

# From Low Probability to High Confidence in Stochastic Convex Optimization

**Damek Davis**

*School of Operations Research and Information Engineering  
Cornell University  
Ithaca, NY 14850, USA*

DSD95@CORNELL.EDU

**Dmitriy Drusvyatskiy**

*Department of Mathematics, University of Washington  
Seattle, WA 98195, USA*

DDRUSV@UW.EDU

**Lin Xiao**

*Facebook AI Research  
Seattle, WA 98019, USA*

LINX@FB.COM

**Junyu Zhang**

*Department of Electrical Engineering, Princeton University  
Princeton, NJ 08544, USA*

JUNYUZ@PRINCETON.EDU

**Editor:** Animashree Anandkumar

## Abstract

Standard results in stochastic convex optimization bound the number of samples that an algorithm needs to generate a point with small function value in expectation. More nuanced *high probability* guarantees are rare, and typically either rely on “light-tail” noise assumptions or exhibit worse sample complexity. In this work, we show that a wide class of stochastic optimization algorithms for strongly convex problems can be augmented with high confidence bounds at an overhead cost that is only logarithmic in the confidence level and polylogarithmic in the condition number. The procedure we propose, called **proxBoost**, is elementary and builds on two well-known ingredients: robust distance estimation and the proximal point method. We discuss consequences for both streaming (online) algorithms and offline algorithms based on empirical risk minimization.

**Keywords:** proximal point method, robust distance estimation, stochastic approximation, empirical risk minimization, composite optimization

## 1. Introduction

Stochastic convex optimization lies at the core of modern statistical and machine learning. Standard results in the subject bound the number of samples that an algorithm needs to generate a point with small function value in *expectation*. Specifically, consider the problem

$$\min_x f(x) := \mathbb{E}_{z \sim \mathcal{P}}[f(x, z)], \quad (1)$$

where the random variable  $z$  follows a fixed unknown distribution  $\mathcal{P}$  and  $f(\cdot, z)$  is convex for almost every  $z \sim \mathcal{P}$ . Given a small tolerance  $\epsilon > 0$ , stochastic gradient methods typically

produce a point  $x_\epsilon$  satisfying

$$\mathbb{E}[f(x_\epsilon)] - \min f \leq \epsilon.$$

The cost of the algorithms, measured by the required number of stochastic (sub-)gradient evaluations, is  $\mathcal{O}(1/\epsilon^2)$  or  $\mathcal{O}(1/\epsilon)$  if  $f$  is strongly convex (see, e.g., Nemirovsky and Yudin, 1983; Polyak and Juditsky, 1992; Ghadimi and Lan, 2013; Hazan and Kale, 2014).

In this paper, we are interested in procedures that can produce an approximate solution with *high probability*, meaning a point  $x_{\epsilon,p}$  satisfying

$$\mathbb{P}(f(x_{\epsilon,p}) - \min f \leq \epsilon) \geq 1 - p, \tag{2}$$

where  $p > 0$  can be arbitrarily small. By Markov’s inequality, one can guarantee (2) by generating a point  $x_{\epsilon,p}$  satisfying  $\mathbb{E}[f(x_{\epsilon,p})] - f^* \leq p\epsilon$ , e.g., by using standard stochastic gradient methods. However, the resulting sample complexity can be very high for small  $p$  with the typical scaling of  $\mathcal{O}(1/(p\epsilon))$  or  $\mathcal{O}(1/(p\epsilon)^2)$ . Existing literature does provide a path to reducing the dependence of the sample complexity on  $p$  to  $\log(1/p)$ , but this usually comes with cost of either worse dependence on  $\epsilon$  (e.g., Bousquet and Elisseeff, 2002; Nesterov and Vial, 2008; Shalev-Shwartz et al., 2009) or more restrictive sub-Gaussian assumptions on the stochastic gradient noise (e.g., Nemirovski et al., 2008; Juditsky and Nesterov, 2014; Ghadimi and Lan, 2012, 2013; Harvey et al., 2019a,b).

**The goal of this work.** We aim to develop *generic* low-cost procedures that equip stochastic optimization algorithms with high confidence guarantees, without making restrictive noise assumptions. Consequently, it will be convenient to treat such algorithms as black boxes. More formally, suppose that the function  $f$  may only be accessed through a *minimization oracle*  $\mathcal{M}(f, \epsilon)$ , which on input  $\epsilon > 0$ , returns a point  $x_\epsilon$  satisfying the low confidence bound

$$\mathbb{P}(f(x_\epsilon) - \min f \leq \epsilon) \geq \frac{2}{3}. \tag{3}$$

By Markov’s inequality, minimization oracles arise from any algorithm that can generate  $x_\epsilon$  satisfying  $\mathbb{E}f(x_\epsilon) - \min f \leq \epsilon/3$ . Let  $\mathcal{C}_{\mathcal{M}}(f, \epsilon)$  denote the cost of the oracle call  $\mathcal{M}(f, \epsilon)$ . Given a minimization oracle and its cost, we investigate the following question:

Is there a procedure within this oracle model of computation that returns a point  $x_{\epsilon,p}$  satisfying the high confidence bound (2) at a total cost that is only a “small” multiple of  $\mathcal{C}_{\mathcal{M}}(f, \epsilon) \cdot \log(\frac{1}{p})$ ?

We will see that when  $f$  is strongly convex, the answer is yes for a wide class of oracles  $\mathcal{M}(f, \epsilon)$ . To simplify discussion, suppose  $f$  is  $\mu$ -strongly convex and  $L$ -smooth (differentiable with  $L$ -Lipschitz continuous gradient). Then the cost  $\mathcal{C}_{\mathcal{M}}(f, \epsilon)$  typically depends on the condition number  $\kappa := L/\mu \gg 1$ , as well as scale sensitive quantities such as initialization quality and upper bound on the gradient variances, etc. The procedures introduced in this paper execute the minimization oracle multiple times in order to boost its confidence, with the total cost on the order of

$$\log\left(\frac{\log(\kappa)}{p}\right) \log(\kappa) \cdot \mathcal{C}_{\mathcal{M}}\left(f, \frac{\epsilon}{\log(\kappa)}\right).$$

Thus, high probability bounds are achieved with a small cost increase, which depends only logarithmically on  $1/p$  and polylogarithmically on the condition number  $\kappa$ .

**Known techniques and limitations.** Before introducing our approach, we discuss two techniques for boosting the confidence of a minimization oracle, both of which have limitations. As a first approach, one may query the oracle  $\mathcal{M}(f, \epsilon)$  multiple times and pick the “best” iterate from the batch. This is a flawed strategy since often one cannot test which iterate is “best” without increasing sample complexity. To illustrate, consider estimating the expectation  $f(x) = \mathbb{E}_z[f(x, z)]$  to  $\epsilon$ -accuracy for a fixed point  $x$ . This task amounts to mean estimation, which requires on the order of  $1/\epsilon^2$  samples, even under sub-Gaussian assumptions (Catoni, 2012). In this paper, the cost  $\mathcal{C}_{\mathcal{M}}(f, \epsilon)$  typically scales at worst as  $1/\epsilon$ , and therefore mean estimation would significantly degrade the overall sample complexity.

The second approach leverages the fact that, with strong convexity, (3) implies

$$\mathbb{P}(\|x_\epsilon - \bar{x}\| \leq \sqrt{2\epsilon/\mu}) \geq \frac{2}{3},$$

where  $\bar{x}$  is the minimizer of  $f$ . Given this bound, one may apply the *robust distance estimation* technique of Nemirovsky and Yudin (1983, p. 243) and Hsu and Sabato (2016) to choose a point near  $\bar{x}$ : Run  $m$  trials of  $\mathcal{M}(f, \epsilon)$  and find one iterate  $x_{i^*}$  around which the other points “cluster”. Then the point  $x_{i^*}$  will be within a distance of  $\sqrt{18\epsilon/\mu}$  from  $\bar{x}$  with probability  $1 - \exp(-m/18)$ . Similar guarantees also hold for the geometric median estimator of Minsker (2015). The downside of this strategy is that when converting naively back to function values, the suboptimality gap becomes  $f(x_{i^*}) - \min f \leq \frac{L}{2}\|x_{i^*} - \bar{x}\|^2 \leq 9\kappa\epsilon$ . Thus the function gap at  $x_{i^*}$  may be significantly larger than the expected function gap at  $x_\epsilon$ , by a factor of the condition number. Therefore, robust distance estimation exhibits a trade-off between robustness and efficiency.

## 1.1 Contribution: The proxBoost algorithm

The procedure we introduce, called **proxBoost**, is based on the following simple observation: although robust distance estimation induces a trade-off between robustness and efficiency, this trade-off disappears for smooth losses that are perfectly conditioned. Leveraging this fact, we design a continuation procedure that links together a short sequence of robust distance estimators for nearby problems with rapidly improving condition numbers. More specifically, the **proxBoost** algorithm generates a sequence of iterates  $x_0, \dots, x_T$ . The first iterate,  $x_0$ , is simply the output of the robust distance estimator for minimizing  $f$ . Subsequently, given an iterate  $x_t$ , the procedure forms the better conditioned function  $f^t(x) := f(x) + \frac{\mu^2}{2}\|x - x_t\|^2$  and declares the next iterate  $x_{t+1}$  to be the output of the robust distance estimator for minimizing  $f^t$ . One may in principle apply **proxBoost** with any minimization oracle  $\mathcal{M}(f^t, \epsilon)$ . The real benefit arises for concrete oracles, such as those based on streaming algorithms (e.g., stochastic gradient) or offline methods (e.g., empirical risk minimization), for which the cost of computing the robust distance estimator rapidly decreases as  $t$  increases and conditioning improves. When used within the **proxBoost** method, these oracles benefit from new high confidence guarantees with only a modest logarithmic and polylogarithmic cost increase in  $1/p$  and  $\kappa$ , respectively. We now illustrate this claim.

