{1}{N} \sum_{n=0}^{N-1} (f(x_n) - f(x_{n+1})) \leq \frac{f(x_0) - f_\star}{N}. \quad \square$$

We summarize the results for GD in Table 1.

Assumptions	Criterion	Iterations
$\alpha$ -convex, $\beta$ -smooth	$\ x_N - x_\star\  \leq \varepsilon$	$O(\kappa \log(R/\varepsilon))$
$\alpha$ -convex, $\beta$ -smooth	$f(x_N) - f_\star \leq \varepsilon$	$O(\kappa \log(\alpha R^2/\varepsilon))$
convex, $\beta$ -smooth	$f(x_N) - f_\star \leq \varepsilon$	$O(\beta R^2/\varepsilon)$
$\alpha$ -( <a href="#">PL</a> ), $\beta$ -smooth	$f(x_N) - f_\star \leq \varepsilon$	$O(\kappa \log(\Delta_0/\varepsilon))$
$\beta$ -smooth	$\min_{n=0,1,\dots,N-1} \ \nabla f(x_n)\  \leq \varepsilon$	$O(\beta \Delta_0/\varepsilon^2)$

Table 1: Rates for [GD](#) with step size  $1/\beta$ . Here,  $R := \|x_0 - x_\star\|$  and  $\Delta_0 := f(x_0) - f_\star$ .

**Example 3.8 (logistic regression revisited).** For fun, let us revisit logistic regression ([Example 1.2](#)) from a statistical lens. For concreteness, we consider Gaussian design,  $X_i \stackrel{\text{i.i.d.}}{\sim} \text{normal}(0, I)$ , and assume that the data is generated from the model with a true parameter  $\theta^\star$ . Let  $\widehat{\mathcal{L}}$  denote the MLE objective, let  $\mathcal{L} := \mathbb{E} \widehat{\mathcal{L}}$  denote the population risk, and let  $R := \|\theta^\star\| \geq 1$ . The state-of-the-art result [[CLM24](#)] shows that if  $n \gtrsim Rd$  for a sufficiently large implied constant,  $\widehat{\theta}_{\text{MLE}}$  exists with probability  $\geq 1 - \exp(-d)$  and satisfies the optimal risk bound  $\mathcal{L}(\widehat{\theta}_{\text{MLE}}) - \mathcal{L}(\theta^\star) \lesssim d/n$ .

In practice, we cannot compute  $\widehat{\theta}_{\text{MLE}}$  exactly, so we use optimization. From [[CLM24](#)], any estimator  $\widehat{\theta}$  satisfying  $\widehat{\mathcal{L}}(\widehat{\theta}) - \widehat{\mathcal{L}}(\widehat{\theta}_{\text{MLE}}) \lesssim d/n$  satisfies the same statistical risk bound as  $\widehat{\theta}_{\text{MLE}}$ , up to a universal constant. We take  $\widehat{\theta} = \widehat{\theta}_{\text{GD}}$  to be the output of [GD](#) after  $N$  steps, and check how large  $N$  must be in order for this to hold. As justified in [Exercise 3.3](#), we can expect an iteration complexity of  $N \asymp R^2 n/d$ .

## Bibliographical notes

My understanding of [Theorem 3.4](#) is again based on extensive discussions with Jason M. Altschuler, Adil Salim, Andre Wibisono, and Ashia Wilson.

## Exercises

**Exercise 3.1.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and  $\beta$ -smooth. Apply [Lemma 3.1](#) to the function  $y \mapsto f(y) - \langle \nabla f(x), y \rangle$  and observe that this function is minimized at  $x$  in order to prove

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2\beta} \|\nabla f(y) - \nabla f(x)\|^2. \quad (3.4)$$

From this, deduce that

$$\|\nabla f(y) - \nabla f(x)\|^2 \leq \beta \langle \nabla f(y) - \nabla f(x), y - x \rangle. \quad (3.5)$$

Finally, use the Cauchy–Schwarz inequality to show that  $\nabla f$  is  $\beta$ -Lipschitz, i.e., that (1.8) holds. Note that this proof that convexity and  $\beta$ -smoothness together imply (1.8) does not require  $f$  to be twice differentiable.

**Exercise 3.2.** Let  $f$  be  $\alpha$ -convex and  $\beta$ -smooth. Let  $T := \text{id} - h \nabla f$  denote the one-step GD mapping. By the fundamental theorem of calculus,

$$\begin{aligned} \|y^+ - x^+\| &= \|T(y) - T(x)\| = \left\| \int_0^1 \nabla T((1-t)x + ty) (y-x) dt \right\| \\ &\leq \left( \int_0^1 \|\nabla T((1-t)x + ty)\|_{\text{op}} dt \right) \|y-x\|. \end{aligned}$$

For any  $z \in \mathbb{R}^d$ , bound the eigenvalues of  $\nabla T(z)$  and show that the choice of step size  $h$  which minimizes the bound on  $\|\nabla T(z)\|_{\text{op}}$  is  $h = 2/(\alpha + \beta)$ . Deduce the sharp rate

$$\|y^+ - x^+\| \leq \frac{\kappa - 1}{\kappa + 1} \|y - x\|.$$

Note that for large  $\kappa$ , the contraction factor is approximately  $\exp(-2/\kappa)$ , so this improves upon the iteration complexity implied by Theorem 3.3 by a factor of nearly 4.

**Exercise 3.3.** What does Theorem 3.4 imply for logistic regression (Example 1.2)? In the setting of Example 3.8, use the fact that  $\lambda_{\max}(\frac{1}{n} \sum_{i=1}^n X_i X_i^\top) \lesssim 1$  with high probability<sup>3</sup> to justify the claimed  $R^2 n/d$  iteration complexity.

## 4 [1/23] Lower bounds for smooth optimization

The goal of this section is to establish lower complexity bounds for convex smooth optimization. Refer to §1.1 for a conceptual first discussion of the oracle model.

Before doing so, we present some reductions between the convex and strongly convex settings which save us some effort.

### 4.1 Reductions between the convex and strongly convex settings

For brevity, let us say that an algorithm *successfully optimizes* a function class  $\mathcal{F}$  in  $\phi(\mathcal{F}, R, \varepsilon)$  iterations if, given any  $f \in \mathcal{F}$  and  $x_0 \in \mathbb{R}^d$  with  $\|x_0 - x_\star\| \leq R$ , it outputs  $x$  with  $f(x) - f_\star \leq \varepsilon$  using no more than  $\phi(\mathcal{F}, R, \varepsilon)$  queries to a first-order oracle for  $f$ .

---

<sup>3</sup>This is a standard fact about the Wishart distribution; see, e.g., [Ver18, Theorem 4.4.5].

**Lemma 4.1.** Assume there is an algorithm which successfully optimizes the class of convex and  $\beta$ -smooth functions in  $\phi(\beta R^2/\varepsilon)$  iterations.

Then, there is an explicit algorithm which successfully optimizes the class of  $\alpha$ -convex and  $\beta$ -smooth functions in  $O(\phi(8\kappa) \log(\alpha R^2/\varepsilon))$  iterations.

*Proof.* Let  $f$  be  $\alpha$ -strongly convex and  $\beta$ -smooth, and apply the given algorithm to  $f$  to obtain a new point  $x_1$  with tolerance  $\varepsilon_1$ . By (QG), we have

$$\frac{\alpha}{2} \|x_1 - x_\star\|^2 \leq f(x_1) - f_\star \leq \varepsilon_1.$$

Set  $\varepsilon_1 = \alpha R^2/8$ , so that

$$\|x_1 - x_\star\| \leq \frac{1}{2} R = \frac{1}{2} \|x_0 - x_\star\|. \quad (4.1)$$

For  $\kappa := \beta/\alpha$ , this requires  $\phi(8\kappa)$  iterations. From (4.1), if we now repeat this procedure  $O(\log(\alpha R^2/\varepsilon))$  times, we can reach a point  $\tilde{x}$  satisfying  $\tilde{R} := \|\tilde{x} - x_\star\| \leq \sqrt{\varepsilon/\alpha}$ . Finally, apply the given algorithm one more time starting from  $\tilde{x}$  with target accuracy  $\varepsilon$  to obtain a point  $x$  with  $f(x) - f_\star \leq \varepsilon$ . The complexity of this final step is  $\phi(\beta \tilde{R}^2/\varepsilon) = \phi(\kappa)$ .  $\square$

For example, if we combine the  $\alpha = 0$  case of Theorem 3.4 with Lemma 4.1, taking  $\phi(x) = O(x)$ , we recover the  $\alpha > 0$  case of Theorem 3.4, up to constants.

**Lemma 4.2.** Assume there is an algorithm which successfully optimizes the class of  $\alpha$ -convex and  $\beta$ -smooth functions in  $\phi(\kappa) \log(\alpha R^2/\varepsilon)$  iterations.

Then, there is an explicit algorithm which successfully optimizes the class of convex and  $\beta$ -smooth functions in  $O(\phi(2\beta R^2/\varepsilon))$  iterations.

*Proof.* Let  $f$  be convex and  $\beta$ -smooth. We apply the given algorithm to the regularized function  $f_\delta := f + \frac{\delta}{2} \|\cdot - x_0\|^2$ , obtaining a point  $x$  such that  $f_\delta(x) \leq \min f_\delta + \varepsilon/2$ . If  $x_{\delta,\star}$  denotes the minimizer of  $f_\delta$ , then

$$f(x) \leq f_\delta(x) \leq f_\delta(x_{\delta,\star}) + \frac{\varepsilon}{2} \leq f_\delta(x_\star) + \frac{\varepsilon}{2} = f_\star + \frac{\delta}{2} \|x_0 - x_\star\|^2 + \frac{\varepsilon}{2}.$$

We now set  $\delta = \varepsilon/R^2$ , so that  $f(x) - f_\star \leq \varepsilon$ .

It remains to estimate the complexity. We first note that  $f_\delta(x_{\delta,\star}) \leq f_\delta(x_\star)$  implies  $\|x_0 - x_{\delta,\star}\| \leq \|x_0 - x_\star\|$ , so the initial distance to the minimizer of  $f_\delta$  is also bounded by  $R$ . We can assume that  $\varepsilon \leq \beta R^2$  (or else the minimization problem is trivial). Then, the smoothness of  $f_\delta$  is bounded by  $\beta + \delta \leq 2\beta$ , and the condition number of  $f_\delta$  is bounded by  $2\beta R^2/\varepsilon$ . Substitute these quantities into the complexity of the given algorithm.  $\square$



Thus, the  $\alpha > 0$  case of [Theorem 3.4](#) and [Lemma 4.2](#) recover the  $\alpha = 0$  case of [Theorem 3.4](#) up to constants.

Taken together, [Lemma 4.1](#) and [Lemma 4.2](#) show that the 0-convex and strongly convex settings are essentially equivalent to each other, in that an optimal method for one class yields an optimal method for the other class. Thus, we now aim to address the following question: what is the smallest possible  $\phi(\cdot)$ ?

## 4.2 Lower bounds

According to the discussion in §1.1, establishing a lower complexity bound requires showing that *any* algorithm which interacts with the first-order oracle using at most a prescribed number of queries cannot have performance better than the lower bound. Actually, although this is possible (see [\[NY83\]](#)), it is not especially easy. It was shown by Nesterov in an earlier edition of [\[Nes18\]](#) that by imposing natural restrictions on the class of algorithms under consideration, it is possible to establish the lower bounds in a more transparent way. Accordingly, his approach has become standard in the field, and it is the approach we adopt here as well. It does, however, have the drawback of not applying to general query algorithms; for example, it does not apply against randomized algorithms.

The class of algorithms we consider is the following one.

**Definition 4.3.** An algorithm is called a **gradient span** algorithm if it deterministically generates a sequence of points  $\{x_n\}_{n \in \mathbb{N}}$  such that for all  $n \in \mathbb{N}$ ,

$$x_{n+1} \in x_0 + \text{span}\{\nabla f(x_0), \dots, \nabla f(x_n)\}.$$

For example, [GD](#) is a gradient span algorithm. On the basis of this assumption, we now establish the following result; recall the asymptotic notation  $\gtrsim$ , which only hides a universal constant.

**Theorem 4.4 (lower bound for convex, smooth minimization).** For any  $1 \leq N \leq \frac{d-1}{2}$ ,  $\beta > 0$ , and  $x_0 \in \mathbb{R}^d$ , there exists a convex and  $\beta$ -smooth function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  such that for any gradient span algorithm,

$$f(x_N) - f_\star \gtrsim \frac{\beta \|x_0 - x_\star\|^2}{N^2}.$$

In other words, in order to obtain  $f(x_N) - f_\star \leq \varepsilon$ , the number of iterations must satisfy

$$N \gtrsim \sqrt{\frac{\beta \|x_0 - x_\star\|^2}{\varepsilon}}.$$

Before proving this result, we observe that by applying [Lemma 4.2](#) with  $\phi(x) \asymp \sqrt{x}$ , it yields the following corollary.

**Theorem 4.5** (lower bound for strongly convex, smooth minimization). For any  $0 < \alpha < \beta$ , any  $\varepsilon > 0$ , any  $d$  sufficiently large, and any  $x_0 \in \mathbb{R}^d$ , there exists an  $\alpha$ -convex and  $\beta$ -smooth function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  such that for any gradient span algorithm, in order to obtain  $f(x_N) - f_\star \leq \varepsilon$ , the number of iterations must satisfy

$$N \gtrsim \sqrt{d} \log \frac{\alpha \|x_0 - x_\star\|^2}{\varepsilon}.$$

*Proof of Theorem 4.4.* By translating the problem, we may assume  $x_0 = 0$ . The construction is based on the following function:

$$f_n : \mathbb{R}^d \rightarrow \mathbb{R}, \quad f_n(x) := \frac{\beta}{4} \left\{ \frac{1}{2} \left( x[1]^2 + \sum_{k=1}^{n-1} (x[k] - x[k+1])^2 + x[n]^2 \right) - x[1] \right\}.$$

For any  $v \in \mathbb{R}^d$ ,

$$\langle v, \nabla^2 f_n(x) v \rangle = \frac{\beta}{4} \left( v[1]^2 + \sum_{k=1}^{n-1} (v[k] - v[k+1])^2 + v[n]^2 \right) \leq \beta \|v\|^2,$$

so each  $f_n$  is convex and  $\beta$ -smooth.

We prove by induction that when we apply a gradient span algorithm to  $f_d$ , the  $n$ -th iterate  $x_n$  belongs to the subspace

$$\mathcal{V}_n := \{x \in \mathbb{R}^d : x[k] = 0 \text{ for all } k = n+1, \dots, d\}.$$

Clearly,  $x_0 \in \mathcal{V}_0$ . Inductively, suppose that  $x_k \in \mathcal{V}_k$  for all  $k \leq n$ . Then,

$$\nabla f_d(x_k) = \frac{\beta}{4} (x_k[1] e_1 + \sum_{j=1}^k (x_k[j] - x_k[j+1]) (e_j - e_{j+1})) - \frac{\beta}{4} e_1 \in \mathcal{V}_{k+1},$$

hence

$$x_{n+1} \in \text{span}\{\nabla f_d(x_0), \dots, \nabla f_d(x_n)\} \subseteq \mathcal{V}_{n+1}.$$

This completes the induction. Also, since  $f_N = f_d$  on  $\mathcal{V}_N$ , it follows that

$$f_d(x_N) = f_N(x_N) \geq (f_N)_\star.$$

The next step is to estimate  $(f_n)_\star := \min f_n$  for all  $n$ . By setting the gradient to zero,  $\nabla f_n(x_{n,\star}) = 0$ , we obtain the following system of equations:

$$\begin{aligned} 2x_{n,\star}[1] - x_{n,\star}[2] &= 1, \\ x_{n,\star}[k-1] - 2x_{n,\star}[k] + x_{n,\star}[k+1] &= 0, \quad \text{for } k = 1, \dots, n, \\ -x_{n,\star}[n-1] + 2x_{n,\star}[n] &= 0. \end{aligned}$$

The solution is  $x_{n,\star}[k] = 1 - \frac{k}{n+1}$  for all  $k \in [n]$ . Writing  $f_n(x) = \frac{\beta}{4} \{ \frac{1}{2} \langle x, A_n x \rangle - \langle e_1, x \rangle \}$ , the system above reads  $A_n x_{n,\star} = e_1$ , hence

$$(f_n)_\star = f_n(x_{n,\star}) = -\frac{\beta}{8} \langle e_1, x_{n,\star} \rangle = -\frac{\beta}{8} \left(1 - \frac{1}{n+1}\right).$$

Moreover,  $\|x_0 - x_{n,\star}\|^2 = \|x_{n,\star}\|^2 \leq n$ . Finally, it yields

$$\begin{aligned} f_d(x_N) - (f_d)_\star &\geq (f_N)_\star - (f_d)_\star = \frac{\beta}{8} \left( \frac{1}{N+1} - \frac{1}{d+1} \right) \\ &\geq \frac{\beta \|x_0 - x_{d,\star}\|^2}{8d} \left( \frac{1}{N+1} - \frac{1}{d+1} \right). \end{aligned}$$

Choosing  $d \asymp N$ , e.g.,  $d = 2N + 1$ , yields the stated lower bound.  $\square$

Notably, the iteration complexity lower bounds [Theorem 4.4](#) and [Theorem 4.5](#) are smaller than the bounds attained by [GD](#) in [Theorem 3.4](#) by a square root. As developed in the next sections, in fact the lower bounds are tight and [GD](#) is suboptimal.

We make two further remarks. First, it is perhaps surprising that the lower bound construction is a *quadratic* function; in some sense, quadratics are the hardest convex and smooth functions to optimize. Second, the lower bound requires the ambient dimension to be larger than the iteration count; this is crucial for the proof technique, which relies on the algorithm discovering one new dimension per iteration. This turns out to be fundamental because there are better methods in low dimension, for quadratics and even for general convex functions.

## Exercises

**Exercise 4.1.** In the setting of [Theorem 4.4](#) and using the same construction as in the proof, show that  $\|x_N - x_\star\|^2 \gtrsim \|x_0 - x_\star\|^2$ . In other words, in the 0-convex case, it is not possible to make progress in the sense of distance to the minimizer by more than a constant factor.

**Exercise 4.2.** We used the reductions from §4.1 to reduce the strongly convex lower bound to the 0-convex lower bound for the sake of brevity, but it is of course possible to develop the strongly convex lower bound directly. Consider the function

$$f : \mathbb{R}^\infty \rightarrow \mathbb{R}, \quad f(x) := \frac{\beta - \alpha}{8} \left\{ x[1]^2 + \sum_{n=1}^{\infty} (x[n] - x[n+1])^2 - 2x[1] \right\} + \frac{\alpha}{2} \|x\|^2.$$

By adapting the proof of Theorem 4.4, show that any gradient span algorithm satisfies

$$f(x_N) - f_\star \geq \frac{\alpha}{2} \|x_N - x_\star\|^2 \geq \frac{\alpha}{2} \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{2N} \|x_0 - x_\star\|^2.$$

## 5 [1/28–1/30] Acceleration

We now show that the lower bounds of Theorem 4.4 and Theorem 4.5 can be attained via algorithms which improve upon GD. This is known as the *acceleration* phenomenon in optimization. We begin with the quadratic case.

### 5.1 Quadratic case: the conjugate gradient method

In this section, the objective function is quadratic:

$$f : \mathbb{R}^d \rightarrow \mathbb{R}, \quad f(x) = \frac{1}{2} \langle x, Ax \rangle - \langle b, x \rangle,$$

where  $A$  is a symmetric matrix,  $A \succ 0$ . Note also that minimizing  $f$  corresponds to solving the system of equations  $Ax_\star = b$ . We now introduce the *conjugate gradient* method.

The method is succinctly described as follows:

$$x_{n+1} := \arg \min \{ f(x) \mid x \in x_0 + \text{span}\{\nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_n)\} \}. \quad (\text{CG})$$

This scheme is very natural in light of the definition of a gradient span algorithm (Definition 4.3) that we encountered for the lower bounds. However, it is not yet clear that (CG) can be implemented cheaply. Using the fact that  $f$  is quadratic, our aim is to show that (CG) can be rewritten as a simple iteration that uses one gradient query per step.

As is usually the case in linear algebra, instead of working with the set of vectors  $\{\nabla f(x_0), \nabla f(x_1), \dots, \nabla f(x_n)\}$ , it is more convenient to work with an *orthogonal* set  $\{p_0, p_1, \dots, p_n\}$ . Here, orthogonality is with respect to the inner product  $\langle \cdot, \cdot \rangle_A$ , i.e., we will require  $\langle p_i, A p_j \rangle = 0$  for all  $i \neq j$ . We start with  $p_0 := \nabla f(x_0)$ , and we write  $\mathcal{K}_n := \text{span}\{p_0, p_1, \dots, p_n\}$ . We must address the following two questions:

- Given  $\mathcal{K}_n$  and  $x_n$ , how can we compute  $x_{n+1} = \arg \min_{x_0 + \mathcal{K}_n} f$ ?
- Given  $\mathcal{K}_n$  and  $\nabla f(x_{n+1})$ , how can we compute  $p_{n+1}$  and thus  $\mathcal{K}_{n+1}$ ?

For the first question, we may assume inductively that  $x_n = \arg \min_{x_0 + \mathcal{K}_{n-1}} f$ , which means that  $\langle \nabla f(x_n), p_k \rangle = 0$  for all  $k < n$ . The next point is taken to be  $x_{n+1} = x_n + h_n p_n$ , chosen so that  $\langle \nabla f(x_{n+1}), p_k \rangle = 0$  for all  $k \leq n$ . Since  $\nabla f$  is linear,

$$\langle \nabla f(x_{n+1}), p_k \rangle = \langle \nabla f(x_n) + h_n A p_n, p_k \rangle.$$

For  $k < n$ , this equals zero by the inductive hypothesis on  $x_n$ , and the orthogonality of  $\{p_0, p_1, \dots, p_n\}$ . We choose  $h_n$  to ensure that this equals zero for  $k = n$  too:

$$h_n = -\frac{\langle \nabla f(x_n), p_n \rangle}{\|p_n\|_A^2}.$$

For the second question, we want to compute the Gram–Schmidt orthogonalization of  $\nabla f(x_{n+1})$  w.r.t.  $\{p_0, p_1, \dots, p_n\}$  in the  $\langle \cdot, \cdot \rangle_A$  inner product. We claim that  $\nabla f(x_{n+1})$  is already  $A$ -orthogonal to  $p_k$  for  $k < n$ , so that

$$p_{n+1} = \nabla f(x_{n+1}) - \langle \nabla f(x_{n+1}), p_n \rangle_A \frac{p_n}{\|p_n\|_A^2}. \quad (5.1)$$

To justify this, we show that for  $k < n$ ,  $\boxed{A p_k \in \mathcal{K}_{k+1}}$ , hence

$$\langle \nabla f(x_{n+1}), p_k \rangle_A = \langle \nabla f(x_{n+1}), A p_k \rangle = 0$$

using the fact shown above that  $\nabla f(x_{n+1})$  is orthogonal (in the usual inner product) to  $\mathcal{K}_n$ . Finally, the boxed equation is shown through the following lemma.

**Lemma 5.1.** For all  $n \in \mathbb{N}$ ,

$$\mathcal{K}_n = \text{span}\{p_0, A p_0, \dots, A^n p_0\}.$$

*Proof.* We proceed via induction, where the case  $n = 0$  is obvious. Assuming it holds at iteration  $n$ , let us show that  $\widetilde{\mathcal{K}}_{n+1} \in \widetilde{\mathcal{K}}_{n+1} := \text{span}\{p_0, A p_0, \dots, A^{n+1} p_0\}$ . By (5.1), it suffices to show that  $\nabla f(x_{n+1}) \in \widetilde{\mathcal{K}}_{n+1}$ . However, as discussed above,  $\nabla f(x_{n+1}) = \nabla f(x_n) + h_n A p_n = p_0 + h_0 A p_0 + \dots + h_n A p_n \in \widetilde{\mathcal{K}}_{n+1}$ .

Conversely, we must show that  $A^{n+1} p_0 \in \mathcal{K}_{n+1}$ . Since  $A^n p_0 \in \mathcal{K}_n$ , we can write  $A^n p_0 = \sum_{k=0}^n c_k p_k$ , thus  $A^{n+1} p_0 = \sum_{k=0}^n c_k A p_k$ . By the inductive hypothesis, each  $A p_k$  for  $k < n$  belongs to  $\mathcal{K}_n$ , so it suffices to have  $A p_n \in \mathcal{K}_{n+1}$ . However, we can observe that  $A p_n = h_n^{-1} (\nabla f(x_{n+1}) - \nabla f(x_n)) \in \mathcal{K}_{n+1}$  by (5.1).  $\square$

**Definition 5.2.** The subspaces  $\{\mathcal{K}_n\}_{n \in \mathbb{N}}$  are called **Krylov subspaces**.

Finally, let us write the iterations in a form which is convenient for implementation. Note first that  $\langle \nabla f(x_n), \nabla f(x_{n+1}) \rangle = 0$  (indeed,  $\nabla f(x_{n+1})$  is orthogonal to all of  $\mathcal{K}_n$ ). So,

$$\frac{\langle \nabla f(x_{n+1}), p_n \rangle_A}{\|p_n\|_A^2} = \frac{\langle \nabla f(x_{n+1}), \nabla f(x_{n+1}) - \nabla f(x_n) \rangle}{h_n \|p_n\|_A^2} = -\frac{\|\nabla f(x_{n+1})\|^2}{\langle \nabla f(x_n), p_n \rangle}$$

and  $\|\nabla f(x_n)\|^2 = \langle \nabla f(x_n), \nabla f(x_n) \rangle = \langle \nabla f(x_n), p_n \rangle$  using (5.1) and the fact that  $\nabla f(x_n)$  is orthogonal to  $\mathcal{K}_{n-1}$ . This yields the following iteration, where we write  $r_n := Ax_n - b = \nabla f(x_n)$  for the residual.

$$x_{n+1} = x_n - \frac{\|r_n\|^2}{\langle p_n, A p_n \rangle} p_n, \quad r_{n+1} = r_n - \frac{\|r_n\|^2}{\langle p_n, A p_n \rangle} A p_n, \quad p_{n+1} = r_{n+1} + \frac{\|r_{n+1}\|^2}{\|r_n\|^2} p_n.$$

This algorithm requires one matrix-vector multiplication per iteration, namely, the computation of  $A p_n$ .

Note that if  $p_{n+1} = 0$ , then  $\nabla f(x_{n+1}) \in \mathcal{K}_n$ , yet  $\nabla f(x_{n+1}) \perp \mathcal{K}_n$  and thus  $\nabla f(x_{n+1}) = 0$ ,  $x_{n+1} = x_\star$ . Since  $p_{d+1} = 0$  (an orthogonal set in  $\mathbb{R}^d$  cannot have more than  $d$  non-zero elements), we arrive at the following conclusion.

**Theorem 5.3 (termination of CG).** The CG algorithm returns the exact minimizer in at most  $d$  iterations.

Let us now show that CG can find an approximate minimizer at the accelerated rate.

**Theorem 5.4 (accelerated convergence for CG).** Let  $0 < \alpha I \leq A \leq \beta I$ . Then, CG outputs  $x_N$  satisfying  $f(x_N) - f_\star \leq \varepsilon$  in  $N = O(\sqrt{\kappa} \log \frac{f(x_0) - f_\star}{\varepsilon})$  iterations.

*Proof.* By the descent lemma (Lemma 3.1) and the defining property of CG,

$$f(x_{n+1}) \leq f\left(x_n - \frac{1}{\beta} \nabla f(x_n)\right) \leq f(x_n) - \frac{1}{2\beta} \|\nabla f(x_n)\|^2,$$

so that

$$f(x_0) - f_\star \geq \frac{1}{2\beta} \sum_{n=0}^{N-1} \|\nabla f(x_n)\|^2.$$

On the other hand, since  $\nabla f(x_n) \perp x_{k+1} - x_k$  for  $k < n$ ,

$$f_\star - f(x_n) \geq \langle \nabla f(x_n), x_\star - x_n \rangle = \langle \nabla f(x_n), x_\star - x_0 \rangle.$$

If we sum these inequalities and use orthogonality of the gradients,

$$\begin{aligned} N(f(x_N) - f_\star) &\leq \sum_{n=0}^{N-1} (f(x_n) - f_\star) \leq \left\langle \sum_{n=0}^{N-1} \nabla f(x_n), x_0 - x_\star \right\rangle \leq \left\| \sum_{n=0}^{N-1} \nabla f(x_n) \right\| \|x_0 - x_\star\| \\ &\leq \left( \sum_{n=0}^{N-1} \|\nabla f(x_n)\|^2 \right)^{1/2} \sqrt{\frac{2(f(x_0) - f_\star)}{\alpha}} \leq 2\sqrt{\kappa} (f(x_0) - f_\star). \end{aligned}$$

Let  $N$  be such that  $f(x_N) - f_\star \geq (f(x_0) - f_\star)/2$ . The inequality above then implies that  $N \leq 4\sqrt{\kappa}$ . Thus, every  $4\sqrt{\kappa}$  iterations, the objective gap decreases by a factor of 2.  $\square$

By applying the restart strategy as in [Lemma 4.1](#), one can also show an iteration complexity scaling with  $\sqrt{\kappa}$  in the strongly convex case. However, we instead give a different proof in order to explain the classical link with polynomial approximation.

Due to [Lemma 5.1](#),  $x_N - x_0 \in \mathcal{K}_{N-1}$  can be written in the form  $x_N - x_0 = \sum_{n=0}^{N-1} c_n A^n p_0$ , so  $x_N - x_\star = x_0 - x_\star + \sum_{n=0}^{N-1} c_n A^{n+1} (x_0 - x_\star) = P_N(A) (x_0 - x_\star)$  where  $P_N$  is a polynomial of degree at most  $N$  satisfying  $P_N(0) = 1$ . Conversely, if  $Q_N$  is any other degree- $N$  polynomial with  $Q_N(0) = 1$ , then  $\tilde{x}_N := x_0 + A^{-1} (Q_N(A) - I) p_0 \in x_0 + \mathcal{K}_{N-1}$  satisfies  $\tilde{x}_N - x_\star = x_0 - x_\star + A^{-1} (Q_N(A) - I) p_0 = Q_N(A) (x_0 - x_\star)$ .

This equivalence, together with the fact that the output  $x_N$  of [CG](#) minimizes  $f$  over  $x_0 + \mathcal{K}_{N-1}$ , shows that

$$f(x_N) - f_\star \leq \frac{1}{2} \min \{ \|Q_N(A) (x_0 - x_\star)\|_A^2 : Q_N \in \mathbb{R}_{\leq N}[X], Q_N(0) = 1 \},$$

where  $\mathbb{R}_{\leq N}[X]$  denotes the set of polynomials with real-valued coefficients and with degree at most  $N$ . Furthermore, since  $A$  and  $Q_N(A)$  commute,

$$\|Q_N(A) (x_0 - x_\star)\|_A^2 \leq \|Q_N(A)\|_{\text{op}}^2 \|x_0 - x_\star\|_A^2 \leq \left( \max_{[\lambda_{\min}(A), \lambda_{\max}(A)]} |Q_N(\lambda)|^2 \right) \|x_0 - x_\star\|_A^2.$$

We have arrived at the following result.

**Lemma 5.5 (CG and polynomial approximation).** Assume that  $0 < \alpha I \leq A \leq \beta I$ . Then, the output  $x_N$  of [CG](#) satisfies

$$f(x_N) - f_\star \leq \min \left\{ \max_{\lambda \in [\alpha, \beta]} |Q_N(\lambda)|^2 : Q_N \in \mathbb{R}_{\leq N}[X], Q_N(0) = 1 \right\} (f(x_0) - f_\star).$$

Informally, this result states that **CG** performs as well as the best possible degree- $N$  polynomial in  $A$ . To bound the rate of convergence of **CG**, it therefore remains to exhibit a judicious polynomial  $Q_N$ . This is accomplished by the family of Chebyshev polynomials, on which many volumes have been written.

**Definition 5.6.** The degree- $n$  **Chebyshev polynomial**  $T_n$  is defined so that  $\cos(n\theta) = T_n(\cos \theta)$  for all  $\theta \in \mathbb{R}$ .

It is not obvious at first glance that  $T_n$  is indeed a degree- $n$  polynomial, but this can be established via trigonometric identities. The use of the Chebyshev polynomials to establish a rate of convergence for **CG** is explored in [Exercise 5.1](#).

Here, we point out another interesting fact that arises from this connection. Recall from the proof of [Lemma 5.5](#) that if we can compute  $\tilde{x}_N := x_0 + A^{-1} (Q_N(A) - I) p_0$ , then it incurs error at most  $f(\tilde{x}_N) - f_\star \leq (\max_{\lambda \in [\alpha, \beta]} |Q_N(\lambda)|^2) (f(x_0) - f_\star)$ . In particular, rather than using **CG**, we can try to compute the polynomial  $x \mapsto (Q_N(x) - 1)/x$  directly, where  $Q_N$  is the polynomial in [Exercise 5.1](#) which witnesses the fast convergence of **CG**. Although we omit the details, it is worth noting that the family of Chebyshev polynomials satisfies a so-called three-term recurrence:

$$T_{n+1}(x) = 2x T_n(x) - T_{n-1}(x), \quad x \in \mathbb{R}.$$

In fact, orthogonal families of polynomials usually do.<sup>4</sup> From an algorithmic standpoint, it leads to an optimization algorithm of the form

$$x_{n+1} = c_0 A x_n + c_1 x_{n-1} + c_2 b,$$

where  $c_0, c_1, c_2 \in \mathbb{R}$  are fixed coefficients. Note that unlike **GD**,  $x_{n+1}$  depends on the previous *two* iterates. This is often referred to as *momentum*, and also forms the basis for acceleration for general convex functions.

**Remark 5.7 (practicality of **CG**).** Solving the linear system  $Ax = b$  via Gaussian elimination requires  $O(d^3)$  operations and is numerically unstable, whereas for well-conditioned matrices  $A$ , **CG** returns an approximate solution in  $\tilde{O}(\sqrt{\kappa})$  iterations, each of which requires a matrix-vector multiplication. A matrix-vector multiplication requires  $O(d^2)$  time in the worst case, but can be faster if  $A$  is sparse. In practice, **CG** is widely used, especially when combined with other strategies such as preconditioning.

<sup>4</sup>This arises in connection with second-order differential operators.



## 5.2 General case: continuous time

Although it does not follow the historical development of events, we begin our treatment of acceleration for general convex smooth functions in continuous time. As identified in [SBC16], the continuous-time ODE is

$$\begin{aligned}\dot{x}_t &= p_t, \\ \dot{p}_t &= -\nabla f(x_t) - \gamma_t p_t.\end{aligned}\tag{AGF}$$

We refer to (AGF) as the *accelerated gradient flow*, and the variable  $p_t$  admits the physical interpretation of momentum (for a particle with unit mass). The dynamics consists of two parts: the equations

$$\begin{aligned}\dot{x}_t &= p_t, \\ \dot{p}_t &= -\nabla f(x_t)\end{aligned}$$

are known as Hamilton's equations, and they are the standard first-order reformulation of Newton's law of motion  $\ddot{x}_t = -\nabla f(x_t)$  with potential energy  $f$ . Hamilton's equations conserve the energy (or Hamiltonian)  $H(x, p) := f(x) + \frac{1}{2} \|p\|^2$ , and this conservation property is perhaps undesirable for an optimization algorithm which seeks to minimize  $f$ . Thus, the second part of the dynamics,  $\dot{p}_t = -\gamma_t p_t$  adds a dissipative *friction* force, where  $\gamma_t \geq 0$  is a possibly time-varying coefficient of friction.

In the case where  $f$  is merely assumed to be convex, it turns out that the right choice of friction coefficient is  $\gamma_t = 3/t$ . This is mysterious at first sight and was obtained by taking the continuous-time limit of Nesterov's discrete algorithm in the next subsection. We begin with a convergence analysis in this setting. (Similar caveats as for §2 apply here; we assume that  $f$  is smooth, that it admits a minimizer  $x_\star$ , and that (AGF) is well-posed.)

**Theorem 5.8 (convergence of AGF under convexity).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and let  $(x_t)_{t \geq 0}$  evolve along AGF with  $\gamma_t = 3/t$  and  $p_0 = 0$ . Then, for all  $t \geq 0$ ,

$$f(x_t) - f_\star \leq \frac{2 \|x_0 - x_\star\|^2}{t^2}.$$

*Proof.* Consider the auxiliary point  $z_t := x_t + \frac{t}{2} p_t$ , and the Lyapunov function

$$\mathcal{L}_t := \frac{t^2}{2} (f(x_t) - f_\star) + \|z_t - x_\star\|^2.$$

The computation below shows that  $\dot{\mathcal{L}}_t \leq 0$ , which implies the result. The choice of Lyapunov function is mysterious, so we partially demystify it after the proof.

Straightforward differentiation and convexity yield

$$\begin{aligned}\dot{\mathcal{L}}_t &= t (f(x_t) - f_\star) + \frac{t^2}{2} \langle \nabla f(x_t), p_t \rangle - t \langle \nabla f(x_t), z_t - x_\star \rangle \\ &= t (f(x_t) - f_\star) - t \langle \nabla f(x_t), x_t - x_\star \rangle \leq 0.\end{aligned}\quad \square$$

Although the Lyapunov function above appears fortuitous, it can be derived in a reasonably systematic manner; see [Exercise 5.2](#). The strongly convex case is similar, and is left as [Exercise 5.3](#).

**Theorem 5.9** (convergence of AGF under strong convexity). Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be  $\alpha$ -convex and let  $(x_t)_{t \geq 0}$  evolve along AGF with  $\gamma_t = 2\sqrt{\alpha}$  and  $p_0 = 0$ . For all  $t \geq 0$ ,

$$f(x_t) - f_\star \leq 2 \exp(-\sqrt{\alpha} t) (f(x_0) - f_\star).$$

Recall that under convexity and  $\alpha$ -convexity, the objective gap  $f(x_t) - f_\star$  for GF converges at the rates  $O(1/t)$  and  $O(\exp(-2\alpha t))$  respectively. On the other hand, for AGF, the convergence happens at the rates  $O(1/t^2)$  and  $O(\exp(-\sqrt{\alpha} t))$  respectively. This is strongly suggestive of the square root factor speed-up, that is, *acceleration*. However, we caution that it is dangerous to deduce conclusions from continuous-time analysis alone. For example, we can run any ODE faster, which can make the continuous-time convergence rate arbitrarily fast; however, this does not translate into a better discrete-time algorithm, since speeding up time makes the ODE more unstable and therefore requires a smaller step size for discretization.

So how, then, can we discretize AGF? Part of the subtlety of acceleration is that not all discretizations work. For example, we could consider

$$\begin{aligned}x_{n+1} &\approx x_n + h p_{n+1}, \\ p_{n+1} &\approx p_n - h \nabla f(x_n) - \gamma_n h p_n\end{aligned}$$

which is equivalent to the update

$$x_{n+1} = x_n - h^2 \nabla f(x_n) + (1 - \gamma_n h) (x_n - x_{n-1}).$$

Or, if we do not presume to know the coefficients for the discrete-time scheme in advance, we could write the update as

$$x_{n+1} = x_n - \eta_n \nabla f(x_n) + \theta_n (x_n - x_{n-1}).$$

In other words, we take a gradient step and then apply momentum. This is known as Polyak's heavy ball method, and although it can be tuned to converge at the rate

of [CG](#) for quadratic objectives, this same tuning leads to divergence for general convex functions [\[LRP16\]](#). On the other hand, the optimal method in the next subsection can be written in the form

$$x_{n+1} = x_n + \theta_n (x_n - x_{n-1}) - \eta_n \nabla f(x_n + \theta_n (x_n - x_{n-1})) .$$

In other words, we add momentum and then take a gradient step.

### 5.3 General case: discrete time

The acceleration phenomenon is undoubtedly one of the most elusive and fascinating aspects of optimization, so it is no surprise that it has been explored through many different angles over the course of countless research papers. At this junction, we must choose how to present the method and in what level of detail.

Having explored acceleration carefully in the quadratic case and in continuous time, here we follow the expedient route by giving perhaps the most direct and shortest proof, at the cost of generality and intuition.<sup>5</sup>

We analyze the following method with  $x_{-1} = x_0$ :

$$x_{n+1} := x_n + \theta_n (x_n - x_{n-1}) - \frac{1}{\beta} \nabla f(x_n + \theta_n (x_n - x_{n-1})) . \quad (\text{AGD})$$

**Theorem 5.10 (convergence of AGD).** Let  $f$  be convex and  $\beta$ -smooth. Define the sequence:  $\lambda_0 := 0$  and  $\lambda_{n+1} := \frac{1}{2} (1 + \sqrt{1 + 4\lambda_n^2})$  for  $n \in \mathbb{N}$ . Set  $\theta_n := (\lambda_n - 1)/\lambda_{n+1}$ . Then, [AGD](#) satisfies

$$f(x_N) - f_\star \leq \frac{2\beta \|x_0 - x^\star\|^2}{N^2} .$$

*Proof.* Let  $y_n := x_n + \theta_n (x_n - x_{n-1})$ , so that  $x_{n+1} = y_n - \frac{1}{\beta} \nabla f(y_n)$ . Recall from [\(3.3\)](#) that for any  $z \in \mathbb{R}^d$ , it holds that

$$\|x_{n+1} - z\|^2 \leq \|y_n - z\|^2 - \frac{2}{\beta} (f(x_{n+1}) - f(z)) .$$

Rearranging, it yields

$$f(x_{n+1}) - f(z) \leq \frac{\beta}{2} (\|y_n - z\|^2 - \|x_{n+1} - z\|^2) = -\frac{\beta}{2} \|x_{n+1} - y_n\|^2 - \beta \langle x_{n+1} - y_n, y_n - z \rangle .$$

---

<sup>5</sup>Perhaps I will change my mind in a future edition of these notes.

We apply this inequality with two points,  $z = x_n$  and  $z = x_\star$ . By multiplying the first inequality by  $\lambda_{n+1} - 1 \geq 0$  and adding it to the second inequality, it implies

$$\begin{aligned} & (\lambda_{n+1} - 1) (f(x_{n+1}) - f(x_n)) + f(x_{n+1}) - f_\star \\ & \leq -\frac{\beta\lambda_{n+1}}{2} \|x_{n+1} - y_n\|^2 - \beta \langle x_{n+1} - y_n, \lambda_{n+1} y_n - (\lambda_{n+1} - 1) x_n - x_\star \rangle \\ & = \frac{\beta}{2\lambda_{n+1}} (\|\lambda_{n+1} y_n - (\lambda_{n+1} - 1) x_n - x_\star\|^2 - \|\lambda_{n+1} x_{n+1} - (\lambda_{n+1} - 1) x_n - x_\star\|^2), \end{aligned}$$

where the last line uses the identity  $\|a\|^2 + 2\langle a, b \rangle = \|a + b\|^2 - \|b\|^2$ . Our goal is to produce a telescoping sum, which is the case if we ensure that

$$\lambda_{n+1} x_{n+1} - (\lambda_{n+1} - 1) x_n = \lambda_{n+2} y_{n+1} - (\lambda_{n+2} - 1) x_{n+1}.$$

By substituting in  $y_{n+1} = x_{n+1} + \theta_{n+1} (x_{n+1} - x_n)$ , some algebra shows that it suffices to take  $\theta_{n+1} = (\lambda_{n+1} - 1)/\lambda_{n+2}$ .

After multiplying the above inequality by  $\lambda_{n+1}$  and summing, we find that

$$\frac{\beta}{2} \|\lambda_1 y_0 - (\lambda_1 - 1) x_0 - x_\star\|^2 \geq \sum_{n=0}^{N-1} \{\lambda_{n+1}^2 (f(x_{n+1}) - f_\star) - \lambda_{n+1} (\lambda_{n+1} - 1) (f(x_n) - f_\star)\}.$$

We also want the right-hand side to telescope, so we set  $\lambda_{n+1} (\lambda_{n+1} - 1) = \lambda_n^2$ , which yields the recursion  $\lambda_{n+1} = \frac{1}{2} (1 + \sqrt{1 + 4\lambda_n^2})$ . With  $\lambda_0 = 0$ , it yields

$$f(x_N) - f_\star \leq \frac{\beta \|y_0 - x_\star\|^2}{2\lambda_N^2} = \frac{\beta \|x_0 - x_\star\|^2}{2\lambda_N^2}.$$

Finally, it is straightforward to show by induction that  $\lambda_N \geq N/2$ . □

By applying the reduction in [Lemma 4.1](#), it also yields an accelerated algorithm for the strongly convex case, i.e., an algorithm that achieves  $f(x_N) - f_\star \leq \varepsilon$  in  $O(\sqrt{\kappa} \log \frac{\alpha R^2}{\varepsilon})$  iterations, where  $R := \|x_0 - x_\star\|$ .

**Example 5.11.** If we apply the accelerated method to logistic regression (see [Example 3.8](#)), it improves the iteration complexity from  $O(R^2 n/d)$  to  $O(R\sqrt{n/d})$ .

## Bibliographical notes

The simple proof of [Theorem 5.4](#) is taken from [\[NY83\]](#). The discussion on Chebyshev polynomials follows [\[Vis12\]](#).

The literature on acceleration is too large to be surveyed here, but we mention a recent result in a somewhat different direction: what is the best rate of [GD](#) just by changing the step sizes? Thus, we consider the iteration  $x_{n+1} = x_n - h_n \nabla f(x_n)$ , with the only freedom being to choose the sequence  $\{h_n\}_{n \in \mathbb{N}}$ . It turns out a constant step size schedule is not optimal, and as established in [\[AP24a; AP24b\]](#), the so-called silver step size schedule achieves the rates of [Lemma 4.1](#) and [Lemma 4.2](#) with  $\phi(x) = x^{\log_\rho 2} \approx x^{0.786}$  with  $\rho := 1 + \sqrt{2}$ . This is a rate intermediate between the unaccelerated rate of [GD](#) and the accelerated rate of [AGD](#).

## Exercises

**Exercise 5.1.** Define the polynomial  $Q_n(x) = T_n(\frac{\alpha+\beta-2x}{\beta-\alpha})/T_n(\frac{\alpha+\beta}{\beta-\alpha})$ . Show that  $Q_n(0) = 1$  and use [Definition 5.6](#) to establish the identity

$$T_n(x) = \frac{1}{2} \left( (x - \sqrt{x^2 - 1})^n + (x + \sqrt{x^2 - 1})^n \right) \quad \text{for } x \in [-1, 1].$$

One can show that this identity actually holds for all  $x \in \mathbb{R}$ . Use this to show that

$$\max_{x \in [\alpha, \beta]} |Q_n(x)| \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^n.$$

Note that by combining this with [Lemma 5.5](#), it yields an exponential rate of convergence for [CG](#) matching the lower bound of [Exercise 4.2](#).

**Exercise 5.2.** To better understand the proof of [Theorem 5.8](#), consider a Lyapunov function of the form

$$\mathcal{L}_t = \|x_t - x_\star\|^2 + a_t \langle x_t - x_\star, p_t \rangle + b_t \|p_t\|^2 + c_t (f(x_t) - f_\star).$$

Note that this is the most general Lyapunov function consisting of a combination of a quadratic function in  $x_t - x_\star$  and  $p_t$ , as well as the objective gap; here, it is crucial that we include the mixed term  $a_t \langle x_t - x_\star, p_t \rangle$ . Our goal is to choose the coefficients  $a_t, b_t, c_t$  so that  $\dot{\mathcal{L}}_t \leq 0$ .

Compute the derivative in time of  $\mathcal{L}_t$  along [AGF](#) with  $\gamma_t = 3/t$ , and apply convexity to the term  $\langle \nabla f(x_t), x_t - x_\star \rangle$ . In the resulting expression, since the terms  $\langle x_t - x_\star, p_t \rangle$  and  $\langle \nabla f(x_t), p_t \rangle$  do not have definite signs, ensure that the coefficients in front of these terms

vanish through a suitable choice of  $a_t, b_t, c_t$ . Show that this leads to  $a_t = t + \bar{a}t^3$  for some  $\bar{a} \geq 0$ . Next, from the remaining terms, obtain the condition  $\dot{b}_t \leq \min\{\frac{a_t}{2}, \frac{6b_t}{t} - a_t\}$ , which implies  $3\dot{b}_t \leq 6b_t/t$ , hence we consider  $b_t = b_0 + \bar{b}t^2$  for some  $b_0, \bar{b} \geq 0$ . Furthermore, argue that we must take  $\bar{a} = 0$  and  $\bar{b} = \frac{1}{4}$ . To ensure that  $\mathcal{L}_0$  only depends on  $\|x_0 - x_\star\|$ , we set  $b_0 = c_0 = 0$ . Finally, check that with these choices, we have  $b_t \geq a_t^2/4$ , which is necessary to ensure that  $\mathcal{L}_t \geq c_t (f(x_t) - f_\star)$ .

Show that the Lyapunov function derived in this way coincides with the one used in [Theorem 5.8](#).

**Exercise 5.3.** Prove [Theorem 5.9](#).

*Hint:* Let  $z_t := x_t + \frac{2}{\gamma} p_t$  and consider

$$\mathcal{L}_t := f(x_t) - f_\star + \frac{\alpha}{2} \|z_t - x_\star\|^2.$$

## 6 [2/4–2/13] Non-smooth convex optimization

Thus far, we have considered the *unconstrained* minimization of convex and *smooth* functions  $f$ . The next step is to consider a far more general class of problems by allowing for constraints and non-smoothness.

The two issues are related. To minimize  $f$  over a convex set  $\mathcal{C}$ , it is equivalent to minimize  $f + \chi_{\mathcal{C}}$  over all of  $\mathbb{R}^d$ , where  $\chi_{\mathcal{C}}$  is the convex indicator function for  $\mathcal{C}$ :

$$\chi_{\mathcal{C}}(x) := \begin{cases} 0, & x \in \mathcal{C}, \\ +\infty, & x \notin \mathcal{C}. \end{cases} \quad (6.1)$$

In this reformulation, the objective function is allowed to take the value  $+\infty$  and is certainly non-smooth. Even if we do not reformulate the problem in this way, convex constraint sets often arise as the intersection of primitive constraints:  $\mathcal{C} = \{f_i \leq 0 \text{ for all } i \in [m]\}$ . This is equivalent to  $\mathcal{C} = \{\max_{i \in [m]} f_i \leq 0\}$ , and the function  $\max_{i \in [m]} f_i$  is non-smooth.

On the other hand, without strong convexity, it is not guaranteed that  $f$  admits a minimizer over all of  $\mathbb{R}^d$  (e.g.,  $f$  is a linear function, or consider the exponential function over  $\mathbb{R}$ ). It often makes sense to consider non-smooth minimization over bounded sets. Thus, we tackle constraints and non-smoothness together.

Although we do not assume smoothness, we still need some minimal regularity for the function  $f$ . As justified in [Lemma 6.7](#), convex functions are actually Lipschitz continuous in the interior of their domains, so it is natural to take as our new function class under consideration the class of convex and Lipschitz functions over bounded convex sets.

## 6.1 Convex analysis

We now work with convex functions  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ . The fact that  $f$  can now take on the value  $+\infty$  leads to some technical issues, but it allows us to seamlessly handle constraints. Convexity can be defined in the usual way, but it is sometimes convenient to instead work with the epigraph.

**Definition 6.1.** The **epigraph** of  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is the following subset of  $\mathbb{R}^d \times \mathbb{R}$ :

$$\text{epi } f := \{(x, t) \in \mathbb{R}^d \times \mathbb{R} : f(x) \leq t\}.$$

**Definition 6.2.** A function  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is **convex** if for all  $x, y \in \mathbb{R}^d$  and all  $t \in [0, 1]$ , it holds that

$$f((1-t)x + ty) \leq (1-t)f(x) + tf(y).$$

Equivalently,  $f$  is convex if and only if  $\text{epi } f$  is a convex set.

**Definition 6.3.** The **domain** of a function  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is the set

$$\text{dom } f := \{x \in \mathbb{R}^d : f(x) < \infty\}.$$

The first point to emphasize is that at this level of generality,  $f$  can still be quite pathological. Indeed, consider the following function:

$$f(x) := \begin{cases} 0, & \|x\| < 1, \\ \phi(x), & \|x\| = 1, \\ +\infty, & \|x\| > 1, \end{cases} \quad (6.2)$$

where  $\phi$  is an *arbitrary* non-negative function defined on the sphere  $\{\|\cdot\| = 1\}$ . Then, one can check that  $f$  is convex. However,  $\phi$  need not be continuous or be coherent in any way whatsoever. To avoid these types of situations, the basic regularity property that we impose is that  $f$  is lower semicontinuous.

**Definition 6.4.** A function  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is **lower semicontinuous** if for all sequences  $\{x_n\}_{n \in \mathbb{N}}$  converging to a point  $x \in \mathbb{R}^d$ , it holds that

$$f(x) \leq \liminf_{n \rightarrow \infty} f(x_n).$$

In other words, when we pass to the limit of a convergent sequence, the value of  $f$  can only drop down. One way to motivate the relevance of this condition for convex optimization is that we often consider suprema  $f = \sup_{\omega \in \Omega} f_\omega$  where  $\{f_\omega\}_{\omega \in \Omega}$  is a collection of continuous functions; in fact, in many cases, we consider suprema of affine functions. When  $\Omega$  is finite, we know that the maximum of finitely many continuous functions is continuous. But when  $\Omega$  is infinite, the suprema of infinitely many continuous functions need not be continuous. The class of lower semicontinuous functions is the smallest class of functions which contains all continuous functions and is closed under taking arbitrary suprema. Further properties are explored in [Exercise 6.1](#).

It follows from that exercise that  $f$  is convex and lower semicontinuous if and only if its epigraph is closed and convex. So, when it comes to functions, we impose convexity and lower semicontinuity; and when it comes to sets, we impose convexity and closedness. For example, one can also check that the convex indicator  $\chi_C$  is lower semicontinuous if and only if  $C$  is closed. We use the following terminology.<sup>6</sup>

**Definition 6.5.** A convex function  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is **regular** if: it is not identically equal to  $+\infty$ , it is lower semicontinuous, and its domain has non-empty interior.

Note that the definition excludes one more pathological case, the function  $f(x) = +\infty$  for all  $x \in \mathbb{R}^d$ , which is of no interest to us. Since the domain of a convex function is a convex set, if it has empty interior then it must be contained in a lower-dimensional affine space, and when we restrict to that space, the domain then has a non-empty interior; this is usually summarized by saying that any non-empty convex set has a non-empty *relative interior*. We do not delve into the details here, but this is why we regard the condition that the domain has non-empty interior as “without loss of generality”.

We also note that in the proof of existence of a minimizer, it is really only lower semicontinuity that matters.

**Lemma 6.6 (existence of minimizer).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be lower semicontinuous and its level sets be bounded. Then, there exists a global minimizer of  $f$ .

*Proof.* The proof is the same as for [Lemma 1.7](#), except that lower semicontinuity substitutes for continuity. □

---

<sup>6</sup>This is not standard terminology but it is convenient.



**Regularity.** Our next order of business is to establish properties of regular convex functions which allow us to manipulate them in proofs. In particular, we show that they are “almost” differentiable, even though we did not assume it a priori; the source of this regularity is the convexity condition.

**Lemma 6.7 (Lipschitz continuity).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be convex and let  $x_0 \in \text{int dom } f$ . Then,  $f$  is locally Lipschitz continuous around  $x_0$ .

*Proof.* We may assume that  $x_0 = 0$ . Since 0 belongs to the interior of  $\text{dom } f$ , we can fit a simplex centered at the origin inside the domain: namely, there exists  $\varepsilon > 0$  such that  $\mathcal{C} := \text{conv}\{\pm \varepsilon e_k : k \in [d]\}$  belongs to  $\text{dom } f$ . First, we show that  $f$  is bounded on  $\mathcal{C}$ : the upper bound follows because  $f(\pm \varepsilon e_k) < \infty$  for all  $k \in [d]$  and the maximum of  $f$  over  $\mathcal{C}$  is attained at one of the vertices (why?). For the lower bound, by convexity we have  $f(x) \geq 2f(0) - f(-x) \geq 2f(0) - \max_{\mathcal{C}} f$  for all  $x \in \mathcal{C}$ .

Next, we show that  $f$  is Lipschitz on the smaller set  $\mathcal{C}' := \text{conv}\{\pm \frac{\varepsilon}{2} e_k : k \in [d]\}$ . The point is that there is a constant  $c_{d,\varepsilon} > 0$  such that for all  $x, y \in \mathcal{C}'$ , there is a point  $y^+ \in \mathcal{C}$  such that the line segment from  $x$  to  $y$  is contained in the line segment from  $x$  to  $y^+$ , and the extension is not too short:  $\|y^+ - x\| \geq c_{d,\varepsilon}$ . Then, by convexity,

$$f(y) = f\left(\frac{\|y^+ - y\|}{\|y^+ - x\|} x + \frac{\|y - x\|}{\|y^+ - x\|} y^+\right) \leq \frac{\|y^+ - y\|}{\|y^+ - x\|} f(x) + \frac{\|y - x\|}{\|y^+ - x\|} f(y^+),$$

hence

$$f(y) - f(x) \leq \frac{\|y - x\|}{\|y^+ - x\|} (f(y^+) - f(x)) \leq \frac{\sup_{\mathcal{C}} f - \inf_{\mathcal{C}} f}{c_{d,\varepsilon}} \|y - x\|.$$

Interchanging  $x$  and  $y$  proves the Lipschitz bound.  $\square$

This lemma shows that locally near  $x_0$ ,  $f(x)$  grows at most linearly in the distance  $\|x - x_0\|$  (as opposed to, say,  $\sqrt{\|x - x_0\|}$ ). This suggests that  $f$  may be differentiable at  $x_0$ . This is not quite right, because  $f$  may have a kink at  $x_0$ , but nevertheless we can find an appropriate substitute for differentiability.

**Definition 6.8.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be convex. We say that  $p \in \mathbb{R}^d$  is a **subgradient** of  $f$  at  $x$  if for all  $y \in \mathbb{R}^d$ , it holds that

$$f(y) \geq f(x) + \langle p, y - x \rangle. \quad (6.3)$$

We denote the set of subgradients of  $f$  at  $x$  as  $\partial f(x)$ , and we refer to this set as the **subdifferential** of  $f$  at  $x$ . Also, we set

$$\partial f := \{(x, p) \in \mathbb{R}^d \times \mathbb{R}^d : p \in \partial f(x)\}.$$

Note that by definition, if  $0 \in \partial f(x)$ , then  $x$  is a global minimizer of  $f$ .

If  $f$  is differentiable at  $x_0 \in \text{int dom } f$ , then  $\partial f(x_0)$  is a singleton:  $\partial f(x_0) = \{\nabla f(x_0)\}$  (Exercise 6.2). However, the subdifferential can be multi-valued. A key example is the absolute value function,  $f : x \mapsto |x|$ , for which  $\partial f(0) = [-1, 1]$ .

For the purpose of optimization, it is enough to have at least one subgradient, which is the content of the following theorem.

**Theorem 6.9 (subdifferential).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be a regular convex function. If  $x_0 \in \text{int dom } f$ , then  $\partial f(x_0)$  is **non-empty**, bounded, convex, and closed.

We follow a traditional route of deducing the non-emptiness from a separation theorem. The proof of the following result is deferred.

**Theorem 6.10 (supporting hyperplane).** Let  $\mathcal{C}$  be a closed and convex set, and let  $x \in \partial \mathcal{C}$ . Then, there exists a non-zero  $p \in \mathbb{R}^d$  such that

$$\langle p, x \rangle \leq \inf_{\mathcal{C}} \langle p, \cdot \rangle.$$

*Proof of Theorem 6.9.* Since  $(x_0, f(x_0)) \in \partial \text{epi } f$ , and  $\text{epi } f$  is closed and convex (by regularity of  $f$ ), there is a supporting hyperplane  $(p, q)$ :

$$\langle p, x_0 \rangle + q f(x_0) \leq \inf_{(x,t) \in \text{epi } f} \{ \langle p, x \rangle + q t \}.$$

We can normalize the coefficients so that  $\|p\|^2 + q^2 = 1$ , and we note that  $q \geq 0$ .

If  $x$  is sufficiently close to  $x_0$ , then

$$\langle p, x_0 - x \rangle \leq q (f(x) - f(x_0)) \leq Lq \|x - x_0\|,$$

where  $L$  is the Lipschitz constant of  $f$  near  $x_0$ . Taking  $x = x_0 - \varepsilon p$  for small  $\varepsilon > 0$ , we deduce that  $\|p\| \leq Lq$ , hence from the normalization condition,  $q \neq 0$ . Thus, for any  $x \in \text{dom } f$ , we deduce that

$$f(x) \geq f(x_0) - \frac{1}{q} \langle p, x - x_0 \rangle,$$

thus,  $-p/q \in \partial f(x_0)$ .

The set  $\partial f(x_0)$  is closed and convex as an intersection of the constraints in (6.3). Boundedness follows from Exercise 6.3.  $\square$

**Constraints.** When the constraint set  $\mathcal{C}$  is simple, it is reasonable to suppose that we can compute the projection onto  $\mathcal{C}$ . We study some properties of this projection operator.

**Definition 6.11.** Let  $\mathcal{C}$  be closed and convex. The **projection onto  $\mathcal{C}$**  is the mapping  $\Pi_{\mathcal{C}} : \mathbb{R}^d \rightarrow \mathcal{C}$  defined by

$$\Pi_{\mathcal{C}}(x) := \arg \min_{y \in \mathcal{C}} \|y - x\|^2.$$

The “arg min” is non-empty because  $\mathcal{C}$  is closed, and the uniqueness of the minimizer follows from a strict convexity argument as in [Lemma 1.10](#). When  $\mathcal{C}$  is a linear subspace, then  $\Pi_{\mathcal{C}}$  coincides with the linear algebra definition of projection, and in this case  $\Pi_{\mathcal{C}}$  is linear. In general, however,  $\Pi_{\mathcal{C}}$  is a *non-linear* operator.

The following lemma characterizes the projection.

**Lemma 6.12 (characterization of projection).** Let  $\mathcal{C}$  be closed and convex, and let  $x \notin \mathcal{C}$ . Then,  $\Pi_{\mathcal{C}}(x)$  is the unique point satisfying the following condition:

$$\langle \Pi_{\mathcal{C}}(x) - x, x' - \Pi_{\mathcal{C}}(x) \rangle \geq 0 \quad \text{for all } x' \in \mathcal{C}. \quad (6.4)$$

*Proof.* As in the proof of [Lemma 1.8](#), the first-order necessary condition for optimality reads  $\langle \Pi_{\mathcal{C}}(x) - x, v \rangle \geq 0$ . However, because the optimization problem is constrained to lie in  $\mathcal{C}$ , this time we do not have the inequality for all  $v$ , but only for  $v$  of the form  $x' - \Pi_{\mathcal{C}}(x)$  where  $x' \in \mathcal{C}$ .  $\square$

This lemma furnishes the following important property.

**Lemma 6.13 (convex projections are non-expansive).** Let  $\mathcal{C}$  be closed and convex. Then, for all  $x, y \in \mathbb{R}^d$ ,

$$\|\Pi_{\mathcal{C}}(y) - \Pi_{\mathcal{C}}(x)\| \leq \|y - x\|.$$

*Proof.* By (6.4),

$$\begin{aligned} \langle \Pi_{\mathcal{C}}(x) - x, \Pi_{\mathcal{C}}(y) - \Pi_{\mathcal{C}}(x) \rangle &\geq 0, \\ \langle \Pi_{\mathcal{C}}(y) - y, \Pi_{\mathcal{C}}(x) - \Pi_{\mathcal{C}}(y) \rangle &\geq 0. \end{aligned}$$

Adding these inequalities yields

$$\|\Pi_{\mathcal{C}}(y) - \Pi_{\mathcal{C}}(x)\|^2 \leq \langle \Pi_{\mathcal{C}}(y) - \Pi_{\mathcal{C}}(x), y - x \rangle \leq \|\Pi_{\mathcal{C}}(y) - \Pi_{\mathcal{C}}(x)\| \|y - x\|. \quad \square$$

Actually, we can now return to prove the supporting hyperplane theorem.

*Proof of Theorem 6.10.* First, we show that if  $\mathcal{C}$  is a closed convex set and  $x \notin \mathcal{C}$ , then we can separate  $\mathcal{C}$  from  $x$ . Namely, by (6.4), the vector  $p := \Pi_{\mathcal{C}}(x) - x$  is non-zero and satisfies

$$\inf_{x' \in \mathcal{C}} \langle p, x' \rangle \geq \langle p, \Pi_{\mathcal{C}}(x) \rangle = \|\Pi_{\mathcal{C}}(x) - x\|^2 + \langle p, x \rangle \geq \langle p, x \rangle.$$

To prove the supporting hyperplane theorem, note that since  $x \in \partial\mathcal{C}$ , there is a sequence of points  $\{x_n\}_{n \in \mathbb{N}}$  which lies outside of  $\mathcal{C}$ , such that  $x_n \rightarrow x$ . For each  $n$ , let  $p_n$  be a hyperplane that separates  $\mathcal{C}$  from  $x_n$ , and by normalizing we may assume that  $\|p_n\| = 1$ . Since  $\{p_n\}_{n \in \mathbb{N}}$  is a bounded sequence, it contains a subsequence which converges to some unit vector  $p$ . By taking limits, it is easy to see that  $p$  is a supporting hyperplane.  $\square$

## 6.2 Projected subgradient methods

Methods for constrained optimization differ based on what they assume about the constraint set. The first method we study assumes access to the projection mapping  $\Pi_{\mathcal{C}}$  for the set  $\mathcal{C}$ . This assumption is appropriate when the set  $\mathcal{C}$  is particularly “simple”, e.g.,  $\mathcal{C}$  is the ball  $\mathcal{C} = \{\|\cdot\| \leq R\}$ , in which case the projection can be computed in closed form. When  $\mathcal{C}$  is more complex, e.g.,  $\mathcal{C}$  is a polytope, we need more sophisticated methods.

Projected subgradient descent is the following method:

$$x_{n+1} := \Pi_{\mathcal{C}}\left(x_n - h \frac{p_n}{\|p_n\|}\right), \quad p_n \in \partial f(x_n). \quad (\text{PSD})$$

Note that we use the normalized subgradient  $p_n/\|p_n\|$ . If we think about the example of the absolute value function  $|\cdot|$  with subdifferential  $[-1, 1]$  at the origin, we see that the magnitude of an arbitrary element of the subdifferential need not be informative. Instead, the intuition behind non-smooth optimization is to use the subgradients as separating directions: in particular, by convexity,  $f(x) - f(x_n) \geq \langle p_n, x - x_n \rangle$ , so any minimizer must lie on one side of the hyperplane defined by  $p_n$ .

We let  $x_{\star}$  denote a minimizer of  $f$  over the closed convex set  $\mathcal{C}$ , and  $f_{\star} := f(x_{\star})$ .

**Theorem 6.14 (convergence of PSD).** Let  $f$  be convex and  $L$ -Lipschitz continuous on the closed convex set  $\mathcal{C}$ . Then, PSD satisfies

$$f\left(\frac{1}{N} \sum_{n=0}^{N-1} x_n\right) - f_\star \leq \frac{1}{N} \sum_{n=0}^{N-1} (f(x_n) - f_\star) \leq \frac{L}{2Nh} \|x_0 - x_\star\|^2 + \frac{Lh}{2}.$$

In particular, by setting  $h = R/\sqrt{N}$ , where  $R$  is an upper bound on  $\|x_0 - x_\star\|$ , it yields the convergence rate

$$f\left(\frac{1}{N} \sum_{n=0}^{N-1} x_n\right) - f_\star \leq \frac{LR}{\sqrt{N}}.$$

*Proof.* The first inequality holds by convexity, so we focus on the second. The idea is similar to the proof of Theorem 3.4, except that instead of using smoothness to handle the error term, we use Lipschitzness. By expanding the squared distance to the minimizer,

$$\begin{aligned} \|x_{n+1} - x_\star\|^2 &= \left\| \Pi_{\mathcal{C}}\left(x_n - h \frac{p_n}{\|p_n\|}\right) - \Pi_{\mathcal{C}}(x_\star) \right\|^2 \leq \left\| x_n - h \frac{p_n}{\|p_n\|} - x_\star \right\|^2 \\ &= \|x_n - x_\star\|^2 - \frac{2h}{\|p_n\|} \langle p_n, x_n - x_\star \rangle + h^2 \\ &\leq \|x_n - x_\star\|^2 - \frac{2h}{\|p_n\|} (f(x_n) - f_\star) + h^2, \end{aligned}$$

where we used Lemma 6.13. Since  $\|p_n\| \leq L$  for all  $n$  (Exercise 6.3), we sum the inequalities:

$$\frac{1}{N} \sum_{n=0}^{N-1} (f(x_n) - f_\star) \leq \frac{L}{2Nh} \|x_0 - x_\star\|^2 + \frac{Lh}{2}. \quad \square$$

Thus, the averaged iterate  $\bar{x}_N$  satisfies  $f(\bar{x}_N) - f_\star \leq \varepsilon$  provided  $N \geq L^2 R^2 / \varepsilon^2$ . Note that this convergence rate is substantially worse than the one for the smooth case (Theorem 3.4). Another difference is that the descent lemma (Lemma 3.1) is available in the smooth case which implies monotonic decrease of the objective value; here, there is no descent lemma, so the guarantee only holds for the averaged iterate. The analysis can also be performed under strong convexity, see Exercise 6.5.

Interestingly, if we only assume that  $f$  is  $L$ -Lipschitz continuous over  $B(x_\star, R)$ , rather than on all of  $\mathcal{C}$ , it is still possible to show that  $\min_{n=0, \dots, N-1} f(x_n) - f_\star \leq LR/\sqrt{N}$ , although the proof becomes more involved [Nes18, §3.2.3].

The analysis above shows that when the projection operator is cheap to compute, optimization under constraints is straightforward provided that we interleave the gradient steps with projection steps. We next tackle a more general setting in which we separate out the constraints into a “simple” set  $\mathcal{C}$  for which we can compute the projection operator, and additional functional constraints  $\{f_i \leq 0 \text{ for all } i \in [m]\}$ . Thus, we consider

$$\min\{f(x) \mid x \in \mathcal{C}, f_i(x) \leq 0 \text{ for all } i \in [m]\}.$$

We assume that  $f, f_1, \dots, f_m$  are all regular convex functions, and write  $f_{\max} := \max_{i \in [m]} f_i$ . The next algorithm is known as the projected subgradient method with functional constraints. For  $n = 0, 1, \dots, N - 1$ :

- If  $f_{\max}(x_n) \leq \varepsilon$ , set

$$x_{n+1} := \Pi_{\mathcal{C}}\left(x_n - \frac{\varepsilon}{\|p_n\|^2} p_n\right), \quad p_n \in \partial f(x_n).$$

- Otherwise, set

$$x_{n+1} := \Pi_{\mathcal{C}}\left(x_n - \frac{f_{\max}(x_n)}{\|p_n\|^2} p_n\right), \quad p_n \in \partial f_{\max}(x_n).$$

The algorithm requires computing elements of the subdifferential for the function  $\max_{i \in [m]} f_i$ . We therefore first identify this subdifferential.

**Lemma 6.15 (subdifferential of a maximum).** Let  $f_1, \dots, f_m$  be regular convex functions. Then, for all  $x \in \mathbb{R}^d$ ,

$$\partial\left(\max_{i \in [m]} f_i\right)(x) = \text{conv}\left\{\partial f_i(x) \mid i \in [m], f_i(x) = \max_{j \in [m]} f_j(x)\right\}.$$

*Proof.* ( $\supseteq$ ) Let  $f_{\max} := \max_{i \in [m]} f_i$  and  $I_{\star}(x) := \{i \in [m] : f_i(x) = f_{\max}(x)\}$ . If  $\lambda$  is a probability vector and  $p_i \in \partial f_i(x)$  for all  $i \in I_{\star}(x)$ , then

$$f_{\max}(y) \geq \sum_{i \in I_{\star}(x)} \lambda_i f_i(y) \geq \sum_{i \in I_{\star}(x)} \lambda_i (f_i(x) + \langle p_i, y - x \rangle) = f_{\max}(x) + \left\langle \sum_{i \in I_{\star}(x)} \lambda_i p_i, y - x \right\rangle.$$

Hence,  $\sum_{i \in I_{\star}(x)} \lambda_i p_i \in \partial f_{\max}(x)$ .

( $\subseteq$ ) Since the purpose of this lemma from the perspective of these notes is simply to compute an element of  $\partial f_{\max}(x)$ , we omit the proof of this direction. It can be proven, e.g., via Lagrangian duality or via more subdifferential theory.  $\square$

The next theorem provides the convergence rate for the method.

**Theorem 6.16** (convergence of PSD under functional constraints). Let  $f, f_1, \dots, f_m$  be convex and  $L$ -Lipschitz on the closed convex set  $\mathcal{C}$ . Then, PSD under functional constraints satisfies

$$\min\{f(x_n) \mid n = 0, 1, \dots, N-1, f_{\max}(x_n) \leq \varepsilon\} - f_\star \leq \varepsilon \quad (6.5)$$

provided that

$$N \geq \frac{L^2 \|x_0 - x_\star\|^2}{\varepsilon^2}.$$

The theorem says that after  $N$  iterations, we can find a point  $\hat{x}_N$  which almost satisfies the functional constraints, in the sense that  $f_{\max}(\hat{x}_N) \leq \varepsilon$ , and moreover  $f(\hat{x}_N) - f_\star \leq \varepsilon$ . The number of iterations is no more than the case without functional constraints.

*Proof of Theorem 6.16.* There are two cases for the algorithm. If the iteration  $n$  belongs to the first case, then as we saw in the proof of Theorem 6.14,

$$\|x_{n+1} - x_\star\|^2 \leq \|x_n - x_\star\|^2 - \frac{2\varepsilon}{\|p_n\|^2} (f(x_n) - f_\star) + \frac{\varepsilon^2}{\|p_n\|^2}.$$

If  $f(x_n) - f_\star \leq \varepsilon$ , then since  $f_{\max}(x_n) \leq \varepsilon$  (by the definition of the first case), we have met the success condition (6.5). Otherwise,  $f(x_n) - f_\star > \varepsilon$ , and the inequality above implies

$$\|x_{n+1} - x_\star\|^2 < \|x_n - x_\star\|^2 - \frac{\varepsilon^2}{\|p_n\|^2} \leq \|x_n - x_\star\|^2 - \frac{\varepsilon^2}{L^2}.$$

What happens in the second case? Here, we *also* show that  $\|x_{n+1} - x_\star\| < \|x_n - x_\star\|$ : since  $x_\star$  satisfies the functional constraints and  $x_n$  does not, the subgradient  $p_n \in \partial f_{\max}(x_n)$  still acts as a separating hyperplane. Indeed,

$$\begin{aligned} \|x_{n+1} - x_\star\|^2 &= \left\| \Pi_{\mathcal{C}} \left( x_n - \frac{f_{\max}(x_n)}{\|p_n\|^2} p_n \right) - \Pi_{\mathcal{C}}(x_\star) \right\|^2 \leq \left\| x_n - \frac{f_{\max}(x_n)}{\|p_n\|^2} p_n - x_\star \right\|^2 \\ &= \|x_n - x_\star\|^2 - \frac{2f_{\max}(x_n)}{\|p_n\|^2} \langle p_n, x_n - x_\star \rangle + \frac{f_{\max}(x_n)^2}{\|p_n\|^2} \\ &\leq \|x_n - x_\star\|^2 - \frac{2f_{\max}(x_n)}{\|p_n\|^2} f_{\max}(x_n) + \frac{f_{\max}(x_n)^2}{\|p_n\|^2} < \|x_n - x_\star\|^2 - \frac{\varepsilon^2}{L^2}. \end{aligned}$$

Summing these inequalities across the iterations yields

$$\|x_N - x_\star\|^2 < \|x_0 - x_\star\|^2 - \frac{N\varepsilon^2}{L^2}.$$

For  $N \geq L^2 \|x_0 - x_\star\|^2 / \varepsilon^2$ , this is not possible unless we reach the success condition (6.5) by iteration  $N$ .  $\square$

**Example 6.17 (soft-margin SVM).** An example of a problem that can be tackled via projected subgradient methods is soft-margin support vector machine (SVM) classification. Suppose that we have a dataset  $\{(x_i, y_i)\}_{i \in [n]}$ , where  $x_i \in \mathbb{R}^d$  and  $y_i \in \{\pm 1\}$ . The output of the soft-margin SVM is the classifier  $x \mapsto \text{sgn}(\langle \theta^\star, x \rangle)$ , where  $\theta^\star$  minimizes

$$\theta \mapsto \frac{1}{n} \sum_{i=1}^n \ell_{\text{hinge}}(y_i, \langle \theta, x_i \rangle) + \frac{\lambda}{2} \|\theta\|^2.$$

Here,  $\ell_{\text{hinge}}(y, \hat{y}) := \max\{0, 1 - y\hat{y}\}$  is the hinge loss,  $\lambda > 0$  is a regularization parameter, and we have omitted the bias term (which can be handled by augmenting the feature vector  $x$  as usual). This objective is strongly convex and Lipschitz over bounded sets, so we can apply projected subgradient descent (projecting onto, e.g., a Euclidean ball).

### 6.3 Cutting plane methods

Non-smooth optimization uses subgradient directions in order to “localize” the solution set. Pursuing this line of reasoning further leads to the family of cutting plane methods.

Suppose that we wish to minimize  $f$  over a bounded, closed, convex set  $\mathcal{C}$ . Let  $\mathcal{C}_\star$  denote the set of minimizers. The idea is to construct a sequence of convex sets  $\mathcal{C} = \mathcal{C}_0, \mathcal{C}_1, \mathcal{C}_2, \dots$ , which shrink toward  $\mathcal{C}_\star$ . The set  $\mathcal{C}_n$  represents possible candidates for the solution to the problem at iteration  $n$ .

If  $x_n \in \mathcal{C}_n$  and  $p_n \in \partial f(x_n)$ , then the subgradient inequality reads

$$0 \geq f(x_\star) - f(x_n) \geq \langle p_n, x_\star - x_n \rangle \quad \text{for all } x_\star \in \mathcal{C}_\star.$$

Thus,

$$\mathcal{C}_\star \subseteq \mathcal{C}_n \cap \{x \in \mathbb{R}^d : \langle p_n, x \rangle \leq \langle p_n, x_n \rangle\}.$$

We can take  $\mathcal{C}_{n+1}$  to be any superset of the right-hand side above.