### 1.1.1 STREAMING ORACLES

Stochastic gradient methods may serve as minimization oracles  $\mathcal{M}(f, \epsilon)$  within **proxBoost**. For these oracles, the cost  $\mathcal{C}_{\mathcal{M}}(f, \epsilon)$  is measured by the number of stochastic gradient esti-

mates that the algorithm must generate in order to reach functional accuracy  $\epsilon$  in expectation. Although many such oracles exist and may be used within `proxBoost`, our goal is to use the *optimal* algorithm of Ghadimi and Lan (2013) as an oracle and equip it with high confidence guarantees. This algorithm is optimal in the sense that it has minimal cost among stochastic gradient methods within a standard oracle model of computation. More specifically, the method generates a point  $x_\epsilon$  satisfying  $\mathbb{E}[f(x_\epsilon) - \min f] \leq \epsilon$  with

$$\mathcal{O}\left(\sqrt{\kappa} \ln\left(\frac{\Delta_{\text{in}}}{\epsilon}\right) + \frac{\sigma^2}{\mu\epsilon}\right) \quad (4)$$

stochastic gradient evaluations, where the quantity  $\sigma^2$  is an upper bound on the variance of the stochastic gradient estimator  $\nabla f(x, z)$  and  $\Delta_{\text{in}}$  is a known upper bound on the initial function gap  $\Delta_{\text{in}} \geq f(x_0) - f^*$ . A simpler algorithm with a similar efficiency estimate was recently presented by Kulunchakov and Mairal (2019), and was based on estimate sequences. Aybat et al. (2019) developed an algorithm with similar efficiency, but in contrast to previous work, it does not require the variance  $\sigma^2$  and the initial gap  $\Delta_{\text{in}}$  as inputs.

It is intriguing to ask if one can equip the stochastic gradient method and its accelerated variant with high confidence guarantees. In their original work, Ghadimi and Lan (2013, 2012) provide an affirmative answer under the additional assumption that the stochastic gradient estimator has light tails. The very recent work of Juditsky et al. (2019) shows that one can avoid the light tail assumption for the basic stochastic gradient method, and for mirror descent more generally, by truncating the gradient estimators. High confidence bounds for the accelerated method, without light tail assumptions, remain open.

In this work, the optimal method of Ghadimi and Lan (2013) will be used as a minimization oracle within `proxBoost`, allowing us to nearly match the efficiency estimate (4) without “light-tail” assumptions. Equipped with this oracle, `proxBoost` returns a point  $x_{\epsilon,p}$  satisfying (2) and the overall cost of the procedure is

$$\tilde{\mathcal{O}}\left(\log\left(\frac{1}{p}\right)\left(\sqrt{\kappa} \ln\left(\frac{\Delta_{\text{in}}}{\epsilon} \vee \kappa\right) + \frac{\sigma^2}{\mu\epsilon}\right)\right).$$

Here,  $\tilde{\mathcal{O}}(\cdot)$  only suppresses logarithmic dependencies in  $\kappa$ ; see Section 5 for a precise guarantee. Thus for small  $\epsilon$ , the sample complexity of the robust procedure is roughly  $\log(1/p)$  times the efficiency estimate (4) of the low-confidence algorithm. Subsequent to this work, Gorbunov et al. (2020) develop an accelerated clipped gradient method whose efficiency is superior by a logarithmic factor in the condition number.

It is worthwhile to note that `proxBoost` seeded with the stochastic gradient method resembles the weight decay schedule, which is commonly used in practice; see e.g. Ge et al. (2019); Yang et al. (2018). The procedure is also related to, but distinct from, the SGD3 algorithm of Allen-Zhu (2018), which rapidly drives the gradient of the objective function to zero in expectation.

### 1.1.2 EMPIRICAL RISK MINIMIZATION ORACLES

An alternative approach to streaming algorithms, such as the stochastic gradient method, is based on empirical risk minimization (ERM) or sample average approximation (SAA)

(Shapiro and Nemirovski, 2005). Namely, we draw i.i.d. samples  $z_1, \dots, z_n \sim \mathcal{P}$  and minimize the empirical average

$$\min_x f_S(x) := \frac{1}{n} \sum_{i=1}^n f(x, z_i). \quad (5)$$

A key question is to determine the number  $n$  of samples that would ensure that the minimizer  $x_S$  of the empirical risk  $f_S$  has low generalization error  $f(x_S) - \min f$ , with high probability. There is a vast literature on this subject; see for example Hsu and Sabato (2016); Bartlett and Mendelson (2002); Shalev-Shwartz et al. (2009); Shalev-Shwartz and Ben-David (2014). Minimax optimal guarantees are available in some special cases. Notably, for ridge regression, a minimax estimator that succeeds with high probability appears in Audibert et al. (2011), although it is unknown how to implement the estimator efficiently.

We build here on the work of Hsu and Sabato (2016), who focused on high confidence guarantees for nonnegative losses  $f(x, z)$ . They showed that the empirical risk minimizer  $x_S$  yields a robust distance estimator of the true minimizer of  $f$ . As a consequence they deduced that ERM can find a point  $x_S$  satisfying the relative error guarantee

$$\mathbb{P}[f(x_S) \leq (1 + \gamma)f^*] \geq 1 - p,$$

with the sample complexity  $n$  on the order of

$$\mathcal{O} \left( \log \left( \frac{1}{p} \right) \cdot \frac{\hat{\kappa} \kappa}{\gamma} \right).$$

Loosely speaking, here  $\kappa$  and  $\hat{\kappa}$  are the condition numbers of  $f$  and  $f_S$ , respectively. It is unknown whether this sample complexity can be improved, but the appearance of a squared “condition number” makes the sample complexity cost of ERM much larger than that of streaming algorithms. In this work, we provide such an improvement by embedding ERM within proxBoost, yielding an order of magnitude better complexity

$$\tilde{\mathcal{O}} \left( \log \left( \frac{1}{p} \right) \left( \frac{\hat{\kappa}}{\gamma} + \hat{\kappa} \right) \right),$$

where the symbol  $\tilde{\mathcal{O}}$  only suppresses polylogarithmic dependence on  $\kappa$  and  $\hat{\kappa}$ . See Section 4 for the precise sample complexity guarantee.

### 1.1.3 CONVEX COMPOSITE OPTIMIZATION

The results, previewed so far rely on the assumption that  $f$  is strongly convex and smooth. These techniques can not directly accommodate constraints or nonsmooth regularizers. To illustrate the difficulty, consider the convex composite optimization problem

$$\min_x f(x) = g(x) + h(x), \quad (6)$$

where  $g: \mathbf{R}^d \rightarrow \mathbf{R}$  is smooth and strongly convex and  $h: \mathbf{R}^d \rightarrow \mathbf{R} \cup \{+\infty\}$  is an arbitrary closed convex function. For example, a constrained optimization problem can be modeled by setting  $h$  to be zero on the feasible region and plus infinity elsewhere. The approach for

the unconstrained problems, outlined previously, heavily relies on the fact that the function gap  $f(x) - \min f$  and the squared distance to the solution  $\|x - \bar{x}\|^2$  are proportional up to multiplication by the condition number. The analogous statement for the composite setting (6) is decisively false. In particular, it is unclear how to turn low probability guarantees on the function gap  $f(x) - \min f$  to high probability outcomes, even if one was willing to degrade the accuracy by the condition number of  $g$ .

In the last section of the paper, we resolve this apparent difficulty and thereby generalize the `proxBoost` framework to the entire composite problem class (6). The key tool is a new robust distance estimation technique for convex composite problems, which may be of independent interest. Consequences for regularized empirical risk minimization and proximal streaming algorithms, in the spirit of Sections 1.1.1 and 1.1.2, follow immediately.

## 1.2 Related literature

Our paper rests on two pillars: the proximal point method and robust distance estimation. The two techniques have been well studied in the optimization and statistics literature respectively. The proximal point method was introduced by Martinet (1972, 1970) and further popularized by Rockafellar (1976). This construction is also closely related to the smoothing function of Moreau (1965). Recently, there has been a renewed interest in the proximal point method, most notably due to its uses in accelerating variance-reduction methods for minimizing finite sums of convex functions (Lin et al., 2015; Frostig et al., 2015; Lan and Zhou, 2018; Shalev-Shwartz and Zhang, 2016). The proximal point method has also featured prominently as a guiding principle in nonconvex optimization, with the works of Asi and Duchi (2019, 2018); Duchi and Ruan (2018); Davis and Drusvyatskiy (2019); Davis and Grimmer (2017). The stepsize schedule we use within the proximal point method is geometrically decaying, in contrast to the more conventional polynomially decaying schemes. Geometrically decaying schedules for subgradient methods were first used by Goffin (1977) and have regained some attention recently due to their close connection to the popular step-decay schedule in stochastic optimization (Ge et al., 2019; Aybat et al., 2019; Xu et al., 2016; Yang et al., 2018).