To finish specifying the scheme, we need a rule for choosing the points  $x_n$  and the sets  $\mathcal{C}_n$ , with the goal of  $\mathcal{C}_n$  shrinking as fast as possible. The key is the following lemma from convex geometry, which we do not prove.



**Lemma 6.18 (Grünbaum).** Let  $\mathcal{C} \subseteq \mathbb{R}^d$  be a convex body (i.e., a compact convex set with non-empty interior) and let  $x_{\mathcal{C}}$  denote the *centroid* of  $\mathcal{C}$ :  $x_{\mathcal{C}} := (\text{vol } \mathcal{C})^{-1} \int_{\mathcal{C}} x \, dx$ . Then, for any half-space  $\mathcal{H}$  containing  $x_{\mathcal{C}}$ ,

$$\frac{\text{vol}(\mathcal{C} \cap \mathcal{H})}{\text{vol}(\mathcal{C})} \geq \left(\frac{d}{d+1}\right)^d \geq \frac{1}{e},$$

where  $e \approx 2.72$  is a numerical constant.

Consequently, if we choose  $x_n$  to be the centroid of  $\mathcal{C}_n$  and set

$$\mathcal{C}_{n+1} = \mathcal{C}_n \cap \{\langle p_n, \cdot \rangle \leq \langle p_n, x_n \rangle\}, \quad x_n = x_{\mathcal{C}_n}, \quad (\text{CoGM})$$

then Grünbaum's inequality shows that  $\text{vol}(\mathcal{C}_n \setminus \mathcal{C}_{n+1})/\text{vol}(\mathcal{C}_n) \leq 1/e$ , or

$$\frac{\text{vol}(\mathcal{C}_{n+1})}{\text{vol}(\mathcal{C}_n)} \leq 1 - \frac{1}{e}.$$

Thus, we cut away a constant fraction of the volume at each iteration. This is known as the *center of gravity* method.

As stated, **CoGM** is not a practical method. The feasible set  $\mathcal{C}_n$  at iteration  $n$  can be quite complicated, making it prohibitively expensive to compute its centroid. Centroids can be computed via Markov chain Monte Carlo (MCMC) methods for numerical integration, with guarantees available due to recent advances in log-concave sampling, but it is generally understood that this is a more difficult computational problem than the original convex optimization problem we set out to solve. Nevertheless, **CoGM** achieves the optimal complexity bound in the oracle model, so let us analyze its efficiency.

**Theorem 6.19 (center of gravity).** Let  $D := \text{diam } \mathcal{C}$  and let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and  $L$ -Lipschitz on  $\mathcal{C}$ . Then, **CoGM** satisfies

$$f(x_{N-1}) - f_{\star} \leq DL \left(1 - \frac{1}{e}\right)^{N/d}.$$

*Proof.* By the argument above, at iteration  $N$ ,  $\text{vol}(\mathcal{C}_N)/\text{vol}(\mathcal{C}) \leq \lambda^N$ , where we can take  $\lambda = 1 - 1/e$ . Now consider the set  $\hat{\mathcal{C}} := (1-t)x_{\star} + t\mathcal{C}$ , where we choose  $t$  so that  $\text{vol}(\hat{\mathcal{C}}) > \text{vol}(\mathcal{C}_N)$ ; since  $\text{vol}(\hat{\mathcal{C}}) = t^d \text{vol}(\mathcal{C})$ , we can take any  $t > \lambda^{N/d}$ . With this choice, there exists  $\hat{x} \in \hat{\mathcal{C}} \setminus \mathcal{C}_N$ . By the definition of  $\mathcal{C}_N$ ,

$$f(x_{N-1}) - f_{\star} \leq f(\hat{x}) - f_{\star} \leq t \left( \sup_{\mathcal{C}} f - f_{\star} \right) \leq tDL.$$

The result follows by letting  $t \searrow \lambda^{N/d}$ . □

Thus, in principle, we can achieve  $f(x_{N-1}) - f_\star \leq \varepsilon$  in  $O(d \log(DL/\varepsilon))$  iterations. Compared to [Theorem 6.14](#), this result incurs only a logarithmic dependence on the ratio  $DL/\varepsilon$ , i.e., we can output a high-accuracy solution even for poorly conditioned convex sets. On the other hand, it incurs dependence on the dimension.

Recall that the lower bound for convex smooth optimization ([Theorem 4.4](#)) only applies in dimension  $d \gtrsim \sqrt{\beta R^2/\varepsilon}$ . The center of gravity method explains why: a  $\beta$ -smooth function over a ball of radius  $R$  is also  $\beta R$ -Lipschitz, so [Theorem 6.19](#) yields an oracle complexity of  $O(d \log(\beta R^2/\varepsilon))$  in this case. This is smaller than the lower bound of  $\Omega(\sqrt{\beta R^2/\varepsilon})$  in [Theorem 4.4](#) when  $d \ll \sqrt{\beta R^2/\varepsilon}/\log(\beta R^2/\varepsilon)$ , so a lower bound construction cannot exist in any smaller dimension.<sup>7</sup> Note also that for convex quadratic minimization, there are methods which find the minimizer in  $d$  queries (e.g., [Theorem 5.3](#) for [CG](#)); the center of gravity method almost achieves this guarantee for general convex optimization.

Toward making cutting plane methods more practical, a famous example is the *ellipsoid method*. In this scheme, we take each set  $\mathcal{C}_n$  to be an ellipsoid,

$$\mathcal{C}_n = \{x \in \mathbb{R}^d : \langle x - x_n, \Sigma_n^{-1} (x - x_n) \rangle \leq 1\}. \quad (6.6)$$

At the next iteration, we must find a new ellipsoid  $\mathcal{C}_{n+1}$  such that

$$\mathcal{C}_{n+1} \supseteq \mathcal{C}_n \cap \{x \in \mathbb{R}^d : \langle p_n, x \rangle \leq \langle p_n, x_n \rangle\}. \quad (6.7)$$

Here, we use the following geometric lemma ([Exercise 6.7](#)).

**Lemma 6.20 (ellipsoid).** Let  $\mathcal{C}_n$  be the ellipsoid (6.6) and let  $p_n \in \mathbb{R}^d$  be a non-zero vector. Define  $\mathcal{C}_{n+1} := \{x \in \mathbb{R}^d : \langle x - x_{n+1}, \Sigma_{n+1}^{-1} (x - x_{n+1}) \rangle \leq 1\}$ , where

$$\begin{aligned} x_{n+1} &:= x_n - \frac{1}{d+1} \frac{\Sigma_n p_n}{\sqrt{\langle p_n, \Sigma_n p_n \rangle}}, \\ \Sigma_{n+1} &:= \frac{d^2}{d^2-1} \left( \Sigma_n - \frac{2}{d+1} \frac{\Sigma_n p_n p_n^\top \Sigma_n}{\langle p_n, \Sigma_n p_n \rangle} \right). \end{aligned}$$

Then, for  $d > 1$ ,  $\mathcal{C}_{n+1}$  satisfies (6.7) and

$$\frac{\text{vol}(\mathcal{C}_{n+1})}{\text{vol}(\mathcal{C}_n)} = \sqrt{\frac{d-1}{d+1}} \left( \frac{d^2}{d^2-1} \right)^d = 1 - \Omega\left(\frac{1}{d}\right).$$

<sup>7</sup>This discussion is not entirely correct since [Theorem 4.4](#) only applies to gradient span algorithms, which does not cover [CoGM](#). However, the moral of the discussion is true for bona fide oracle lower bounds.

By following the proof of [Theorem 6.19](#), replacing  $\lambda$  by  $1 - \Omega(1/d)$ , one obtains the same guarantee as for [CoGM](#) but with iteration count  $O(d^2 \log(LD/\varepsilon))$ . (See [Exercise 6.6](#) for details.) Thus, the cost of obtaining an implementable version of the center of gravity method is a larger query complexity. Naturally, there have been numerous follow-up works in the field which aim at achieving the best of both worlds.

## 6.4 Lower bounds

In this section, we study lower bounds for convex non-smooth optimization.

**Theorem 6.21 (lower bound for convex, non-smooth minimization).** For any  $x_0 \in \mathbb{R}^d$ ,  $d > N$ , and  $L, R > 0$ , there exists a convex and  $L$ -Lipschitz function  $f$  over  $B(x_\star, R)$  such that  $x_0 \in B(x_\star, R)$  and for any gradient span algorithm,

$$f(x_N) - f_\star \gtrsim \frac{LR}{\sqrt{N}}.$$

*Proof.* Assume  $x_0 = 0$  and define the function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  by

$$f(x) := \gamma \max_{i \in [d]} x[i] + \frac{\alpha}{2} \|x\|^2,$$

where  $\alpha, \gamma > 0$  are to be chosen. Note that this function is Lipschitz with constant  $\gamma + \alpha(\|x_\star\| + R)$ . Also, if  $I_\star(x) := \{i \in [d] : x[i] = \max_{j \in [d]} x[j]\}$ , then from [Lemma 6.15](#),

$$\partial f(x) = \alpha x + \gamma \text{conv}\{e_i : i \in I_\star(x)\}.$$

The optimal point is  $x_\star[k] = -\gamma/(\alpha d)$  for  $k \in [d]$ , by checking that  $0 \in \partial f(x_\star)$ . Thus,  $\|x_\star\| = \gamma/(\alpha\sqrt{d})$  and the Lipschitz constant is at most  $2\gamma + \alpha R$ .

We take a subgradient oracle which, given a point  $x$ , outputs  $\alpha x + \gamma e_i \in \partial f(x)$ , where  $i = \min I_\star(x)$  is the *first* coordinate of  $x$  that achieves the maximum. From this property, it is straightforward to show via induction that  $x_n \in \mathcal{V}_n$  for all  $n$ , where  $\mathcal{V}_n$  is the subspace from the proof of [Theorem 4.4](#).

Since  $d > N$ , it follows that  $f(x_N) \geq 0$ . On the other hand,

$$f_\star = f(x_\star) = -\frac{\gamma^2}{\alpha d} + \frac{\gamma^2}{2\alpha d} = -\frac{\gamma^2}{2\alpha d}.$$

We set  $d = N + 1$ ,  $\gamma = L/4$ ,  $\alpha = \gamma/(R\sqrt{d})$  (to ensure that  $\|x_0 - x_\star\| \leq R$ ), which leads to a Lipschitz constant of  $L/2 + L/(4\sqrt{d}) \leq L$ . It yields

$$f(x_N) - f_\star \geq -f(x_\star) \gtrsim \frac{LR}{\sqrt{N}}. \quad \square$$

Note that this matches the guarantee of PSD (Theorem 6.14), so projected subgradient descent is *optimal* in the non-smooth setting. In other words, without smoothness, there is no acceleration phenomenon.

There is a version of Theorem 6.21 in the strongly convex case (Exercise 6.8).

**Theorem 6.22** (lower bound for strongly convex, non-smooth minimization). For any  $x_0 \in \mathbb{R}^d$ ,  $d > N$ , and  $\alpha, L > 0$ , there exists  $R > 0$  and an  $\alpha$ -convex and  $L$ -Lipschitz function  $f$  over  $B(x_\star, R)$  such that  $x_0 \in B(x_\star, R)$  and for any gradient span algorithm,

$$f(x_N) - f_\star \gtrsim \frac{L^2}{\alpha N}.$$

Next, in the low-dimensional setting, the following lower bound holds.

**Theorem 6.23** (lower bound for convex, non-smooth minimization II). The oracle complexity of minimizing convex,  $L$ -Lipschitz functions over  $[-R, R]^d$  to accuracy  $\varepsilon$  is at least  $\Omega(d \log(LR/\varepsilon))$ .

This shows that CoGM is optimal as well. Actually, we do not prove Theorem 6.23; instead, we focus on the related but harder task of feasibility.

**Definition 6.24.** Let  $0 < \delta < R$ . Let  $\mathcal{C} \subseteq [-R, R]^d$  be a closed convex set such that there exists a ball  $B(x_\star, \delta) \subseteq \mathcal{C}$ . The **feasibility problem** with parameters  $(\delta, R)$  is the problem of outputting a point in  $\text{int } \mathcal{C}$ , given access to a separation oracle. Namely, given a point  $x \in \mathbb{R}^d$ , the separation oracle either reports that  $x \in \mathcal{C}$ , or it outputs a non-zero vector  $p \in \mathbb{R}^d$  such that  $\sup_{\mathcal{C}} \langle p, \cdot \rangle \leq \langle p, x \rangle$ .

If one can solve the feasibility problem, then one can solve the convex Lipschitz minimization problem. Indeed, given a convex,  $L$ -Lipschitz function  $f$  over  $[-R, R]^d$ , suppose for the sake of argument that we know the optimal value  $f_\star$ . Consider the feasibility problem for set  $\mathcal{C} := \{f - f_\star \leq \varepsilon\}$ . For  $x_\star := \arg \min_{[-R, R]^d} f$ , we claim that  $B(x_\star, \varepsilon/L) \subseteq \mathcal{C}$ ; indeed this follows from  $L$ -Lipschitzness.<sup>8</sup> Also, the subgradient oracle

<sup>8</sup>Actually this is not exactly true because  $x_\star$  could lie near the boundary of  $[-R, R]^d$ . To fix this, one could instead look for a minimizer of  $f$  over  $\mathcal{C}' := [-R + \delta, R - \delta]^d$ , i.e., define  $x_{\delta, \star}$  to be a minimizer over this smaller cube and set  $\mathcal{C} := \mathcal{C}' \cap \{f - f(x_{\delta, \star}) \leq \varepsilon\}$ . If we take  $\delta = \varepsilon/(L\sqrt{d})$ , then by  $L$ -Lipschitzness we see that any point in  $\mathcal{C}$  is a  $2\varepsilon$ -minimizer of  $f$  over  $[-R, R]^d$ , and now  $B(x_{\delta, \star}, \delta) \subseteq \mathcal{C}$ . This does not really change the argument.

for  $f$  yields a separation oracle for  $\mathcal{C}$ . Thus, solving the feasibility problem for  $\mathcal{C}$  with parameters  $(\varepsilon/L, R)$  yields an  $\varepsilon$ -solution to the problem of minimizing  $f$ .

Since the feasibility problem is harder, the following theorem is weaker than [Theorem 6.23](#). However, it is easier to prove, and it contains most of the main ideas.

**Theorem 6.25 (lower bound for feasibility).** For any deterministic algorithm, the feasibility problem with parameters  $(\varepsilon, R)$  requires  $\Omega(d \log(R/\varepsilon))$  queries.

*Proof.* We play a game with the algorithm. Suppose that the algorithm has chosen query points  $x_1, \dots, x_n$  thus far. Our goal is to choose a vector  $p_n$ —which is supposed to correspond to the output of a separation oracle—and we provide the algorithm with this vector, which it then uses to produce a new point  $x_{n+1}$  and so on. Simultaneously, we also maintain a sequence of convex bodies (actually, boxes)  $\mathcal{C}_0, \mathcal{C}_1, \dots, \mathcal{C}_N$ .

At the end of the game, the algorithm has produced points  $x_1, \dots, x_N$ , and we have produced vectors  $p_1, \dots, p_N$ . By itself, this is not yet meaningful; the algorithm is not designed to produce useful results, *unless*  $p_1, \dots, p_N$  are valid outputs from a separation oracle corresponding to a convex body  $\mathcal{C}$  satisfying the assumptions of the feasibility problem. So, we aim to choose  $p_1, \dots, p_N$  so that this holds with  $\mathcal{C} = \mathcal{C}_N$ . Now, we can use the following post hoc reasoning: *had we* run the algorithm with the separation oracle for  $\mathcal{C}_N$  from the outset, then the algorithm would have output the same sequence of points  $x_1, \dots, x_N$ , because it is deterministic, so this construction yields a valid lower bound (i.e., it requires more than  $N$  iterations to solve the feasibility problem). This proof technique is known as the method of *resisting oracles*, and its main drawback is that it does not apply to randomized algorithms.<sup>9</sup>

Let us instantiate the resisting oracle for the feasibility problem. At each iteration  $n$ , the convex body  $\mathcal{C}_n$  is the box  $\{x \in \mathbb{R}^d : a_n \leq x \leq b_n\}$ ; here,  $a_n, b_n \in \mathbb{R}^d$  and the inequality is interpreted pointwise. We start with  $a_0 = -R\mathbf{1}_d$ ,  $b_0 = +R\mathbf{1}_d$ , where  $\mathbf{1}_d$  is the all-ones vector; thus,  $\mathcal{C}_0 = [-R, R]^d$ .

When the algorithm makes the first query  $x_1$ , we update the box by cutting it in half, based on the first coordinate of  $x_1$ . Namely, if  $x_1[1] \leq 0$ , we set  $a_1[1] = 0$ , and  $a_1[k] = a_0[k]$  for all  $k > 1$ ; we output the separating vector  $-e_1$ . If  $x_1[1] \geq 0$ , we set  $b_1[1] = 0$  and  $b_1[k] = b_0[k]$  for all  $k > 1$ ; we output the separating vector  $+e_1$ . In either case,  $\text{vol}(\mathcal{C}_1) = \frac{1}{2} \text{vol}(\mathcal{C}_0)$  and  $x_1 \notin \text{int } \mathcal{C}_1$ .

When the algorithm makes the second query  $x_2$ , we repeat this procedure except that we cut the box in half along the second coordinate. We continue in this fashion, cycling through the coordinates.

<sup>9</sup>Lower bounds for randomized algorithms require the use of information theory.

Let  $c_n$  denote the center of  $\mathcal{C}_n$ . We now claim that for each  $n$ ,  $B(c_n, r_n) \subseteq \mathcal{C}_n$ , where  $r_n = (R/2) (1/2)^{n/d}$ . Indeed, this is true for  $n = 0$ . Also, for  $n = ad$  for integer  $a$ , each side of the box has length  $R (1/2)^a$ , so the result is true in this case too. Finally, for  $n = ad + b$ , we have  $B(c_{(a+1)d}, R/2^{a+1}) \subseteq \mathcal{C}_{(a+1)d} \subseteq \mathcal{C}_n$  hence  $B(c_n, R/2^{a+1}) \subseteq \mathcal{C}_n$ , and we note that  $R/2^{a+1} \leq (R/2) (1/2)^{n/d}$ .

The resisting oracle construction succeeds up to iteration  $N$  provided that  $\mathcal{C}_N$  contains a ball of radius  $\varepsilon$ . It therefore suffices to have  $(R/2) (1/2)^{N/d} \geq \varepsilon$ , i.e.,  $N \gtrsim d \log(R/\varepsilon)$ .  $\square$

## Exercises

### Exercise 6.1.

1. Prove that a function  $f$  is lower semicontinuous if and only if for all  $c \in \mathbb{R}$ , the level set  $\{f \leq c\}$  is closed.
2. Prove that a supremum of lower semicontinuous functions is lower semicontinuous.
3. Show that the function defined in (6.2) is lower semicontinuous if and only if  $\phi = 0$ .

**Exercise 6.2.** Prove that if  $f$  is differentiable at  $x_0 \in \text{int dom } f$ , then  $\partial f(x_0) = \{\nabla f(x_0)\}$ .

**Exercise 6.3.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be continuous and convex on a convex set  $\mathcal{C}$ . Prove that  $f$  is Lipschitz continuous over  $\mathcal{C}$  with constant  $L$  if and only if for every  $x_0 \in \text{int } \mathcal{C}$  and every  $p \in \partial f(x_0)$ , we have  $\|p\| \leq L$ .

**Exercise 6.4.** Compute the subdifferential of the Euclidean norm  $\|\cdot\|$ .

**Exercise 6.5.** Assume that  $f$  is  $\alpha$ -strongly convex and  $L$ -Lipschitz continuous over the closed convex set  $\mathcal{C}$ . Prove that for PSD,

$$f(\bar{x}_N) - f_\star \leq \frac{\alpha}{2 \{(1 - \alpha h/L)^{-N} - 1\}} \|x_0 - x_\star\|^2 + \frac{Lh}{2},$$

where  $\bar{x}_N$  is a suitable averaged iterate. Deduce that by setting  $h = \varepsilon/L$ , one can achieve  $f(\bar{x}_N) - f_\star \leq \varepsilon$  in  $O(\frac{L^2}{\alpha\varepsilon} \log(\frac{\alpha R^2}{\varepsilon}))$  iterations (compared with  $O(L^2 R^2/\varepsilon^2)$  iterations, as implied by Theorem 6.14).

Also, show that under these assumptions,  $\|x_0 - x_\star\| \leq 2L/\alpha$ .

**Exercise 6.6.** The analysis of the ellipsoid method (and general cutting plane schemes) presents an additional difficulty: since the next set  $\mathcal{C}_{n+1}$  is only chosen to be a superset of  $\mathcal{C}_n \cap \{\langle p_n, \cdot \rangle \leq \langle p_n, x_n \rangle\}$ , it is not guaranteed that  $\mathcal{C} \subseteq \mathcal{C}_n$  for all  $n$ ; in particular, the chosen point  $x_n$  may lie outside of  $\mathcal{C}$ .

Assume that we have access to a separation oracle for  $\mathcal{C}$ : given a point  $x \notin \mathcal{C}$ , the oracle outputs a non-zero vector  $p \in \mathbb{R}^d$  such that  $\sup_{\mathcal{C}} \langle p, \cdot \rangle \leq \langle p, x \rangle$ . Modify the cutting plane method as follows: if a chosen point  $x_n$  does not lie in  $\mathcal{C}$ , then let  $p_n$  be vector that separates  $x_n$  from  $\mathcal{C}$  and instead update  $\mathcal{C}_{n+1}$  to be a superset of  $\mathcal{C}_n \cap \{\langle p_n, \cdot \rangle \leq \langle p_n, x_n \rangle\}$ . We also allow  $\mathcal{C}_0 \supseteq \mathcal{C}$ , so that  $x_0$  is not necessarily feasible either. Prove that if the sets are chosen so that  $\text{vol}(\mathcal{C}_{n+1})/\text{vol}(\mathcal{C}_n) \leq \lambda < 1$  for all  $n$ , then the following assertions hold.

1. If  $\text{vol}(\mathcal{C}_N) < \text{vol}(\mathcal{C})$ , then there exists  $n < N$  with  $x_n \in \mathcal{C}$ .
2. If  $\text{vol}(\mathcal{C}_N) < \text{vol}(\mathcal{C})$ , then there exists  $n < N$  with  $x_n \in \mathcal{C}$  and

$$f(x_n) - f_\star \leq DL\lambda^{N/d} \left( \frac{\text{vol } \mathcal{C}_0}{\text{vol } \mathcal{C}} \right)^{1/d}.$$

*Hint:* Define a sequence of sets  $\mathcal{C}'_0, \mathcal{C}'_1, \mathcal{C}'_2, \dots$  as follows. Start with  $\mathcal{C}'_0 = \mathcal{C}$  and  $n_{-1} := 0$ . For each  $k \in \mathbb{N}$ , let  $n_k$  denote the first integer greater than  $n_{k-1}$  for which  $x_{n_k} \in \mathcal{C}$  and set  $\mathcal{C}'_{k+1} := \mathcal{C}'_k \cap \{\langle p_{n_k}, \cdot \rangle \leq \langle p_{n_k}, x_{n_k} \rangle\}$ . Prove via induction that if  $k(N)$  is the largest integer such that  $n_{k(N)} \leq N$ , then  $\mathcal{C}'_{n_{k(N)}} \subseteq \mathcal{C}_N$ .

**Exercise 6.7.** Prove [Lemma 6.20](#).

**Exercise 6.8.** Prove [Theorem 6.22](#). (Use the same construction as in the proof of [Theorem 6.21](#), but choose the parameters  $\alpha$  and  $\gamma$  differently.)

## 7 [2/18] Frank–Wolfe

In order to overcome the lower bounds in the black-box setting, we must take advantage of additional structure in the problem. The first method we study in this vein is the *Frank–Wolfe* or *conditional gradient* method. Instead of assuming access to a projection oracle for the constraint set  $\mathcal{C}$ , it instead assumes access to a *linear optimization oracle* (LOO) over the set  $\mathcal{C}$ :

$$\text{Given } p \in \mathbb{R}^d, \text{ output } \arg \min_{\mathcal{C}} \langle p, \cdot \rangle. \quad (\text{LOO})$$

Here, we assume that  $\mathcal{C}$  is compact (bounded and closed).

The oracle equivalently maximizes the convex function  $-\langle p, \cdot \rangle$  over  $\mathcal{C}$ , so the arg min is attained at a vertex of  $\mathcal{C}$ . Let us define these concepts properly.

**Definition 7.1.** A point  $x \in \mathcal{C}$  is called an **extreme point** or a **vertex** of  $\mathcal{C}$  if there do not exist  $x_0, x_1 \in \mathcal{C}$  and  $t \in (0, 1)$  such that  $x = (1 - t)x_0 + tx_1$ .

**Theorem 7.2.** Every compact convex set is the convex hull of its extreme points.

For example, the set of vertices of the closed unit ball  $\overline{B(0, 1)}$  is the sphere  $\partial B(0, 1)$ . It follows that to implement (LOO), it suffices to solve  $\arg \min_{\text{vertices of } \mathcal{C}} \langle p, \cdot \rangle$ .

We now present the Frank–Wolfe method for minimizing  $f$  over  $\mathcal{C}$ :

$$x_{n+1} := (1 - h_n) x_n + h_n \text{LOO}(\nabla f(x_n)). \quad (\text{FW})$$

**Theorem 7.3 (convergence of FW).** Let  $f$  be convex and  $\beta$ -smooth over  $\mathcal{C}$ . Let  $D := \text{diam } \mathcal{C}$  and  $h_n = 2/(n + 2)$ . Then, for any  $N \geq 1$ , FW satisfies

$$f(x_N) - f_\star \leq \frac{2\beta D^2}{N + 1}.$$

*Proof.* Let  $y_n := \text{LOO}(\nabla f(x_n))$ . Using  $\beta$ -smoothness,

$$\begin{aligned} f(x_{n+1}) - f(x_n) &\leq \langle \nabla f(x_n), x_{n+1} - x_n \rangle + \frac{\beta}{2} \|x_{n+1} - x_n\|^2 \\ &\leq h_n \langle \nabla f(x_n), y_n - x_n \rangle + \frac{\beta D^2 h_n^2}{2} \leq h_n \langle \nabla f(x_n), x_\star - x_n \rangle + \frac{\beta D^2 h_n^2}{2} \\ &\leq -h_n (f(x_n) - f_\star) + \frac{\beta D^2 h_n^2}{2}. \end{aligned}$$

Rearranging,

$$f(x_{n+1}) - f_\star \leq (1 - h_n) (f(x_n) - f_\star) + \frac{\beta D^2 h_n^2}{2}.$$

For  $h_n = 2/(n + 2)$ , we now prove the error bound by induction on  $n$ , where the base case  $n = 0$  follows from the inequality above. If the error bound holds at iteration  $n$ , then

$$f(x_{n+1}) - f_\star \leq \frac{n}{n+2} \frac{2\beta D^2}{n+1} + \frac{2\beta D^2}{(n+2)^2} \leq \frac{2\beta D^2}{n+2}. \quad \square$$

The analysis above is actually not the most natural one, since it fails to capture the affine invariance of the Frank–Wolfe algorithm (Exercise 7.1).

Besides positing different oracle access than projected gradient methods, the Frank–Wolfe method has the appealing property of producing sparse solutions. This connects with results known as approximate Carathéodory theorems. First, let us recall the classical statement of Carathéodory’s theorem.



**Theorem 7.4 (Carathéodory).** Let  $\mathcal{C} \subseteq \mathbb{R}^d$  be a compact convex set and let  $x \in \mathcal{C}$ . Then,  $x$  can be written as a convex combination of  $d + 1$  vertices of  $\mathcal{C}$ .

Caution: in this theorem, the choice of  $d + 1$  vertices of course depends on  $x$  itself. If every point in  $\mathcal{C}$  could be written as a convex combination of the *same*  $d + 1$  vertices, this would say that  $\mathcal{C}$  only has  $d + 1$  vertices at all.

Carathéodory's theorem says that even if a convex body has exponentially many vertices, such as the cube  $[-1, 1]^d$ , any given point has a succinct representation using only  $d + 1$  vertices. However, the size of the representation grows with the ambient dimension. What happens if we relax the requirement that the representation is exact? The following simple argument, often attributed to B. Maurey, shows that the size of the representation is *dimension-free*, and the convex combination even uses equal weights.

**Theorem 7.5 (approximate Carathéodory).** Let  $\mathcal{C} \subseteq \mathbb{R}^d$  be a compact convex set with diameter  $D$ , let  $0 < \varepsilon < 1$ , and let  $x \in \mathcal{C}$ . Then, there exist vertices  $y_1, \dots, y_N \in \mathcal{C}$  with

$$\left\| x - \frac{1}{N} \sum_{i=1}^N y_i \right\| \leq \varepsilon D, \quad N \leq \frac{1}{\varepsilon^2}.$$

*Proof.* By [Theorem 7.4](#), there exist vertices  $\bar{y}_1, \dots, \bar{y}_{d+1} \in \mathcal{C}$  and a probability distribution  $\lambda$  over  $[d + 1]$  such that  $x = \sum_{j=1}^{d+1} \lambda_j \bar{y}_j$ . Now consider the distribution  $\mu = \sum_{j=1}^{d+1} \lambda_j \delta_{\bar{y}_j}$  and sample points  $Y_1, \dots, Y_N \stackrel{\text{i.i.d.}}{\sim} \mu$ . Note that each  $Y_i$  is a vertex of  $\mathcal{C}$ . Then, since the mean of  $\mu$  is  $x$ , the usual variance calculation shows that

$$\mathbb{E} \left[ \left\| x - \frac{1}{N} \sum_{i=1}^N Y_i \right\|^2 \right] \leq \frac{\sum_{j=1}^{d+1} \lambda_j \|x - \bar{y}_j\|^2}{N} \leq \frac{D^2}{N}.$$

Choose  $N$  to make the right-hand side at most  $\varepsilon^2 D^2$ . □

The approximate Carathéodory theorem has implications, e.g., for controlling the covering numbers of polytopes. But more broadly, the proof technique is quite influential and is at the root of other important developments, e.g., the existence of neural networks of small width which approximate functions in the Barron class [\[Bar93\]](#).

Now comes the punchline: Frank-Wolfe renders the approximate Carathéodory theorem constructive. Indeed, suppose that the [LOO](#) always outputs a vertex. After  $N - 1$  iterations of [FW](#) starting from a vertex, the iterate  $x_{N-1}$  is a convex combination of at most  $N$  vertices. At the same time, if we apply [Theorem 7.3](#) to the 2-smooth function  $f : z \mapsto \|x - z\|^2$ , where  $x \in \mathcal{C}$  and  $f_\star = 0$ , we see that  $\|x_{N-1} - x\|^2 \leq 4D^2/N$ .

The full statement of [Theorem 7.3](#) can therefore be seen as a generalization of the approximate Carathéodory principle: the iterate of [FW](#) is a sparse combination of vertices which is approximately optimal. We next demonstrate an example in which this sparsity property is crucial.

**Example 7.6 (low-rank estimation).** Consider the nuclear norm ball

$$\mathcal{C} = \left\{ X \in \mathbb{R}^{d \times d} : \|X\|_* = \sum_{i=1}^d \sigma_i(X) \leq 1 \right\}.$$

This constraint set often arises in low-rank matrix recovery as a convex relaxation of a rank constraint. Projection onto the set  $\mathcal{C}$  requires projecting the singular values onto the simplex; this requires computing a full SVD, which uses  $O(d^3)$  arithmetic operations. On the other hand, since

$$\mathcal{C} = \text{conv}\{uv^\top : u, v \in \mathbb{R}^d, \|u\| = \|v\| = 1\},$$

the [LOO](#) for  $\mathcal{C}$  involves solving, for any  $P \in \mathbb{R}^{d \times d}$ ,

$$\arg \min_{X \in \mathcal{C}} \langle P, X \rangle = \arg \min_{u, v \in \mathbb{R}^d, \|u\| = \|v\| = 1} \langle P, uv^\top \rangle.$$

Solving this amounts to computing the top singular vector of  $P$ , which is often implemented via power iteration at cost  $O(d^2)$  per step. Moreover, [FW](#) yields an  $\varepsilon$ -accurate solution with rank  $O(1/\varepsilon)$ .

## Exercises

**Exercise 7.1.** Show that [FW](#) is affine-invariant in the following sense. Let  $A \in \mathbb{R}^{d \times d}$  be an invertible matrix. Show that the iterates  $\{\hat{x}_n\}_{n \in \mathbb{N}}$  of [FW](#) applied to the problem of minimizing  $\hat{x} \mapsto f(A\hat{x})$  over the set  $A^{-1}\mathcal{C}$  are related to the iterates  $\{x_n\}_{n \in \mathbb{N}}$  of [FW](#) on the original problem via  $x_n = A\hat{x}_n$ .

## 8 [2/20] Proximal methods

Can we solve non-smooth problems at the same rate as smooth problems? The black-box lower bounds say *no* in general, but if the non-smooth part is “simple” in the sense that it admits an implementable proximal oracle, the answer becomes yes.

**Definition 8.1.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ . The **proximal oracle** for  $f$  is the mapping  $\text{prox}_f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  given by

$$\text{prox}_f(y) := \arg \min_{x \in \mathbb{R}^d} \left\{ f(x) + \frac{1}{2} \|y - x\|^2 \right\}.$$

If  $f$  is a regular convex function, then the optimization problem defining the proximal oracle is strongly convex, so it admits a unique minimizer by [Lemma 1.10](#) and [Lemma 6.6](#). Note also that

$$\text{prox}_{hf}(y) = \arg \min_{x \in \mathbb{R}^d} \left\{ hf(x) + \frac{1}{2} \|y - x\|^2 \right\} = \arg \min_{x \in \mathbb{R}^d} \left\{ f(x) + \frac{1}{2h} \|y - x\|^2 \right\},$$

where  $h > 0$  plays the role of a step size.

The value of the optimization problem defining  $\text{prox}_f$  also has a name.

**Definition 8.2.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ . The **Moreau–Yosida envelope** of  $f$  with parameter  $h > 0$  is the mapping  $f_h : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  given by

$$f_h(y) := \inf_{x \in \mathbb{R}^d} \left\{ f(x) + \frac{1}{2h} \|y - x\|^2 \right\}.$$

## 8.1 Algorithms and examples

The proximal oracle is a regularized version of the original optimization problem. Assuming for the moment that we can compute the proximal oracle easily, let us explore its uses for algorithm design.

The simplest algorithm is to repeatedly iterate the proximal mapping. This is known as the *proximal point method*.

$$x_{n+1} := \text{prox}_{hf}(x_n). \quad (\text{PPM})$$

Assume for the moment that  $f$  is smooth and that the next point  $x_{n+1}$  can be obtained from the first-order optimality condition for  $\text{prox}_{hf}$ . This leads to

$$0 = \nabla f(x_{n+1}) + \frac{1}{h} (x_{n+1} - x_n) \iff x_{n+1} = x_n - h \nabla f(x_{n+1}).$$

Note that this is similar to the [GD](#) update, except that the gradient is evaluated at the subsequent point  $x_{n+1}$ . In numerical analysis, we say that [GD](#) is an *explicit* discretization of

the gradient flow, whereas **PPM** is an *implicit* discretization. The advantage of an explicit method is easy of implementation; it does not require solving a (non-linear) system in order to perform an update. The advantage of an implicit method is stability.

Recall that the results in §2 for **GF** do not require smoothness of  $f$ , whereas the results in §3 for **GD** do. (We studied the non-smooth case for **GD** in §6.2, but it requires decreasing step sizes and averaging.) Shortly, we shall see that **PPM** is similar to **GF**, in that it also does not require smoothness.

The most powerful results using the proximal oracle, however, are for the problem of *composite optimization*. Here, the goal is to minimize a sum of functions:

$$\text{minimize} \quad F := f + g.$$

We assume that  $f$  is smooth and that  $g$  is non-smooth.

**Example 8.3 (LASSO as composite optimization).** The computation of the LASSO estimator from Example 1.3 is the canonical example of composite optimization, where

$$f : \theta \mapsto \frac{1}{2n} \sum_{i=1}^n (Y_i - \langle \theta, X_i \rangle)^2, \quad g : \theta \mapsto \lambda \|\theta\|_1.$$

In this example, the non-smooth part is particularly simple, so we can compute its proximal oracle in closed form. First, note that it is coordinate-wise decomposable:

$$\begin{aligned} \text{prox}_{\lambda \|\cdot\|_1}(y) &= \arg \min_{x \in \mathbb{R}^d} \left\{ \lambda \|x\|_1 + \frac{1}{2} \|y - x\|^2 \right\} \\ &= \sum_{i=1}^d \left( \arg \min_{x[i] \in \mathbb{R}} \left\{ \lambda |x[i]| + \frac{1}{2} (y[i] - x[i])^2 \right\} \right) e_i. \end{aligned}$$

Therefore, it suffices to solve the problem in dimension one. A direct computation (see Exercise 8.1) then yields

$$\text{prox}_{\lambda \|\cdot\|_1}(y) = (|y| - \lambda)_+ \text{sgn } y =: \text{thesh}_\lambda(y)$$

where  $(\cdot)_+ := \max\{0, \cdot\}$  denotes the positive part. The operator  $\text{thesh}_\lambda$ , known as the *soft thresholding operator*, reduces the magnitude of its input by  $\lambda$ , or to 0 if the original magnitude is less than  $\lambda$ . The proximal operator for  $\lambda \|\cdot\|_1$  simply applies  $\text{thesh}_\lambda$  to each coordinate.

**Example 8.4 (constrained optimization as composite optimization).** Consider the problem of minimizing a smooth function  $f$  over a closed convex set  $\mathcal{C}$ . We can also treat this as composite optimization with

$$g = \chi_{\mathcal{C}}.$$

(Recall the convex indicator defined in (6.1).) In this case, the proximal oracle for  $g$  is

$$\text{prox}_{h\chi_{\mathcal{C}}}(y) = \arg \min_{x \in \mathbb{R}^d} \left\{ \chi_{\mathcal{C}}(x) + \frac{1}{2h} \|y - x\|^2 \right\} = \arg \min_{x \in \mathcal{C}} \left\{ \frac{1}{2h} \|y - x\|^2 \right\} = \Pi_{\mathcal{C}}(y).$$

So, the proximal oracle for  $\chi_{\mathcal{C}}$  is the projection oracle for  $\mathcal{C}$ .

The above examples motivate the assumption that we have access to the proximal oracle for the non-smooth part  $g$ . Further examples of computable proximal oracles can be found on the website [proximity-operator.net](http://proximity-operator.net).