Robust distance estimation has a long history. The estimator we use was first introduced by Nemirovsky and Yudin (1983, p. 243), and can be viewed as a multivariate generalization of the median of means estimator (Alon et al., 1999; Jerrum et al., 1986). Robust distance estimation was further investigated by Hsu and Sabato (2016) with a focus on high probability guarantees for empirical risk minimization. A different generalization based on the geometric median was studied by Minsker (2015). Other recent articles related to the subject include median of means tournaments (Lugosi and Mendelson, 2016), robust multivariate mean estimators (Joly et al., 2017; Lugosi and Mendelson, 2019), and bandits with heavy tails (Bubeck et al., 2013).

One of the main applications of our techniques is to streaming algorithms. Most currently available results that establish high confidence convergence guarantees make sub-Gaussian assumptions on the stochastic gradient estimator (Nemirovski et al., 2008; Juditsky and Nesterov, 2014; Ghadimi and Lan, 2012, 2013; Harvey et al., 2019a,b). More recently, there has been renewed interest in obtaining robust guarantees without the light-tails assumption. For example, Chen et al. (2017) and Yin et al. (2018) make use of the

geometric median of means technique to robustly estimate the gradient in distributed optimization. A different technique was recently developed by Juditsky et al. (2019), where the authors establish high confidence guarantees for mirror descent type algorithms by truncating the gradient.

A preliminary version of the current paper, written by the first two authors, appeared in the conference proceedings (Davis and Drusvyatskiy, 2020). The earlier version covered the material in Sections 1-5. The current paper extends the `proxBoost` algorithm to constrained and regularized settings, develops consequences for both streaming and offline algorithms, and develops a smoothing technique that enables application of `proxBoost` for nonsmooth problems.

The outline of the paper is as follows. Section 2 presents the problem setting and robust distance estimation. Section 3 develops the `proxBoost` procedure. Section 4 presents consequences for empirical risk minimization, while Section 5 discusses consequences for streaming algorithms, both in the strongly convex and smooth setting. The final Section 6 extends the aforementioned techniques to convex composite problems. It is worth mentioning that the approach we take here depends on the smoothness of the loss functions. For problems with nonsmooth losses, either smoothing (see Section 6.3) or some new techniques should be developed, which we leave for future work.

## 2. Preliminary: a robust distance estimator

Throughout, we follow standard notation of convex optimization, as set out for example in the monographs of Nesterov (2018) and Beck (2017). We let  $\mathbf{R}^d$  denote an Euclidean space with inner product  $\langle \cdot, \cdot \rangle$  and the induced norm  $\|x\| = \sqrt{\langle x, x \rangle}$ . The symbol  $B_\varepsilon(x)$  will stand for the closed ball around  $x$  of radius  $\varepsilon > 0$ . We will use the shorthand interval notation  $[1, m] := \{1, \dots, m\}$  for any number  $m \in \mathbb{N}$ . Abusing notation slightly, for any set of real numbers  $\{r_i\}_{i=1}^m$  we will let  $\text{median}(r_1, r_2, \dots, r_m)$  denote the  $\lceil \frac{m}{2} \rceil$ 'th entry in the ordered list  $r_{[1]} \leq r_{[2]} \leq \dots \leq r_{[m]}$ .

Consider a function  $f: \mathbf{R}^d \rightarrow \mathbf{R} \cup \{+\infty\}$ . The effective domain of  $f$ , denoted  $\text{dom } f$ , consists of all points where  $f$  is finite. The function  $f$  is called  $\mu$ -strongly convex if the perturbed function  $f - \frac{\mu}{2} \|\cdot\|^2$  is convex. We say that  $f$  is  $L$ -smooth if it is differentiable with  $L$ -Lipschitz continuous gradient. If  $f$  is both  $\mu$ -strongly convex and  $L$ -smooth, then standard results in convex optimization (e.g., Nesterov, 2018, § 2.1) imply for all  $x, y \in \mathbf{R}^d$  the bound

$$\langle \nabla f(y), x - y \rangle + \frac{\mu}{2} \|x - y\|^2 \leq f(x) - f(y) \leq \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2.$$

In particular, if  $\bar{x}$  is the minimizer of  $f$ , we have  $\nabla f(\bar{x}) = 0$  and thus the two-sided bound:

$$\frac{\mu}{2} \|x - \bar{x}\|^2 \leq f(x) - f(\bar{x}) \leq \frac{L}{2} \|x - \bar{x}\|^2 \quad \text{for all } x \in \mathbf{R}^d. \quad (7)$$

The ratio  $\kappa := L/\mu$  is called the condition number of  $f$ .

**Assumption 1** Throughout this work, we consider the optimization problem

$$\min_{x \in \mathbf{R}^d} f(x) \quad (8)$$

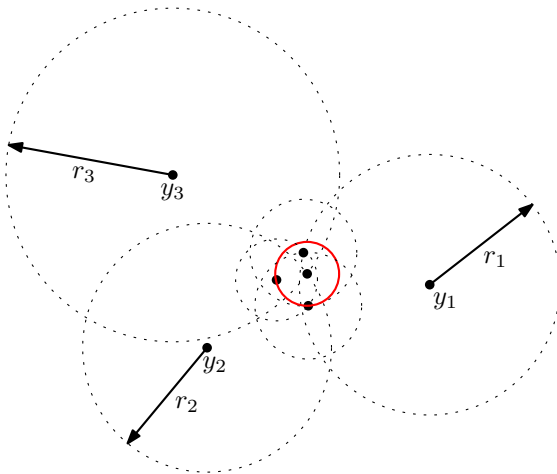


Figure 1: Illustration of the robust distance estimator  $\mathcal{D}(\varepsilon, m)$ .

where the function  $f: \mathbf{R}^d \rightarrow \mathbf{R} \cup \{+\infty\}$  is closed and  $\mu$ -strongly convex. We denote the minimizer of  $f$  by  $\bar{x}$  and its minimal value by  $f^* := \min f$ .

Let us suppose for the moment that the only access to  $f$  is by querying a black-box procedure that estimates  $\bar{x}$ . Namely following Hsu and Sabato (2016) we will call a procedure  $\mathcal{D}(\varepsilon)$  a *weak distance oracle* for the problem (8) if it returns a point  $x$  satisfying

$$\mathbb{P}[\|x - \bar{x}\| \leq \varepsilon] \geq \frac{2}{3}. \tag{9}$$

We will moreover assume that when querying  $\mathcal{D}(\varepsilon)$  multiple times, the returned vectors are all statistically independent. Weak distance oracles arise naturally in stochastic optimization both in streaming and offline settings. We will discuss specific examples in Sections 4 and 5. The numerical value  $2/3$  plays no real significance and can be replaced by any fraction greater than a half.

It is well known from Nemirovsky and Yudin (1983, p. 243) and Hsu and Sabato (2016) that the low-confidence estimate (9) can be improved to a high confidence guarantee by a clustering technique. Following Hsu and Sabato (2016), we define the *robust distance estimator*  $\mathcal{D}(\varepsilon, m)$  to be the following procedure (Algorithm 1).

<p><b>Algorithm 1:</b> Robust Distance Estimation <math>\mathcal{D}(\varepsilon, m)</math></p> <p><b>Input:</b> access to a weak distance oracle <math>\mathcal{D}(\varepsilon)</math> and trial count <math>m</math>.          Query <math>m</math> times the oracle <math>\mathcal{D}(\varepsilon)</math> and let <math>Y = \{y_1, \dots, y_m\}</math> consist of the responses.  <b>Step</b> <math>i = 1, \dots, m</math>:              Compute <math>r_i = \min\{r \geq 0 :  B_r(y_i) \cap Y  &gt; \frac{m}{2}\}</math>.          Set <math>i^* = \operatorname{argmin}_{i \in [1, m]} r_i</math>  <b>Return</b> <math>y_{i^*}</math></p>
---

Thus the robust distance estimator  $\mathcal{D}(\varepsilon, m)$  first generates  $m$  statistically independent random points  $y_1, \dots, y_m$  by querying  $m$  times the weak distance oracle  $\mathcal{D}(\varepsilon)$ . Then the



procedure computes the smallest radius ball around each point  $y_i$  that contains more than half of the generated points  $\{y_1, \dots, y_m\}$ . Finally, the point  $y_{i^*}$  corresponding to the smallest such ball is returned. See Figure 1 for an illustration.

The intuition underlying the algorithm is that by Chernoff’s bound, with high confidence, the ball  $B_\varepsilon(\bar{x})$  will contain strictly more than  $m/2$  of the generated points. Therefore in this event, the estimate  $r_{i^*} < 2\varepsilon$  holds. Moreover since the two sets,  $B_\varepsilon(\bar{x})$  and  $B_{r_{i^*}}(y_{i^*})$  intersect, it follows that  $\bar{x}$  and  $y_{i^*}$  are within a distance of  $3\varepsilon$  of each other. For a complete argument, see Nemirovsky and Yudin (1983, p. 243) or Hsu and Sabato (2016, Propositions 8,9).

**Lemma 1 (Robust Distance Estimator)** *The point  $x$  returned by  $\mathcal{D}(\varepsilon, m)$  satisfies*

$$\mathbb{P}(\|x - \bar{x}\| \leq 3\varepsilon) \geq 1 - \exp\left(-\frac{m}{18}\right).$$

It is worthwhile to mention that there are other estimation procedures, besides Algorithm 1, that equip weak distance oracles with the same guarantees as in Lemma 1. The most notable example is the geometric median of Minsker (2015). For the sake of simplicity, we focus only on Algorithm 1, though other robust distance estimators can be used instead.

We seek to understand how one may use a robust distance estimator  $\mathcal{D}(\varepsilon, m)$  to compute a point  $x$  satisfying  $f(x) - \min f \leq \delta$  with high probability, where  $\delta > 0$  is a specified accuracy. As motivation, consider the case when  $f$  is also  $L$ -smooth. Then one immediate approach is to appeal to the upper bound in (7). Hence by Lemma 1, the point  $x = \mathcal{D}(\varepsilon, m)$ , with  $\varepsilon = \sqrt{\frac{2\delta}{9L}}$ , satisfies the guarantee

$$\mathbb{P}(f(x) - f^* \leq \delta) \geq \mathbb{P}(\|x - \bar{x}\| \leq 3\varepsilon) \geq 1 - \exp\left(-\frac{m}{18}\right).$$

We will follow an alternative approach, which can significantly decrease the overall cost in the regime  $\kappa \gg 1$ . The optimistic goal is to replace the accuracy  $\varepsilon \approx \sqrt{\frac{\delta}{L}}$  used in the call to  $\mathcal{D}(\varepsilon, m)$  by the potentially much larger quantity  $\sqrt{\frac{\delta}{\mu}}$ . The strategy we propose will apply a robust distance estimator  $\mathcal{D}$  to a sequence of optimization problems that are better and better conditioned, thereby amortizing the overall cost. In the initial step, we will simply apply  $\mathcal{D}$  to  $f$  with the low accuracy  $\sqrt{\frac{\delta}{\mu}}$ . In step  $i$ , we will apply  $\mathcal{D}$  to a new function  $f^i$ , which has condition number  $\kappa_i \approx \frac{L+\mu 2^i}{\mu+\mu 2^i}$ , with accuracy  $\varepsilon_i \approx \sqrt{\frac{\delta}{\mu+\mu 2^i}}$ . Continuing this process for  $T \approx \log_2\left(\frac{L}{\mu}\right)$  rounds, we arrive at accuracy  $\varepsilon_T \approx \sqrt{\frac{\delta}{\mu+L}}$  and a function  $f^T$  that is nearly perfectly conditioned with  $\kappa_T \leq 2$ . In this way, the total cost is amortized over the sequence of optimization problems. The key of course is to control the error incurred by varying the optimization problems along the iterations.

### 3. The proxBoost method

The continuation procedure outlined at the end of the previous section can be succinctly described within the framework of an *inexact proximal point method*. Henceforth, fix an

increasing sequence of penalties  $\lambda_0, \dots, \lambda_T$  and a sequence of centers  $x_0, \dots, x_T$ . For each index  $i = 0, \dots, T$ , define the quadratically perturbed functions and their minimizers:

$$f^i(x) := f(x) + \frac{\lambda_i}{2} \|x - x_i\|^2, \quad \bar{x}_{i+1} := \operatorname{argmin}_x f^i(x).$$

The exact proximal point method (Martinet, 1972, 1970; Rockafellar, 1976) proceeds by inductively declaring  $x_i = \bar{x}_i$  for  $i \geq 1$ . Since computing  $\bar{x}_i$  exactly is in general impossible, we will instead monitor the error  $\|\bar{x}_i - x_i\|$ . The following elementary result will form the basis for the rest of the paper. To simplify notation, we will set  $\bar{x}_0 := \operatorname{argmin} f$  and  $\lambda_{-1} := 0$ , throughout.

**Theorem 2 (Inexact proximal point method)** *For all  $j \geq 0$ , the following estimate holds:*

$$f^j(\bar{x}_{j+1}) - f^* \leq \sum_{i=0}^j \frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2. \quad (10)$$

Consequently, we have the error decomposition:

$$\boxed{f(x_{j+1}) - f^* \leq (f^j(x_{j+1}) - f^j(\bar{x}_{j+1})) + \sum_{i=0}^j \frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2.} \quad (11)$$

Moreover, if  $f$  is  $L$ -smooth, then for all  $j \geq 0$  the estimate holds:

$$f(x_j) - f^* \leq \frac{L + \lambda_{j-1}}{2} \|\bar{x}_j - x_j\|^2 + \sum_{i=0}^{j-1} \frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2. \quad (12)$$

**Proof** We first establish (10) by induction. For the base case  $j = 0$ , observe  $\lambda_- = 0$  and

$$f^0(\bar{x}_1) = \min_x f^0(x) \leq f^0(\bar{x}_0) = f^* + \frac{\lambda_0}{2} \|\bar{x}_0 - x_0\|^2.$$

As the inductive assumption, suppose (10) holds up to iteration  $j - 1$ . We then conclude

$$\begin{aligned} f^j(\bar{x}_{j+1}) &\leq f^j(\bar{x}_j) = f(\bar{x}_j) + \frac{\lambda_j}{2} \|\bar{x}_j - x_j\|^2 \\ &\leq f^{j-1}(\bar{x}_j) + \frac{\lambda_j}{2} \|\bar{x}_j - x_j\|^2 \leq f^* + \sum_{i=0}^j \frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2, \end{aligned}$$

where the last inequality follows by the inductive assumption. This completes the proof of (10). To see (11), we observe using (10) the estimate

$$\begin{aligned} f(x_{j+1}) - f^* &\leq f^j(x_{j+1}) - f^* = (f^j(x_{j+1}) - f^j(\bar{x}_{j+1})) + f^j(\bar{x}_{j+1}) - f^* \\ &\leq (f^j(x_{j+1}) - f^j(\bar{x}_{j+1})) + \sum_{i=0}^j \frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2. \end{aligned}$$

Finally, if  $f$  is  $L$ -smooth, then  $f^j$  is  $(L + \lambda_j)$ -smooth. An analogous result to (7) yields

$$f^j(x_{j+1}) - f^j(\bar{x}_{j+1}) \leq \frac{L + \lambda_j}{2} \|\bar{x}_{j+1} - x_{j+1}\|^2.$$

Inequality (12) follows from applying this bound in (11). ■

The main conclusion of Theorem 2 is the decomposition of the functional error described in (11). Namely, the estimate (11) upper bounds the error  $f(x_{j+1}) - \min f$  as the sum of the suboptimality in the last step  $f^T(x_{T+1}) - f^T(\bar{x}_{T+1})$  and the errors  $\frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2$  incurred along the way. By choosing  $T$  sufficiently large, we can be sure that the function  $f^T$  is well-conditioned. Moreover in order to ensure that each term in the sum  $\frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2$  is of order  $\delta$ , it suffices to guarantee  $\|\bar{x}_i - x_i\| \leq \sqrt{\frac{2\delta}{\lambda_i}}$  for each index  $i$ . Since  $\lambda_i$  is an increasing sequence, it follows that we may gradually decrease the tolerance on the errors  $\|\bar{x}_i - x_i\|$ , all the while improving the conditioning of the functions we encounter. With this intuition in mind, we introduce the **proxBoost** procedure (Algorithm 2). The algorithm, and its latter modifications, depend on the amplitude sequence  $\{\lambda_j\}_{j=1}^T$  governing the proximal regularization terms. To simplify notation, we will omit this sequence from the algorithm input and instead treat it as a global parameter specified in theorems.

<p><b>Algorithm 2:</b> <b>proxBoost</b>(<math>\delta, p, T</math>)</p> <p><b>Input:</b> <math>\delta \geq 0, p \in (0, 1), T \in \mathbb{N}</math></p> <p>Set <math>\lambda_{-1} = 0, \varepsilon_{-1} = \sqrt{\frac{2\delta}{\mu}}</math></p> <p>Generate a point <math>x_0</math> satisfying <math>\ x_0 - \bar{x}_0\  \leq \varepsilon_{-1}</math> with probability <math>1 - p</math>.</p> <p><b>for</b> <math>j = 0, \dots, T - 1</math> <b>do</b></p> <div style="margin-left: 20px;"> <p>Set <math>\varepsilon_j = \sqrt{\frac{2\delta}{\mu + \lambda_j}}</math></p> <p>Generate a point <math>x_{j+1}</math> satisfying</p> <math display="block">\mathbb{P} [\ x_{j+1} - \bar{x}_{j+1}\  \leq \varepsilon_j \mid E_j] \geq 1 - p, \tag{13}</math> <p style="margin-left: 40px;">where <math>E_j</math> denotes the event <math>E_j := \{x_i \in B_{\varepsilon_{i-1}}(\bar{x}_i) \text{ for all } i \in [0, j]\}</math>.</p> </div> <p><b>end</b></p> <p>Generate a point <math>x_{T+1}</math> satisfying</p> $\mathbb{P} [f^T(x_{T+1}) - \min f^T \leq \delta \mid E_T] \geq 1 - p. \tag{14}$ <p><b>Return</b> <math>x_{T+1}</math></p>
---

Thus **proxBoost** consists of three stages, which we now examine in detail.