The algorithm we consider in this context is known as *proximal gradient descent*.

$$x_{n+1} := \arg \min_{x \in \mathbb{R}^d} \left\{ f(x_n) + \langle \nabla f(x_n), x - x_n \rangle + g(x) + \frac{1}{2h} \|x - x_n\|^2 \right\}. \quad (\text{PGD})$$

In other words, we take the objective function  $F = f + g$  and linearize only the smooth part. The update can be rewritten as follows. By completing the square,

$$x_{n+1} = \arg \min_{x \in \mathbb{R}^d} \left\{ g(x) + \frac{1}{2h} \|x - x_n + h \nabla f(x_n)\|^2 \right\} = \text{prox}_{hg}(x_n - h \nabla f(x_n)).$$

This corresponds to taking an explicit step on  $f$ , followed by an implicit step on  $g$ . It is not obvious that this algorithm converges to  $x_{\star}$ , the minimizer of  $F = f + g$ . However, note that if  $g$  is differentiable, then

$$x_{n+1} = x_n - h \nabla f(x_n) - h \nabla g(x_{n+1}).$$

If  $x_n = x_{\star}$ , then  $x_{n+1} = x_{\star}$  is the solution since  $0 = \nabla F(x_{\star}) = \nabla f(x_{\star}) + \nabla g(x_{\star})$ . Thus, provided that  $f$  and  $g$  are convex and differentiable,  $x_{\star}$  is the unique fixed point.

For the LASSO problem, the iteration reads

$$x_{n+1} = \text{thresh}_{\lambda h}(x_n - h \nabla f(x_n)).$$

In the literature, this is known as the *iterative shrinking-thresholding algorithm* (ISTA). For constrained optimization, proximal gradient descent is projected gradient descent.

## 8.2 Convergence analysis

We study the convergence of [PGD](#), since it includes [PPM](#) as a special case (take  $f = 0$ ).

**Theorem 8.5 (convergence of PGD).** Let  $f$  be  $\alpha_f$ -convex and  $\beta_f$ -smooth, and let  $g$  be  $\alpha_g$ -convex. Let the step size  $h$  satisfy  $h \leq 1/\beta_f$ , let  $x^+$  denote the next iterate of [PGD](#) started from  $x$ , and let  $y \in \mathbb{R}^d$ . Then,

$$(1 + \alpha_g h) \|y - x^+\|^2 \leq (1 - \alpha_f h) \|y - x\|^2 - 2h (F(x^+) - F(y)). \quad (8.1)$$

In particular, if we set  $y = x_\star$  and iterate, it yields

$$F(x_N) - F_\star \leq \frac{\alpha_f + \alpha_g}{2(\lambda_h^{-N} - 1)} \|x_0 - x_\star\|^2,$$

where  $\lambda_h := (1 - \alpha_f h)/(1 + \alpha_g h)$ .

*Proof.* Let  $\psi_x$  denote the objective function in the definition of [PGD](#). Then,  $\psi_x$  is  $(\alpha_g + 1/h)$ -strongly convex with minimizer  $x^+$ , so by the quadratic growth inequality,

$$\psi_x(y) \geq \psi_x(x^+) + \frac{\alpha_g + 1/h}{2} \|y - x^+\|^2.$$

On one hand, by  $\alpha_f$ -convexity,

$$\psi_x(y) = f(x) + \langle \nabla f(x), y - x \rangle + g(y) + \frac{1}{2h} \|y - x\|^2 \leq F(y) + \frac{1/h - \alpha_f}{2} \|y - x\|^2.$$

On the other hand, by  $\beta_f$ -smoothness,

$$\begin{aligned} \psi_x(x^+) &= f(x) + \langle \nabla f(x), x^+ - x \rangle + g(x^+) + \frac{1}{2h} \|x^+ - x\|^2 \\ &\geq F(x^+) + \frac{1/h - \beta_f}{2} \|x^+ - x\|^2 \geq F(x^+). \end{aligned}$$

Combining these inequalities and rearranging,

$$(1 + \alpha_g h) \|y - x^+\|^2 \leq (1 - \alpha_f h) \|y - x\|^2 - 2h (F(x^+) - F(y)).$$

Note that by taking  $y = x$ , it yields the descent property

$$F(x^+) - F(x) \leq -\frac{1 + \alpha_g h}{2h} \|x - x^+\|^2 \leq 0.$$

The final bound follows from [Lemma 3.5](#) and algebra. □

Refined analyses of PPM are presented in [Exercise 8.3](#), [Corollary 9.13](#), and [Exercise 9.2](#).

The key feature of [Theorem 8.5](#) is that it essentially recovers the *smooth* rate for GD despite the presence of non-smoothness in the objective. Thus, for the LASSO problem ([Example 8.3](#)), we can solve it as quickly as if it were a smooth problem via ISTA.

Moreover, the one-step inequality (8.1) is the PGD analogue of the inequality (3.3) which holds for GD, and in turn, (3.3) is the only property of GD which plays a role in the proof of Nesterov acceleration ([Theorem 5.10](#)); the remainder of the proof is purely algebraic. This naturally leads to an accelerated algorithm for composite optimization.

Starting with  $x_{-1} = x_0$ , consider

$$x_{n+1} := x_n + \theta_n (x_n - x_{n-1}) - \text{PGD}_{F,1/\beta}(x_n + \theta_n (x_n - x_{n-1})), \quad (\text{APGD})$$

where  $\text{PGD}_{F,1/\beta}$  denotes one step of PGD on  $F = f + g$  with step size  $h = 1/\beta$ .

**Theorem 8.6 (convergence of APGD).** Let  $f$  be convex and  $\beta$ -smooth, and let  $g$  be convex. Define the sequence:  $\lambda_0 := 0$  and  $\lambda_{n+1} := \frac{1}{2} (1 + \sqrt{1 + 4\lambda_n^2})$  for  $n \in \mathbb{N}$ . Set  $\theta_n := (\lambda_n - 1)/\lambda_{n+1}$ . Then, APGD satisfies

$$F(x_N) - F_\star \leq \frac{2\beta \|x_0 - x_\star\|^2}{N^2}.$$

When applied to LASSO, this algorithm is known as *fast ISTA* or *FISTA*. Rates in the strongly convex setting can be obtained from the reduction in [Lemma 4.1](#).

## Exercises

**Exercise 8.1.** Verify the computation of  $\text{prox}_{\lambda|\cdot|}$  in [Example 8.3](#).

**Exercise 8.2.** Prove that even for non-convex  $f$ , as long as  $x^+ := \text{prox}_{hf}(x)$  is well-defined,

$$f(x^+) - f_\star \leq \frac{1}{2h} \|x - x_\star\|^2.$$

Thus, if we can implement PPM for arbitrarily large step sizes  $h > 0$ , we can solve non-convex optimization.

**Exercise 8.3.** To avoid technical difficulties, assume that  $f$  is convex and differentiable everywhere. Show that for the PPM, (8.1) can be refined as follows: for all  $y \in \mathbb{R}^d$ ,

$$\|x^+ - y\|^2 \leq \|x - y\|^2 - 2h (f(x^+) - f(y)) - h^2 \|\nabla f(x^+)\|^2.$$

Next, in analogy to [Exercise 2.1](#), define the Lyapunov function

$$\mathcal{L}_n := n^2 h^2 \|\nabla f(x_n)\|^2 + 2nh (f(x_n) - f_\star) + \|x_n - x_\star\|^2,$$

where  $\{x_n\}_{n \in \mathbb{N}}$  are the iterates of [PPM](#), and show that  $\mathcal{L}_{n+1} \leq \mathcal{L}_n$ . (Use the fact that [PPM](#) is contractive, see [Corollary 9.13](#).) Deduce the bounds

$$\|\nabla f(x_N)\| \leq \frac{\|x_0 - x_\star\|}{Nh}, \quad f(x_N) - f_\star \leq \frac{\|x_0 - x_\star\|^2}{4Nh}.$$

Observe that if  $h \searrow 0$  while  $Nh \rightarrow t$ , it recovers the results of [Exercise 2.1](#).

## 9 [2/25–2/27] Fenchel duality

In this section, we study a notion of duality for convex functions.

**Definition 9.1.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be proper ( $\text{dom } f \neq \emptyset$ ). The **convex conjugate** or **Fenchel–Legendre conjugate** of  $f$  is the function  $f^* : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  defined by

$$f^*(y) := \sup_{x \in \mathbb{R}^d} \{\langle x, y \rangle - f(x)\}.$$

For any proper function  $f$ , the conjugate  $f^*$  is always *convex* and *lower semicontinuous*, since it is a supremum of affine functions. Conversely, if  $f$  is a regular convex function (thus: proper, convex, and lower semicontinuous), then  $f = f^{**}$  ([Theorem 9.7](#)).

**Example 9.2.** The verification of these examples is left as [Exercise 9.1](#).

1. If  $f(x) = \frac{1}{2} \langle x, Ax \rangle$  where  $A \succ 0$ , then  $f^*(y) = \frac{1}{2} \langle y, A^{-1}y \rangle$ .
2. If  $f(x) = |x|^p/p$  for  $p > 1$  and  $x \in \mathbb{R}$ , then  $f^*(y) = |y|^q/q$  where  $1/p + 1/q = 1$ .
3. Let  $\|\cdot\|$  denote a norm over  $\mathbb{R}^d$  (not necessarily Euclidean), and let  $\|\cdot\|_*$  denote the dual norm:  $\|y\|_* := \sup\{\langle x, y \rangle : x \in \mathbb{R}^d, \|x\| \leq 1\}$ .  
If  $f(x) = \|x\|$ , then  $f^*(y) = \chi_{\mathcal{C}}(y)$  where  $\mathcal{C} := \{y \in \mathbb{R}^d : \|y\|_* \leq 1\}$  is the closed unit ball in the dual norm.

Before formally establishing further properties of this duality, we take a detour to explain the origin of this concept in classical mechanics.



## 9.1 (Optional) Connection with classical mechanics

**Disclaimer:** The material in this subsection is not necessarily the most relevant for optimization, and it is included for the sake of broader historical context. We make no attempt to be rigorous: assume all functions are smooth, etc.

Newton's law of motion states that the trajectory  $(x_t)_{t \geq 0}$  of a particle of mass  $m$  obeys the differential equation  $m\ddot{x}_t = F(x_t)$ , where  $F$  is the force. The force is typically given as the gradient of a potential:  $F = -\nabla\phi$ .

In 1662, Pierre de Fermat proposed an explanation for the law of refraction via his principle of least action: light takes the path which minimizes the total travel time. Is there such a principle for classical mechanics as well? In 1760, Joseph-Louis Lagrange found such a variational principle: let  $L(x, v) := \frac{1}{2} m \|v\|^2 - \phi(x)$  denote the *Lagrangian*, where  $v$  denotes the velocity of the particle. Note that the Lagrangian is the difference of the kinetic energy and the potential energy. The action functional is

$$\mathcal{A}((x_t)_{t \in [0, T]}) := \int_0^T L(x_t, \dot{x}_t) dt.$$

Lagrangian mechanics states that if a particle starts at  $x_0$  at time 0, and ends at  $x_T$  at time  $T$ , then the path it takes in between is a stationary point of the action functional subject to the endpoint constraints.

We solve for the path using calculus of variations. Let  $x_{[0, T]} := (x_t)_{t \in [0, T]}$  be a shorthand for the path. If  $x_{[0, T]}$  is a stationary point, it means that for any perturbation  $\delta x_{[0, T]}$ , the difference  $\mathcal{A}(x_{[0, T]} + \delta x_{[0, T]}) - \mathcal{A}(x_{[0, T]})$  should vanish to first order in  $\delta x_{[0, T]}$ . The endpoint constraints require that  $\delta x_0 = \delta x_T = 0$ . Thus,

$$\begin{aligned} \mathcal{A}(x_{[0, T]} + \delta x_{[0, T]}) - \mathcal{A}(x_{[0, T]}) &= \int_0^T \{L(x_t + \delta x_t, \dot{x}_t + \delta \dot{x}_t) - L(x_t, \dot{x}_t)\} dt \\ &= \int_0^T \{\langle \nabla_x L(x_t, \dot{x}_t), \delta x_t \rangle + \langle \nabla_v L(x_t, \dot{x}_t), \delta \dot{x}_t \rangle\} dt + o(\|\delta x\|) \\ &= \int_0^T \langle \nabla_x L(x_t, \dot{x}_t) - \partial_t \nabla_v L(x_t, \dot{x}_t), \delta x_t \rangle dt + o(\|\delta x\|). \end{aligned}$$

The stationary point therefore satisfies the Euler–Lagrange equation

$$\partial_t \nabla_v L(x_t, \dot{x}_t) = \nabla_x L(x_t, \dot{x}_t).$$

For  $L(x, v) = \frac{1}{2} m \|v\|^2 - \phi(x)$ , it recovers Newton's equation.

We now introduce the Legendre transform. Define the *Hamiltonian*  $H$  to be the convex conjugate of  $L$  with respect to the  $v$ -variable, i.e.,

$$H(x, p) := \sup_{v \in \mathbb{R}^d} \{\langle p, v \rangle - L(x, v)\}.$$

The first-order condition reveals that

$$p = \nabla_v L(x, v) .$$

Instead of working with the variables  $(x, v)$ , we now work with the variables  $(x, p)$ . The inverse of the transformation is given by

$$v = \nabla_p H(x, p) . \tag{9.1}$$

Indeed, we will argue that a regular convex function  $f$  satisfies  $f = f^{**}$  (Theorem 9.7). Assuming that  $v \mapsto L(x, v)$  is regular convex, it yields the dual representation

$$L(x, v) = \sup_{p \in \mathbb{R}^d} \{ \langle p, v \rangle - H(x, p) \} ,$$

and the first-order condition for this problem yields (9.1).

Thus, if we define  $p_t := \nabla_v L(x_t, \dot{x}_t)$ , we can reformulate the Euler–Lagrange equation as follows. First,  $\dot{x}_t = v_t = \nabla_p H(x_t, p_t)$  by (9.1). Also,  $\nabla_x H(x, p) = -\nabla_x L(x, v)$  by the envelope theorem, so  $\dot{p}_t = \partial_t \nabla_v L(x_t, \dot{x}_t) = \nabla_x L(x_t, \dot{x}_t) = -\nabla_x H(x_t, p_t)$ . In summary,

$$\dot{x}_t = \nabla_p H(x_t, p_t) , \quad \dot{p}_t = -\nabla_x H(x_t, p_t) .$$

These are known as *Hamilton's equations*, and it is easy to verify that they conserve the Hamiltonian:  $\partial_t H(x_t, p_t) = 0$ . Compared to Newton's law, which is a second-order differential equation for the trajectory, Hamilton's equations are a system of coupled first-order differential equations evolving in phase space.

For our running example,  $p = mv$  is interpreted as the *momentum*, and

$$H(x, p) = \left\langle p, \frac{p}{m} \right\rangle - \frac{1}{2} m \left\| \frac{p}{m} \right\|^2 + \phi(x) = \frac{1}{2m} \|p\|^2 + \phi(x) ,$$

which is the total energy (kinetic plus potential). Hamilton's equations read

$$m\dot{x}_t = p_t , \quad \dot{p}_t = -\nabla \phi(x_t) .$$

What does duality say about the action functional? Surprisingly, it relates back to other concepts we have already seen. Specialize now to the case where the Lagrangian only depends on  $v$  ( $\phi = 0$ ; no external potential, so we expect particles to move in straight lines). Define the following function of space and time:

$$u(t, x) := \inf \left\{ \int_0^t L(\dot{x}_s) \, ds + f(x_0) \mid x : [0, t] \rightarrow \mathbb{R}^d , \, x_t = x \right\} .$$

In words, we minimize the action functional up to time  $t$ , subject to the constraint that we hit  $x$  at time  $t$ . We also add an initial cost  $f(x_0)$ . The function  $u$  resembles the notion of the *value function* or *cost-to-go function* in dynamic programming, and indeed it satisfies a dynamic programming principle: for  $0 \leq s < t$ ,

$$u(t, y) = \inf_{x \in \mathbb{R}^d} \left\{ (t-s) L\left(\frac{y-x}{t-s}\right) + u(s, x) \right\}. \quad (9.2)$$

The heuristic derivation of this identity is as follows: consider a potential candidate  $x$  for the value of the path at time  $s$ . Given  $x$ , the best possible value of  $\int_0^s L(\dot{x}_r) dr + f(x_0)$  is  $u(s, x)$ . For the remaining part, by convexity,

$$\int_s^t L(\dot{x}_r) dr \geq (t-s) L\left(\frac{1}{t-s} \int_s^t \dot{x}_r dr\right) = (t-s) L\left(\frac{y-x}{t-s}\right).$$

The lower bound is achieved if  $\dot{x}_r$  is constant for  $r \in [s, t]$ , i.e.,  $x_{[s,t]}$  is a straight line.

In particular, since  $u(0, \cdot) = f$ , we see that

$$u(t, y) := \inf_{x \in \mathbb{R}^d} \left\{ tL\left(\frac{y-x}{t}\right) + f(x) \right\}. \quad (9.3)$$

**Definition 9.3.** The **Hopf–Lax semigroup**  $(Q_t)_{t \geq 0}$  is a family of operators which maps functions to functions, such that  $Q_t f(y)$  is defined to be the right-hand side of (9.3).

The dynamic programming principle (9.2) shows that  $Q_t f = Q_{t-s}(Q_s f)$ . Thus, we have the properties  $Q_0 = \text{id}$ ,  $Q_{s+t} = Q_s Q_t = Q_t Q_s$  for all  $s, t \geq 0$ , which are the defining properties of a semigroup.

These concepts are fundamental, so it is unsurprising that they have been rediscovered in different contexts. In the context of convex analysis, the corresponding operation is known as the infimal convolution.

**Definition 9.4.** Let  $f, g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ . The **infimal convolution** of  $f$  and  $g$ , denoted  $f \square g$ , is the function defined by

$$(f \square g)(y) := \inf_{x \in \mathbb{R}^d} \{f(x) + g(y-x)\}.$$

In this notation,  $Q_t f = tL(\cdot/t) \square f$ . Interestingly, the operation of convex conjugation turns addition into infimal convolution and vice versa.

**Theorem 9.5** (convex conjugation and infimal convolution). Let  $f, g$  be regular convex functions. Then,

$$(f \square g)^* = f^* + g^* .$$

Conversely, if  $\text{int dom } f \cap \text{int dom } g \neq \emptyset$ , then

$$(f + g)^* = f^* \square g^* .$$

*Proof.* For the first statement, note that

$$\begin{aligned} (f \square g)^*(y) &= \sup_{x \in \mathbb{R}^d} \{ \langle x, y \rangle - (f \square g)(x) \} = \sup_{x \in \mathbb{R}^d} \{ \langle x, y \rangle - \inf_{z \in \mathbb{R}^d} \{ f(z) + g(x - z) \} \} \\ &= \sup_{x, z \in \mathbb{R}^d} \{ \langle z, y \rangle - f(z) + \langle x - z, y \rangle - g(x - z) \} = f^*(y) + g^*(y) . \end{aligned}$$

The first statement also implies that  $(f^* \square g^*)^* = f^{**} + g^{**} = f + g$  by [Theorem 9.7](#). By applying convex conjugation to both sides,  $(f^* \square g^*)^{**} = (f + g)^*$ , which implies the second statement if  $f^* \square g^*$  equals its double conjugate. For this, we need to know that  $f^* \square g^*$  is regular convex, which follows from the condition on the domains (see [\[Roc97, Theorem 16.4\]](#)).  $\square$

There is a surprising analogy with the Fourier transform, which transforms convolution into multiplication. Recall that for  $f, g : \mathbb{R}^d \rightarrow \mathbb{C}$ , the Fourier transform is given by  $\mathcal{F}f(\xi) := \int f(x) \exp(-2\pi i \langle \xi, x \rangle) dx$ , the convolution is given by  $(f * g)(y) := \int f(x) g(y - x) dx$ , and we have the key property  $\mathcal{F}(f * g) = \mathcal{F}f \mathcal{F}g$ .

To see a connection more precisely, note that we usually work with the algebra  $(+, \cdot)$  with its familiar properties: there is an additive identity  $0$  such that  $x + 0 = 0 + x = x$  for all  $x$ ; every  $x$  has an additive inverse  $-x$  satisfying  $x + (-x) = 0$ ; multiplication distributes over addition; etc. Now introduce a new structure, consisting of the operations  $(\min, +)$ . This shares some properties with the usual algebra: the identity element for  $\min$  is  $+\infty$ , and  $+$  distributes over  $\min$ , i.e.,  $x + \min(y, z) = \min(x + y, x + z)$ . However, we also lose some properties: e.g., not every element has an inverse for the  $\min$  operation. This is sometimes known as the *min-plus algebra* despite the fact that it is not technically an algebra; more accurately, it is called the *tropical semiring*.

If we think of integrals as continuous summations, then convolution is a sum of products; infimal convolution is a min of sums. Hence, infimal convolution is the tropical analogue of convolution. The following table summarizes further analogies.

$(+, \times)$	$(\min, +)$
convolution	infimal convolution
Fourier transform	convex conjugate
Gaussians	convex quadratics
diffusion processes	gradient flow
heat equation	Hamilton–Jacobi equation
heat semigroup	Hopf–Lax semigroup

We conclude this discussion by using this perspective to show that the Hopf–Lax semigroup solves the following PDE, known as the *Hamilton–Jacobi equation*:

$$\partial_t u + H(\nabla_x u) = 0. \quad (9.4)$$

The proof is patterned on the following derivation of the solution to the heat equation  $\partial_t u = \Delta u$  with initial condition  $u(0, \cdot) = f$ ; here  $\Delta u = \sum_{i=1}^d \partial_i^2 u$  is the Laplacian. If we take the Fourier transform of both sides of the equation, then  $\partial_t \mathcal{F}u = \mathcal{F} \Delta u = -4\pi^2 \|\cdot\|^2 \mathcal{F}u$ , where the last equality follows from differentiating the Fourier transform under the integral. This implies that  $\partial_t \log \mathcal{F}u = -4\pi^2 \|\cdot\|^2$ , or  $\mathcal{F}u(t, \cdot) = \mathcal{F}f \exp(-4\pi^2 t \|\cdot\|^2)$ . Using the fact that the inverse Fourier transform transforms multiplication into convolution, one can then show that  $u(t, \cdot) = f * \mathcal{F} \exp(-4\pi^2 t \|\cdot\|^2) = f * \text{normal}(0, 4tI)$ .

In the same way, we start with (9.4) and take the convex conjugate of both sides. Using the shorthand notation  $f_t := u(t, \cdot)$ , and since  $f_t^*(p) = \sup_{v \in \mathbb{R}^d} \{\langle p, v \rangle - f_t(v)\}$  with the supremum attained at  $v = \nabla f_t^*(p)$ ,

$$\partial_t f_t^*(p) = -\partial_t f_t(\nabla f_t^*(p)) = H(\nabla f_t(\nabla f_t^*(p))) = H(p).$$

Hence,  $f_t^* = tH + f$  and  $f_t = (tH)^* \square f = tL(\cdot/t) \square f$ . Thus, the solution to (9.4) is given by the Hopf–Lax semigroup as claimed.

When  $L = H = \frac{1}{2} \|\cdot\|^2$ , (9.4) becomes  $\partial_t u + \frac{1}{2} \|\nabla_x u\|^2 = 0$  and the Hopf–Lax semigroup  $Q_t f$  coincides with the Moreau–Yosida envelope (Definition 8.2). This yields an unexpected connection between the Hamilton–Jacobi equation and the PPM.

## 9.2 Duality correspondences

**Theorem 9.6 (Fenchel–Young).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be regular and convex. Then,

$$f(x) + f^*(p) \geq \langle p, x \rangle \quad \text{for all } p, x \in \mathbb{R}^d.$$

Moreover, equality holds if and only if  $p \in \partial f(x)$ , if and only if  $x \in \partial f^*(p)$ .

*Proof.* The inequality is trivial from the definition of  $f^*$ . If equality holds, then for any  $p', x' \in \mathbb{R}^d$ ,

$$\begin{aligned} f(x') &\geq \langle p, x' \rangle - f^*(p) = f(x) + \langle p, x' - x \rangle, \\ f^*(p') &\geq \langle p', x \rangle - f(x) = f^*(p) + \langle x, p' - p \rangle, \end{aligned}$$

i.e.,  $p \in \partial f(x)$  and  $x \in \partial f^*(p)$ . Conversely, if  $p \in \partial f(x)$ , then

$$f^*(p) = \sup_{x' \in \mathbb{R}^d} \{\langle p, x' \rangle - f(x')\} \leq \langle p, x \rangle - f(x). \quad \square$$

**Theorem 9.7 (double conjugation).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ . Then,  $f \geq f^{**}$ . Moreover, if  $f$  is regular and convex, then equality holds:  $f = f^{**}$ .

*Proof.* For the first statement,

$$f^{**}(z) = \sup_{y \in \mathbb{R}^d} \left\{ \langle y, z \rangle - \sup_{x \in \mathbb{R}^d} \{\langle x, y \rangle - f(x)\} \right\} = \sup_{y \in \mathbb{R}^d} \inf_{x \in \mathbb{R}^d} \{\langle y, z - x \rangle + f(x)\} \leq f(z) \quad (9.5)$$

by choosing  $x = z$ .

Now assume that  $f$  is regular and convex. If  $z \in \text{int dom } f$ , then by [Theorem 6.9](#) there exists  $p \in \partial f(z)$ , so that  $f(x) \geq f(z) + \langle p, x - z \rangle$  for all  $x \in \mathbb{R}^d$ . By taking  $y = p$ ,

$$f^{**}(z) \geq \inf_{x \in \mathbb{R}^d} \{\langle p, z - x \rangle + f(x)\} \geq f(z),$$

which proves the equality for such  $z$ . For brevity, we omit the proof for  $z \notin \text{int dom } f$  (see, e.g., [Roc97](#), Theorem 12.2).  $\square$

This result implies that in general, if  $f_*$  is the largest convex and lower semicontinuous function which is smaller than  $f$ , then  $f_* = f^{**}$ . Indeed,  $f_* \geq f^{**}$  by definition, whereas  $f \geq f_*$  implies  $f^{**} \geq (f_*)^{**} = f_*$ . The proof above also shows that whenever  $\partial f(x) \neq \emptyset$ , then  $f(x) = f^{**}(x)$ . In particular, if  $x_*$  is a minimizer of  $f$ , then  $0 \in \partial f(x_*)$  and  $f(x_*) = f^{**}(x_*)$ ; moreover, by taking  $y = 0$  in (9.5) we see that  $\inf f = \inf f^{**}$ . Thus, we can start with a non-convex function  $f$  and “convexify” it by replacing it with  $f^{**}$  while preserving the optimal value, although this is seldom useful in practice.

Properties of  $f$  are often reflected as “dual” properties for  $f^*$ . For example, if  $f$  is regular convex, the following assertions hold (see [Roc97](#)):

- $f$  is Lipschitz if and only if  $\text{dom } f^*$  is bounded.
- $\text{epi } f$  contains no non-vertical half-lines if and only if  $\text{dom } f^* = \mathbb{R}^d$ .

- $f$  has no lines along which it is affine if and only if  $\text{int dom } f^* \neq \emptyset$ .
- $f$  has bounded level sets if and only if  $0 \in \text{int dom } f^*$ .
- $f$  is differentiable at  $x$  with  $\nabla f(x) = p$  if and only if  $(p, f^*(p))$  is an exposed point of  $\text{epi } f^*$ . (An *exposed point* of a convex set is a point at which some linear function attains its strict maximum over the convex set.)

For our purposes, we are most interested in conditions under which  $\nabla f$  is a well-defined bijection from an open convex set  $\mathcal{C}$  to its image  $\nabla f(\mathcal{C})$ , with inverse given by  $(\nabla f)^{-1} = \nabla f^*$ . In this case, the correspondence between  $f$  and  $f^*$  is known as the Legendre transformation and we informally discussed it in the previous subsection. We accept the results in the following discussion without proof; see [Roc97, §26] for details.

**Definition 9.8.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be regular convex.

- We say that  $f$  is **essentially smooth** if  $f$  is differentiable on  $\mathcal{C} := \text{int dom } f$  and  $\lim_{n \rightarrow \infty} \|\nabla f(x_n)\| \rightarrow \infty$  whenever  $\{x_n\}_{n \in \mathbb{N}}$  is a sequence in  $\mathcal{C}$  converging to  $\partial \mathcal{C}$ .
- We say that  $f$  is **essentially strictly convex** if  $f$  is strictly convex on every convex subset of  $\text{dom } \partial f := \{x \in \mathbb{R}^d : \partial f(x) \neq \emptyset\}$ .

**Lemma 9.9.** A regular convex function  $f$  is essentially smooth if and only if  $f^*$  is essentially strictly convex.

**Theorem 9.10.** Let  $f$  be regular, strictly convex, and essentially smooth over  $\mathcal{C} = \text{int dom } f$ . Then,  $f^*$  is regular, strictly convex, and essentially smooth over  $\mathcal{C}^* := \text{int dom } f^*$ . Moreover,  $\nabla f : \mathcal{C} \rightarrow \mathcal{C}^*$  is a continuous bijection with  $(\nabla f)^{-1} = \nabla f^*$ .

**Definition 9.11.** We say that a function  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is of **Legendre type** if it satisfies the assumptions of Theorem 9.10.

To summarize, the condition that  $f$  is regular convex ensures duality at the level of  $f = f^{**}$ . The condition that  $f$  is of Legendre type ensures duality at the level of  $(\nabla f)^{-1} = \nabla f^*$ . Note also that if  $f, f^*$  are sufficiently smooth, then by differentiating the equality  $\nabla f(\nabla f^*) = \text{id}$  one obtains the identity

$$\nabla^2 f \circ \nabla f^* = [\nabla^2 f^*]^{-1}.$$

In particular,  $\nabla^2 f \geq \alpha I$  is equivalent to  $[\nabla^2 f^*]^{-1} \leq \alpha^{-1} I$ , i.e., there is a duality between the properties of strong convexity and smoothness. Let us prove this last fact without assuming differentiability.

**Lemma 9.12 (convexity–smoothness duality).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be regular and  $\alpha$ -convex for some  $\alpha > 0$ . Then,  $f^*$  is  $\alpha^{-1}$ -smooth.

*Proof.* By the duality correspondences (including Lemma 9.9),  $\text{dom } f^* = \mathbb{R}^d$  and  $f^*$  is differentiable everywhere. For two points  $y, y' \in \mathbb{R}^d$ , let  $x, x' \in \mathbb{R}^d$  achieve the suprema in the definitions of  $f^*(y), f^*(y')$  respectively. By Theorem 9.6,  $x = \nabla f^*(y)$  and  $x' = \nabla f^*(y')$ . Then, by strong convexity of  $f - \langle \cdot, y \rangle$ ,

$$f(x') - \langle x', y \rangle \geq f(x) - \langle x, y \rangle + \frac{\alpha}{2} \|x' - x\|^2.$$

Adding this to the analogous inequality with  $x$  and  $x'$  swapped,

$$\begin{aligned} \alpha \|\nabla f^*(y') - \nabla f^*(y)\|^2 &= \alpha \|x' - x\|^2 \leq \langle x, y \rangle + \langle x', y' \rangle - \langle x', y \rangle - \langle x, y' \rangle \\ &= \langle x' - x, y' - y \rangle. \end{aligned}$$

Rearranging this after Cauchy–Schwarz proves  $\|\nabla f^*(y') - \nabla f^*(y)\| \leq \alpha^{-1} \|y' - y\|$ .  $\square$

**Corollary 9.13 (contractivity of the proximal operator).** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be  $\alpha$ -convex. Then,  $\text{prox}_f$  is  $1/(1 + \alpha)$ -Lipschitz.

*Proof.* We can write

$$\text{prox}_f(y) := \arg \min_{x \in \mathbb{R}^d} \left\{ f(x) + \frac{1}{2} \|y - x\|^2 \right\} = - \arg \max_{x \in \mathbb{R}^d} \left\{ \langle x, y \rangle - f(x) - \frac{1}{2} \|x\|^2 \right\}.$$

This shows that  $-\text{prox}_f$  is the gradient of the convex conjugate of the function  $f + \frac{1}{2} \|\cdot\|^2$ , which is  $(1 + \alpha)$ -convex.  $\square$

For a closed convex set  $\mathcal{C}$ ,  $\Pi_{\mathcal{C}} = \text{prox}_{\chi_{\mathcal{C}}}$ , so this corollary recovers Lemma 6.13. It also shows that PPM with an  $\alpha$ -convex function contracts with rate  $1/(1 + \alpha h)$ .



## Bibliographical notes

The treatment of classical mechanics is based on [Eva10, §3.3]. Variational principles also lead to an influential perspective on acceleration, see [WWJ16].

The problem of minimizing an action functional can be generalized to the problem of optimal control, and in that context, the corresponding Hamilton–Jacobi equation is known as the *Hamilton–Jacobi–Bellman equation*. The analogy between dynamic programming and diffusions can be pushed even further to obtain “laws of large numbers” and “central limit theorems” for the former, see [Bac+92, §9.4].

The result of [Exercise 9.2](#) is from [Che+22].

## Exercises

**Exercise 9.1.** Verify the assertions in [Example 9.2](#).

**Exercise 9.2.** To avoid technical difficulties, assume that  $f$  is differentiable everywhere and satisfies (PL) with constant  $\alpha > 0$ . The goal of this exercise is to derive the sharp rate of convergence of the PPM in this setting (which turns out to be non-trivial).

For any  $x \in \mathbb{R}^d$ , let  $(Q_t)_{t \geq 0}$  denote the Hopf–Lax semigroup and  $x_t := \text{prox}_{tf}(x)$ . From the Hamilton–Jacobi equation (9.4) or via direct computation, show that  $\partial_t Q_t f(x) = -\|x_t - x\|^2 / (2t^2)$ . Use this to deduce that

$$\partial_t \{Q_t f(x) - f_\star\} \leq -\frac{\alpha}{1 + \alpha t} \{Q_t f(x) - f_\star\}. \quad (9.6)$$

Finally, by adapting Grönwall’s lemma ([Lemma 2.3](#)), use this to prove the *sharp* rate

$$f(x_h) - f_\star \leq \frac{f(x) - f_\star}{(1 + \alpha h)^2}.$$

**Exercise 9.3.** The inequality (9.6) implies that

$$Q_h f(x) - f_\star \leq \frac{1}{1 + \alpha h} (f(x) - f_\star). \quad (9.7)$$

In this exercise, we give an alternative proof of this fact under the assumption that  $f$  is  $\alpha$ -convex. As a consequence, we also obtain a result for  $\alpha = 0$ .

Write out the definition of  $Q_h f(x)$  as an infimum, and choose as a test point the interpolant  $(1 - t)x_\star + tx$ . Pick  $t > 0$  to derive (9.7). Also, in the case  $\alpha = 0$ , show that if  $f(x) - f_\star \leq \|x - x_\star\|^2 / h$ ,

$$Q_h f(x) - f_\star \leq \left(1 - \frac{f(x) - f_\star}{2\|x - x_\star\|^2} h\right) (f(x) - f_\star).$$

## 10 [3/4–3/6] Mirror methods

Consider the following situation.

**Example 10.1 (optimization in a different norm).** Suppose that  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is convex and that we wish to minimize it over the simplex  $\Delta_d := \{x \in \mathbb{R}_+^d : \sum_{i=1}^d x[i] = 1\}$ . If  $f$  is Lipschitz, then we can apply PSD and obtain an  $\varepsilon$ -approximate solution in  $O(L^2 R^2 / \varepsilon^2)$  iterations. Here,  $L$  is the Lipschitz constant, and  $R \leq 2$  is the radius.

For example, suppose that we have  $d$  actions and that the loss of the  $i$ -th action is  $\ell[i]$ , where the losses are bounded:  $|\ell[i]| \leq 1$ . If we choose an action randomly according to a probability distribution  $x \in \Delta_d$ , the expected loss is  $\langle \ell, x \rangle =: f(x)$ . We then seek to minimize the expected loss over  $\Delta_d$ .<sup>†</sup> The Lipschitz constant of  $f$  is  $\|\ell\|$ , which could be as large as  $\sqrt{d}$  in the worst case; the resulting complexity estimate of  $O(d/\varepsilon^2)$  is poor in high dimension.

Implicit in this discussion, however, is that we are measuring the Lipschitz constant and the radius with respect to the usual Euclidean norm. In this setting, however, it may make more sense to use the  $\ell_1$  norm, in which case the Lipschitz constant is  $\|\ell\|_\infty \leq 1$ .

<sup>†</sup>Trivially, the solution to the optimization problem is given by the distribution which puts all of its mass on  $\arg \min_{i \in [d]} \ell[i]$ . This problem is simply meant to illustrate the pitfalls of naively using the Euclidean norm, but it also forms the basis for the more interesting setting of Example 10.14.

Until now, we have identified points  $x$  and gradients  $\nabla f(x)$  as part of the same space  $\mathbb{R}^d$ , but this is because of the self-dual nature of the Euclidean norm. Suppose now that  $(\mathcal{X}, \|\cdot\|)$  is a general (finite-dimensional) normed vector space and  $f : \mathcal{X} \rightarrow \mathbb{R} \cup \{\infty\}$ . The dual space is  $(\mathcal{X}^*, \|\cdot\|_*)$ , where  $\mathcal{X}^*$  is the space of linear functionals  $\ell : \mathcal{X} \rightarrow \mathbb{R}$ , equipped with the dual norm  $\|\ell\|_* := \sup\{|\ell(x)| : \|x\| \leq 1\}$ . The derivative of  $f$  at  $x$  is defined to be the linearization at  $x$ : if there exists an element  $\ell \in \mathcal{X}^*$  such that

$$|f(x+v) - f(x) - \ell(v)| = o(\|v\|) \quad \text{as } v \rightarrow 0,$$

we say that  $f$  is differentiable<sup>10</sup> at  $x$  and we write  $Df(x)$  for the functional  $\ell$ . Note that in this formalism, the derivative  $Df(x)$  is an element of the dual space.

Above, we wrote  $Df(x)$  instead of  $\nabla f(x)$  to emphasize that in this context, we should no longer think of  $Df(x)$  as belonging to the original space  $\mathcal{X}$ . However, when  $\mathcal{X} = \mathbb{R}^d$ , it is still convenient to identify  $Df(x)$  as a vector in  $\mathbb{R}^d$ , and we therefore continue to use the notation  $\nabla f(x)$ . This is fine as long as we remember the following two points:

- It does not make sense to add a point  $x \in \mathcal{X}$  to a gradient  $\nabla f(x) \in \mathcal{X}^*$ .

<sup>10</sup>Strictly speaking, this is the Fréchet derivative.

- The size of  $\nabla f(x)$  should be measured in the dual norm  $|||\cdot|||_*$ .

In the context of [Example 10.1](#), the dual norm for  $\|\cdot\|_1$  is the  $\ell_\infty$  norm  $\|\cdot\|_\infty$ .

Immediately, the first point above rules out [GD](#) and [PSD](#) as sensible algorithms. A first attempt to remedy this issue is to somehow develop analogues of these algorithms in different norms, but this is seriously complicated by the fact that non-Euclidean norms lack crucial properties, e.g., we cannot “expand the square” as we did in previous proofs.

Instead, the idea of [\[NY83\]](#) is to use the Fenchel–Legendre duality.

Throughout this section,  $\phi : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  is a convex function of Legendre type. We refer to it as the *mirror map*.

The idea is to use the auxiliary function  $\phi$  to map the iterate  $x_n$  into the dual space via  $x_n^* = \nabla \phi(x_n)$ . Now that we are in the dual space, it makes sense to take a gradient step:  $x_{n+1}^* = x_n^* - h \nabla f(x_n)$ . Then, we use  $\nabla \phi^*$  to return:  $x_{n+1} = \nabla \phi^*(x_{n+1}^*)$ . The goal of this section is to formalize this idea and its analysis.

## 10.1 Bregman divergences and relative convexity/smoothness

We introduce a key definition which substitutes for the squared Euclidean norm.

**Definition 10.2.** Given a function  $\phi : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  of Legendre type over  $\mathcal{C}_\phi$ , the corresponding **Bregman divergence** associated with  $\phi$  is the map  $D_\phi : \mathbb{R}^d \times \mathcal{C}_\phi \rightarrow \mathbb{R} \cup \{\infty\}$  defined by

$$D_\phi(x, y) := \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle.$$

In words,  $D_\phi(\cdot, y)$  is defined by subtracting from  $\phi$  its linearization at  $y$ . We can observe the following properties:

- $D_\phi \geq 0$ : this is equivalent to the convexity of  $\phi$ .
- $D_\phi$  is convex with respect to its first argument.
- If  $\phi$  is twice continuously differentiable, then

$$D_\phi(x, y) \sim \frac{1}{2} \langle x - y, \nabla^2 \phi(y) (x - y) \rangle \quad \text{as } x \rightarrow y. \quad (10.1)$$

The last property indicates that  $D_\phi$  should behave as a squared distance between  $x$  and  $y$ . In some respects this is true, e.g.,  $D_\phi$  satisfies a Pythagorean inequality ([Exercise 10.1](#)). However, (10.1) is a purely local statement, and a priori there does not seem to be a reason for  $D_\phi$  to have useful global properties. For example,  $D_\phi$  is *asymmetric*, and  $\sqrt{D_\phi}$  does not in general satisfy a triangle inequality. Nevertheless, it turns out that  $D_\phi$  is a powerful global measure of progress, which is arguably the greatest surprise of mirror methods.

**Example 10.3 (mirror maps).**

1. Let  $\phi(x) = \frac{1}{2} \|x\|^2$ . Then,  $\nabla\phi$  is the identity mapping and  $D_\phi$  is one-half times the squared Euclidean distance. So, our study of mirror methods subsumes the preceding Euclidean methods.
2. Let  $\phi(x) = \sum_{i=1}^d \{x[i] \log x[i] - x[i]\}$  for  $x \in \mathbb{R}_+^d$ . Then,  $\nabla\phi(x) = \log x$ , where  $\log$  is applied coordinate-wise. The associated Bregman divergence is the Kullback–Leibler divergence  $D_\phi(x, y) = \sum_{i=1}^d \{x[i] \log(x[i]/y[i]) - x[i] + y[i]\}$ .
3. Let  $\phi(X) = \text{tr}(X \log X - X)$  for  $X \succ 0$ ; this is known as the von Neumann entropy. The associated Bregman divergence is the quantum relative entropy  $D_\phi(X, Y) = \text{tr}(X (\log X - \log Y) - X + Y)$ .

Let us define notions of convexity and smoothness *relative to*  $\phi$ .

**Definition 10.4.** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be differentiable on  $\text{int dom } f \subseteq \mathcal{C}_\phi$ .

- We say that  $f$  is  $\alpha$ -**convex relative to**  $\phi$  if  $D_f \geq \alpha D_\phi$ .
- We say that  $f$  is  $\beta$ -**smooth relative to**  $\phi$  if  $D_f \leq \beta D_\phi$ .

Similarly to §1.2, there are equivalent reformulations of these definitions; see [LFN18, Proposition 1.1] for details.

**Proposition 10.5 (relative convexity).** For any  $\alpha \geq 0$ , the following are equivalent.

- $f$  is  $\alpha$ -convex relative to  $\phi$ .
- $f - \alpha\phi$  is convex.
- $\langle \nabla f(y) - \nabla f(x), y - x \rangle \geq \alpha \langle \nabla\phi(y) - \nabla\phi(x), y - x \rangle$  for all  $x, y \in \text{int dom } f$ .

If  $f$  is twice continuously differentiable on  $\text{int dom } f$ , the above are also equivalent to:

- $\nabla^2 f \geq \alpha \nabla^2 \phi$  on  $\text{int dom } f$ .

**Proposition 10.6 (relative smoothness).** For any  $\beta \geq 0$ , the following are equivalent.

- $f$  is  $\beta$ -smooth relative to  $\phi$ .
- $\beta\phi - f$  is convex.
- $\langle \nabla f(y) - \nabla f(x), y - x \rangle \leq \beta \langle \nabla \phi(y) - \nabla \phi(x), y - x \rangle$  for all  $x, y \in \text{int dom } f$ .

If  $f$  is twice continuously differentiable on  $\text{int dom } f$ , the above are also equivalent to:

- $\nabla^2 f \leq \beta \nabla^2 \phi$  on  $\text{int dom } f$ .

For the case of  $\phi = \frac{1}{2} \|\cdot\|^2$ , we recover the usual notions of convexity and smoothness described in §1.2. These relative definitions satisfy similar properties as convexity/smoothness, e.g., if  $f_1, f_2$  are respectively  $\alpha_1$ - and  $\alpha_2$ -convex relative to  $\phi$  and  $\lambda_1, \lambda_2 > 0$ , then  $\lambda_1 f_1 + \lambda_2 f_2$  is  $(\lambda_1 \alpha_1 + \lambda_2 \alpha_2)$ -convex relative to  $\phi$ . Also, we have a growth bound.

**Lemma 10.7 (relative growth).** Suppose that  $f$  is  $\alpha$ -convex relative to  $\phi$  for some  $\alpha > 0$ , and that  $f$  is minimized at an interior point  $x_\star$  of its domain. Then, for all  $x \in \mathbb{R}^d$ ,

$$f(x) - f_\star \geq \alpha D_\phi(x, x_\star).$$

*Proof.* The left-hand side is  $D_f(x, x_\star)$ . □

Other useful properties of Bregman divergences are explored in [Exercise 10.1](#).

## 10.2 Algorithms and convergence analysis

Briefly, let us first consider the continuous-time picture. Since we add the gradient of  $f$  in the dual space, the dynamics we consider evolve according to

$$\partial_t \nabla \phi(x_t) = -\nabla f(x_t). \quad (10.2)$$

By the chain rule, this is equivalent to the following evolution in the primal space:

$$\dot{x}_t = -[\nabla^2 \phi(x_t)]^{-1} \nabla f(x_t). \quad (10.3)$$

This can be interpreted as a preconditioned gradient flow.

Despite the fact that (10.2) and (10.3) are equivalent in continuous time, they lead to different discretizations. The discretization of (10.3) is usually called natural gradient descent and it is related to the subject of information geometry [AN00]. In fact, one can

view the use of the mirror map  $\phi$  as equipping the space  $\mathcal{C}_\phi$  with a Riemannian metric, namely, a local inner product  $\langle u, v \rangle_x := \langle u, \nabla^2 \phi(x) v \rangle$ . This turns  $\mathcal{C}_\phi$  into a so-called Hessian manifold. From this perspective, the natural objects are geometric in nature: geodesics, length, curvature, etc.

However, this is **not** what we consider; mirror descent is obtained from discretization of (10.2) in the dual. This conceptual point is so important that we isolate it into a remark.

**Remark 10.8.** The key distinguishing feature of mirror methods from preconditioned or Riemannian gradient methods is the existence of the *global* progress measure given by the Bregman divergence  $D_\phi$ . In contrast, preconditioned/Riemannian gradient methods are purely *local* in nature.

Now that we have emphasized the conceptual underpinnings of the methods, let us now turn to concrete algorithms. We begin with the smooth case, and we consider the following *mirror proximal gradient descent* method:

$$x_{n+1} := \arg \min_{x \in \mathbb{R}^d} \left\{ f(x_n) + \langle \nabla f(x_n), x - x_n \rangle + g(x) + \frac{1}{h} D_\phi(x, x_n) \right\}. \quad (\text{MPGD})$$

Note that this incorporates the proximal splitting considered in §8, except that we replace  $\frac{1}{2} \|x - x_n\|^2$  with the more general  $D_\phi(x, x_n)$ . We consider this iteration for the sake of generality, since it encompasses the following algorithms.

- When  $g = 0$ , since  $\nabla_1 D_\phi(x, x_n) = \nabla \phi(x) - \nabla \phi(x_n)$ , the first-order optimality condition reads

$$\nabla \phi(x_{n+1}) = \nabla \phi(x_n) - h \nabla f(x_n).$$

This is known as *mirror descent*.

- When  $f = 0$ , we obtain

$$x_{n+1} = \arg \min_{x \in \mathbb{R}^d} \left\{ g(x) + \frac{1}{h} D_\phi(x, x_n) \right\} =: \text{prox}_{hg}^\phi(x_n),$$

which is the *mirror proximal point method*.

- When  $g = \chi_{\mathcal{C}}$ , where  $\mathcal{C} \subseteq \mathcal{C}_\phi$  is a closed convex set,

$$x_{n+1} = \arg \min_{x \in \mathcal{C}} \left\{ \langle \nabla f(x_n), x - x_n \rangle + \frac{1}{h} (\phi(x) - \langle \nabla \phi(x_n), x - x_n \rangle) \right\}$$

$$\begin{aligned}
&= \arg \min_{x \in \mathcal{C}} \{ \phi(x) - \langle \nabla \phi(x_n) - h \nabla f(x_n), x - x_n \rangle \} \\
&= \Pi_{\mathcal{C}}^{\phi}(\nabla \phi^*(\nabla \phi(x_n) - h \nabla f(x_n))),
\end{aligned}$$

where  $\Pi_{\mathcal{C}}^{\phi}$  is the Bregman projection (see [Exercise 10.1](#)). This is the mirror analogue of projected gradient descent.

**Theorem 10.9 (convergence of MPGD).** Let  $f$  be  $\alpha_f$ -convex and  $\beta_f$ -smooth, and let  $g$  be  $\alpha_g$ -convex, all relative to  $\phi$ . Let the step size  $h$  satisfy  $h \leq 1/\beta_f$ , let  $x^+$  denote the next iterate of [MPGD](#) started from  $x$ , and let  $y \in \mathbb{R}^d$ . Then,

$$(1 + \alpha_g h) D_{\phi}(y, x^+) \leq (1 - \alpha_f h) D_{\phi}(y, x) - h (F(x^+) - F(y)).$$

In particular, if we set  $y = x_{\star}$  and iterate, it yields

$$F(x_N) - F_{\star} \leq \frac{\alpha_f + \alpha_g}{\lambda_h^{-N} - 1} D_{\phi}(x_{\star}, x_0),$$

where  $\lambda_h := (1 - \alpha_f h)/(1 + \alpha_g h)$ .

*Proof.* The proof is patterned upon the proof of [Theorem 8.5](#). Let  $\psi_x$  denote the objective in ([MPGD](#)) starting from  $x$  (rather than  $x_n$ ). Then,  $\psi_x$  is  $(\alpha_g + 1/h)$ -convex relative to  $\phi$  with minimizer  $x^+$ , so by the growth inequality ([Lemma 10.7](#)),

$$\psi_x(y) \geq \psi_x(x^+) + \left(\alpha_g + \frac{1}{h}\right) D_{\phi}(y, x^+).$$

On one hand, by  $\alpha_f$ -convexity,

$$\begin{aligned}
\psi_x(y) &= f(x) + \langle \nabla f(x), y - x \rangle + g(y) + \frac{1}{h} D_{\phi}(y, x) \\
&= f(y) - D_f(y, x) + g(y) + \frac{1}{h} D_{\phi}(y, x) \leq F(y) + \left(\frac{1}{h} - \alpha_f\right) D_{\phi}(y, x).
\end{aligned}$$

On the other hand, by  $\beta_f$ -smoothness,

$$\begin{aligned}
\psi_x(x^+) &= f(x) + \langle \nabla f(x), x^+ - x \rangle + g(x^+) + \frac{1}{h} D_{\phi}(x^+, x) \\
&= f(x^+) - D_f(x^+, x) + g(x^+) + \frac{1}{h} D_{\phi}(x^+, x) \geq F(x^+) + \left(\frac{1}{h} - \beta_f\right) D_{\phi}(x^+, x).
\end{aligned}$$

Drop the term  $(1/h - \beta_f) D_\phi(x^+, x)$  and combine the inequalities to prove the one-step bound. If we set  $y = x$  in the one-step bound, it yields the descent lemma  $F(x^+) - F(x) \leq -h^{-1} (1 + \alpha_g h) D_\phi(x, x^+) \leq 0$ , so we can iterate the one-step bound using the discrete Grönwall lemma ([Lemma 3.5](#)).  $\square$

Although this result is the analogue of the smooth convergence rate for [GD](#) ([Theorem 3.4](#)), since  $\nabla\phi$  necessarily blows up at the boundary  $\partial\mathcal{C}_\phi$ , so can  $\nabla f$ . Therefore, this theorem actually covers examples in which  $f$  is not at all smooth in the usual sense.

To relate this back to [Example 10.1](#), consider convexity/smoothness relative to a norm.

**Definition 10.10.** A function  $f$  is  $\alpha$ -**convex** (resp.  $\beta$ -**smooth**) **relative to a norm**  $\|\cdot\|$  if for all  $x, y \in \text{int dom } f$ ,

$$D_f(x, y) \geq \frac{\alpha}{2} \|y - x\|^2 \quad (\text{resp. } D_f \leq \frac{\beta}{2} \|y - x\|^2).$$

Suppose that  $\phi$  is strongly convex relative to a norm  $\|\cdot\|$ . Then, to check that  $f$  is smooth relative to  $\phi$ , it suffices to check that  $f$  is smooth relative to  $\|\cdot\|$ , so the norm can act as a useful intermediary. Moreover, whereas the Bregman structure is crucial for carrying out the iterative analysis of [MPGD](#), the norm structure is often convenient too, e.g., for the use of tools such as Cauchy–Schwarz.

To illustrate this, we now consider the non-smooth case. Here, we assume that  $f$  is Lipschitz with respect to  $\|\cdot\|$ :

$$|f(x) - f(y)| \leq L \|x - y\| \quad \text{for all } x, y \in \mathcal{C}_\phi.$$

We again consider [MPGD](#), except that  $\nabla f(x_n)$  should be interpreted as a subgradient; we leave the notation unchanged because it should not cause confusion. The Lipschitz condition is then equivalent to the subgradient bound

$$\|\nabla f(x)\|_* \leq L \quad \text{for all } x \in \mathcal{C}_\phi.$$



**Theorem 10.11** (convergence of MPGD, non-smooth case). Let  $f$  and  $g$  be convex, and let  $f$  be  $L$ -Lipschitz with respect to a norm  $\|\cdot\|$ . Let  $\phi$  be  $\alpha_\phi$ -convex relative to  $\|\cdot\|$ . Then, for MPGD, it holds that

$$F\left(\frac{1}{N} \sum_{n=1}^N x_n\right) - F_\star \leq \frac{1}{N} \sum_{n=1}^N (F(x_n) - F_\star) \leq \frac{D_\phi(x_\star, x_0)}{Nh} + \frac{2L^2h}{\alpha_\phi}.$$

In particular, if  $R_\phi^2 \geq D_\phi(x_\star, x_0)$  and we choose step size  $h^2 = \alpha_\phi R_\phi^2 / (2L^2N)$ , then

$$F\left(\frac{1}{N} \sum_{n=1}^N x_n\right) - F_\star \leq LR_\phi \sqrt{\frac{8}{\alpha_\phi N}}.$$

*Proof.* Following the proof of Theorem 10.9, we still have

$$\psi_x(x^+) + \frac{1}{h} D_\phi(x_\star, x^+) \leq \psi_x(x_\star) \leq F(x_\star) + \frac{1}{h} D_\phi(x_\star, x).$$

In the lower bound for  $\psi_x(x^+)$ , we originally used smoothness to upper bound  $D_f(x^+, x)$ , which is no longer available to us. Instead, by Cauchy–Schwarz,

$$\begin{aligned} D_f(x^+, x) &= f(x^+) - f(x) - \langle \nabla f(x), x^+ - x \rangle \leq L \|x^+ - x\| + \|\nabla f(x)\|_* \|x^+ - x\| \\ &\leq 2L \|x^+ - x\|. \end{aligned}$$

Thus,

$$\begin{aligned} \psi_x(x^+) &= F(x^+) - D_f(x^+, x) + \frac{1}{h} D_\phi(x^+, x) \geq F(x^+) - 2L \|x^+ - x\| + \frac{\alpha_\phi}{2h} \|x^+ - x\|^2 \\ &\geq F(x^+) - \frac{2L^2h}{\alpha_\phi}. \end{aligned}$$

This leads to the one-step bound

$$D_\phi(x_\star, x^+) \leq D_\phi(x_\star, x) - h(F(x^+) - F_\star) + \frac{2L^2h^2}{\alpha_\phi}.$$

Iterating this inequality finishes the proof.  $\square$

**Example 10.12 (optimization over the simplex).** We return to [Example 10.1](#), and we use the entropic mirror map  $\phi(x) = \sum_{i=1}^d \{x[i] \log x[i] - x[i]\}$ . Then,  $\phi$  is 1-convex relative to the  $\ell_1$ -norm  $\|\cdot\|_1$  over the probability simplex  $\Delta_d$ ; this is known as Pinsker's inequality ([Exercise 10.3](#)).

To minimize  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  over  $\Delta_d$ , we apply [MPGD](#) with  $g = \chi_{\Delta_d}$ . Then,

$$\nabla \phi^*(\nabla \phi(x_n) - h \nabla f(x_n)) = x_n \odot \exp(-h \nabla f(x_n)),$$

where  $\exp$  is applied pointwise and  $\odot$  is the Hadamard (or pointwise) product. Also, one can check that  $\Pi_{\Delta_d}^\phi(x) = x/\|x\|_1$  simply normalizes the vector ([Exercise 10.4](#)). Hence, the algorithm reads

$$x_{n+1} = \frac{x_n \odot \exp(-h \nabla f(x_n))}{\|x_n \odot \exp(-h \nabla f(x_n))\|_1}.$$

Consider initializing at the uniform distribution  $x_0 = \mathbf{1}_d/d$ . Then, for any  $x_\star \in \Delta_d$ ,

$$D_\phi(x_\star, x_0) = \text{KL}(x_\star \parallel x_0) = \log d - \sum_{i=1}^d x_\star[i] \log \frac{1}{x_0[i]} \leq \log d,$$

by Jensen's inequality. Consequently, we can take  $R_\phi = \sqrt{\log d}$ , and

$$f\left(\frac{1}{N} \sum_{n=1}^N x_n\right) - f_\star \leq L_1 \sqrt{\frac{8 \log d}{N}},$$

where  $L_1$  is the Lipschitz constant of  $f$  in the  $\ell_1$  norm. This estimate is far better than the one described in [Example 10.1](#) for the Euclidean norm; we only pay an overhead which is logarithmic in the dimension.

### 10.3 Online algorithms and multiplicative weights

Let us examine the proof of [Theorem 10.11](#) once more. In that proof, we start with

$$\psi_x(x^+) + \frac{1}{h} D_\phi(y, x^+) \leq \psi_x(y),$$

which holds for all  $y \in \mathbb{R}^d$ . If we expand out the terms, this is equivalent to

$$\langle \nabla f(x), x^+ - x \rangle + g(x^+) + \frac{1}{h} D_\phi(x^+, x) + \frac{1}{h} D_\phi(y, x^+)$$

$$\leq \langle \nabla f(x), y - x \rangle + g(y) + \frac{1}{h} D_\phi(y, x).$$

On the left-hand side, if we apply Lipschitzness,

$$\langle \nabla f(x), x^+ - x \rangle + \frac{1}{h} D_\phi(x^+, x) \geq -\|\nabla f(x)\|_* \|x^+ - x\| + \frac{\alpha_\phi}{2h} \|x^+ - x\|^2 \geq -\frac{L^2 h}{2\alpha_\phi}.$$

If we now specialize to the case  $g = \chi_{\mathcal{C}}$ , then for any  $y \in \mathcal{C}$ ,

$$\langle \nabla f(x), x - y \rangle \leq \frac{1}{h} (D_\phi(y, x) - D_\phi(y, x^+)) + \frac{L^2 h}{2\alpha_\phi}.$$

Normally, we apply convexity to further lower bound the left-hand side, but let us now refrain from doing so. We make the observation that in the derivation thus far, *we have not used any property of  $f$* ; in fact, the same inequality holds even if  $\nabla f(x)$  is replaced by an *arbitrary* vector  $p \in \mathbb{R}^d$  with bounded dual norm, provided that we redefine the update in [MPGD](#) accordingly.

We now define the online version of mirror descent. Let  $\mathcal{C} \subseteq \mathcal{C}_\phi$  be a closed convex set, and let  $\{p_n\}_{n \in \mathbb{N}}$  be an arbitrary sequence of vectors. Define the updates

$$x_{n+1} := \arg \min_{x \in \mathcal{C}} \left\{ \langle p_n, x - x_n \rangle + \frac{1}{h} D_\phi(x, x_n) \right\} = \Pi_{\mathcal{C}}^\phi(\nabla \phi^*(\nabla \phi(x_n) - h p_n)). \quad (\text{OMD})$$

We immediately obtain the following theorem.

**Theorem 10.13 (regret guarantee for OMD).** Let  $\mathcal{C} \subseteq \mathcal{C}_\phi$  be a closed convex set, let  $\phi$  be  $\alpha_\phi$ -convex relative to a norm  $\|\cdot\|$ , and suppose that  $\{p_n\}_{n=0}^{N-1}$  are bounded in dual norm by  $L$ , i.e.,  $\|p_n\|_* \leq L$  for all  $n$ . Then, [OMD](#) satisfies

$$\sum_{n=0}^{T-1} \langle p_n, x_n \rangle \leq \inf_{y \in \mathcal{C}} \left\{ \sum_{n=0}^{T-1} \langle p_n, y \rangle + \frac{D_\phi(y, x_0)}{h} \right\} + \frac{L^2 T h}{2\alpha_\phi}.$$

In particular, if  $R_\phi^2 \geq \sup_{y \in \mathcal{C}} D_\phi(y, x_0)$  and  $h = R_\phi \sqrt{2\alpha_\phi} / (L\sqrt{T})$ , then

$$\sum_{n=0}^{T-1} \langle p_n, x_n \rangle \leq \inf_{y \in \mathcal{C}} \sum_{n=0}^{T-1} \langle p_n, y \rangle + L R_\phi \sqrt{2T/\alpha_\phi}.$$

In the setting of online learning (with full information feedback), at each round  $n$ , the player must play an action  $x_n$  belonging to some set  $\mathcal{C}$  of actions. An adversary then

chooses a loss function  $\ell_n$  belonging to some class of losses, the player incurs the loss  $\ell_n(x_n)$ , and the function  $\ell_n(\cdot)$  is revealed to the player. Thus, the total loss incurred by the player after  $T$  rounds is  $\sum_{n=0}^{T-1} \ell_n(x_n)$ . Since the losses are chosen in an adversarial fashion, one cannot hope to compete with a changing benchmark, so the measure of progress is to compare against the best *fixed* point that one could have played in hindsight, which incurs loss  $\inf_{y \in \mathcal{C}} \sum_{n=0}^{T-1} \ell_n(y)$ . The difference  $\sum_{n=0}^{T-1} \ell_n(x_n) - \inf_{y \in \mathcal{C}} \sum_{n=0}^{T-1} \ell_n(y)$  is called the *regret*, and the goal is to minimize it. In particular, regret bounds that scale linearly with  $T$  are often considered “trivial”, whereas regret bounds that scale as  $o(T)$  indicate that the algorithm has learned from its past mistakes.

With this context in mind, [Theorem 10.13](#) is a regret guarantee for the [OMD](#) algorithm for the *linear bandit* problem in which the loss functions are linear,  $\ell_n(\cdot) = \langle p_n, \cdot \rangle$ , and the vectors belong to the dual norm ball  $\{\|\cdot\|_* \leq L\}$ . This result is already interesting in the Euclidean case  $\phi = \frac{1}{2} \|\cdot\|^2$ , but the simplex setting is of particular interest to its connection with a well-established algorithm.

**Example 10.14 (learning from expert advice).** On each day  $n$ , an investor seeks to predict the price of a stock. There are  $d$  so-called “experts” who give daily predictions. On the following day, the investor compares their predictions with reality and assigns them losses  $\ell_n[1], \dots, \ell_n[d] \in [-1, 1]$ . (For example, we could set  $\ell_n[i] = +1$  if expert  $i$  incorrectly predicted the direction of change of the stock price on day  $n$ , and  $\ell_n[i] = -1$  otherwise.) Not all of the experts are necessarily reliable, but some might be. Can we aggregate the expert forecasts and compete with the best of them in hindsight, i.e., incur small regret?

The algorithm maintains a vector  $x_n \in \Delta_d$  in the probability simplex. On each day  $n$ , the algorithm picks an expert  $i_n \sim x_n$  and trusts the advice of the  $i_n$ -th expert. Note that the expected loss incurred by the algorithm is  $\mathbb{E}_{i_n \sim x_n} \ell_n[i_n] = \langle \ell_n, x_n \rangle$ , where  $\ell_n \in \mathbb{R}^d$  is the vector of losses. (This is the online version of [Example 10.1](#).) The regret is

$$\sum_{n=0}^{T-1} \langle \ell_n, x_n \rangle - \inf_{x \in \Delta_d} \sum_{n=0}^{T-1} \langle \ell_n, x \rangle = \sum_{n=0}^{T-1} \langle \ell_n, x_n \rangle - \min_{i \in [d]} \sum_{n=0}^{T-1} \ell_n[i].$$

We update the vector  $x_n$  using [OMD](#) with  $p_n = \ell_n$  and the entropic mirror map  $\phi$ . Note that  $\|\ell_n\|_* = \|\ell_n\|_\infty \leq 1$ , and by [Example 10.12](#), we can take  $x_0 = \mathbf{1}_d/d$  for which  $R_\phi \leq \sqrt{\log d}$ . Therefore, [Theorem 10.13](#) implies

$$\text{Regret}_T(\text{OMD}) \leq \sqrt{2T \log d}.$$

The corresponding algorithm, with updates

$$x_{n+1} = \frac{x_n \odot \exp(-h\ell_n)}{\|x_n \odot \exp(-h\ell_n)\|_1},$$

is known as the *multiplicative weights* algorithm.

## Bibliographical notes

The definitions and usage of relative convexity and smoothness are from [\[BBT17; LFN18\]](#). An interesting discussion of the various ways to discretize [\(10.2\)](#) and [\(10.3\)](#) can be found in the paper [\[GWS21\]](#). The example in [Exercise 10.6](#) is taken from [\[BBT17\]](#).

This section provides an introduction to online learning, although it should be noted that many of the interesting questions revolve around the more challenging setting of *bandit feedback*, i.e., after each round  $n$ , the player only receives the value  $\ell_n(x_n)$  of the incurred loss rather than the full loss function  $\ell_n(\cdot)$ . Tackling this setting requires

significant new ideas; see, e.g., [BC12] for an exposition.

## Exercises

### Exercise 10.1.

1. Prove that for all  $x, x' \in \mathcal{C}_\phi$ ,  $D_\phi(x, x') = D_{\phi^*}(\nabla\phi(x'), \nabla\phi(x))$ .
2. Let  $\mathcal{C} \subseteq \mathcal{C}_\phi$  be a closed convex set and let  $\Pi_{\mathcal{C}}^\phi : \mathcal{C}_\phi \rightarrow \mathcal{C}$  denote the Bregman projection operator:

$$\Pi_{\mathcal{C}}^\phi(x) := \arg \min_{\mathcal{C} \cap \mathcal{C}_\phi} D_\phi(\cdot, x).$$

Show that  $\langle \nabla\phi(\Pi_{\mathcal{C}}^\phi(x)) - \nabla\phi(x), \Pi_{\mathcal{C}}^\phi(x) - z \rangle \leq 0$  for all  $z \in \mathcal{C}$ . Use this to justify the Pythagorean inequality

$$D_\phi(z, x) \geq D_\phi(z, \Pi_{\mathcal{C}}^\phi(x)) + D_\phi(\Pi_{\mathcal{C}}^\phi(x), x).$$

3. Let  $X$  be a random variable with  $\mathbb{E}|\phi(X)| < \infty$ . For any  $v \in \mathcal{C}_\phi$ , establish the identity  $\mathbb{E} D_\phi(X, v) - \mathbb{E} D_\phi(X, \mathbb{E} X) = D_\phi(\mathbb{E} X, v)$ . From this, deduce that the Bregman barycenter coincides with the usual mean:  $\arg \min_{v \in \mathcal{C}_\phi} \mathbb{E} D_\phi(X, v) = \mathbb{E} X$ .

**Exercise 10.2.** Show that if  $\phi$  is  $\alpha$ -convex relative to a norm  $\|\cdot\|$ , then  $\phi^*$  is  $\alpha^{-1}$ -smooth relative to the dual norm  $\|\cdot\|_*$ .

**Exercise 10.3.** Prove that the entropic mirror map is 1-convex relative to  $\|\cdot\|_1$  over the probability simplex  $\Delta_d$ .

**Exercise 10.4.** For the entropic mirror map  $\phi$ , prove that the Bregman projection  $\Pi_{\Delta_d}^\phi$  onto the probability simplex simply normalizes the vector:  $x \mapsto x/\|x\|_1$ .

### Exercise 10.5.

1. More generally, show that **MPGD** can be rewritten as the update

$$x_{n+1} = \text{prox}_{h\|\cdot\|_1}^\phi(\nabla\phi^*(\nabla\phi(x_n - h \nabla f(x_n)))) .$$

2. Consider the mirror map  $\phi : x \mapsto -\sum_{i=1}^d \log x[i]$ , defined over  $\mathbb{R}_+^d$ . Compute  $\text{prox}_{h\|\cdot\|_1}^\phi$ , the Bregman proximal operator for  $\|\cdot\|_1$ .

**Exercise 10.6.** Consider the problem of recovering an image  $x \in \mathbb{R}_{++}^d$  from a noisy observation  $y \approx Ax$ , where  $y \in \mathbb{R}_+^n$  and  $A \in \mathbb{R}_{++}^{n \times d}$  is a matrix with positive entries. To solve this problem, we can set up the problem of minimizing

$$x \mapsto D_{\phi_{\text{ent}}}(y, Ax) + \lambda \|x\|_1,$$

where  $\phi_{\text{ent}}$  is the entropic mirror map. We apply [MPGD](#), using the negative logarithm as a mirror map, i.e.,  $\phi : x \mapsto -\sum_{i=1}^d \log x[i]$ . Show that the first term in the objective is relatively convex and smooth, with smoothness constant bounded by  $\|y\|_1$ . Deduce that we can obtain an  $\varepsilon$ -approximate solution in  $O(\|y\|_1 D_{\phi}(x_{\star}, x_0)/\varepsilon)$  iterations. Also, write out the algorithm iterates explicitly.

## 11 [3/25–4/3] Alternating minimization

In this section, we study the method of alternating minimization. The goal is to minimize a function  $f$  by decomposing the optimization variable  $x$  into  $D$  variables  $x^1, \dots, x^D$ . In this decomposition, the individual variables do not have to be one-dimensional, so we let  $x^i \in \mathbb{R}^{d_i}$ . The method is defined as follows:

$$x_{n+1}^i := \arg \min_{x^i \in \mathbb{R}^{d_i}} f(x_{n+1}^1, \dots, x_{n+1}^{i-1}, x^i, x_n^{i+1}, \dots, x_n^D).$$

In other words, we iterate through the variables cyclically and minimize  $f$  over the  $i$ -th variable  $x^i$ , holding the other variables fixed. The decomposition is chosen so that it is cheap to compute the minimizer over each individual variable.

**Example 11.1 (low-rank matrix recovery).** Suppose that we want to recover an unknown matrix  $X_\star \in \mathbb{R}^{p_1 \times p_2}$  which is observed through noisy observations  $y_i \approx \langle A_i, X_\star \rangle$ ; here, the matrices  $A_i \in \mathbb{R}^{p_1 \times p_2}$  are part of the design and are known. If we further posit that  $X_\star$  is low-rank, say of rank at most  $r$ , then we aim to solve

$$\underset{X \in \mathbb{R}^{p_1 \times p_2}}{\text{minimize}} \quad \sum_{i=1}^n (y_i - \langle A_i, X \rangle)^2 \quad \text{such that} \quad \text{rank } X \leq r.$$

The rank constraint is difficult to deal with, so we instead factorize the matrix as  $X = UV^\top$  where  $U \in \mathbb{R}^{p_1 \times r}$  and  $V \in \mathbb{R}^{p_2 \times r}$ . This factorization is known as the *Burer–Monteiro factorization*, after [BM03; BM05]. The problem becomes

$$\underset{U \in \mathbb{R}^{p_1 \times r}, V \in \mathbb{R}^{p_2 \times r}}{\text{minimize}} \quad \sum_{i=1}^n (y_i - \langle A_i, UV^\top \rangle)^2.$$

This is a non-convex problem, but at least it is now amenable to gradient-based methods. Alternatively, we can apply alternating minimization. In words, we minimize over  $U$  while holding  $V$  fixed, and then minimize over  $V$  while holding  $U$  fixed, and so on. Each iteration corresponds to solving an unconstrained least-squares problem and admits a closed-form solution.

Although we present Example 11.1 as motivation for the design of alternating minimization methods in practice, as usual in these notes, we focus on guarantees in the convex case. Nevertheless, the analysis of the convex case still applies to relevant problems; see the bibliographical notes for examples.

## 11.1 Alternating projections

We can use alternating minimization to find a point in the intersection of two closed convex sets  $\mathcal{C}_1$  and  $\mathcal{C}_2$ . In this case, we take

$$f(x, y) = \chi_{\mathcal{C}_1}(x) + \chi_{\mathcal{C}_2}(y) + \|y - x\|^2.$$

If there exists  $x_\star \in \mathcal{C}_1 \cap \mathcal{C}_2$ , then  $(x_\star, x_\star)$  is a minimizer for  $f$ , and the alternating minimization algorithm reads

$$\begin{aligned} x_{n+1} &:= \arg \min_{x \in \mathbb{R}^d} f(x, y_n) = \Pi_{\mathcal{C}_1}(y_n), \\ y_{n+1} &:= \arg \min_{y \in \mathbb{R}^d} f(x_{n+1}, y) = \Pi_{\mathcal{C}_2}(x_{n+1}). \end{aligned}$$



Thus, we alternate projecting onto  $\mathcal{C}_1$  and onto  $\mathcal{C}_2$ . This method is quite useful when projections onto  $\mathcal{C}_1, \mathcal{C}_2$  individually are cheap, but the projection onto  $\mathcal{C}_1 \cap \mathcal{C}_2$  is expensive. The method easily generalizes to the intersection of more than two convex sets.

As a “warm up”, we first study the method of alternating projections. Actually, we consider a generalization to alternating Bregman projections, which is needed for §11.3:

$$x_{n+1} := \Pi_{\mathcal{C}_1}^\phi(y_n), \quad y_{n+1} := \Pi_{\mathcal{C}_2}^\phi(x_{n+1}), \quad (\text{ABP})$$

where  $\Pi_{\mathcal{C}}^\phi$  is the Bregman projection from [Exercise 10.1](#). We assume that  $\mathcal{C}_1 \cap \mathcal{C}_2 \neq \emptyset$  and that  $\mathcal{C}_1, \mathcal{C}_2 \subseteq \mathcal{C}_\phi$ .

**Lemma 11.2.** For any  $x_\star \in \mathcal{C}_1 \cap \mathcal{C}_2$ , the iterates of [ABP](#) satisfy

$$\sum_{n=0}^{\infty} \{D_\phi(x_n, y_{n-1}) + D_\phi(y_n, x_n)\} \leq D_\phi(x_\star, y_0).$$

Also, monotonicity holds:

$$D_\phi(x_\star, y_0) \geq D_\phi(x_\star, x_1) \geq D_\phi(x_\star, y_1) \geq \dots$$

*Proof.* By the Pythagorean inequality ([Exercise 10.1](#)),

$$D_\phi(x_\star, y_n) \geq D_\phi(x_\star, \Pi_{\mathcal{C}_1}^\phi(y_n)) + D_\phi(\Pi_{\mathcal{C}_1}^\phi(y_n), y_n) = D_\phi(x_\star, x_{n+1}) + D_\phi(x_{n+1}, y_n).$$

By adding this to a similar inequality for the other projection and summing,

$$D_\phi(x_\star, y_0) \geq \sum_{n=1}^{\infty} \{D_\phi(x_n, y_{n-1}) + D_\phi(y_n, x_n)\}. \quad \square$$

We can use the preceding lemma to prove a convergence result for [ABP](#). The following corollary relies on two additional technical assumptions for  $\phi$  which must be checked, but note that they hold for the Euclidean case  $\phi = \frac{\|\cdot\|^2}{2}$ .

**Corollary 11.3 (convergence of ABP).** Assume that the following conditions hold:

1. For any  $x \in \mathcal{C}_\phi$ , the sublevel sets of  $D_\phi(x, \cdot)$  are compact.
2. If  $\{z_n\}_{n \in \mathbb{N}}, \{z'_n\}_{n \in \mathbb{N}} \subseteq \mathcal{C}_\phi$  are such that  $D_\phi(z_n, z'_n) \rightarrow 0$ , then  $z_n - z'_n \rightarrow 0$ .

Then, the iterates of [ABP](#) satisfy  $x_n \rightarrow x_\star$  and  $y_n \rightarrow x_\star$  for some  $x_\star \in \mathcal{C}_1 \cap \mathcal{C}_2$ .

*Proof.* The first assumption ensures that there is a convergent subsequence  $\{x_{n_k}\}_{k \in \mathbb{N}}$  that converges to some  $x_\star \in \mathcal{C}_\phi$ . Since  $x_n \in \mathcal{C}_1$  for all  $n$  and  $\mathcal{C}_1$  is closed, then  $x_\star \in \mathcal{C}_1$ . Moreover, by [Lemma 11.2](#),  $D_\phi(\Pi_{\mathcal{C}_2}^\phi(x_n), x_n) = D_\phi(y_n, x_n) \rightarrow 0$ , so the second property shows that  $\Pi_{\mathcal{C}_2}^\phi(x_n) - x_n \rightarrow 0$ . Since  $\mathcal{C}_2$  is closed,  $x_\star \in \mathcal{C}_2$  as well.

To upgrade the subsequential convergence to full convergence, we observe that  $D_\phi(x_\star, x_{n_k}) \rightarrow 0$ , whence the monotonicity statement in [Lemma 11.2](#) implies  $D_\phi(x_\star, x_n) \rightarrow 0$  and  $D_\phi(x_\star, y_n) \rightarrow 0$ . By the second assumption,  $x_n \rightarrow x_\star$  and  $y_n \rightarrow x_\star$ .  $\square$

Furthermore, [Lemma 11.2](#) implies

$$\min_{n=1, \dots, N} \{D_\phi(x_n, y_{n-1}) + D_\phi(y_n, x_n)\} \leq \frac{D_\phi(x_\star, y_0)}{N}. \quad (11.1)$$

This does not, however, imply a rate of convergence for  $x_n$  to  $x_\star$ . For example, if  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are two lines that meet each other at a very small angle, then the successive projections can lie very close to each other even though they are both very far from the common point of intersection.