**Stage I: Initialization.** Algorithm 2 begins by generating a point  $x_0$  that is a distance of  $\sqrt{\frac{2\delta}{\mu}}$  away from the minimizer of  $f$  with probability  $1 - p$ . This task can be achieved by applying a robust distance estimator on  $f$ , as discussed in Section 2.

**Stage II: Proximal iterations.** In each subsequent iteration,  $x_{j+1}$  is defined to be a point that is within a radius of  $\varepsilon_j = \sqrt{\frac{2\delta}{\mu + \lambda_j}}$  from the minimizer of  $f^j$  with probability

$1 - p$  conditioned on the event  $E_j$ . The event  $E_j$  encodes that each previous iteration was successful in the sense that the point  $x_i$  indeed lies inside the ball  $B_{\varepsilon_{i-1}}(\bar{x}_i)$  for all  $i = 0, \dots, j$ . Thus  $x_{j+1}$  can be determined by a procedure that conditioned on the event  $E_j$  is a robust distance estimator on the function  $f^j$ .

**Stage III: Cleanup.** In the final step, the algorithm outputs a  $\delta$ -minimizer of  $f^T$  with probability  $1 - p$  conditioned on the event  $E_T$ . In particular, if  $f$  is  $L$ -smooth then we may use a robust distance estimator on  $f^T$  directly. Namely, taking into account the upper bound in (7), we may declare  $x_{T+1}$  to be any point satisfying

$$\mathbb{P} \left[ \|x_{T+1} - \bar{x}_{T+1}\| \leq \sqrt{\frac{2\delta}{L+\lambda_T}} \mid E_T \right] \geq 1 - p.$$

Notice that by choosing  $\lambda_T$  sufficiently large, we may ensure that the condition number  $\frac{\mu+\lambda_T}{L+\lambda_T}$  of  $f^T$  is arbitrarily close to one. If  $f$  is not smooth, such as when constraints or additional regularizers are present, we can not use a robust distance estimator in the cleanup stage. We will see in Section 6 a different approach for convex composite problems, based on a modified robust distance estimation technique.

The following theorem summarizes the guarantees of the `proxBoost` procedure.

**Theorem 3 (Proximal Boost)** *Fix a constant  $\delta > 0$ , a probability of failure  $p \in (0, 1)$  and a natural number  $T \in \mathbb{N}$ . Then with probability at least  $1 - (T + 2)p$ , the point  $x_{T+1} = \text{proxBoost}(\delta, p, T)$  satisfies*

$$f(x_{T+1}) - \min f \leq \delta \left( 1 + \sum_{i=0}^T \frac{\lambda_i}{\mu + \lambda_{i-1}} \right). \quad (15)$$

**Proof** We first prove by induction the estimate

$$\mathbb{P}[E_t] \geq 1 - (t + 1)p \quad \text{for all } t = 0, \dots, T. \quad (16)$$

The base case  $t = 0$  is immediate from the definition of  $x_0$ . Suppose now that (16) holds for some index  $t - 1$ . Then the inductive assumption and the definition of  $x_t$  yield

$$\mathbb{P}[E_t] = \mathbb{P}[E_t \mid E_{t-1}] \mathbb{P}[E_{t-1}] \geq (1 - p)(1 - tp) \geq 1 - (t + 1)p,$$

thereby completing the induction. Thus the inequalities (16) hold. Define the event

$$F = \{f^T(x_{T+1}) - \min f^T \leq \delta\}.$$

We therefore deduce

$$\mathbb{P}[F \cap E_T] = \mathbb{P}[F \mid E_T] \cdot \mathbb{P}[E_T] \geq (1 - (T + 1)p)(1 - p) \geq 1 - (T + 2)p.$$

Suppose now that the event  $F \cap E_T$  occurs. Then using the estimate (11), we conclude

$$f(x_{T+1}) - \min f \leq (f^T(x_{T+1}) - f^T(\bar{x}_{T+1})) + \sum_{i=0}^T \frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2 \leq \delta + \sum_{i=0}^T \frac{\delta \lambda_i}{\mu + \lambda_{i-1}},$$

where the last inequality uses the definitions of  $x_{T+1}$  and  $\varepsilon_j$ . This completes the proof.  $\blacksquare$

Looking at the estimate (15), we see that the final error  $f(x_{T+1}) - \min f$  is controlled by the sum  $\sum_{i=0}^T \frac{\lambda_i}{\mu + \lambda_{i-1}}$ . A moment of thought yields an appealing choice  $\lambda_i = \mu 2^i$  for the proximal parameters. Indeed, then every element in the sum  $\frac{\lambda_i}{\mu + \lambda_{i-1}}$  is upper bounded by two. Moreover, if  $f$  is  $L$ -smooth, then the condition number  $\frac{L + \lambda_T}{\mu + \lambda_T}$  of  $f^T$  is upper bounded by two after only  $T = \lceil \log(L/\mu) \rceil$  rounds.

**Corollary 4 (Proximal boost with geometric decay)** *Fix an iteration count  $T$ , a target accuracy  $\epsilon > 0$ , and a probability of failure  $p \in (0, 1)$ . Define the algorithm parameters:*

$$\delta = \frac{\epsilon}{2 + 2T} \quad \text{and} \quad \lambda_i = \mu 2^i \quad \forall i \in [0, T].$$

*Then the point  $x_{T+1} = \text{proxBoost}(\delta, p, T)$  satisfies*

$$\mathbb{P}(f(x_{T+1}) - \min f \leq \epsilon) \geq 1 - (T + 2)p.$$

In the next two sections, we seed the `proxBoost` procedure with (accelerated) stochastic gradient algorithms and methods based on empirical risk minimization. The reader, however, should keep in mind that `proxBoost` is entirely agnostic to the inner workings of the robust distance estimators it uses. The only point to be careful about is that some distance estimators (e.g., when using stochastic gradient methods) require auxiliary quantities as input, such as an upper estimate on the function gap at the initial point. Therefore, we may have to update such estimates along the iterations of `proxBoost`.

#### 4. Consequences for empirical risk minimization

In this section, we explore the consequences of the `proxBoost` algorithm for empirical risk minimization. Setting the stage, fix a probability space  $(\Omega, \mathcal{F}, \mathcal{P})$  and equip  $\mathbf{R}^d$  with the Borel  $\sigma$ -algebra. Consider the optimization problem

$$\min_x f(x) = \mathbb{E}_{z \sim \mathcal{P}} [f(x, z)], \tag{17}$$

where  $f: \mathbf{R}^d \times \Omega \rightarrow \mathbf{R}_+$  is a measurable nonnegative function. A common approach to problems of the form (17) is based on empirical risk minimization. Namely, one collects i.i.d. samples  $z_1, \dots, z_n \sim \mathcal{P}$  and minimizes the empirical average

$$\min_x f_S(x) := \frac{1}{n} \sum_{i=1}^n f(x, z_i). \tag{18}$$

A central question is to determine the number  $n$  of samples that would ensure that the minimizer  $x_S$  of the empirical risk has low generalization error  $f(x_S) - \min f$ , with reasonably high probability. There is a vast literature on this subject; some representative works include Hsu and Sabato (2016); Bartlett and Mendelson (2002); Shalev-Shwartz et al. (2009); Shalev-Shwartz and Ben-David (2014). We build here on the work of Hsu-Sabato Hsu and Sabato (2016), who specifically focused on high confidence guarantees for smooth strongly convex minimization. As in the previous sections, we let  $\bar{x}$  be a minimizer of  $f$  and define the shorthand  $f^* = \min f$ .

**Assumption 2** Following Hsu and Sabato (2016), we make the following assumptions on the loss.

1. **(Strong convexity)** There exist a real  $\mu > 0$  and a natural number  $N \in \mathbb{N}$  such that:
  - (a) the population loss  $f$  is  $\mu$ -strongly convex,
  - (b) the empirical loss  $x \mapsto f_S(x)$  is  $\mu$ -strongly convex with probability at least  $5/6$ , whenever  $|S| \geq N$ .
2. **(Smoothness)** There exist constants  $L, \hat{L} > 0$  such that:
  - (a) for a.e.  $z \sim \mathcal{P}$ , the loss  $x \mapsto f(x, z)$  is nonnegative and  $\hat{L}$ -smooth,
  - (b) the population objective  $x \mapsto f(x)$  is  $L$ -smooth. (It holds that  $L \leq \hat{L}$ .)

In addition, we assume  $f^* := \min f > 0$ .

The following result of Hsu and Sabato (2016, Theorem 15) shows that the empirical risk minimizer is a weak distance oracle for the problem (17).