## 11.2 Convergence analysis for alternating minimization

We now return to the alternating minimization method. We use the shorthand  $x^S$  to denote the components in  $S$ ,  $x^S := \{x^i\}_{i \in S}$ , where we abbreviate consecutive indices  $\{i, \dots, j\}$  as  $i:j$ . We perform an analysis in the smooth case. However, similarly to how gradient-based methods do not suffer from non-smoothness provided that one has access to a proximal oracle, it turns out that coordinate-based methods do not suffer from non-smoothness provided that the non-smooth part respects the coordinate decomposition. Hence, we consider the slightly more general problem of minimizing

$$F : \mathbb{R}^{d_1 \times \dots \times d_D} \rightarrow \mathbb{R}, \quad F(x^{1:D}) := f(x^{1:D}) + \sum_{i=1}^D g_i(x^i),$$

where  $f$  is convex and smooth, and each  $g_i$  is convex. For shorthand, write  $g := \bigoplus_{i=1}^D g_i$ , that is,  $g(x^{1:D}) := \sum_{i=1}^D g_i(x^i)$ . The algorithm reads

$$x_{n+1}^i \in \arg \min_{x^i \in \mathbb{R}^{d_i}} \{f(x_{n+1}^{1:i-1}, x^i, x_n^{i+1:D}) + g_i(x^i)\}. \quad (\text{AM})$$

**Theorem 11.4 (convergence of AM).** Assume that  $f$  is convex and  $\beta$ -smooth, and that each  $g_i$  is convex. Then, AM achieves  $F(x_N^{1:D}) - F_\star \leq \varepsilon$  if

$$N \geq \left( \log_{1/2} \frac{F(x_0^{1:D}) - F_\star}{4\beta D^2 R^2} \right)_+ + \frac{8\beta D^2 R^2}{\varepsilon},$$

where  $R := \sup_{n \in \mathbb{N}} \|x_n^{1:D} - x_\star^{1:D}\|$ .

*Proof.* By (3.4),

$$\begin{aligned} f(x_n^{1:D}) &\geq f(x_{n+1}^1, x_n^{2:D}) + \langle \nabla_1 f(x_{n+1}^1, x_n^{2:D}), x_n^1 - x_{n+1}^1 \rangle \\ &\quad + \frac{1}{2\beta} \|\nabla f(x_{n+1}^1, x_n^{2:D}) - \nabla f(x_n^{1:D})\|^2. \end{aligned}$$

On the other hand, since  $\nabla_1 f(x_{n+1}^1, x_n^{2:D}) \in -\partial g_1(x_{n+1}^1)$ ,

$$g_1(x_n^1) + \langle \nabla_1 f(x_{n+1}^1, x_n^{2:D}), x_n^1 - x_{n+1}^1 \rangle \geq g_1(x_{n+1}^1).$$

By summing these two inequalities together with the corresponding ones for the other coordinates, we obtain a “descent lemma”

$$F(x_n^{1:D}) - F(x_{n+1}^{1:D}) \geq \frac{1}{2\beta} \sum_{i=1}^D \|\nabla f(x_{n+1}^{1:i}, x_n^{i+1:D}) - \nabla f(x_{n+1}^{1:i-1}, x_n^{i:D})\|^2.$$

Next,

$$\begin{aligned} F(x_{n+1}^{1:D}) - F_\star &\leq \langle \nabla f(x_{n+1}^{1:D}), x_{n+1}^{1:D} - x_\star^{1:D} \rangle + g(x_{n+1}^{1:D}) - g(x_\star^{1:D}) \\ &= \sum_{i=1}^D \{ \langle \nabla_i f(x_{n+1}^{1:D}), x_{n+1}^i - x_\star^i \rangle + g_i(x_{n+1}^i) - g_i(x_\star^i) \} \\ &\leq \sum_{i=1}^D \langle \nabla_i f(x_{n+1}^{1:D}) - \nabla_i f(x_{n+1}^{1:i}, x_n^{i+1:D}), x_{n+1}^i - x_\star^i \rangle \\ &\leq \sum_{i=1}^D \|\nabla f(x_{n+1}^{1:D}) - \nabla f(x_{n+1}^{1:i}, x_n^{i+1:D})\| \|x_{n+1}^i - x_\star^i\| \\ &\leq \sum_{i=1}^D \left( \sum_{j=i}^{D-1} \|\nabla f(x_{n+1}^{1:j+1}, x_n^{j+2:D}) - \nabla f(x_{n+1}^{1:j}, x_n^{j+1:D})\| \right) \|x_{n+1}^i - x_\star^i\| \end{aligned}$$

$$\begin{aligned}
&\leq \left( \sum_{i=0}^{D-1} \|\nabla f(x_{n+1}^{1:i+1}, x_n^{i+2:D}) - \nabla f(x_{n+1}^{1:i}, x_n^{n+1:D})\| \right) \left( \sum_{i=1}^D \|x_{n+1}^i - x_\star^i\| \right) \\
&\leq D \sqrt{\left( \sum_{i=0}^{D-1} \|\nabla f(x_{n+1}^{1:i+1}, x_n^{i+2:D}) - \nabla f(x_{n+1}^{1:i}, x_n^{n+1:D})\|^2 \right) \left( \sum_{i=1}^D \|x_{n+1}^i - x_\star^i\|^2 \right)} \\
&\leq DR \sqrt{\sum_{i=0}^{D-1} \|\nabla f(x_{n+1}^{1:i+1}, x_n^{i+2:D}) - \nabla f(x_{n+1}^{1:i}, x_n^{n+1:D})\|^2}.
\end{aligned}$$

Combined with the previous inequality, it yields, for  $\Delta_n := F(x_n^{1:D}) - F_\star$ ,

$$\Delta_{n+1} - \Delta_n \leq -\frac{1}{2\beta} \sum_{i=1}^D \|\nabla f(x_{n+1}^{1:i}, x_n^{i+1:D}) - \nabla f(x_{n+1}^{1:i-1}, x_n^{i:D})\|^2 \leq -\frac{1}{2\beta D^2 R^2} \Delta_{n+1}^2.$$

If  $\Delta_{n+1} \geq \Delta_n/2$ , this yields  $\Delta_{n+1} - \Delta_n \leq -\Delta_n^2/(8\beta D^2 R^2)$ , so in general

$$\Delta_{n+1} \leq \max\left\{\frac{1}{2}, \left(1 - \frac{\Delta_n}{8\beta D^2 R^2}\right)\right\} \Delta_n.$$

This implies that the error decays exponentially fast until iteration  $n_0$  which satisfies  $\Delta_{n_0} \leq 4\beta D^2 R^2$ . Thereafter,

$$\frac{1}{\Delta_n} - \frac{1}{\Delta_{n+1}} = \frac{\Delta_{n+1} - \Delta_n}{\Delta_n \Delta_{n+1}} \leq -\frac{1}{8\beta D^2 R^2},$$

which yields

$$\Delta_N \leq \frac{8\beta D^2 R^2 \Delta_{n_0}}{8\beta D^2 R^2 + (N - n_0) \Delta_{n_0}} \leq \frac{8\beta D^2 R^2}{N - n_0}.$$

□

Although [Theorem 11.4](#) provides a convergence guarantee for [AM](#), it incurs a poor dependence on the number of blocks  $D$ —the complexity scales as  $D^3$ . It turns out that this can be alleviated by randomly choosing a block at each iteration. More precisely, define the following randomized alternating minimization algorithm:

$$x_{n+1} := \arg \min_{x^{i(n)} \in \mathbb{R}^{d_{i(n)}}} \{f(x_n^{1:i(n)-1}, x^{i(n)}, x_n^{i(n)+1:D}) + g_{i(n)}(x^{i(n)})\}, \quad i(n) \sim \text{uniform}([D]).$$

(RAM)

The analysis below also handles anisotropic smoothness: we assume that for each  $i$ ,

$$f(x^{1:i-1}, \cdot, x^{i+1:D}) \text{ is } \beta_i\text{-smooth for each } x^{1:D} \in \mathbb{R}^{d_1 \times \dots \times d_D}.$$

We refer to this condition succinctly by saying that  $f$  is  $\beta$ -smooth, where  $\beta = (\beta_1, \dots, \beta_D)$ . It induces the norm

$$\|x^{1:D}\|_{\beta} := \left( \sum_{i=1}^D \beta_i \|x^i\|^2 \right)^{1/2}.$$

**Theorem 11.5 (convergence of RAM).** Assume that  $f$  is  $\beta$ -smooth and  $\alpha_f$ -convex w.r.t.  $\|\cdot\|_{\beta}$ . Also, assume that  $g$  is  $\alpha_g$ -convex w.r.t.  $\|\cdot\|_{\beta}$ . Then, the iterates of RAM satisfy the following bounds. If  $\alpha_f + \alpha_g > 0$ , then

$$\mathbb{E} F(x_N^{1:D}) - F_{\star} \leq \left( 1 - \frac{\alpha_f + \alpha_g}{(1 + \alpha_g) D} \right)^N (F(x_0^{1:D}) - F_{\star}).$$

Otherwise, if  $\alpha_f + \alpha_g = 0$ , then

$$\mathbb{E} F(x_N^{1:D}) - F_{\star} \leq \frac{2DR_{\beta}^2}{N},$$

where  $R_{\beta} \geq \sup_{n \in \mathbb{N}} \max\{\sqrt{F(x_n^{1:D}) - F_{\star}}, \|x_n^{1:D} - x_{\star}^{1:D}\|_{\beta}\}$  almost surely.

Before proving the theorem, let us compare the computational costs implied by these various results in the weakly convex case. Suppose, for the sake of argument, that computing a full gradient  $\nabla f$  is  $O(D)$ . If the proximal oracle for  $g$  is available, we can run PGD to obtain an  $\varepsilon$ -solution at cost  $O(\beta D R^2 / \varepsilon)$ . On the other hand, each iteration of AM requires minimization with respect to one of the variables, leading to a total cost of roughly  $O(\beta D^3 R^2 / \varepsilon)$ , assuming that minimization over one variable is comparable in cost to computing a partial gradient.

For RAM, let  $\beta_{\max} := \max_{i \in [D]} \beta_i$ . The overall smoothness of  $f$  satisfies  $\beta_{\max} \leq \beta \leq \sum_{i=1}^D \beta_i \leq D\beta_{\max}$  and, ignoring the first term in the definition of  $R_{\beta}$ ,  $R_{\beta}^2 \leq \beta_{\max} R^2$ .<sup>11</sup> The cost for RAM is therefore  $O(\beta_{\max} D R^2 / \varepsilon)$ , which is “never worse” than the rate for PGD, but could be substantially better when  $\beta$  is closer to the upper bound  $D\beta_{\max}$ . This is the case when the directions of high smoothness are not aligned with the coordinate directions (e.g., imagine that the Hessian matrix looks like the all-ones matrix). Thus, RAM can “adapt” to directional smoothness, which is particularly appealing since RAM has no tuning parameters (not even a step size!).

<sup>11</sup>The result for RAM requires an almost sure bound on the iterates, but let us ignore this detail for the sake of this discussion.

*Proof of Theorem 11.5.* Let  $y^{1:D} \in \mathbb{R}^{d_1 \times \dots \times d_D}$  and let  $\mathbb{E}_n$  denote the expectation over  $i(n)$  only. Then,

$$\begin{aligned}
\mathbb{E}_n F(x_{n+1}^{1:D}) &\leq \mathbb{E}_n F(x_n^{1:i(n)-1}, y^{i(n)}, x_n^{i(n)+1:D}) \\
&\leq \mathbb{E}_n \left[ f(x_n^{1:D}) + \langle \nabla_{i(n)} f(x_n^{1:D}), y^{i(n)} - x_n^{i(n)} \rangle + \frac{\beta_{i(n)}}{2} \|y^{i(n)} - x_n^{i(n)}\|^2 \right] \\
&\quad + \mathbb{E}_n \left[ \sum_{i \neq i(n)} g_i(x_n^i) + g_{i(n)}(y^{i(n)}) \right] \\
&= f(x_n^{1:D}) + \frac{1}{D} \langle \nabla f(x_n^{1:D}), y^{1:D} - x_n^{1:D} \rangle + \frac{1}{2D} \|y^{1:D} - x_n^{1:D}\|_\beta^2 \\
&\quad + \frac{D-1}{D} g(x_n^{1:D}) + \frac{1}{D} g(y^{1:D}) \\
&\leq \frac{D-1}{D} f(x_n^{1:D}) + \frac{1}{D} f(y^{1:D}) + \frac{1-\alpha_f}{2D} \|y^{1:D} - x_n^{1:D}\|_\beta^2 \\
&\quad + \frac{D-1}{D} g(x_n^{1:D}) + \frac{1}{D} g(y^{1:D}).
\end{aligned}$$

By taking the infimum over  $y^{1:D}$ , we have shown that

$$\mathbb{E}_n F(x_{n+1}^{1:D}) \leq \frac{D-1}{D} F(x_n^{1:D}) + \frac{1}{D} Q_{1/(1-\alpha_f)} F(x_n^{1:D}),$$

where  $(Q_t)_{t \geq 0}$  denotes the Hopf–Lax semigroup (Definition 9.3) with respect to  $\|\cdot\|_\beta$ . By Exercise 9.3, we have

$$\frac{Q_{1/(1-\alpha_f)} F(x_n^{1:D}) - F_\star}{F(x_n^{1:D}) - F_\star} \leq \begin{cases} (1-\alpha_f)/(1+\alpha_g), & \text{if } \alpha_f + \alpha_g > 0, \\ 1 - (F(x_n^{1:D}) - F_\star)/(2R_\beta^2), & \text{if } \alpha_f + \alpha_g = 0. \end{cases}$$

By taking the expectation again, in the first case it yields

$$\mathbb{E} F(x_{n+1}^{1:D}) - F_\star \leq \left(1 - \frac{\alpha_f + \alpha_g}{(1+\alpha_g)D}\right) \{\mathbb{E} F(x_n^{1:D}) - F_\star\},$$

and in the second case, by Jensen's inequality,

$$\begin{aligned}
\mathbb{E} F(x_{n+1}^{1:D}) - F_\star &\leq \mathbb{E} \left[ \left(1 - \frac{F(x_n^{1:D}) - F_\star}{2DR_\beta^2}\right) \{F(x_n^{1:D}) - F_\star\} \right] \\
&\leq \left(1 - \frac{\mathbb{E} F(x_n^{1:D}) - F_\star}{2DR_\beta^2}\right) \{\mathbb{E} F(x_n^{1:D}) - F_\star\}.
\end{aligned}$$

The results follow by iterating these inequalities.  $\square$

### 11.3 Case study: entropic optimal transport

As a case study, we apply these ideas to a concrete problem of modern interest. Namely, over the past decade, there has been considerable interest in applications of optimal transport to machine learning, among other domains. In this problem, we are given two probability distributions  $\mu, \nu$  over spaces  $\mathcal{X}$  and  $\mathcal{Y}$  respectively, as well as a cost function  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ . In this section, we focus on the case where  $\mathcal{X}, \mathcal{Y}$  are finite sets, although the ideas presented here readily generalize. The optimal transport cost between  $\mu$  and  $\nu$  with cost  $c$  is

$$\text{OT}(\mu, \nu) := \inf \{ \mathbb{E} c(X, Y) : X \sim \mu, Y \sim \nu \},$$

where the infimum is taken over all *couplings*  $(X, Y)$ , i.e., jointly defined random variables with marginal laws  $\mu$  and  $\nu$  respectively. A particularly common choice is the Euclidean cost,  $c(x, y) = \|y - x\|^2$ , but other choices are common too. Since the structure of the cost does not play any role here, we leave it general.

Focus on the case where  $\mathcal{X}, \mathcal{Y}$  are finite sets. If we write  $\gamma$  for the joint distribution of  $(X, Y)$ , the optimal transport problem can be written

$$\underset{\gamma \in \mathbb{R}_+^{\mathcal{X} \times \mathcal{Y}}}{\text{minimize}} \quad \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} c_{x,y} \gamma_{x,y} \quad \text{such that} \quad \begin{cases} \sum_{y \in \mathcal{Y}} \gamma_{x,y} = \mu_x \text{ for all } x \in \mathcal{X}, \\ \sum_{x \in \mathcal{X}} \gamma_{x,y} = \nu_y \text{ for all } y \in \mathcal{Y}. \end{cases}$$

More compactly, if we write  $C = (c_{x,y})_{x \in \mathcal{X}, y \in \mathcal{Y}}$  for the cost matrix and  $\mathbf{1}_{\mathcal{X}}, \mathbf{1}_{\mathcal{Y}}$  for the all-ones vectors on  $\mathcal{X}$  and on  $\mathcal{Y}$  respectively, this can be written

$$\underset{\gamma \in \mathbb{R}_+^{\mathcal{X} \times \mathcal{Y}}}{\text{minimize}} \quad \langle C, \gamma \rangle \quad \text{such that} \quad \begin{cases} \gamma \mathbf{1}_{\mathcal{Y}} = \mu, \\ \gamma^T \mathbf{1}_{\mathcal{X}} = \nu. \end{cases}$$

This is readily recognized as a linear program, but solving it as such is expensive. There are specialized combinatorial algorithms—see [PC19]—but the computational cost scales at least as  $d^3$  if  $d = |\mathcal{X}| = |\mathcal{Y}|$  (for simplicity, the input dimensions are the same).

On the other hand, the size of the input matrix  $C$  is  $d^2$ , and optimistically we ask if we can solve the problem in  $\tilde{O}(d^2)$  time—that is, nearly linear time in the size of the input. We shall see that this is the case, provided that we add some entropic regularization to the problem, as popularized in machine learning by Cuturi [Cut13].

Given a regularization parameter  $\varepsilon_{\text{reg}} > 0$ , the goal is to instead solve

$$\underset{\gamma \in \mathbb{R}_+^{\mathcal{X} \times \mathcal{Y}}}{\text{minimize}} \quad \langle C, \gamma \rangle + \varepsilon_{\text{reg}} \text{KL}(\gamma \parallel \mu \otimes \nu) \quad \text{such that} \quad \begin{cases} \gamma \mathbf{1}_{\mathcal{Y}} = \mu, \\ \gamma^T \mathbf{1}_{\mathcal{X}} = \nu. \end{cases} \quad (11.2)$$

Call the value of this problem  $\text{OT}_{\varepsilon_{\text{reg}}}(\mu, \nu)$ . Why does this make the problem so much easier to solve? The answer is that by Kantorovich duality, the dual to the unregularized problem turns out to be

$$\underset{f \in \mathbb{R}^{\mathcal{X}}, g \in \mathbb{R}^{\mathcal{Y}}}{\text{maximize}} \quad \sum_{x \in \mathcal{X}} f_x \mu_x + \sum_{y \in \mathcal{Y}} g_y \nu_y \quad \text{such that} \quad f_x + g_y \leq c_{x,y} \text{ for all } x \in \mathcal{X}, y \in \mathcal{Y},$$

which is still a constrained problem. On the other hand, the dual to the regularized problem is unconstrained.

**Theorem 11.6 (EOT duality).** Consider the following dual problem:

$$\underset{f \in \mathbb{R}^{\mathcal{X}}, g \in \mathbb{R}^{\mathcal{Y}}}{\text{maximize}} \quad \sum_{x \in \mathcal{X}} f_x \mu_x + \sum_{y \in \mathcal{Y}} g_y \nu_y - \varepsilon_{\text{reg}} \left( \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \exp\left(\frac{f_x + g_y - c_{x,y}}{\varepsilon_{\text{reg}}}\right) \mu_x \nu_y - 1 \right) \quad (11.3)$$

Let  $f^*, g^*$  solve the dual problem. Then,  $\gamma^*$  defined by

$$\gamma_{x,y}^* := \exp\left(\frac{f_x^* + g_y^* - c_{x,y}}{\varepsilon_{\text{reg}}}\right) \mu_x \nu_y \quad (11.4)$$

is the unique solution to the entropic optimal transport problem.

Moreover,  $\gamma^*$  is characterized as the unique distribution of the form (11.4) for some  $f^*, g^*$  with the correct marginals.

For a proof, see, e.g., [CNR25, Proposition 4.3], although in this discrete setting it can be proven via Lagrange multipliers (Exercise 11.2).

By replacing  $c_{x,y}$  with  $c_{x,y}/\varepsilon_{\text{reg}}$  and rescaling  $f_x, g_y$  accordingly, we may set  $\varepsilon_{\text{reg}} = 1$  without loss of generality, so we adopt this convention henceforth.

Let  $\mathcal{D}(f, g)$  denote the dual functional, i.e., the objective of (11.3). Since the dual is unconstrained, we propose to solve it by alternating maximization. Namely, given iterates  $f^n, g^n$ , we set

$$f^{n+1} := \arg \max_{f \in \mathbb{R}^{\mathcal{X}}} \mathcal{D}(f, g^n), \quad g^{n+1} := \arg \max_{g \in \mathbb{R}^{\mathcal{Y}}} \mathcal{D}(f^{n+1}, g).$$

By solving for the first-order conditions, the updates can be written explicitly:

$$f_x^{n+1} = -\log \sum_{y \in \mathcal{Y}} \exp(g_y^n - c_{x,y}) \nu_y, \quad g_y^{n+1} = -\log \sum_{x \in \mathcal{X}} \exp(f_x^{n+1} - c_{x,y}) \mu_x. \quad (11.5)$$



At this point, one can try to apply [Theorem 11.4](#), but the correct geometry for this problem is not Euclidean.

Instead, consider what happens to the matrices

$$\gamma_{x,y}^n := \exp(f_x^n + g_y^n - c_{x,y}) \mu_x \nu_y, \quad \gamma_{x,y}^{n+1/2} := \exp(f_x^{n+1} + g_y^n - c_{x,y}) \mu_x \nu_y.$$

Performing the update  $f^n \mapsto f^{n+1}$  implicitly performs the update  $\gamma^n \mapsto \gamma^{n+1/2}$ . The  $\mathcal{X}$ -marginal of  $\gamma^{n+1/2}$  is computed as follows:

$$\sum_{y \in \mathcal{Y}} \gamma_{x,y}^{n+1/2} = \mu_x \sum_{x \in \mathcal{X}} \exp(f_x^{n+1} + g_y^n - c_{x,y}) \nu_y = \mu_x,$$

by (11.5). Thus,  $\gamma^{n+1/2}$  has the correct  $\mathcal{X}$ -marginal  $\mu$ , although its  $\mathcal{Y}$ -marginal may not be correct. In fact, one can see that  $\gamma^{n+1/2}$  is obtained from  $\gamma^n$  by normalizing its rows to fix its  $\mathcal{X}$ -marginal. Similarly, the update  $g^n \mapsto g^{n+1}$  implicitly performs the update  $\gamma^{n+1/2} \mapsto \gamma^{n+1}$ , and  $\gamma^{n+1}$  has the correct  $\mathcal{Y}$ -marginal  $\nu$ . In this form, this is known as *Sinkhorn's matrix scaling algorithm* [[Sin64](#)]. In words, Sinkhorn's algorithm iteratively “fixes the rows, and then fixes the columns, and then fixes the rows...”.

We can therefore identify Sinkhorn's algorithm as an instance of alternating Bregman projections. Indeed, we define the constraint sets

$$\mathcal{C}^\mu := \{\gamma \in \mathbb{R}_+^{\mathcal{X} \times \mathcal{Y}} : \gamma \mathbf{1}_{\mathcal{Y}} = \mu\}, \quad \mathcal{C}_\nu := \{\gamma \in \mathbb{R}_+^{\mathcal{X} \times \mathcal{Y}} : \gamma^\top \mathbf{1}_{\mathcal{X}} = \nu\},$$

and we let  $\phi : \gamma \mapsto \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} (\gamma_{x,y} \log \gamma_{x,y} - \gamma_{x,y})$  denote the entropic mirror map, then similarly to [Exercise 10.4](#) one can show that the Bregman projections onto  $\mathcal{C}^\mu$  and  $\mathcal{C}_\nu$  normalize the rows and columns respectively. This yields

$$\gamma^{n+1/2} = \Pi_{\mathcal{C}^\mu}^\phi(\gamma^n), \quad \gamma^{n+1} = \Pi_{\mathcal{C}_\nu}^\phi(\gamma^{n+1/2}).$$

From this perspective, Sinkhorn's algorithm aims to find a point in the intersection  $\mathcal{C}^\mu \cap \mathcal{C}_\nu$ . Does this mean that the intersection is a singleton, which is the solution to the entropic optimal transport problem? No! Note that the intersection  $\mathcal{C}^\mu \cap \mathcal{C}_\nu$  only encodes the *constraints* of the original problem, not the objective which depends on the cost function  $c$ . In fact, by [Theorem 11.6](#), different choices of  $c$  give rise to different entropic optimal couplings, so  $\mathcal{C}^\mu \cap \mathcal{C}_\nu$  is certainly not a singleton.

What is true, however, is that if  $\Gamma_c$  denotes the set of joint distributions of the form (11.4), then the unique element of  $\mathcal{C}^\mu \cap \mathcal{C}_\nu \cap \Gamma_c$  solves the entropic optimal transport problem. This is the last statement of [Theorem 11.6](#). Moreover, Sinkhorn's algorithm maintains the property that if we initialize at an element of  $\Gamma_c$ , e.g., by taking  $f^0 = g^0 = 0$ , then the

algorithm iterates all remain in  $\Gamma_c$ . So, remarkably, the alternating Bregman projections do indeed solve our problem.

Let us now see what [Lemma 11.2](#) implies for Sinkhorn's algorithm. In the following, we assume that we initialize at a probability distribution  $\gamma^0 \in \Gamma_c$ , e.g., we can take

$$\gamma^0 = \frac{\exp(-c) (\mu \otimes \nu)}{\|\exp(-c) (\mu \otimes \nu)\|_1}.$$

**Theorem 11.7 (convergence of Sinkhorn's algorithm).** Initialize Sinkhorn's algorithm at a probability distribution  $\gamma^0 \in \Gamma_c$ . Suppose that the number of iterations  $N$  satisfies

$$N \geq \frac{2 \text{KL}(\gamma^\star \parallel \gamma^0)}{\varepsilon^2}.$$

Then, there exists an iteration  $n \in \{0, 1, \dots, N-1\}$  and  $\hat{\gamma} \in \{\gamma^n, \gamma^{n+1/2}\}$  such that if  $\hat{\mu}, \hat{\nu}$  denote the marginals of  $\hat{\gamma}$ , then

$$\|\hat{\mu} - \mu\|_1 + \|\hat{\nu} - \nu\|_1 \leq \varepsilon.$$

*Proof.* By [\(11.1\)](#), we know that

$$\min_{n=0,1,\dots,N-1} \{\text{KL}(\gamma^{n+1/2} \parallel \gamma^n) + \text{KL}(\gamma^{n+1} \parallel \gamma^{n+1/2})\} \leq \frac{\text{KL}(\gamma^\star \parallel \gamma^0)}{N}.$$

Let  $\mu^n$  denote the  $\mathcal{X}$ -marginal of  $\gamma^n$ :

$$\mu_x^n = \sum_{y \in \mathcal{Y}} \gamma_{x,y}^n = \mu_x \exp(f_x^n) \sum_{y \in \mathcal{Y}} \exp(g_y^n - c_{x,y}) \nu_y = \mu_x \exp(f_x^n - f_x^{n+1}).$$

Therefore, since the  $\mathcal{X}$ -marginal of  $\gamma^{n+1/2}$  is  $\mu$ ,

$$\begin{aligned} \text{KL}(\gamma^{n+1/2} \parallel \gamma^n) &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \gamma_{x,y}^{n+1/2} \log \frac{\exp(f_x^{n+1} + g_y^n - c_{x,y})}{\exp(f_x^n + g_y^n - c_{x,y})} = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \gamma_{x,y}^{n+1/2} (f_x^{n+1} - f_x^n) \\ &= \sum_{x \in \mathcal{X}} \mu_x (f_x^{n+1} - f_x^n) = \sum_{x \in \mathcal{X}} \mu_x \log \frac{\mu_x}{\mu_x^n} = \text{KL}(\mu \parallel \mu^n) \geq \frac{1}{2} \|\mu^n - \mu\|_1^2, \end{aligned}$$

where the last inequality is Pinsker's inequality. The result follows.  $\square$

However, this is not the last word on Sinkhorn's algorithm. For instance, it does not provide convergence of the last iterate. It turns out that Sinkhorn's algorithm admits a third interpretation: as an instantiation of mirror descent ([Exercise 11.3](#)). Using this, one can prove the following theorem.

**Theorem 11.8** (convergence of Sinkhorn’s algorithm, II). Initialize Sinkhorn’s algorithm at a probability distribution  $\gamma^0 \in \Gamma_c$ . Then, if  $\mu^N$  denotes the  $\mathcal{X}$ -marginal of  $\gamma^N$ ,

$$\text{KL}(\mu^N \parallel \mu) \leq \frac{\text{KL}(\gamma^\star \parallel \gamma^0)}{N}.$$

## Bibliographical notes

The analysis of alternating minimization has recently inspired analyses of the coordinate ascent variational inference (CAVI) algorithm [AL24; LZ24], the expectation maximization (EM) algorithm [CJ24], and Gibbs sampling [ALZ24]. The proof of Theorem 11.5 is also taken from [LZ24].

For an introduction to optimal transport for statisticians, see [CNR25]. Other treatments of optimal transport, aimed at more mathematical audiences, include [Vil03; Vil09; San15]. The literature on (entropic) optimal transport is vast, so we only mention a few relevant references: the proof of Theorem 11.7 is similar in spirit to [ANR17]; Sinkhorn’s algorithm as interpreted as mirror descent in [Lég21]; and the interpretation in Exercise 11.3 is from [AKL22].

## Exercises

**Exercise 11.1.** Here, we present a simple convergence analysis of alternating minimization in the case where there are only two blocks,  $f$  satisfies (PL) with constant  $\alpha$  and is  $\beta$ -smooth, and  $g = 0$ .

For any  $h > 0$ , by definition of alternating minimization,

$$f(x_{n+1}^1, x_n^2) \leq f(x_n^1 - h \nabla_1 f(x_n^1, x_n^2), x_n^2).$$

Apply the descent lemma for GD (Lemma 3.1), the fact that  $\nabla_2 f(x_n^1, x_n^2) = 0$  (why?), and the PL inequality to deduce a one-step inequality  $f(x_{n+1}^1, x_n^2) - f_\star \leq (1 - \alpha/\beta) (f(x_n^1, x_n^2) - f_\star)$ . Deduce that

$$f(x_N^1, x_N^2) - f_\star \leq \left(1 - \frac{\alpha}{\beta}\right)^{2N} (f(x_0^1, x_0^2) - f_\star).$$

### Exercise 11.2.

1. Setting  $\varepsilon_{\text{reg}} = 1$  in (11.2), show that the objective is equivalent to minimizing  $\gamma \mapsto \text{KL}(\gamma \parallel \gamma^0)$ , where  $\gamma_{x,y}^0 = \mu_x \nu_y \exp(-c_{x,y})$ .

2. Introduce two Lagrange multipliers  $\kappa, \lambda \in \mathbb{R}$  and argue that (11.2) is equivalent to

$$\min_{\gamma \in \mathbb{R}_+^{\mathcal{X} \times \mathcal{Y}}} \max_{\kappa, \lambda \in \mathbb{R}} \left\{ \text{KL}(\gamma \parallel \gamma^0) + \kappa (\gamma \mathbf{1}_{\mathcal{Y}} - \mu) + \lambda (\gamma^T \mathbf{1}_{\mathcal{X}} - \nu) \right\}.$$

Without justification, assume that we can switch the min and max, so that the above problem is equivalent to

$$\max_{\kappa, \lambda \in \mathbb{R}} \min_{\gamma \in \mathbb{R}_+^{\mathcal{X} \times \mathcal{Y}}} \left\{ \text{KL}(\gamma \parallel \gamma^0) + \kappa (\gamma \mathbf{1}_{\mathcal{Y}} - \mu) + \lambda (\gamma^T \mathbf{1}_{\mathcal{X}} - \nu) \right\}.$$

Argue that the solution  $\gamma$  to this problem is of the form  $\gamma_{x,y} = \exp(f_x + g_y - c_{x,y}) \mu_x \nu_y$  for some  $f \in \mathbb{R}^{\mathcal{X}}, g \in \mathbb{R}^{\mathcal{Y}}$ .

**Exercise 11.3.** For a joint distribution  $\gamma$ , let  $(\Pi_{\mathcal{X}})_{\#} \gamma$  denote its  $\mathcal{X}$ -marginal. Consider the objective functional  $\mathcal{F} : \gamma \mapsto \text{KL}((\Pi_{\mathcal{X}})_{\#} \gamma \parallel \mu)$ . Show that the iteration  $\gamma^n \mapsto \gamma^{n+1}$  of Sinkhorn's algorithm can be viewed as one step of mirror descent on  $\mathcal{F}$  with the entropic mirror map  $\phi$ , constraint set  $\mathcal{C}_{\nu}$ , and step size 1. By checking relative convexity and smoothness, prove Theorem 11.8.

## 12 [4/8–4/17] Stochastic optimization

Our next topic is optimization with stochastic gradients. Besides its relevance in situations where the gradient cannot be computed exactly, stochastic optimization is particularly important for machine learning and statistics for at least two reasons. First, it can be viewed as a method for directly minimizing the population risk, and we can establish generalization bounds provided that we perform a single pass over our data. Second, it is routinely used to minimize empirical risk functions by approximating the full gradient by mini-batches over the data. Our treatment therefore centers around these applications.

### 12.1 Stochastic mirror proximal gradient descent

We start with the fundamental convergence result. Suppose that we wish to minimize  $F = f + g$ , where we only have access to stochastic gradients for  $f$ . More precisely, we assume that at each  $x \in \text{int dom } f$ , we can compute a random vector  $\hat{\nabla} f(x)$  which is unbiased:  $\mathbb{E} \hat{\nabla} f(x) = \nabla f(x)$ . Actually, we can also let  $f$  be non-smooth, in which case we require that  $\mathbb{E} \hat{\nabla} f(x) \in \partial f(x)$ . The analysis below can also handle the case where the stochastic gradient is biased, at the expense of an additional error term.

Let  $\phi : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$  be a mirror map. We consider the following iteration:

$$x_{n+1} := \arg \min_{x \in \mathbb{R}^d} \left\{ f(x_n) + \langle \hat{\nabla} f(x_n), x - x_n \rangle + g(x) + \frac{1}{h} D_{\phi}(x, x_n) \right\}. \quad (\text{SMPGD})$$

This is the *stochastic mirror proximal gradient descent* algorithm.<sup>12</sup> For most applications, we do not need all of these aspects (stochastic, mirror, proximal) simultaneously, but we may as well include them to emphasize that a unified proof is possible. Anyway, it is helpful to include a proximal term since it allows for projections, and the use of a Bregman divergence is a bonus since it covers stochastic mirror descent.

**Theorem 12.1 (convergence of SMPGD).** Let  $\phi$  be a mirror map which is  $\alpha_\phi$ -convex relative to a norm  $\|\cdot\|$ . We assume that  $f$  is  $\alpha_f$ -convex and  $g$  is  $\alpha_g$ -convex, relative to  $\phi$ ; we let  $\lambda_h := (1 - \alpha_f h)/(1 + \alpha_g h)$ . We assume that the stochastic gradient is unbiased.

- (smooth case) Assume that  $f$  is  $\beta_f$ -smooth relative to  $\phi$ , that  $h \leq 1/(2\beta_f)$ , and that we have a variance bound for the stochastic gradient:

$$\mathbb{E}[\|\hat{\nabla}f(x) - \nabla f(x)\|_*^2] \leq \sigma^2 d \quad \text{for all } x \in \mathcal{C}_\phi. \quad (12.1)$$

Then, for a suitably averaged iterate  $\bar{x}_N$ ,

$$\mathbb{E} F(\bar{x}_N) - F_\star \leq \frac{\alpha_f + \alpha_g}{\lambda_h^{-N} - 1} D_\phi(x_\star, x_0) + \frac{(1 + \alpha_g h) \sigma^2 d h}{\alpha_\phi}.$$

- (non-smooth case) Assume that the stochastic gradients are  $L^2$ -bounded,

$$\mathbb{E}[\|\hat{\nabla}f(x)\|_*^2] \leq L^2 \quad \text{for all } x \in \mathcal{C}_\phi. \quad (12.2)$$

Then, for a suitably averaged iterate  $\bar{x}_N$ ,

$$\mathbb{E} F(\bar{x}_N) - F_\star \leq \frac{\alpha_f + \alpha_g}{\lambda_h^{-N} - 1} D_\phi(x_\star, x_0) + \frac{2(1 + \alpha_g h) L^2 h}{\alpha_\phi}.$$

*Proof.* We prove a one-step inequality for

$$x^+ := \arg \min \psi_x, \quad \psi_x(y) := f(x) + \langle \hat{\nabla}f(x), y - x \rangle + g(y) + \frac{1}{h} D_\phi(y, x).$$

By the relative growth inequality (Lemma 10.7), for any  $y \in \mathcal{C}_\phi$ ,

$$\left(\alpha_g + \frac{1}{h}\right) D_\phi(y, x^+) + \psi_x(x^+) \leq \psi_x(y).$$

---

<sup>12</sup>Yes, the name is comically long.