**Lemma 5** Fix an i.i.d. sample  $z_1, \dots, z_n \sim \mathcal{P}$  of size  $n \geq N$ . Suppose Assumption 2 holds. Then the minimizer  $x_S$  of the empirical risk (18) satisfies the bound:

$$\mathbb{P} \left[ \|x_S - \bar{x}\| \leq \sqrt{\frac{96\hat{L}f^*}{n\mu^2}} \right] \geq 2/3.$$

In particular, using Algorithm 1 one may turn empirical risk minimization into a robust distance estimator for the problem (4) using a total of  $mn$  samples. Let us estimate the function value at the generated point by a direct application of smoothness. Appealing to Lemma 1 and the two-sided bound (7), we deduce that with probability  $1 - \exp(-m/18)$  the procedure will return a point  $x$  satisfying

$$f(x) \leq \left( 1 + \frac{432\hat{L}L}{n\mu^2} \right) f^*.$$

Observe that this is an estimate of *relative error*. In particular, let  $p \in (0, 1)$  be some acceptable probability of failure and let  $\gamma > 0$  be a desired level of relative accuracy. Then setting  $m = \lceil 18 \ln(1/p) \rceil$  and  $n \geq \max\{\frac{432\hat{\kappa}\kappa}{\gamma}, N\}$ , we conclude that  $x$  satisfies

$$\mathbb{P}[f(x) \leq (1 + \gamma)f^*] \geq 1 - p, \tag{19}$$

while the overall sample complexity is

$$\left\lceil 18 \ln \left( \frac{1}{p} \right) \right\rceil \cdot \max \left\{ \left\lceil \frac{432\hat{\kappa}\kappa}{\gamma} \right\rceil, N \right\}, \tag{20}$$

where  $\hat{\kappa} = \hat{L}/\mu$  and  $\kappa = L/\mu$ . This is exactly the result of Hsu and Sabato (2016, Corollary 16).

We will now see how to find a point  $x$  satisfying (19) with significantly fewer samples by embedding empirical risk minimization within `proxBoost`. Algorithm 3 encodes the empirical risk minimization process on a quadratically regularized problem. Algorithm 4 is the robust distance estimator induced by Algorithm 3. Finally, Algorithm 5 is the `proxBoost` algorithm specialized to empirical risk minimization.

<p><b>Algorithm 3:</b> <math>\text{ERM}(n, \lambda, x)</math></p> <p><b>Input:</b> sample count <math>n \in \mathbb{N}</math>, center <math>x \in \mathbf{R}^d</math>, amplitude <math>\lambda &gt; 0</math>. Generate i.i.d. samples <math>z_1, \dots, z_n \sim \mathcal{P}</math> and compute the minimizer <math>\bar{y}</math> of</p> $\min_y \frac{1}{n} \sum_{i=1}^n f(y, z_i) + \frac{\lambda}{2} \ y - x\ ^2.$ <p><b>Return</b> <math>\bar{y}</math></p>
--

<p><b>Algorithm 4:</b> <math>\text{ERM-R}(n, m, \lambda, x)</math></p> <p><b>Input:</b> sample count <math>n \in \mathbb{N}</math>, trial count <math>m \in \mathbb{N}</math>, center <math>x \in \mathbf{R}^d</math>, amplitude <math>\lambda &gt; 0</math>. Query <math>m</math> times <math>\text{ERM}(n, \lambda, x)</math> and let <math>Y = \{y_1, \dots, y_m\}</math> consist of the responses. <b>Step</b> <math>j = 1, \dots, m</math>:     Compute <math>r_j = \min\{r \geq 0 :  B_r(y_j) \cap Y  &gt; \frac{m}{2}\}</math>.     Set <math>i^* = \text{argmin}_{i \in [1, m]} r_i</math> <b>Return</b> <math>y_{i^*}</math></p>
---

<p><b>Algorithm 5:</b> <math>\text{BoostERM}(\gamma, T, m)</math></p> <p><b>Input:</b> <math>T, m \in \mathbb{N}</math>, <math>\gamma &gt; 0</math> Set <math>\lambda_{-1} = 0</math>, <math>x_{-1} = 0</math>, <math>n_{-1} = \frac{432\hat{L}}{\gamma\mu}</math> <b>Step</b> <math>j = 0, \dots, T</math>:     <math>x_j = \text{ERM-R}(n_{j-1}, m, \lambda_{j-1}, x_{j-1})</math>     <math>n_j = 432 \left\lceil \frac{\hat{L} + \lambda_j}{\mu + \lambda_j} \left( \frac{1}{\gamma} + \sum_{i=0}^j \frac{\lambda_i}{\mu + \lambda_{i-1}} \right) \right\rceil \vee N</math> <b>Return</b> <math>x_{T+1} = \text{ERM-R}\left(\frac{L + \lambda_T}{\mu + \lambda_T} \cdot n_T, m, \lambda_T, x_T\right)</math></p>
---

Using Theorem 3, we can now prove the following result.

**Theorem 6 (Efficiency of BoostERM)** *Fix a target relative accuracy  $\gamma > 0$  and numbers  $T, m \in \mathbb{N}$ . Then with probability at least  $1 - (T + 2) \exp(-\frac{m}{18})$ , the point  $x_{T+1} = \text{BoostERM}(\gamma, T, m)$  satisfies*

$$f(x_{T+1}) - f^* \leq \left( 1 + \sum_{i=0}^T \frac{\lambda_i}{\mu + \lambda_{i-1}} \right) \gamma f^*.$$

**Proof** We will verify that Algorithm 5 is an instantiation of Algorithm 2 with  $\delta = \gamma f^*$  and  $p = \exp(-\frac{m}{18})$ . More precisely, we will prove by induction that with this choice of  $p$  and

$\delta$ , the iterates  $x_j$  satisfy (13) for each index  $j = 0, \dots, T$  and  $x_{T+1}$  satisfies (14). As the base case, consider the evaluation  $x_0 = \text{ERM-R}(n_{-1}, m, \lambda_{-1}, x_{-1})$  where  $x_{-1}$  can be arbitrary since  $\lambda_{-1} = 0$ . Then Lemma 1 and Theorem 5 guarantee

$$\mathbb{P} \left[ \|x_0 - \bar{x}_0\| \leq 3\sqrt{\frac{96\hat{L}f^*}{n_{-1}\mu^2}} \right] \geq 1 - \exp\left(-\frac{m}{18}\right).$$

Taking into account the definitions of  $n_{-1}$  in Algorithm 5 and  $\epsilon_{-1}$  in Algorithm 2, we deduce

$$\mathbb{P} [\|x_0 - \bar{x}_0\| \leq \epsilon_{-1}] \geq 1 - p,$$

as claimed. As an inductive hypothesis, suppose that (13) holds for  $x_0, x_1, \dots, x_{j-1}$ . We will prove it holds for  $x_j = \text{ERM-R}(n_{j-1}, m, \lambda_{j-1}, x_{j-1})$ . To this end, suppose that the event  $E_{j-1}$  occurs. Then by the same reasoning as in the base case, the point  $x_j$  satisfies

$$\mathbb{P} \left[ \|x_j - \bar{x}_j\| \leq 3\sqrt{\frac{96(\hat{L} + \lambda_{j-1})f^{j-1}(\bar{x}_j)}{n_{j-1}(\mu + \lambda_{j-1})^2}} \right] \geq 1 - \exp\left(-\frac{m}{18}\right). \quad (21)$$

Now, using (10) and the inductive assumption that  $\|x_i - \bar{x}_i\| \leq \epsilon_{i-1} = \sqrt{\frac{2\delta}{\mu + \lambda_{i-1}}}$  for all  $i \in [0, j-1]$  (conditioned on  $E_{j-1}$ ), we have

$$f^{j-1}(\bar{x}_j) - f^* \leq \sum_{i=0}^{j-1} \frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2 \leq \delta \sum_{i=0}^{j-1} \frac{\lambda_i}{\mu + \lambda_{i-1}},$$

which, together with  $\delta = \gamma f^*$ , implies

$$f^{j-1}(\bar{x}_j) \leq f^* + \delta \sum_{i=0}^{j-1} \frac{\lambda_i}{\mu + \lambda_{i-1}} = \left(1 + \gamma \sum_{i=0}^{j-1} \frac{\lambda_i}{\mu + \lambda_{i-1}}\right) f^*.$$

Combining this inequality with (21), we conclude that conditioned on the event  $E_{j-1}$ , we have with probability  $1 - p$  the guarantee

$$\frac{\mu + \lambda_{j-1}}{2} \|x_j - \bar{x}_j\|^2 \leq \frac{432(\hat{L} + \lambda_{j-1})(1 + \gamma \sum_{i=0}^{j-1} \frac{\lambda_i}{\mu + \lambda_{i-1}})}{n_{j-1}(\mu + \lambda_{j-1})} \cdot f^* \leq \gamma f^* = \delta, \quad (22)$$

where the last inequality follows from the definition of  $n_{j-1}$ . This implies that the estimate (13) holds for  $x_j$  with  $\epsilon_{j-1} = \sqrt{\frac{2\delta}{\mu + \lambda_{j-1}}}$ . Therefore, it holds for all iterates  $x_0, \dots, x_T$ , as needed. Suppose now that that event  $E_T$  occurs. Then by exactly the same reasoning that led to (22), and considering the extra factor  $\frac{L + \lambda_T}{\mu + \lambda_T}$  multiplied to  $n_T$  in the last call of  $\text{ERM-R}$ , we have the estimate

$$\frac{\mu + \lambda_T}{2} \|x_{T+1} - \bar{x}_{T+1}\|^2 \leq \frac{\mu + \lambda_T}{L + \lambda_T} \gamma f^*.$$

Using smoothness, we therefore deduce  $f^T(x_{T+1}) - \min f^T \leq \gamma f^* = \delta$ , as claimed. An application of Theorem 3 completes the proof.  $\blacksquare$

Finally, using the proximal parameters  $\lambda_i = \mu 2^i$  yields the following guarantee.