For the right-hand side, in both cases,

$$\begin{aligned}\mathbb{E} \psi_x(y) &= f(x) + \langle \nabla f(x), y - x \rangle + g(y) + \frac{1}{h} D_\phi(y, x) = F(y) - D_f(y, x) + \frac{1}{h} D_\phi(y, x) \\ &\leq F(y) + \frac{1 - \alpha_f h}{h} D_\phi(y, x).\end{aligned}$$

**Smooth case.** Here, we lower bound  $\mathbb{E} \psi_x(x^+)$  as follows: since  $h \leq 1/(2\beta_f)$ ,

$$\begin{aligned}\mathbb{E} \psi_x(x^+) &= \mathbb{E} \left[ f(x) + \langle \hat{\nabla} f(x), x^+ - x \rangle + g(x^+) + \frac{1}{h} D_\phi(x^+, x) \right] \\ &= \mathbb{E} \left[ f(x) + \langle \nabla f(x), x^+ - x \rangle + g(x^+) + \frac{1}{h} D_\phi(x^+, x) + \langle \hat{\nabla} f(x) - \nabla f(x), x^+ - x \rangle \right] \\ &= \mathbb{E} \left[ F(x^+) - D_f(x^+, x) + \frac{1}{h} D_\phi(x^+, x) + \langle \hat{\nabla} f(x) - \nabla f(x), x^+ - x \rangle \right] \\ &\geq \mathbb{E} \left[ F(x^+) + \frac{1}{2h} D_\phi(x^+, x) - \|\hat{\nabla} f(x) - \nabla f(x)\|_* \|x^+ - x\| \right] \\ &\geq \mathbb{E} \left[ F(x^+) + \frac{\alpha_\phi}{4h} \|x^+ - x\|^2 \right] - \sqrt{\mathbb{E} [\|\hat{\nabla} f(x) - \nabla f(x)\|_*^2] \mathbb{E} [\|x^+ - x\|^2]} \\ &\geq \mathbb{E} \left[ F(x^+) + \frac{\alpha_\phi}{4h} \|x^+ - x\|^2 \right] - \sqrt{\sigma^2 d \mathbb{E} [\|x^+ - x\|^2]} \geq \mathbb{E} F(x^+) - \frac{\sigma^2 d h}{\alpha_\phi}.\end{aligned}$$

This leads to the one-step bound

$$(1 + \alpha_g h) \mathbb{E} D_\phi(y, x^+) \leq (1 - \alpha_f h) D_\phi(y, x) - h (\mathbb{E} F(x^+) - F(y)) + \frac{\sigma^2 d h^2}{\alpha_\phi}.$$

**Non-smooth case.** In this case, note that

$$\|\nabla f(x)\|_* = \|\mathbb{E} \hat{\nabla} f(x)\|_* \leq \sqrt{\mathbb{E} [\|\hat{\nabla} f(x)\|_*^2]} \leq L,$$

so that  $f$  is  $L$ -Lipschitz with respect to  $\|\cdot\|$ . Then,

$$\begin{aligned}\mathbb{E} \psi_x(x^+) &= \mathbb{E} \left[ f(x) + \langle \hat{\nabla} f(x), x^+ - x \rangle + g(x^+) + \frac{1}{h} D_\phi(x^+, x) \right] \\ &= \mathbb{E} \left[ F(x^+) + f(x) - f(x^+) + \langle \hat{\nabla} f(x), x^+ - x \rangle + \frac{1}{h} D_\phi(x^+, x) \right] \\ &\geq \mathbb{E} \left[ F(x^+) - L \|x^+ - x\| - \|\hat{\nabla} f(x)\|_* \|x^+ - x\| + \frac{1}{h} D_\phi(x^+, x) \right]\end{aligned}$$

$$\geq \mathbb{E}\left[F(x^+) + \frac{\alpha_\phi}{2h} \|\|x^+ - x\|\|^2\right] - 2L\sqrt{\mathbb{E}[\|\|x^+ - x\|\|^2]} \geq \mathbb{E}F(x^+) - \frac{2L^2h}{\alpha_\phi}.$$

This leads to the one-step bound

$$(1 + \alpha_g h) \mathbb{E} D_\phi(y, x^+) \leq (1 - \alpha_f h) D_\phi(y, x) - h (\mathbb{E} F(x^+) - F(y)) + \frac{2L^2h^2}{\alpha_\phi}.$$

**Completing the proof.** Observe that the one-step bounds have exactly the same form in both cases. We take  $y = x_\star$  and iterate as usual. Let  $E = \sigma^2 dh^2/\alpha_\phi$  in the first case, and  $E = 2L^2h^2/\alpha_\phi$  in the second case. Then, the discrete Grönwall lemma ([Lemma 3.5](#)) and some computations yield

$$\sum_{n=1}^N \frac{\lambda_h^{N-n}}{\sum_{k=1}^N \lambda_h^k} \{\mathbb{E} F(x_n) - F_\star\} \leq \frac{\alpha_f + \alpha_g}{\lambda_h^{-N} - 1} D_\phi(x_\star, x_0) + \frac{1 + \alpha_g h}{h} E.$$

This yields a convergence rate for a suitably averaged iterate.  $\square$

For simplicity, let us assume that  $\alpha_\phi = 1$  and  $R^2 \geq D_\phi(x_\star, x_0)$ . We let  $\alpha := \alpha_f + \alpha_g$ . To obtain  $\varepsilon$  error, [Theorem 12.1](#) implies the rates in [Table 2](#).

Assumptions	Iterations
convex, smooth	$O(\sigma^2 d R^2 / \varepsilon^2)$
strongly convex, smooth	$O(\sigma^2 d \log(\alpha R^2 / \varepsilon) / (\alpha \varepsilon))$
convex, non-smooth	$O(L^2 R^2 / \varepsilon^2)$
strongly convex, non-smooth	$O(L^2 \log(\alpha R^2 / \varepsilon) / (\alpha \varepsilon))$

Table 2: Rates for [SMPGD](#) with an appropriate step size and averaging.

A few remarks are in order.

1. The effect of the stochasticity of the gradients, at least under our assumptions, is qualitatively the same as the effect of non-smoothness. Indeed, the rates of  $O(1/\varepsilon^2)$  under convexity and  $O(1/\varepsilon)$  under strong convexity reflect the rates for [PSD](#) in [§6.2](#). Similarly, in this setting we do not have a descent lemma, and hence we should average the iterates.
2. In the non-smooth case, stochasticity of the gradients does not hurt the rate at all, provided that the stochastic gradients satisfy an appropriate  $L^2$  bound.

3. By Jensen's inequality, it is not hard to see that (12.2) implies (12.1) with  $\sigma^2 d \leq 4L^2$ . Hence, the variance bound (12.1) is indeed a weaker assumption, although the analysis requires a stronger assumption—smoothness—for  $f$ . At first glance, it may seem that (12.1) and (12.2) are similar, but the former is a variance bound and the latter is an  $L^2$  bound. Suppose, for instance, that  $\|\cdot\|$  is the Euclidean norm, and that  $\hat{\nabla}f$  is computed on the basis of a mini-batch of  $B$  samples. Then, we expect (12.1) to decay as  $1/B$ , whereas (12.2) does not.
4. Since stochastic optimization behaves like non-smooth optimization, we also expect that it is not possible to “accelerate”. Indeed, many of the rates in Theorem 12.1 can be shown to be optimal up to logarithmic terms [Aga+12].

## 12.2 Implications for statistical generalization

Suppose that we have a dataset  $\{Z_i : i \in [n]\} \subseteq \mathcal{Z}$  of i.i.d. samples, a parameter space  $\Theta \subseteq \mathbb{R}^d$ , and a loss  $\ell : \Theta \times \mathcal{Z} \rightarrow \mathbb{R}$ . Define the *empirical risk* and the *population risk*:

$$\mathcal{R}_n(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(\theta; Z_i), \quad \mathcal{R}(\theta) := \mathbb{E} \mathcal{R}_n(\theta).$$

We further assume that  $\theta \mapsto \ell(\theta; z)$  is  $L$ -Lipschitz for all  $z \in \mathcal{Z}$ , and that  $\Theta$  is the ball  $B(0, R)$  of radius  $R$ .

**Example 12.2 (regression).** Suppose that  $Z_i = (X_i, Y_i)$  with  $X_i \in \mathbb{R}^d$  and  $Y_i \in [-1, 1]$ . Assume that for each  $\theta \in \Theta$ , we have a predictor  $f_\theta : \mathbb{R}^d \rightarrow [-1, 1]$ . Then, we can consider the squared loss

$$\ell(\theta; z) = \ell(\theta; (x, y)) = (y - f_\theta(x))^2,$$

If we further consider linear regression,  $f_\theta : x \mapsto \langle \theta, x \rangle$ , then  $f_\theta$  is  $4r$ -Lipschitz for all  $z = (x, y)$  with  $\|x\| \leq r$  and  $|y| \leq 1$ . Linear regression is not as restrictive as it may seem, since we can imagine that each  $X_i$  is actually the output of a high-dimensional feature map, in which case we obtain kernel regression.

The setting we have described is actually the standard one for statistical learning theory, and it covers much more than regression (e.g., classification and density estimation), but regression is a good representative example. Also, since our conclusions below will not involve the dimension  $d$ , we can think of the function class as infinite-dimensional.



Define the minimizers

$$\widehat{\theta}_n \in \arg \min_{\theta \in \Theta} \mathcal{R}_n(\theta), \quad \theta^\star \in \arg \min_{\theta \in \Theta} \mathcal{R}(\theta).$$

We view  $\theta^\star$  as the ground truth value of the parameter that we wish to recover. Since we do not have access to the population risk  $\mathcal{R}$ , we must base our procedures on the samples  $\{Z_i\}_{i \in [n]}$ , so it is natural to use  $\widehat{\theta}_n$ , the *empirical risk minimizer* (ERM), as our estimator. This is the starting point for statistical theory, but it says nothing about how we can actually compute  $\widehat{\theta}_n$ .

The power of stochastic gradient descent (SGD) is that we can view it as a method to *directly* minimize the population risk. We consider the iteration

$$\theta_{k+1} := \theta_k - h \nabla_{\theta} \ell(\theta_k; Z_{k+1}). \quad (\text{SGD})$$

If we denote  $\hat{\nabla} \mathcal{R}(\theta) = \nabla_{\theta} \ell(\theta; Z)$ , then this is indeed an unbiased stochastic gradient:  $\mathbb{E} \hat{\nabla} \mathcal{R}(\theta) = \nabla \mathcal{R}(\theta)$ . Due to our Lipschitz assumption, we also know that

$$\mathbb{E}[\|\hat{\nabla} \mathcal{R}(\theta)\|^2] \leq L^2.$$

However, the implicit assumption in [Theorem 12.1](#) is that the randomness is fresh at each iteration, so we are not allowed to reuse any samples. This limits the total number of iterations of [SGD](#) to the sample size  $n$ , and we refer to this as *one-pass* SGD.

From [Theorem 12.1](#), if we further assume that  $\theta \mapsto \ell(\theta; z)$  is *convex* for every  $z \in \mathcal{Z}$ , then [SGD](#) with an optimized step size and averaging satisfies

$$\mathbb{E} \mathcal{R}(\bar{\theta}_n) - \mathcal{R}(\theta^\star) \lesssim \frac{LR}{\sqrt{n}}.$$

Here, the guarantee is in terms of the difference between the expected risk of our estimator, the averaged iterate of [SGD](#), and the best possible risk. This is known as the *excess risk* and it is the best we can hope for, since  $\theta^\star$  may not be identifiable (unique).

How does this compare to the performance of ERM? Analysis of the ERM estimator starts with the following decomposition:

$$\begin{aligned} \mathcal{R}(\widehat{\theta}_n) - \mathcal{R}(\theta^\star) &= \mathcal{R}(\widehat{\theta}_n) - \mathcal{R}_n(\widehat{\theta}_n) + \underbrace{\mathcal{R}_n(\widehat{\theta}_n) - \mathcal{R}_n(\theta^\star)}_{\leq 0} + \mathcal{R}_n(\theta^\star) - \mathcal{R}(\theta^\star) \\ &\leq 2 \sup_{\theta \in \Theta} |\mathcal{R}_n(\theta) - \mathcal{R}(\theta)|. \end{aligned}$$

We therefore want to show that  $\sup_{\theta \in \Theta} |\mathcal{R}_n(\theta) - \mathcal{R}(\theta)|$  tends to zero at a certain rate, which is known as a uniform convergence argument. This is a type of stochastic process

(since  $\mathcal{R}_n$  is random) known as an empirical process, and sophisticated tools have been developed for its study. After applying a number of them (symmetrization, contraction principle, control of the Rademacher complexity), one can show that

$$\mathbb{E} \mathcal{R}(\hat{\theta}_n) - \mathcal{R}(\theta^\star) \lesssim \frac{LR}{\sqrt{n}}, \quad (12.3)$$

just as for [SGD](#).

Actually, it is worth remarking that the bounds from empirical process theory depend on various notions of the complexity of the class  $\{\ell(\theta; \cdot) : \theta \in \Theta\}$ . A traditional approach measures this complexity essentially by counting the number of free parameters in the class (Vapnik–Chervonenkis or VC theory), and would not match the dimension-free rate attained by SGD. In order to do so, one needs to turn toward “size-based” measures of complexity that take into account the fact that  $\Theta$  lies in a ball, hence the need to control the Rademacher complexity directly.

Anyway, we can show that the ERM estimator satisfies (12.3), and moreover, this argument does not require convexity of the loss. However, when we discuss how to compute the ERM estimator, we need to assume convexity anyway. Let us suppose that we compute it by running [GD](#) (or more specifically, [PSD](#)) on the empirical risk  $\mathcal{R}_n$ . Since the statistical error is already  $LR/\sqrt{n}$ , we only need to optimize to this level of accuracy. Applying [Theorem 6.14](#), we see that the number of iterations of [PSD](#) is roughly  $n$ . This is the same number of iterations as one-pass [SGD](#), except that each iteration of [SGD](#) is roughly  $n$  times cheaper than the corresponding one for [PSD](#).

In conclusion, one-pass [SGD](#) produces an estimator which has comparable statistical performance to the ERM estimator, with a computational cost roughly  $n$  times cheaper than minimizing the empirical risk directly using [PSD](#).

The dependence  $n^{-1/2}$  dependence on the sample size  $n$  is known as a “slow rate”. It can be improved when we additionally assume that for every  $z \in \mathcal{Z}$ ,  $\theta \mapsto \ell(\theta; z)$  is  $\alpha$ -convex for some  $\alpha > 0$ . In this case, [Theorem 12.1](#) yields

$$\mathbb{E} \mathcal{R}(\bar{\theta}_n) - \mathcal{R}(\theta^\star) \lesssim \frac{L^2}{\alpha n}.$$

(Actually, [Theorem 12.1](#) yields a slightly worse result by a logarithmic factor, but this can be fixed with a time-varying choice of step sizes.) Due to strong convexity, it also implies parameter recovery:

$$\sqrt{\mathbb{E}[\|\bar{\theta}_n - \theta^\star\|^2]} \lesssim \frac{L}{\alpha \sqrt{n}}.$$

What about for the ERM estimator? This time, one needs a refined argument based on *localized* Rademacher complexities, and it again eventually yields

$$\mathbb{E} \mathcal{R}(\hat{\theta}_n) - \mathcal{R}(\theta^\star) \lesssim \frac{L^2}{\alpha n}.$$

As for computation, by applying [PSD](#) and the result of [Exercise 6.5](#) (again, omitting the logarithmic factor which can be removed with better step sizes), the conclusion is the same: the number of iterations is the same as for [SGD](#), but each iteration is roughly  $n$  times more expensive.

This discussion suggests that, at least when the risk is convex, averaged one-pass [SGD](#) is expected to be an excellent estimator, both computationally and statistically. The next subsection will further reinforce this point.

### 12.3 Central limit theorem for Polyak–Ruppert averaging

**Disclaimer:** This subsection is somewhat technical, so rather than tracing through all of the details, you are encouraged to follow the high-level ideas.

We now turn to a celebrated result in stochastic optimization: namely, that the iterates of SGD with Polyak–Ruppert averaging obey a central limit theorem. Let

$$\theta_{n+1} := \theta_n - h_{n+1} \hat{\nabla} f(\theta_n), \quad \bar{\theta}_n := \frac{1}{n} \sum_{j=0}^{n-1} \theta_j. \quad (\text{ASGD})$$

Our goal is to show that  $\sqrt{n}(\bar{\theta}_n - \theta^\star) \rightarrow \text{normal}(0, \Sigma)$  for a certain covariance matrix  $\Sigma$ . Throughout, we write

$$\hat{\nabla} f(\theta_n) = \nabla f(\theta_n) + \xi_{n+1},$$

where conditionally on  $\theta_n$ , the random vector  $\xi_{n+1}$  has zero mean. Our condition on the step sizes is

$$h_n = n^{-\gamma}, \quad \text{for some } \frac{1}{2} < \gamma < 1. \quad (12.4)$$

We recall some preliminaries on convergence in distribution.

- We say that a sequence  $\{X_n\}_{n \in \mathbb{N}}$  of random vectors converges in distribution (or converges weakly in law) to a probability distribution  $\mu$ , denoted  $X_n \xrightarrow{d} \mu$ , if  $\mathbb{E} f(X_n) \rightarrow \int f d\mu$  for every bounded, continuous function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ . This is the same notion of convergence as in the classical central limit theorem (CLT).

- We say that  $\{X_n\}_{n \in \mathbb{N}}$  tends to 0 in probability if for all  $\varepsilon > 0$ ,  $\mathbb{P}(\|X_n\| \geq \varepsilon) \rightarrow 0$ . If  $\mathbb{E}\|X_n\| \rightarrow 0$ , then  $X_n \rightarrow 0$  in probability; this follows from Markov's inequality.
- If  $X_n = Y_n + Z_n$ , and  $Z_n \rightarrow 0$  in probability, then  $\{X_n\}_{n \in \mathbb{N}}$  and  $\{Y_n\}_{n \in \mathbb{N}}$  have the same distributional limit.

**Quadratic case.** We begin with the quadratic case with

$$f : \mathbb{R}^d \rightarrow \mathbb{R}, \quad f(\theta) = \frac{1}{2} \langle \theta - \theta^\star, A(\theta - \theta^\star) \rangle.$$

We assume  $A \succ 0$ . We do not treat this case separately merely as a “warm-up”; the analysis here is crucial for the general case as well.

**Theorem 12.3 (CLT for ASGD, quadratic case).** Assume that  $f$  is a strongly convex quadratic function, and that the step sizes satisfy the condition (12.4). Assume that conditionally on  $\theta_n$ , each  $\xi_{n+1}$  is mean zero and has covariance  $\Xi_{n+1}$ , such that  $n^{-1} \sum_{k=1}^n \Xi_k \rightarrow S_\infty$  in probability and  $\sup_{n \geq 1} \mathbb{E} \operatorname{tr} \Xi_n < \infty$ . Then,

$$\sqrt{n} (\bar{\theta}_n - \theta^\star) \xrightarrow{d} \operatorname{normal}(0, A^{-1} S_\infty A^{-1}).$$

Before turning toward the proof, let us consider a simple example.

**Example 12.4 (estimating the mean of a Gaussian).** Suppose that we have samples  $\{X_n\}_{n \in \mathbb{N}} \stackrel{\text{i.i.d.}}{\sim} \operatorname{normal}(\theta^\star, A^{-1})$  for a known covariance matrix  $A \succ 0$ , and we wish to estimate the mean  $\theta^\star$ . We consider

$$\theta_{n+1} = \theta_n - h_{n+1} A (\theta_n - X_{n+1}), \quad \bar{\theta}_n = \frac{1}{n} \sum_{j=0}^{n-1} \theta_j.$$

This corresponds to the quadratic loss function with  $\xi_n = A(\theta^\star - X_n)$ . In this case, the  $\{\xi_n\}_{n \in \mathbb{N}}$  are i.i.d., with mean zero and covariance matrix  $A$ .

If we use the sample mean, then

$$\sqrt{n} \left( \frac{1}{n} \sum_{j=1}^n X_j - \theta^\star \right) \sim \operatorname{normal}(0, A^{-1}).$$

On the other hand, the CLT for ASGD with  $S_\infty = A$  shows that

$$\sqrt{n} (\bar{\theta}_n - \theta^\star) \rightarrow \operatorname{normal}(0, A^{-1}).$$

The first step is to write out the iterates. For  $\delta_n := \theta_n - \theta^\star$ ,

$$\delta_{n+1} = \delta_n - h_{n+1} A \delta_n - h_{n+1} \xi_{n+1} = (I - h_{n+1} A) \delta_n - h_{n+1} \xi_{n+1}.$$

Unrolling,

$$\delta_n = \left[ \prod_{j=1}^n (I - h_j A) \right] \delta_0 - \sum_{j=1}^n \left[ h_j \prod_{k=j+1}^n (I - h_k A) \right] \xi_j.$$

Defining  $\bar{\delta}_n := n^{-1} \sum_{j=0}^{n-1} \delta_j$ , it yields

$$\begin{aligned} \bar{\delta}_n &= \frac{1}{n} \sum_{j=0}^{n-1} \left[ \prod_{k=1}^j (I - h_k A) \right] \delta_0 - \frac{1}{n} \sum_{j=0}^{n-1} \sum_{k=1}^j \left[ h_k \prod_{\ell=k+1}^j (I - h_\ell A) \right] \xi_k \\ &= \frac{1}{n} \sum_{j=0}^{n-1} \left[ \prod_{k=1}^j (I - h_k A) \right] \delta_0 - \underbrace{\frac{1}{n} \sum_{k=1}^{n-1} h_k \sum_{j=k}^{n-1} \left[ \prod_{\ell=k+1}^j (I - h_\ell A) \right] \xi_k}_{=: M_k^n} \\ &= \frac{1}{n} M_0^n \delta_0 - \frac{1}{n} \sum_{k=1}^{n-1} A^{-1} \xi_k - \frac{1}{n} \sum_{k=1}^{n-1} (M_k^n - A^{-1}) \xi_k, \end{aligned} \tag{12.5}$$

where we set  $h_0 := 1$ . The intuition is that if all of the  $h_\ell$  were the same, then  $M_k^n$  would equal  $h \sum_{j=k}^{n-1} (I - hA)^{j-k} \rightarrow A^{-1}$  as  $n \rightarrow \infty$ , so we hope that the last term tends to zero.

**Lemma 12.5.** The  $M_k^n$  matrices are uniformly bounded in operator norm. Also,

$$\frac{1}{n} \sum_{k=1}^{n-1} \|M_k^n - A^{-1}\|_{\text{op}} \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

*Proof sketch.* For intuition, suppose that  $S = \sum_{n=0}^{\infty} (I - hA)^n$  for  $h < 1/\|A\|_{\text{op}}$ . To show that  $hS = A^{-1}$ , one can argue that  $(I - hA)S = S - I$ , which leads to  $hS = A^{-1}$ . We want to replicate this type of proof, but it is significantly complicated due to the time-varying step sizes. Let  $B_j^k := \prod_{\ell=j}^k (I - h_\ell A)$ , so that  $M_k^n = h_k \sum_{j=k}^{n-1} B_{k+1}^j$ . We start with an equation for the  $B$ 's: since

$$B_j^{k+1} = B_j^k - h_{k+1} A B_j^k = \cdots = I - A \sum_{\ell=j-1}^k h_{\ell+1} B_j^\ell,$$

we can write

$$M_k^n = h_k \sum_{j=k}^{n-1} B_{k+1}^j = \sum_{j=k}^{n-1} (h_k - h_{j+1} + h_{j+1}) B_{k+1}^j = \sum_{j=k}^{n-1} (h_k - h_{j+1}) B_{k+1}^j + A^{-1} (I - B_{k+1}^n).$$

Therefore,

$$AM_k^n - I = A \sum_{j=k}^{n-1} (h_k - h_{j+1}) B_{k+1}^j - B_{k+1}^n.$$

Since  $A$  is bounded above and below,

$$\frac{1}{n} \sum_{k=1}^{n-1} \|M_k^n - A^{-1}\|_{\text{op}} \lesssim \frac{1}{n} \sum_{k=1}^{n-1} \sum_{j=k}^{n-1} |h_k - h_{j+1}| \|B_{k+1}^j\|_{\text{op}} + \frac{1}{n} \sum_{k=2}^n \|B_k^n\|_{\text{op}}.$$

Also, it is easy to see that for sufficiently small step sizes (which is all that matters in the asymptotic regime),

$$\|B_k^n\|_{\text{op}} \leq \prod_{\ell=k}^n (I - \alpha h_\ell) \leq \exp\left(-\alpha \sum_{\ell=k}^n h_\ell\right).$$

So, for the second term,

$$\frac{1}{n} \sum_{k=2}^n \|B_k^n\|_{\text{op}} \leq \frac{1}{n} \sum_{k=1}^n \exp\left(-\alpha \sum_{\ell=k}^n h_\ell\right).$$

Let  $\tau_k := \sum_{\ell=1}^k h_\ell \approx (1 - \gamma)^{-1} k^{1-\gamma}$ . If, for any  $t > 0$ , we define  $\tau(t) = t^{1-\gamma}$ , the summation is roughly equivalent to the following integral (up to replacing  $\alpha$  by  $\alpha/(1 - \gamma)$ ):

$$I := \int_1^n \exp(-\alpha (\tau(n) - \tau(k))) dk.$$

For ease of presentation, we focus on bounding the integral instead. By change of variables,  $t = \tau(k)$ ,

$$I \asymp \int_1^{n^{1-\gamma}} t^{\gamma/(1-\gamma)} \exp(-\alpha (n^{1-\gamma} - t)) dt \lesssim n^\gamma \int_1^{n^{1-\gamma}} \exp(-\alpha (n^{1-\gamma} - t)) dt \lesssim n^\gamma.$$

Since  $\gamma < 1$ , the second term tends to zero.

As for the first term, we follow a similar strategy and approximate it via

$$\begin{aligned}
& \frac{1}{n} \iint_{1 \leq k \leq j \leq n} \left( \frac{1}{k^\gamma} - \frac{1}{j^\gamma} \right) \exp(-\alpha(\tau(j) - \tau(k))) \, dj \, dk \\
& \asymp \frac{1}{n} \iint_{1 \leq s \leq t \leq n^{1-\gamma}} \underbrace{(t^{\gamma/(1-\gamma)} - s^{\gamma/(1-\gamma)})}_{\text{apply Taylor expansion}} \exp(-\alpha(t-s)) \, ds \, dt \\
& \lesssim \frac{n^{2\gamma-1}}{n} \iint_{1 \leq s \leq t \leq n^{1-\gamma}} (t-s) \exp(-\alpha(t-s)) \, ds \, dt \lesssim \frac{n^\gamma}{n} \int_1^{n^{1-\gamma}} t \exp(-\alpha t) \, dt \rightarrow 0.
\end{aligned}$$

This completes the heuristic proof of the convergence. The first statement, about the boundedness of the  $M_k^n$ , can be proved via similar arguments.  $\square$

Returning to (12.5), note that

$$\sqrt{n} \bar{\delta}_n = \frac{1}{n^{1/2}} M_0^n \delta_0 - \frac{1}{n^{1/2}} \sum_{k=1}^{n-1} A^{-1} \xi_k - \frac{1}{n^{1/2}} \sum_{k=1}^{n-1} (M_k^n - A^{-1}) \xi_k,$$

where  $\mathbb{E} \|M_0^n \delta_0\| / \sqrt{n} \rightarrow 0$ . For the last term, we use the fact that the  $\xi_k$ 's are orthogonal: for  $k < \ell$ , by conditioning on  $\theta_{1:\ell-1} := (\theta_1, \dots, \theta_{\ell-1})$ ,

$$\mathbb{E} \langle (M_k^n - A^{-1}) \xi_k, (M_\ell^n - A^{-1}) \xi_\ell \rangle = \mathbb{E} \langle (M_k^n - A^{-1}) \xi_k, (M_\ell^n - A^{-1}) \mathbb{E}[\xi_\ell \mid \theta_{1:\ell-1}] \rangle = 0.$$

Therefore,

$$\begin{aligned}
\frac{1}{n} \mathbb{E} \left[ \left\| \sum_{k=1}^{n-1} (M_k^n - A^{-1}) \xi_k \right\|^2 \right] &= \frac{1}{n} \sum_{k=1}^{n-1} \mathbb{E} [\| (M_k^n - A^{-1}) \xi_k \|^2] = \frac{1}{n} \sum_{k=1}^{n-1} \langle (M_k^n - A^{-1})^2, \mathbb{E} \Xi_k \rangle \\
&\lesssim \frac{1}{n} \sum_{k=1}^n \|M_k^n - A^{-1}\|_{\text{op}} \rightarrow 0.
\end{aligned}$$

Hence, to obtain a distributional limit for  $\sqrt{n} \bar{\delta}_n$ , it suffices to obtain one for the second term above. This will be accomplished via martingale theory.

**Martingale CLT.** In the context of stochastic optimization, the noise sequence  $\{\xi_n\}_{n \in \mathbb{N}}$  is not i.i.d.; indeed, we want to consider

$$\xi_n = \hat{\nabla} f(\theta_n) - \nabla f(\theta_n),$$

and since the iterates  $\{\theta_n\}_{n \in \mathbb{N}}$  are random and depend on the noise sequence, it leads to a complicated dependence structure for the noise sequence. Nevertheless, it falls within the framework of martingale theory.

**Definition 12.6.** An increasing sequence of  $\sigma$ -algebras  $\{\mathcal{F}_n\}_{n \in \mathbb{N}}$  is called a **filtration**. We think of  $\mathcal{F}_n$  as the information available to an observer up to iteration  $n$ .

A sequence of random vectors  $\{X_n\}_{n \in \mathbb{N}}$  is called a **martingale** if for all  $n$ ,  $X_n$  is  $\mathcal{F}_n$ -measurable,  $\mathbb{E}\|X_n\| < \infty$ , and

$$\mathbb{E}[X_{n+1} \mid \mathcal{F}_n] = X_n.$$

In other words, the difference  $X_{n+1} - X_n$  is conditionally unbiased given the information  $\mathcal{F}_n$  available at iteration  $n$ . If we set  $X_n := \sum_{k=1}^n \xi_k$ , then  $\{X_n\}_{n \in \mathbb{N}}$  is a martingale; we sometimes refer to the noise sequence  $\{\xi_n\}_{n \in \mathbb{N}}$  as a *martingale difference sequence*.

Our next goal is to establish the following theorem.

**Theorem 12.7 (martingale CLT).** Let  $\{\xi_n\}_{n \in \mathbb{N}}$  be a martingale difference sequence and write  $\Xi_{n+1} := \text{cov}(\xi_{n+1} \mid \mathcal{F}_n)$ . Assume that  $\sup_{n \in \mathbb{N}} \mathbb{E} \text{tr} \Xi_n < \infty$  and that  $n^{-1} \sum_{k=1}^n \Xi_k \rightarrow S_\infty$  in probability. Then,

$$\frac{1}{\sqrt{n}} \sum_{k=1}^n \xi_k \xrightarrow{d} \text{normal}(0, S_\infty) \quad \text{as } n \rightarrow \infty.$$

This is a special case of a more general theorem on triangular arrays of Lindeberg–Feller type. For simplicity, we prove it under the stronger hypothesis that  $\{\xi_n\}_{n \in \mathbb{N}}$  is uniformly bounded.

*Proof.* Let  $X_k := n^{-1/2} \sum_{\ell=1}^k \xi_\ell$ ; note that this depends on  $n$  but we suppress it from the notation. Consider the characteristic function: for  $\mathbf{i} := \sqrt{-1}$  and  $\lambda \in \mathbb{R}^d$ ,

$$\phi_n(\lambda) := \mathbb{E} \exp(\mathbf{i} \langle \lambda, X_n \rangle).$$

Due to standard results on Fourier inversion, it suffices to prove that the characteristic function  $\phi_n$  converges pointwise to the characteristic function of the Gaussian,

$$\phi_\infty(\lambda) := \mathbb{E} \exp(\mathbf{i} \langle \lambda, Z \rangle) = \exp\left(-\frac{1}{2} \langle \lambda, S_\infty \lambda \rangle\right), \quad Z \sim \text{normal}(0, S_\infty).$$

Let  $S_k := n^{-1} \sum_{\ell=1}^k \Xi_\ell$ . We start by writing

$$\begin{aligned} |\phi_n(\lambda) - \phi_\infty(\lambda)| &\leq \left| \mathbb{E} \left[ \exp(\mathbf{i} \langle \lambda, X_n \rangle) \left( 1 - \exp\left(\frac{1}{2} \langle \lambda, S_n \lambda \rangle\right) - \frac{1}{2} \langle \lambda, S_\infty \lambda \rangle \right) \right] \right| \\ &\quad + \left| \exp\left(-\frac{1}{2} \langle \lambda, S_\infty \lambda \rangle\right) \left( \mathbb{E} \exp(\mathbf{i} \langle \lambda, X_n \rangle) + \frac{1}{2} \langle \lambda, S_n \lambda \rangle - 1 \right) \right| \end{aligned}$$



$$\begin{aligned} &\leq \mathbb{E} \left| 1 - \exp\left(\frac{1}{2} \langle \lambda, S_n \lambda \rangle - \frac{1}{2} \langle \lambda, S_\infty \lambda \rangle\right) \right| \\ &\quad + \left| \mathbb{E} \exp(\mathbf{i} \langle \lambda, X_n \rangle + \frac{1}{2} \langle \lambda, S_n \lambda \rangle) - 1 \right|. \end{aligned}$$

Since  $\{\xi_n\}_{n \in \mathbb{N}}$  is bounded, say by  $B$ , then  $\|\Xi_k\|_{\text{op}} \leq B^2$ , so  $\{S_n\}_{n \in \mathbb{N}}$  is bounded. Since  $S_n \rightarrow S_\infty$  in probability, the first term above tends to zero as  $n \rightarrow \infty$ .

For the second term, we peel off the terms in  $X_n$  by conditioning. Indeed,

$$\begin{aligned} \mathbb{E} \exp(\mathbf{i} \langle \lambda, X_n \rangle) &= \mathbb{E} \exp(\mathbf{i} \langle \lambda, X_{n-1} + n^{-1/2} \xi_n \rangle) \\ &= \mathbb{E} [\exp(\mathbf{i} \langle \lambda, X_{n-1} \rangle) \mathbb{E} [\exp(\mathbf{i} \langle \lambda, n^{-1/2} \xi_n \rangle) \mid \mathcal{F}_{n-1}]] . \end{aligned}$$

By Taylor expansion,

$$\mathbb{E} [\exp(\mathbf{i} \langle \lambda, n^{-1/2} \xi_n \rangle) \mid \mathcal{F}_{n-1}] = \mathbb{E} \left[ 1 + \frac{\mathbf{i}}{n^{1/2}} \langle \lambda, \xi_n \rangle - \frac{1}{2n} \langle \lambda, \xi_n \rangle^2 + O(n^{-3/2}) \mid \mathcal{F}_{n-1} \right]$$

where the error term  $O(n^{-3/2})$  is uniform, due to the assumption of boundedness. By the martingale property, this equals

$$1 - \frac{1}{2n} \langle \lambda, \Xi_n \lambda \rangle + O(n^{-3/2}) = \exp\left(-\frac{1}{2n} \langle \lambda, \Xi_n \lambda \rangle + O(n^{-3/2})\right) .$$

Hence,

$$\begin{aligned} \mathbb{E} \exp(\mathbf{i} \langle \lambda, X_n \rangle + \frac{1}{2} \langle \lambda, S_n \lambda \rangle) &= \mathbb{E} \exp(\mathbf{i} \langle \lambda, X_{n-1} \rangle + \frac{1}{2} \langle \lambda, S_n \lambda \rangle - \frac{1}{2n} \langle \lambda, \Xi_n \lambda \rangle + O(n^{-3/2})) \\ &= \mathbb{E} \exp(\mathbf{i} \langle \lambda, X_{n-1} \rangle + \frac{1}{2} \langle \lambda, S_{n-1} \lambda \rangle + O(n^{-3/2})) . \end{aligned}$$

Iterating,

$$\mathbb{E} \exp(\mathbf{i} \langle \lambda, X_n \rangle + \frac{1}{2} \langle \lambda, S_n \lambda \rangle) = \exp(O(n^{-1/2})) ,$$

so the second error term above also tends to zero.  $\square$

This completes the proof of [Theorem 12.3](#) since, if  $n^{-1/2} \sum_{k=1}^n \xi_k \xrightarrow{d} \text{normal}(0, S_\infty)$ , it follows that  $n^{-1/2} \sum_{k=1}^n A \xi_k \xrightarrow{d} \text{normal}(0, A^{-1} S_\infty A^{-1})$ .

**General case.** We now return to [ASGD](#) for a general function  $f$ . In this case, the CLT still holds, where we take  $A = \nabla^2 f(\theta^\star)$  to be the Hessian at the minimizer.

**Theorem 12.8 (CLT for ASGD, general case).** Assume that  $f$  is a strongly convex, smooth, and has a bounded third derivative, and that the step sizes satisfy the condition (12.4). Assume that conditionally on  $\theta_n$ , each  $\xi_{n+1}$  is mean zero and has covariance  $\Xi_{n+1}$ , such that  $n^{-1} \sum_{k=1}^n \Xi_k \rightarrow S_\infty$  in probability and  $\sup_{n \geq 1} \mathbb{E} \operatorname{tr} \Xi_n < \infty$ . Then,

$$\sqrt{n} (\bar{\theta}_n - \theta^\star) \xrightarrow{d} \text{normal}(0, A^{-1} S_\infty A^{-1}), \quad A := \nabla^2 f(\theta^\star).$$

Write the iterate for [ASGD](#) as

$$\begin{aligned} \theta_{n+1} &= \theta_n - h_{n+1} (\nabla f(\theta_n) + \xi_{n+1}) \\ &= \theta_n - h_{n+1} A (\theta_n - \theta^\star) - h_{n+1} \xi_{n+1} - h_{n+1} \underbrace{(\nabla f(\theta_n) - A (\theta_n - \theta^\star))}_{=:\zeta_n}, \end{aligned}$$

which leads to

$$\delta_{n+1} = \delta_n - h_{n+1} A \delta_n - h_{n+1} (\xi_{n+1} + \zeta_n).$$

By applying the derivation of (12.5), replacing  $\xi_{n+1}$  with  $\xi_{n+1} + \zeta_n$ , we obtain

$$\sqrt{n} \bar{\delta}_n = \frac{1}{n^{1/2}} M_0^n \delta_0 - \frac{1}{n^{1/2}} \sum_{k=1}^{n-1} A^{-1} (\xi_k + \zeta_{k-1}) - \frac{1}{n^{1/2}} \sum_{k=1}^{n-1} (M_k^n - A^{-1}) (\xi_k + \zeta_{k-1}).$$

We must show that the extra terms involving the  $\zeta_k$ 's vanish in the limit. Since the  $M_k^n$  matrices are bounded, it suffices to show that

$$\frac{1}{\sqrt{n}} \sum_{k=1}^n \mathbb{E} \|\zeta_k\| \rightarrow 0.$$

Since we assume that the third derivative of  $f$  is bounded, Taylor expansion shows that

$$\|\zeta_n\| = \|\nabla f(\theta_n) - \nabla^2 f(\theta^\star) (\theta_n - \theta^\star)\| \lesssim \|\theta_n - \theta^\star\|^2.$$

By our usual argument, for  $n$  large so that  $h_n$  is small,

$$\begin{aligned} \mathbb{E}[\|\theta_{n+1} - \theta^\star\|^2] &= \mathbb{E}[\|\theta_n - \theta^\star\|^2 - 2h_{n+1} \langle \nabla f(\theta_n), \theta_n - \theta^\star \rangle + h_{n+1}^2 \|\hat{\nabla} f(\theta_n)\|^2] \\ &= \mathbb{E}[\|\theta_n - \theta^\star\|^2 - 2h_{n+1} \langle \nabla f(\theta_n), \theta_n - \theta^\star \rangle] \end{aligned}$$

$$\begin{aligned}
& + \mathbb{E}[h_{n+1}^2 \|\nabla f(\theta_n)\|^2 + h_{n+1}^2 \|\hat{\nabla} f(\theta_n) - \nabla f(\theta_n)\|^2] \\
& \leq \mathbb{E}[(1 - \alpha h_{n+1}) \|\theta_n - \theta^\star\|^2 + h_{n+1}^2 \text{tr } \Xi_n] \\
& = (1 - \alpha h_{n+1}) \mathbb{E}[\|\theta_n - \theta^\star\|^2] + O(h_{n+1}^2).
\end{aligned}$$

Iterating,

$$\mathbb{E}[\|\theta_n - \theta^\star\|^2] \leq \exp\left(-\alpha \sum_{k=1}^n h_k\right) \|\theta_0 - \theta^\star\|^2 + \sum_{k=1}^n O(h_k^2) \exp\left(-\alpha \sum_{\ell=k+1}^n h_\ell\right).$$

The estimate for the summation in the second term is similar to the computation that appears in [Lemma 12.5](#), except that it corresponds to the integral

$$I' := \int_1^n \frac{1}{k^{2\gamma}} \exp(-\alpha(\tau(n) - \tau(k))) dk.$$

A trickier calculation eventually shows that

$$\mathbb{E}[\|\theta_n - \theta^\star\|^2] \leq \exp(-\Omega(n^{1-\gamma})) \|\theta_0 - \theta^\star\|^2 + O(n^{-\gamma}) = O(n^{-\gamma}). \quad (12.6)$$

Hence,

$$\frac{1}{\sqrt{n}} \sum_{k=1}^n \mathbb{E}\|\zeta_k\| \lesssim \frac{1}{\sqrt{n}} \sum_{k=1}^n \mathbb{E}[\|\theta_k - \theta^\star\|^2] \lesssim \frac{1}{\sqrt{n}} \sum_{k=1}^n \frac{1}{k^\gamma} \asymp n^{1/2-\gamma}.$$

This tends to zero provided  $\gamma > 1/2$ , completing the proof of [Theorem 12.8](#).

**Application to parameter recovery.** We now consider the example of parameter recovery in a parametric family of densities  $\{p_\theta\}_{\theta \in \Theta}$ . Write  $\ell(\theta; z) := \log(1/p_\theta(z))$ . In this case, the empirical risk minimizer  $\hat{\theta}_n$  corresponds to the maximum likelihood estimator (MLE), and if  $Z_1, \dots, Z_n \stackrel{\text{i.i.d.}}{\sim} p_{\theta^\star}$ , the population minimizer is indeed  $\theta^\star$  (provided that the model is identifiable).

Consider one-pass averaged SGD, so that

$$\xi_{k+1} = \nabla_\theta \ell(\theta_k; Z_{k+1}) - \int \nabla_\theta \ell(\theta_k; z) p_{\theta^\star}(dz).$$

This is conditionally unbiased, and if we define

$$I(\theta) := \text{cov}_{p_{\theta^\star}} \nabla_\theta \ell(\theta; Z),$$

then  $\Xi_{k+1} = I(\theta_k)$ . Now, adopt the following assumptions:

- For each  $z \in \mathcal{Z}$ , the function  $\theta \mapsto \ell(\theta; z)$  is strongly convex, smooth, and has a bounded third derivative.
- $I(\cdot)$  is Lipschitz continuous.

The second assumption, together with the fact that  $\theta_n \rightarrow \theta^\star$  in probability by (12.6), implies that  $\Xi_n = I(\theta_{n-1}) \rightarrow I(\theta^\star)$  in probability, hence  $n^{-1} \sum_{k=1}^n \Xi_k \rightarrow I(\theta^\star)$  in probability as well. Actually, since  $\mathbb{E}\|I(\theta_{n-1}) - I(\theta^\star)\|_{\text{op}} \lesssim \mathbb{E}\|\theta_{n-1} - \theta^\star\| \rightarrow 0$ , it readily implies that  $\sup_{n \in \mathbb{N}} \mathbb{E} \text{tr} \Xi_n < \infty$ . All of the assumptions of Theorem 12.8 are met.

The value of  $I(\cdot)$  at  $\theta^\star$  is special: it is called the *Fisher information matrix* and we denote it by  $\mathcal{J}$ :

$$\mathcal{J} := I(\theta^\star) = \text{cov}_{p_{\theta^\star}} \nabla_\theta \ell(\theta^\star; Z).$$

Since

$$\begin{aligned} \int \nabla_\theta \ell(\theta; z) p_\theta(dz) &= - \int \nabla_\theta \log p_\theta(z) p_\theta(dz) = - \int \nabla_\theta p_\theta(z) dz \\ &= - \nabla_\theta \int p_\theta(dz) = 0, \end{aligned}$$

it follows that

$$\begin{aligned} 0 &= \nabla_\theta \int \nabla_\theta \ell(\theta^\star; z) p_{\theta^\star}(dz) = \int \nabla_\theta^2 \ell(\theta^\star; z) p_{\theta^\star}(dz) + \int \nabla_\theta \ell(\theta^\star; z) \otimes \nabla_\theta p_{\theta^\star}(z) dz \\ &= \int \nabla_\theta^2 \ell(\theta^\star; z) p_{\theta^\star}(dz) - \int \nabla_\theta \ell(\theta^\star; z) \otimes \nabla_\theta \ell(\theta^\star; z) p_{\theta^\star}(dz). \end{aligned}$$

Combined with the fact that  $\int \nabla_\theta \ell(\theta^\star; z) p_{\theta^\star}(dz) = 0$ , we can identify the second term above as  $\text{cov}_{p_{\theta^\star}} \nabla_\theta \ell(\theta^\star; Z)$ , hence

$$\mathcal{J} = \int \nabla_\theta^2 \ell(\theta^\star; z) p_{\theta^\star}(dz) = \nabla^2 \mathcal{R}(\theta^\star),$$

since  $\mathcal{R}(\theta) = \int \ell(\theta; z) p_{\theta^\star}(dz)$  for every  $\theta \in \Theta$ . Therefore, Theorem 12.8 implies

$$\sqrt{n} (\bar{\theta}_n - \theta^\star) \xrightarrow{d} \text{normal}(0, \mathcal{J}^{-1}).$$

On the other hand, it is classical that under these assumptions, the MLE also has an asymptotically Gaussian limit:

$$\sqrt{n} (\hat{\theta}_n - \theta^\star) \xrightarrow{d} \text{normal}(0, \mathcal{J}^{-1}).$$

This is a celebrated result in statistics because the asymptotic covariance  $\mathcal{F}^{-1}$  is also a lower bound on the covariance of any unbiased estimator of  $\theta^*$ , by the *Cramér–Rao* or *information inequality*. Moreover, via comparison of experiments, it is known that no estimator can perform better than the MLE, in the sense that the MLE is locally asymptotically minimax optimal. We have just shown that this asymptotic optimality property also carries over to Polyak–Ruppert averaging of SGD. Finally, we remark that when  $p_\theta = \text{normal}(\theta, A^{-1})$ , this encompasses [Example 12.4](#).

## 12.4 Variance reduction

In §12.2, we argued that the generalization bounds for [GD](#) and [SGD](#) are comparable, yet the overall computational cost for [GD](#) is roughly  $n$  times larger due to the larger per-iteration cost. In making this comparison, our assumption was that we do not aim to completely minimize the empirical risk; we simply want the optimization error to be comparable to the statistical error. However, if our goal is indeed to fully minimize the empirical risk  $\mathcal{R}_n$ , then [GD](#) can be faster than [SGD](#) for high-accuracy solutions.

For example, in the convex and smooth setting, assume that the cost of computing the full gradient of  $\mathcal{R}_n$  is  $n$  times larger than the cost of computing the gradient of a single term (corresponding to a single sample). Then, in order to obtain an  $\varepsilon$ -approximate minimizer of  $\mathcal{R}_n$ , the overall computational cost for [GD](#) is  $O(n\beta R^2/\varepsilon)$  ([Theorem 3.4](#)), whereas for [SGD](#) it is  $O(\sigma^2 d R^2/\varepsilon^2)$ , where  $\sigma^2 d$  is an upper bound on the variance of the stochastic gradient ([Theorem 12.1](#)). In the strongly convex and smooth setting, the rates are  $\tilde{O}(n\kappa \log(\alpha R^2/\varepsilon))$  and  $\tilde{O}(\sigma^2 d/(\alpha\varepsilon))$  respectively. In general, these rates are incomparable.

In this section, we show that we can improve upon both of these rates through a technique known as *variance reduction*. Namely, the method we develop runs in a number of iterations comparable to [GD](#), but with a per-iteration cost comparable to [SGD](#). The structural assumption we impose on  $f$  is that it is a finite sum of  $n$  functions:

$$f = \frac{1}{n} \sum_{i=1}^n f_i.$$

In this setting, it makes sense to measure the complexity in terms of the number of gradient evaluations of the *individual* functions  $f_i$ .

To see why there is a possibility for variance reduction, note that if we run [SGD](#), where our stochastic gradient  $\hat{\nabla}f(x)$  is taken to be  $\nabla f_i(x)$  for a randomly chosen index  $i \sim \text{uniform}([n])$ , then the variance of the stochastic gradient is

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x)\|^2.$$

Even when we are at the minimizer,  $\nabla f(x_\star) = 0$ , the variance of the stochastic gradient is non-zero. However, if the iterates of the algorithm are converging to the minimizer,  $x_n \rightarrow x_\star$ , then we can hope that the variance of the gradient estimator also tends to zero.

This intuition is carried out the family of variance reduction methods, of which we pick one representative one: *stochastic variance reduced gradient descent* (SVRG) [JZ13]. We generalize our setting to a composite objective:

$$F = f + g = \frac{1}{n} \sum_{i=1}^n f_i + g.$$

The algorithm proceeds via “epochs”, where the  $t$ -th epoch runs for  $N_t$  iterations. In the  $t$ -th epoch, we initialize  $x_0^t := x_{N_{t-1}}^{t-1}$  (that is, we start the  $t$ -epoch at the last iterate of the previous epoch). The algorithm is described as follows:

$$\begin{aligned} x_{n+1}^t &:= \arg \min_{x \in \mathbb{R}^d} \left\{ \langle \hat{\nabla}_n^t f, x - x_n^t \rangle + g(x) + \frac{1}{2h} \|x - x_n^t\|^2 \right\}, \\ \hat{\nabla}_n^t f &:= \nabla_{i_n^t} f(x_n^t) - \nabla_{i_n^t} f(\bar{x}_0^t) + \nabla f(\bar{x}_n^t), \quad i_n^t \sim \text{uniform}([n]). \end{aligned} \quad (\text{SVRG})$$

Here,  $\bar{x}_0^t$  is a certain average of iterates from the previous epoch  $t - 1$ . Note that in each epoch, we compute (and store) the full gradient  $\nabla f(\bar{x}_0^t)$ , which requires  $n$  gradient computations, and then each subsequent iteration requires only one gradient computation. Therefore, the  $t$ -th epoch requires  $n + N_t$  gradient computations, and the total cost after  $T$  epochs is  $Tn + \sum_{t=0}^{T-1} N_t$ .

The intuition here is that if we take the expectation over  $i_n^t$ , then

$$\mathbb{E} \hat{\nabla}_n^t f = \mathbb{E} [\nabla_{i_n^t} f(x_n^t) - \nabla_{i_n^t} f(\bar{x}_0^t) + \nabla f(\bar{x}_n^t)] = \nabla f(x_n^t),$$

so the gradient estimator is indeed unbiased. But the extra centered term that we added to the gradient estimator,  $-\nabla_{i_n^t} f(\bar{x}_0^t) + \nabla f(\bar{x}_n^t)$ , reduces the variance: since  $\bar{x}_0^t, x_n^t \rightarrow x_\star$ , we expect that

$$\hat{\nabla}_n^t f - \nabla f(x_n^t) = \nabla_{i_n^t} f(x_n^t) - \nabla_{i_n^t} f(\bar{x}_0^t) + \nabla f(\bar{x}_n^t) - \nabla f(x_n^t) \rightarrow 0.$$

**Theorem 12.9 (convergence of SVRG).** Assume that  $f$  is  $\alpha_f$ -convex and  $\beta$ -smooth, and that  $g$  is  $\alpha_g$ -convex. Then, the following assertions hold for a suitable choice of step size  $h$  and averaged iterate  $\bar{x}_0^t$ . Let  $\Delta_0 := F(x_0) - F_\star + \beta \|x_0 - x_\star\|^2$ .

- If  $\alpha_f + \alpha_g = 0$ , then SVRG can achieve  $\mathbb{E} F(\bar{x}_0^T) - F_\star \leq \varepsilon$  with a total number of gradient evaluations at most  $O(n \log(\Delta_0/\varepsilon) + \Delta_0/\varepsilon)$ .
- If  $\alpha_f + \alpha_g > 0$ , then SVRG can achieve  $\mathbb{E} F(\bar{x}_0^T) - F_\star \leq \varepsilon$  with a total number of gradient evaluations at most  $O((n + \kappa) \log(\Delta_0/\varepsilon))$ , where  $\kappa := \beta/(\alpha_f + \alpha_g)$ , where  $\kappa := \beta/(\alpha_f + \alpha_g)$ .

*Proof.* We start by analyzing a single epoch; thus, for simplicity of notation, we drop the superscript  $t$ . The one-step inequality for SGD (see [Theorem 12.1](#)) shows that

$$\begin{aligned}\mathbb{E} F(x_{n+1}) - F_\star &\leq \frac{1 - \alpha_f h}{2h} \mathbb{E}[\|x_n - x_\star\|^2] - \frac{1 + \alpha_g h}{2h} \mathbb{E}[\|x_{n+1} - x_\star\|^2] \\ &\quad + h \mathbb{E}[\|\hat{\nabla}_n f - \nabla f(x_n)\|^2].\end{aligned}$$

We upper bound the variance of the stochastic gradient: by (3.4),

$$\begin{aligned}\mathbb{E}[\|\nabla_{i_n} f(x_n) - \nabla_{i_n} f(\bar{x}_0) + \nabla f(\bar{x}_0) - \nabla f(x_n)\|^2] &\leq \mathbb{E}[\|\nabla_{i_n} f(x_n) - \nabla_{i_n} f(\bar{x}_0)\|^2] \\ &\leq 2 \mathbb{E}[\|\nabla_{i_n} f(x_n) - \nabla_{i_n} f(x_\star)\|^2] + 2 \mathbb{E}[\|\nabla_{i_n} f(\bar{x}_0) - \nabla_{i_n} f(x_\star)\|^2] \\ &\leq 2\beta \mathbb{E}[D_{f_{i_n}}(x_n, x_\star) + D_{f_{i_n}}(\bar{x}_0, x_\star)] = 2\beta \mathbb{E}[D_f(x_n, x_\star) + D_f(\bar{x}_0, x_\star)] \\ &\leq 2\beta \mathbb{E}[D_F(x_n, x_\star) + D_F(\bar{x}_0, x_\star)] = 2\beta \mathbb{E}[F(x_n) - F_\star + F(\bar{x}_0) - F_\star].\end{aligned}$$

Note that this already captures the intuition above, namely, the variance decreases with the objective gap. Therefore, we end up with the recursion

$$\begin{aligned}\mathbb{E} F(x_{n+1}) - F_\star &\leq \frac{1 - \alpha_f h}{2h} \mathbb{E}[\|x_n - x_\star\|^2] - \frac{1 + \alpha_g h}{2h} \mathbb{E}[\|x_{n+1} - x_\star\|^2] \\ &\quad + 2\beta h \mathbb{E}[F(x_n) - F_\star + F(\bar{x}_0) - F_\star].\end{aligned}$$

We now choose  $h = 1/(8\beta)$  so that  $2\beta h = 1/4$ . After dividing by  $1 + \alpha_g h$  and iterating using [Lemma 3.5](#), it yields

$$\begin{aligned}\frac{\mathbb{E}[\|x_N - x_\star\|^2]}{2h} &\leq \frac{\lambda_h^N \mathbb{E}[\|x_0 - x_\star\|^2]}{2h} \\ &\quad + \underbrace{\sum_{n=1}^N \lambda_h^{N-n} \left( \frac{1}{4} \{ \mathbb{E} F(x_{n-1}) - F_\star + \mathbb{E} F(\bar{x}_0) - F_\star \} - \{ \mathbb{E} F(x_n) - F_\star \} \right)}_{=(\star)}.\end{aligned}$$

Since by assumption  $1/(4\lambda_h) \leq 1/3$ , the last summation is at most

$$\begin{aligned}(\star) &= - \sum_{n=1}^N \lambda_h^{N-n} \left( 1 - \frac{1}{4\lambda_h} \right) (\mathbb{E} F(x_n) - F_\star) - \frac{1}{4\lambda_h} (\mathbb{E} F(x_N) - F_\star) \\ &\quad + \frac{\lambda_h^{N-1}}{4} (\mathbb{E} F(x_0) - F_\star) + \frac{S}{4} (\mathbb{E} F(\bar{x}_0) - F_\star) \\ &\leq -\frac{2}{3} \sum_{n=1}^N \lambda_h^{N-n} (\mathbb{E} F(x_n) - F_\star) - \frac{1}{4\lambda_h} (\mathbb{E} F(x_N) - F_\star)\end{aligned}$$

$$+ \frac{\lambda_h^{N-1}}{4} (\mathbb{E} F(x_0) - F_\star) + \frac{S}{3} (\mathbb{E} F(\bar{x}_0) - F_\star),$$

where  $S := \sum_{n=0}^{N-1} \lambda_h^n$ . Thus, the above inequality can be rearranged to yield

$$\begin{aligned} & \frac{\lambda_h^N \mathbb{E}[\|x_0 - x_\star\|^2]}{2hS} + \frac{\lambda_h^{N-1}}{4S} (\mathbb{E} F(x_0) - F_\star) + \frac{1}{3} (\mathbb{E} F(\bar{x}_0) - F_\star) \\ & \geq \frac{\mathbb{E}[\|x_N - x_\star\|^2]}{2hS} + \frac{1}{4\lambda_h S} (\mathbb{E} F(x_N) - F_\star) + \frac{2}{3} \sum_{n=1}^N \frac{\lambda_h^{N-n}}{S} (\mathbb{E} F(x_n) - F_\star). \end{aligned}$$

The goal now is to make this inequality telescope across the epochs. We recall that  $x_0^{t+1} = x_{N_t}^t$ , and we define  $\bar{x}_0^{t+1} := S_t^{-1} \sum_{n=1}^{N_t} \lambda_h^{N_t-n} x_n^t$ , where  $S_t := \sum_{n=0}^{N_t} \lambda_h^n$ . By applying convexity to the last term, the inequality can be rewritten

$$\begin{aligned} & \frac{\lambda_h^{N_t} \mathbb{E}[\|x_0^t - x_\star\|^2]}{2hS_t} + \frac{\lambda_h^{N_t-1}}{4S_t} (\mathbb{E} F(x_0^t) - F_\star) + \frac{1}{3} (\mathbb{E} F(\bar{x}_0^t) - F_\star) \\ & \geq \frac{\mathbb{E}[\|x_0^{t+1} - x_\star\|^2]}{2hS_t} + \frac{1}{4\lambda_h S_t} (\mathbb{E} F(x_0^{t+1}) - F_\star) + \frac{2}{3} (\mathbb{E} F(\bar{x}_0^{t+1}) - F_\star). \end{aligned}$$

We now divide the proof up into two cases.

**Convex case.** In this case,  $\lambda_h = 1$ , so  $S_t = N_t$ . Here, we set  $N_{t+1} = 2N_t$ , which leads to

$$\begin{aligned} & \frac{\mathbb{E}[\|x_0^t - x_\star\|^2]}{2hN_t} + \frac{1}{4N_t} (\mathbb{E} F(x_0^t) - F_\star) + \frac{1}{3} (\mathbb{E} F(\bar{x}_0^t) - F_\star) \\ & \geq 2 \left[ \frac{\mathbb{E}[\|x_0^{t+1} - x_\star\|^2]}{2hN_{t+1}} + \frac{1}{4N_{t+1}} (\mathbb{E} F(x_0^{t+1}) - F_\star) + \frac{1}{3} (\mathbb{E} F(\bar{x}_0^{t+1}) - F_\star) \right]. \end{aligned}$$

The inequality clearly telescopes and shows that  $\mathbb{E} F(\bar{x}_0^T) - F_\star \leq \varepsilon$  after  $T$  epochs, where  $T \leq \log_2[O(F(x_0) - F_\star + \beta \|x_0 - x_\star\|^2)/\varepsilon]$  and  $h \asymp 1/\beta$ . The number of gradient evaluations is  $Tn + \sum_{t=0}^{T-1} N_t = Tn + 2^T$ , which yields the final result.

**Strongly convex case.** In this case, we set  $N_t = N$  for all  $t$ , where  $N$  is chosen so that  $\lambda_h^N \leq 1/2$ . With  $h \asymp 1/\beta$ , this leads to  $N \asymp \kappa$  and

$$\begin{aligned} & \frac{\mathbb{E}[\|x_0^t - x_\star\|^2]}{4hS} + \frac{1}{8\lambda_h S} (\mathbb{E} F(x_0^t) - F_\star) + \frac{1}{3} (\mathbb{E} F(\bar{x}_0^t) - F_\star) \\ & \geq 2 \left[ \frac{\mathbb{E}[\|x_0^{t+1} - x_\star\|^2]}{4hS} + \frac{1}{8\lambda_h S} (\mathbb{E} F(x_0^{t+1}) - F_\star) + \frac{1}{3} (\mathbb{E} F(\bar{x}_0^{t+1}) - F_\star) \right]. \end{aligned}$$

Again, this telescopes, and the computational cost is  $Tn + TN = O(T(n + \kappa))$ .  $\square$



The result of [Theorem 12.9](#) indeed improves upon the rates for [GD](#). Before presenting the final rate comparison, however, we note that the rates in [Theorem 12.9](#) are generally incomparable with the ones achieved via acceleration, i.e., for [AGD](#). One can ask whether acceleration can also be combined with variance reduction, and the answer is yes; we state a representative result from [\[Sha+18\]](#).

**Theorem 12.10 (accelerated SVRG).** Assume that each  $f_i$  is convex and  $\beta_i$ -smooth, and that  $g$  is  $\alpha$ -convex.<sup>†</sup> Then, there is an algorithm which achieves the following guarantees. Let  $\Delta_0 := F(x_0) - F_\star + \beta \|x_0 - x_\star\|^2$ .

- If  $\alpha = 0$ , then the algorithm obtains an  $\varepsilon$ -approximate solution with a total number of gradient evaluations at most  $O(n \log(\Delta_0/\varepsilon) + \sqrt{n\Delta_0/\varepsilon})$ .
- If  $\alpha > 0$ , then the algorithm obtains an  $\varepsilon$ -approximate solution with a total number of gradient evaluations at most  $O((n + \sqrt{n\kappa}) \log(\Delta_0/\varepsilon))$ .

<sup>†</sup>The cited paper works under slightly different assumptions compared to [Theorem 12.9](#), but they are broadly comparable.

These accelerated rates are almost the best possible due to nearly matching lower bounds [\[WS16\]](#). Interestingly, in this setting, randomness is crucial for attaining the optimal complexity; otherwise, among the class of deterministic algorithms, [AGD](#) is the best possible (but strictly worse than [Theorem 12.10](#)).

The rates for the finite sum setting are presented in [Table 3](#). Note that the rate for [SGD](#) is not exactly comparable since the assumption of uniformly bounded variance is not met here, but it is included for completeness.

Algorithm	Iterations (Convex)	Iterations (Strongly Convex)
<a href="#">SGD</a> <sup>†</sup>	$O(\sigma^2 d R^2 / \varepsilon^2)$	$O(\sigma^2 d / (\alpha \varepsilon))$
<a href="#">GD</a>	$O(n \Delta_0 / \varepsilon)$	$O(n \kappa \log(\Delta_0 / \varepsilon))$
<a href="#">AGD</a>	$O(n \sqrt{\Delta_0 / \varepsilon})$	$O(n \sqrt{\kappa} \log(\Delta_0 / \varepsilon))$
<a href="#">SVRG</a>	$O(n \log(\Delta_0 / \varepsilon) + \Delta_0 / \varepsilon)$	$O((n + \kappa) \log(\Delta_0 / \varepsilon))$
<a href="#">ASVRG</a>	$O(n \log(\Delta_0 / \varepsilon) + \sqrt{\Delta_0 / \varepsilon})$	$O((n + \sqrt{n\kappa}) \log(\Delta_0 / \varepsilon))$

Table 3: Rates for finite sum minimization. <sup>†</sup>The rates for [SGD](#) are not directly comparable since the stochastic gradients do not have uniformly bounded variance here.

## Bibliographical notes

For more discussion on the statistical performance of **SGD**, see [Bac24]. For an exposition to empirical process theory and statistics, see any standard reference, e.g., [Wai19].

The CLT for **ASGD** was first established in [PJ92]. The treatment of the martingale CLT follows [Bil95]. For an exposition to asymptotic statistics, see [Vaa98].

The proof of **Theorem 12.9** is inspired by [AY16], although care was taken to unify the convex and strongly convex proofs.

## Exercises

**Exercise 12.1.** Often, stochastic gradients do not have uniformly bounded variance. For example, suppose we have the objective function  $f : x \mapsto \frac{1}{2n} \sum_{i=1}^n \langle a_i, x \rangle^2$ , with stochastic gradient  $\hat{\nabla} f(x) = \langle a_i, x \rangle a_i$  with  $i \sim \text{uniform}([n])$ . Then, the variance of the stochastic gradient is

$$\mathbb{E}[\|\hat{\nabla} f(x) - \nabla f(x)\|^2] = \frac{1}{n} \sum_{i=1}^n \left\| \left( a_i a_i^\top - \frac{1}{n} \sum_{j=1}^n a_j a_j^\top \right) x \right\|^2,$$

which grows quadratically with  $\|x\|$ .

Assume therefore that  $f$  is  $\alpha$ -strongly convex and  $\beta$ -smooth with respect to the Euclidean norm, and that the following variance condition holds:

$$\mathbb{E}[\|\hat{\nabla} f(x) - \nabla f(x)\|^2] \leq c_0 + c_1 \|x - x_\star\|^2 \quad \text{for all } x \in \mathbb{R}^d.$$

Show that the iterates of stochastic gradient descent satisfy the following guarantee. If  $\varepsilon$  is sufficiently small and the step size  $h$  is chosen appropriately, then  $\mathbb{E} f(\bar{x}_N) - f_\star \leq \varepsilon$  for a suitably averaged iterate  $\bar{x}_N$  and all

$$N \gtrsim \frac{c_0}{\alpha \varepsilon} \log \frac{\alpha \|x_0 - x_\star\|^2}{\varepsilon}.$$

**Exercise 12.2.** Consider linear regression with fixed design: our dataset is  $\{(X_i, Y_i)\}_{i \in [n]}$ , where the covariates  $X_i$  are deterministic and fixed, and the  $Y_i$  are independent with

$$Y_i = \langle \theta^\star, X_i \rangle + \xi_i, \quad \xi_i \sim \text{normal}(0, \sigma^2 I).$$

The empirical and population risks are

$$\mathcal{R}_n(\theta) := \frac{1}{n} \sum_{i=1}^n (Y_i - \langle \theta, X_i \rangle)^2, \quad \mathcal{R}(\theta) := \mathbb{E} \mathcal{R}_n(\theta).$$

1. Show that the population risk is given by  $\mathcal{R}(\theta) = \sigma^2 + \|X(\theta - \theta^\star)\|^2/n$ , where  $X \in \mathbb{R}^{n \times d}$  is the matrix whose rows are  $\{X_i^\top\}_{i \in [n]}$ .
2. Show that the ERM of minimal norm is the least-squares estimator  $\widehat{\theta}_n = (X^\top X)^\dagger X^\top Y$ , where  $^\dagger$  denotes the Moore–Penrose pseudoinverse. Show that the excess risk of the ERM is given by

$$\mathbb{E} \mathcal{R}(\widehat{\theta}_n) - \mathcal{R}(\theta^\star) = \frac{\sigma^2 \text{rank } X}{n}.$$

3. Consider the iterates of GD on the empirical risk  $\mathcal{R}_n$ . Show that for a step size  $h$  sufficiently small, it holds that

$$\mathbb{E} \mathcal{R}(\theta_N) - \mathcal{R}(\theta^\star) \leq \frac{\sigma^2 \text{rank } X}{n} + O\left(\frac{\|\theta_0 - \theta^\star\|^2}{Nh}\right).$$

*Hints:* For all of these parts, make extensive use of the singular value decomposition of  $X$ . For the third part, write an exact recursion for  $\theta_k - \theta^\star$ , iterate this recursion, and then compute the excess risk. Use the fact that  $\max_{x \in [0,1]} x(1-x)^N \lesssim 1/N$  for  $N \geq 1$ .

## References

- [Aga+12] A. Agarwal, P. L. Bartlett, P. Ravikumar, and M. J. Wainwright. “Information-theoretic lower bounds on the oracle complexity of stochastic convex optimization”. In: *IEEE Trans. Inform. Theory* 58.5 (2012), pp. 3235–3249.
- [AKL22] P.-C. Aubin-Frankowski, A. Korba, and F. Léger. “Mirror descent with relative smoothness in measure spaces, with application to Sinkhorn and EM”. In: *Advances in Neural Information Processing Systems*. Ed. by S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh. Vol. 35. Curran Associates, Inc., 2022, pp. 17263–17275.
- [AL24] M. Arnese and D. Lacker. “Convergence of coordinate ascent variational inference for log-concave measures via optimal transport”. In: *arXiv preprint 2404.08792* (2024).
- [ALZ24] F. Ascolani, H. Lavenant, and G. Zanella. “Entropy contraction of the Gibbs sampler under log-concavity”. In: *arXiv preprint 2410.00858* (2024).
- [AN00] S.-i. Amari and H. Nagaoka. *Methods of information geometry*. Vol. 191. Translations of Mathematical Monographs. Translated from the 1993 Japanese original by Daishi Harada. American Mathematical Society, Providence, RI; Oxford University Press, Oxford, 2000, pp. x+206.

- [ANR17] J. M. Altschuler, J. Niles-Weed, and P. Rigollet. “Near-linear time approximation algorithms for optimal transport via Sinkhorn iteration”. In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett. Curran Associates, Inc., 2017, pp. 1964–1974.
- [AP24a] J. M. Altschuler and P. A. Parrilo. “Acceleration by stepsize hedging: multi-step descent and the silver stepsize schedule”. In: *J. ACM* (Dec. 2024).
- [AP24b] J. M. Altschuler and P. A. Parrilo. “Acceleration by stepsize hedging: silver stepsize schedule for smooth convex optimization”. In: *Mathematical Programming* (2024).
- [AY16] Z. Allen-Zhu and Y. Yuan. “Improved SVRG for non-strongly-convex or sum-of-non-convex objectives”. In: *Proceedings of the 33rd International Conference on Machine Learning*. Ed. by M. F. Balcan and K. Q. Weinberger. Vol. 48. Proceedings of Machine Learning Research. New York, New York, USA: PMLR, June 2016, pp. 1080–1089.
- [Bac+92] F. L. Baccelli, G. Cohen, G. J. Olsder, and J.-P. Quadrat. *Synchronization and linearity*. Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. An algebra for discrete event systems. John Wiley & Sons, Ltd., Chichester, 1992, pp. xx+489.
- [Bac24] F. Bach. *Learning theory from first principles*. MIT Press, 2024.
- [Bar93] A. R. Barron. “Universal approximation bounds for superpositions of a sigmoidal function”. In: *IEEE Trans. Inform. Theory* 39.3 (1993), pp. 930–945.
- [BBT17] H. H. Bauschke, J. Bolte, and M. Teboulle. “A descent lemma beyond Lipschitz gradient continuity: first-order methods revisited and applications”. In: *Math. Oper. Res.* 42.2 (2017), pp. 330–348.
- [BC12] S. Bubeck and N. Cesa-Bianchi. “Regret analysis of stochastic and nonstochastic multi-armed bandit problems”. In: *Foundations and Trends® in Machine Learning* 5.1 (2012), pp. 1–122.
- [Bil95] P. Billingsley. *Probability and measure*. Third. Wiley Series in Probability and Mathematical Statistics. A Wiley-Interscience Publication. John Wiley & Sons, Inc., New York, 1995, pp. xiv+593.
- [BM03] S. Burer and R. D. C. Monteiro. “A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization”. In: vol. 95. 2. Computational semidefinite and second order cone programming: the state of the art. 2003, pp. 329–357.

- [BM05] S. Burer and R. D. C. Monteiro. “Local minima and convergence in low-rank semidefinite programming”. In: *Math. Program.* 103.3 (2005), pp. 427–444.
- [Bub15] S. Bubeck. “Convex optimization: algorithms and complexity”. In: *Foundations and Trends® in Machine Learning* 8.3-4 (2015), pp. 231–357.
- [Che+22] Y. Chen, S. Chewi, A. Salim, and A. Wibisono. “Improved analysis for a proximal algorithm for sampling”. In: *Proceedings of Thirty Fifth Conference on Learning Theory*. Ed. by P.-L. Loh and M. Raginsky. Vol. 178. Proceedings of Machine Learning Research. PMLR, July 2022, pp. 2984–3014.
- [Che25] S. Chewi. *Log-concave sampling*. Available online at [chewisinho.github.io](https://chewisinho.github.io). Forthcoming, 2025.
- [CJ24] R. Caprio and A. M. Johansen. “Fast convergence of the expectation maximization algorithm under a logarithmic Sobolev inequality”. In: *arXiv preprint 2407.17949* (2024).
- [CLM24] H. Chardon, M. Lerasle, and J. Mourtada. “Finite-sample performance of the maximum likelihood estimator in logistic regression”. In: *arXiv preprint 2411.02137* (2024).
- [CNR25] S. Chewi, J. Niles-Weed, and P. Rigollet. *Statistical optimal transport*. Lecture Notes in Mathematics. École d’Été de Probabilités de Saint-Flour XLIX—2019. Springer Cham, 2025, pp. xiv+260.
- [Cut13] M. Cuturi. “Sinkhorn distances: lightspeed computation of optimal transport”. In: *Advances in Neural Information Processing Systems*. Ed. by C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger. Vol. 26. Curran Associates, Inc., 2013.
- [Eva10] L. C. Evans. *Partial differential equations*. Second. Vol. 19. Graduate Studies in Mathematics. American Mathematical Society, Providence, RI, 2010, pp. xxii+749.
- [GWS21] S. Gunasekar, B. Woodworth, and N. Srebro. “Mirrorless mirror descent: a natural derivation of mirror descent”. In: *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics*. Ed. by A. Banerjee and K. Fukumizu. Vol. 130. Proceedings of Machine Learning Research. PMLR, Apr. 2021, pp. 2305–2313.
- [JZ13] R. Johnson and T. Zhang. “Accelerating stochastic gradient descent using predictive variance reduction”. In: *Advances in Neural Information Processing Systems*. Ed. by C. J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger. Vol. 26. Curran Associates, Inc., 2013.

- [KNS16] H. Karimi, J. Nutini, and M. Schmidt. “Linear convergence of gradient and proximal-gradient methods under the Polyak–Łojasiewicz condition”. In: *European Conference on Machine Learning and Knowledge Discovery in Databases—Volume 9851*. ECML PKDD 2016. Riva del Garda, Italy: Springer-Verlag, 2016, pp. 795–811.
- [Lég21] F. Léger. “A gradient descent perspective on Sinkhorn”. In: *Appl. Math. Optim.* 84.2 (2021), pp. 1843–1855.
- [LFN18] H. Lu, R. M. Freund, and Y. Nesterov. “Relatively smooth convex optimization by first-order methods, and applications”. In: *SIAM J. Optim.* 28.1 (2018), pp. 333–354.
- [LMW24] J. Liang, S. Mitra, and A. Wibisono. “On independent samples along the Langevin diffusion and the unadjusted Langevin algorithm”. In: *arXiv preprint 2402.17067* (2024).
- [Łoj63] S. Łojasiewicz. “Une propriété topologique des sous-ensembles analytiques réels”. In: *Les Équations aux Dérivées Partielles (Paris, 1962)*. Éditions du Centre National de la Recherche Scientifique (CNRS), Paris, 1963, pp. 87–89.
- [LRP16] L. Lessard, B. Recht, and A. Packard. “Analysis and design of optimization algorithms via integral quadratic constraints”. In: *SIAM J. Optim.* 26.1 (2016), pp. 57–95.
- [LZ24] H. Lavenant and G. Zanella. “Convergence rate of random scan coordinate ascent variational inference under log-concavity”. In: *SIAM J. Optim.* 34.4 (2024), pp. 3750–3761.
- [Nes18] Y. Nesterov. *Lectures on convex optimization*. Vol. 137. Springer Optimization and Its Applications. Springer, 2018, pp. xxiii+589.
- [NY83] A. S. Nemirovski and D. B. Yudin. *Problem complexity and method efficiency in optimization*. Wiley-Interscience Series in Discrete Mathematics. Translated from the Russian and with a preface by E. R. Dawson. John Wiley & Sons, Inc., New York, 1983, pp. xv+388.
- [OV00] F. Otto and C. Villani. “Generalization of an inequality by Talagrand and links with the logarithmic Sobolev inequality”. In: *J. Funct. Anal.* 173.2 (2000), pp. 361–400.
- [OV01] F. Otto and C. Villani. “Comment on: “Hypercontractivity of Hamilton–Jacobi equations” [J. Math. Pures Appl. (9) 80 (2001), no. 7, 669–696] by S. G. Bobkov, I. Gentil and M. Ledoux”. In: *J. Math. Pures Appl. (9)* 80.7 (2001), pp. 697–700.

- [PC19] G. Peyré and M. Cuturi. *Computational optimal transport: with applications to data science*. Now, 2019.
- [PJ92] B. T. Polyak and A. B. Juditsky. “Acceleration of stochastic approximation by averaging”. In: *SIAM J. Control Optim.* 30.4 (1992), pp. 838–855.
- [Pol63] B. T. Polyak. “Gradient methods for minimizing functionals”. In: *Ž. Vychisl. Mat i Mat. Fiz.* 3 (1963), pp. 643–653.
- [Roc97] R. T. Rockafellar. *Convex analysis*. Princeton Landmarks in Mathematics. Reprint of the 1970 original, Princeton Paperbacks. Princeton University Press, Princeton, NJ, 1997, pp. xviii+451.
- [San15] F. Santambrogio. *Optimal transport for applied mathematicians*. Vol. 87. Progress in Nonlinear Differential Equations and their Applications. Calculus of variations, PDEs, and modeling. Birkhäuser/Springer, Cham, 2015, pp. xxvii+353.
- [SBC16] W. Su, S. Boyd, and E. J. Candès. “A differential equation for modeling Nesterov’s accelerated gradient method: theory and insights”. In: *J. Mach. Learn. Res.* 17 (2016), Paper No. 153, 43.
- [Sha+18] F. Shang, L. Jiao, K. Zhou, J. Cheng, Y. Ren, and Y. Jin. “ASVRG: accelerated proximal SVRG”. In: *Proceedings of the 10th Asian Conference on Machine Learning*. Ed. by J. Zhu and I. Takeuchi. Vol. 95. Proceedings of Machine Learning Research. PMLR, Nov. 2018, pp. 815–830.
- [Sin64] R. Sinkhorn. “A relationship between arbitrary positive matrices and doubly stochastic matrices”. In: *Ann. Math. Statist.* 35 (1964), pp. 876–879.
- [Vaa98] A. W. van der Vaart. *Asymptotic statistics*. Vol. 3. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, Cambridge, 1998, pp. xvi+443.
- [Ver18] R. Vershynin. *High-dimensional probability*. Vol. 47. Cambridge Series in Statistical and Probabilistic Mathematics. An introduction with applications in data science, With a foreword by Sara van de Geer. Cambridge University Press, Cambridge, 2018, pp. xiv+284.
- [Vil03] C. Villani. *Topics in optimal transportation*. Vol. 58. Graduate Studies in Mathematics. American Mathematical Society, Providence, RI, 2003, pp. xvi+370.
- [Vil09] C. Villani. *Optimal transport*. Vol. 338. Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]. Old and new. Springer-Verlag, Berlin, 2009, pp. xxii+973.
- [Vis12] N. K. Vishnoi. “ $Lx = b$  Laplacian solvers and their algorithmic applications”. In: *Found. Trends Theor. Comput. Sci.* 8.1-2 (2012), front matter, 1–141.

- [Wai19] M. J. Wainwright. *High-dimensional statistics*. Vol. 48. Cambridge Series in Statistical and Probabilistic Mathematics. A non-asymptotic viewpoint. Cambridge University Press, Cambridge, 2019, pp. xvii+552.
- [WS16] B. E. Woodworth and N. Srebro. “Tight complexity bounds for optimizing composite objectives”. In: *Advances in Neural Information Processing Systems*. Ed. by D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett. Vol. 29. Curran Associates, Inc., 2016.
- [WWJ16] A. Wibisono, A. C. Wilson, and M. I. Jordan. “A variational perspective on accelerated methods in optimization”. In: *Proc. Natl. Acad. Sci. USA* 113.47 (2016), E7351–E7358.