**Corollary 7 (Efficiency of BoostERM with geometric decay)** *Fix a target relative accuracy  $\gamma' > 0$  and a probability of failure  $p \in (0, 1)$ . Define the algorithm parameters:*

$$T = \lceil \log_2(\kappa) \rceil, \quad m = \left\lceil 18 \ln \left( \frac{T+2}{p} \right) \right\rceil, \quad \gamma = \frac{\gamma'}{2+2T}, \quad \lambda_i = \mu 2^i.$$

*Then with probability of at least  $1-p$ , the point  $x_{T+1} = \text{BoostERM}(\gamma, T, m)$  satisfies  $f(x_{T+1}) \leq (1 + \gamma')f^*$ . Moreover, the total number of samples used by the algorithm is*

$$\mathcal{O} \left( \ln(\kappa) \ln \left( \frac{\ln(\kappa)}{p} \right) \cdot \max \left\{ \left( 1 + \frac{1}{\gamma'} \right) \hat{\kappa} \ln(\kappa), N \right\} \right).$$

Notice that the sample complexity provided by Corollary 7 is an order of magnitude better than (20) in terms of the dependence on the condition numbers  $\hat{\kappa}$  and  $\kappa$ .

## 5. Consequences for stochastic approximation

We next investigate the consequences of `proxBoost` for stochastic approximation. Namely, we will seed `proxBoost` with the robust distance estimator, induced by the stochastic gradient method and its accelerated variant. An important point is that the sample complexity of stochastic gradient methods depends on the initialization quality  $f(x_0) - f^*$ . Consequently, in order to know how many iterations are needed to reach a desired accuracy  $\mathbb{E}[f(x_i)] - f^* \leq \delta$ , we must have available an upper bound on the initialization quality  $\Delta \geq f(x_0) - f^*$ . Therefore, we will have to dynamically update an estimate of the initialization quality for each proximal subproblem along the iterations of `proxBoost`. The following assumption formalizes this idea.

**Assumption 3** Consider the proximal minimization problem

$$\min_y \varphi_x(y) := f(y) + \frac{\lambda}{2} \|y - x\|^2,$$

Let  $\Delta > 0$  be a real number satisfying  $\varphi_x(x) - \min \varphi_x \leq \Delta$ . We will let `Alg`( $\delta, \lambda, \Delta, x$ ) be a procedure that returns a point  $y$  satisfying

$$\mathbb{P}[\varphi_x(y) - \min \varphi_x \leq \delta] \geq \frac{2}{3}.$$

Clearly, `Alg`( $\delta, \lambda, \Delta, x$ ) is a minimization oracle in the sense of (3). We note that specific instantiations of the oracle `Alg`( $\delta, \lambda, \Delta, x$ ) (e.g. SGD or accelerated SGD) require further assumptions on the nature of the noise (e.g bounded gradient variance (24)). Since the proximal function  $\varphi_x$  is  $(\mu + \lambda)$ -strongly convex, it has a unique minimizer  $\bar{y}_x$  and satisfies

$$\frac{\mu + \lambda}{2} \|y - \bar{y}_x\|^2 \leq \varphi_x(y) - \min \varphi_x.$$

Therefore, `Alg`( $\delta, \lambda, \Delta, x$ ) is a weak distance oracle, in the sense that  $\mathbb{P}(\|y - \bar{y}_x\| \leq \varepsilon) \geq \frac{2}{3}$  with  $\varepsilon = \sqrt{\frac{2\delta}{\mu + \lambda}}$ . Following the procedure in Section 2, we may turn it into a robust distance

estimator for minimizing  $\varphi_x$ , as long as  $\Delta$  upper bounds the initialization error. We record the robust distance estimator induced by  $\text{Alg}(\cdot)$  as Algorithm 6.

<p><b>Algorithm 6:</b> <math>\text{Alg-R}(\delta, \lambda, \Delta, x, m)</math></p> <p><b>Input:</b> accuracy <math>\delta &gt; 0</math>, amplitude <math>\lambda &gt; 0</math>, upper bound <math>\Delta &gt; 0</math>, center <math>x \in \mathbf{R}^d</math>, trial count <math>m \in \mathbb{N}</math>.</p> <p>Query <math>m</math> times <math>\text{Alg}(\delta, \lambda, \Delta, x)</math> and let <math>Y = \{y_1, \dots, y_m\}</math> consist of the responses.</p> <p><b>Step</b> <math>j = 1, \dots, m</math>:              Compute <math>r_i = \min\{r \geq 0 :  B_r(y_i) \cap Y  &gt; \frac{m}{2}\}</math>.</p> <p>Set <math>i^* = \text{argmin}_{i \in [1, m]} r_i</math></p> <p><b>Return</b> <math>y_{i^*}</math></p>
---

Henceforth, in addition to Assumptions 1 and 3, we assume that  $f$  is  $L$ -smooth and set  $\kappa = \frac{L}{\mu}$ . It is then straightforward to instantiate `proxBoost` with the robust distance estimator  $\text{Alg-R}$ . We record the resulting procedure as Algorithm 7.

<p><b>Algorithm 7:</b> <math>\text{BoostAlg}(\delta, \Delta_{\text{in}}, x_{\text{in}}, T, m)</math></p> <p><b>Input:</b> accuracy <math>\delta &gt; 0</math>, upper bound <math>\Delta_{\text{in}} &gt; 0</math>, center <math>x_{\text{in}} \in \mathbf{R}^d</math>, and <math>m, T \in \mathbb{N}</math></p> <p>Set <math>\lambda_{-1} = 0, \Delta_{-1} = \Delta_{\text{in}}, x_{-1} = x_{\text{in}}</math></p> <p><b>Step</b> <math>j = 0, \dots, T</math>:  <math>x_j = \text{Alg-R}(\delta/9, \lambda_{j-1}, \Delta_{j-1}, x_{j-1}, m)</math>  <math>\Delta_j = \delta \left( \frac{L + \lambda_{j-1}}{\mu + \lambda_{j-1}} + \sum_{i=0}^{j-1} \frac{\lambda_i}{\mu + \lambda_{i-1}} \right)</math></p> <p><b>Return</b> <math>x_{T+1} = \text{Alg-R}(\frac{\mu + \lambda_T}{L + \lambda_T} \cdot \frac{\delta}{9}, \lambda_T, \Delta_T, x_T, m)</math></p>
---

We can now prove the following theorem on the efficiency of Algorithm 7. The proof is almost a direct application of Theorem 3. The only technical point is to verify that for all indices  $j$ , the quantity  $\Delta_j$  is a valid upper bound on the initialization error  $f^j(x_j) - \min f^j$  in the event  $E_j$  (defined in Algorithm 2).

**Theorem 8 (Efficiency of BoostAlg)** *Fix an arbitrary point  $x_{\text{in}} \in \mathbf{R}^d$  and let  $\Delta_{\text{in}}$  be any constant satisfying  $\Delta_{\text{in}} \geq f(x_{\text{in}}) - \min f$ . Fix natural numbers  $T, m \in \mathbb{N}$ . Then with probability at least  $1 - (T + 2) \exp(-\frac{m}{18})$ , the point  $x_{T+1} = \text{BoostAlg}(\delta, \Delta_{\text{in}}, x_{\text{in}}, T, m)$  satisfies*

$$f(x_{T+1}) - \min f \leq \delta \left( 1 + \sum_{i=0}^T \frac{\lambda_i}{\mu + \lambda_{i-1}} \right).$$

**Proof** We will verify that Algorithm 7 is an instantiation of Algorithm 2 with  $p = \exp(-\frac{m}{18})$ . More precisely, we will prove by induction that with this choice of  $p$ , the iterates  $x_j$  satisfy (13) for each index  $j = 0, \dots, T$  and  $x_{T+1}$  satisfies (14). For the base case  $j = 0$ , Lemma 1 guarantees that with probability  $1 - p$ , the point  $x_0$  produced by the robust distance estimator  $\text{Alg-R}$  satisfies

$$\|x_0 - \bar{x}_0\| \leq 3 \sqrt{\frac{2 \cdot \delta/9}{\mu}} = \varepsilon_{-1}.$$

As an inductive hypothesis, suppose that (13) holds for the iterates  $x_0, \dots, x_{j-1}$  for some  $j \geq 1$ . We will prove it holds for  $x_j$ . To this end, suppose that the event  $E_{j-1}$  occurs. Then using (12) we deduce

$$\begin{aligned} f(x_{j-1}) - f^* &\leq \frac{L + \lambda_{j-2}}{2} \|\bar{x}_{j-1} - x_{j-1}\|^2 + \sum_{i=0}^{j-2} \frac{\lambda_i}{2} \|\bar{x}_i - x_i\|^2 \\ &\leq \frac{\delta(L + \lambda_{j-2})}{\mu + \lambda_{j-2}} + \sum_{i=0}^{j-2} \frac{\delta\lambda_i}{\mu + \lambda_{i-1}} = \Delta_{j-1}, \end{aligned}$$

where the second inequality follows from  $x_i \in B_{\varepsilon_{i-1}}(\bar{x}_i)$  with  $\varepsilon_{i-1} = \sqrt{\frac{2\delta}{\mu + \lambda_{i-1}}}$  for all  $i \in [0, j-1]$ . By examining the definition of  $f^{j-1}$ , we deduce  $f^{j-1}(x_{j-1}) = f(x_{j-1})$  and  $\min f^{j-1} \geq \min f = f^*$ , which imply

$$f^{j-1}(x_{j-1}) - \min f^{j-1} \leq f(x_{j-1}) - f^* \leq \Delta_{j-1}. \quad (23)$$

That is,  $\Delta_{j-1}$  is an upper bound on the initial gap  $f^{j-1}(x_{j-1}) - \min f^{j-1}$  for all  $j$  whenever the event  $E_{j-1}$  occurs. Moreover Lemma 1 guarantees that conditioned on  $E_{j-1}$  with probability  $1 - p$ , the following estimate holds:

$$\|x_j - \bar{x}_j\| \leq 3\sqrt{\frac{2 \cdot \delta/9}{\mu + \lambda_{j-1}}} = \varepsilon_{j-1}.$$

Thus the condition (13) holds for the iterate  $x_j$ , as desired.

Now suppose that the event  $E_T$  holds. Then exactly the same reasoning that led to (23) yields the guarantee  $f^T(x_T) - \min f^T \leq \Delta_T$ . Therefore Lemma 1 guarantees that with probability  $1 - p$  conditioned on  $E_T$ , we have

$$\|x_{T+1} - \bar{x}_{T+1}\| \leq 3\sqrt{\frac{2}{\mu + \lambda_T} \cdot \frac{\delta}{9} \cdot \frac{\mu + \lambda_T}{L + \lambda_T}} = \sqrt{\frac{2\delta}{L + \lambda_T}}.$$

Taking into account the fact that  $f^T$  is  $(L + \lambda_T)$ -smooth, we therefore deduce

$$\mathbb{P}[f^T(x_{T+1}) - \min f^T \leq \delta \mid E_T] \geq 1 - p,$$

thereby establishing (14). An application of Theorem 3 completes the proof.  $\blacksquare$

When using the proximal parameters  $\lambda_i = \mu 2^i$ , we obtain the following guarantee.

**Corollary 9 (Efficiency of BoostAlg with geometric decay)** *Fix an arbitrary point  $x_{\text{in}} \in \mathbf{R}^d$  and let  $\Delta_{\text{in}}$  be any upper bound  $\Delta_{\text{in}} \geq f(x_{\text{in}}) - \min f$ . Fix a target accuracy  $\epsilon > 0$  and probability of failure  $p \in (0, 1)$ , and set the algorithm parameters*

$$T = \lceil \log_2(\kappa) \rceil, \quad m = \left\lceil 18 \ln \left( \frac{2+T}{p} \right) \right\rceil, \quad \delta = \frac{\epsilon}{2+2T}, \quad \lambda_i = \mu 2^i.$$

Then the point  $x_{T+1} = \text{BoostAlg}(\delta, \Delta_{\text{in}}, x_{\text{in}}, T, m)$  satisfies

$$\mathbb{P}(f(x_{T+1}) - \min f \leq \epsilon) \geq 1 - p.$$

Moreover, the total number of calls to  $\text{Alg}(\cdot)$  is

$$\left\lceil 18 \ln \left( \frac{\lceil 2 + \log_2(\kappa) \rceil}{p} \right) \right\rceil \lceil 2 + \log_2(\kappa) \rceil,$$

while the initialization errors satisfy

$$\max_{i=0, \dots, T+1} \Delta_i \leq \frac{\kappa + 1 + 2 \lceil \log_2(\kappa) \rceil}{2 + 2 \lceil \log_2(\kappa) \rceil} \epsilon.$$

We now concretely describe how to use (accelerated) stochastic gradient methods as  $\text{Alg}(\cdot)$  within `proxBoost`.

**Illustration: robust (accelerated) stochastic gradient methods**

Following the standard literature on streaming algorithms, we suppose that the only access to  $f$  is through a stochastic gradient oracle. Namely, fix a probability space  $(\Omega, \mathcal{F}, \mathcal{P})$  and let  $G: \mathbf{R}^d \times \Omega \rightarrow \mathbf{R}$  be a measurable map satisfying

$$\mathbb{E}_z G(x, z) = \nabla f(x) \quad \text{and} \quad \mathbb{E}_z \|G(x, z) - \nabla f(x)\|^2 \leq \sigma^2. \quad (24)$$

We suppose that for any point  $x$ , we may sample  $z \in \Omega$  and compute the vector  $G(x, z)$ , which serves as an unbiased estimator of the gradient  $\nabla f(x)$ . The performance of standard numerical methods within this model of computation is judged by their sample complexity—the number of stochastic gradient evaluations  $G(x, z)$  with  $z \sim \mathcal{P}$  required by the algorithm to produce an approximate minimizer of the problem.

Fix an initial point  $x_{\text{in}}$  and let  $\Delta_{\text{in}} > 0$  satisfy  $\Delta_{\text{in}} \geq f(x_0) - f^*$ . It is well known that an appropriately modified stochastic gradient method can generate a point  $x$  satisfying  $\mathbb{E}f(x) - f^* \leq \epsilon$  with sample complexity

$$\mathcal{O} \left( \kappa \log \left( \frac{\Delta_{\text{in}}}{\epsilon} \right) + \frac{\sigma^2}{\mu \epsilon} \right). \quad (25)$$

The accelerated stochastic gradient method of Ghadimi and Lan (2013, Multi-stage ACSA, Proposition 6) and the simplified optimal algorithm of Kulunchakov and Mairal (2019, Restarted Algorithm C, Corollary 9) have the substantially better sample complexity

$$\mathcal{O} \left( \sqrt{\kappa} \log \left( \frac{\Delta_{\text{in}}}{\epsilon} \right) + \frac{\sigma^2}{\mu \epsilon} \right). \quad (26)$$

Clearly, we may use either of these two procedures as  $\text{Alg}(\cdot)$  within the `proxBoost` framework. Indeed, using Corollary 9, we deduce that the two resulting algorithms will find a point  $x$  satisfying

$$\mathbb{P}[f(x) - f^* \leq \epsilon] \geq 1 - p$$

with sample complexities

$$\mathcal{O} \left( \ln(\kappa) \ln \left( \frac{\ln \kappa}{p} \right) \cdot \left( \kappa \ln \left( \frac{\Delta_{\text{in}} \ln(\kappa)}{\epsilon} \vee \kappa \right) + \frac{\sigma^2 \ln(\kappa)}{\mu \epsilon} \right) \right), \quad (27)$$

and

$$\mathcal{O} \left( \ln(\kappa) \ln \left( \frac{\ln \kappa}{p} \right) \cdot \left( \sqrt{\kappa} \ln \left( \frac{\Delta_{\text{in}} \ln(\kappa)}{\epsilon} \vee \kappa \right) + \frac{\sigma^2 \ln(\kappa)}{\mu \epsilon} \right) \right), \quad (28)$$

for the unaccelerated and accelerated methods, respectively. Thus, `proxBoost` endows the stochastic gradient method and its accelerated variant with high confidence guarantees at an overhead cost that is only polylogarithmic in  $\kappa$  and logarithmic in  $1/p$ .

## 6. Extension to convex composite problems

One limitation of the techniques presented in Sections 4 and 5 is that the function  $f$  to be minimized was assumed to be smooth. In particular, these techniques can not accommodate constraints or nonsmooth regularizers. To illustrate the difficulty, consider the task of minimizing a smooth and strongly convex function  $f$  over a closed convex set  $\mathcal{X}$ . The current approach heavily relies on the two-sided bound (7), which guarantees that the function gap  $f(x) - f^*$  and the squared distance to the solution  $\|x - \bar{x}\|^2$  are proportional up to multiplication by the condition number of  $f$ . When a constraint set  $\mathcal{X}$  is present, the left inequality of (7) still holds, but the right inequality is typically false. In particular, in the clean up stage of `proxBoost`, it is unclear how to turn low probability guarantees on the function gap to high probability guarantees using robust distance estimation. In this section, we show how to overcome this difficulty and extend the aforementioned techniques to convex composite optimization problems.

### 6.1 Geometric intuition in the constrained case

Before delving into the details, it is instructive to first focus on the constrained setting, where no additional regularizers are present. This is the content of this section. In section 6.2, we formally describe the algorithm for regularized convex optimization problems in full generality and prove correctness. Consequently, the reader may safely skip to Section 6.2, without losing continuity.

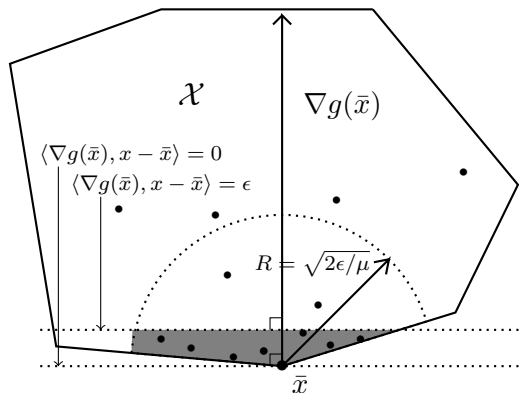
Setting the stage, consider the optimization problem

$$\min_x g(x) \quad \text{subject to} \quad x \in \mathcal{X},$$

where  $g: \mathbf{R}^d \rightarrow \mathbf{R}$  is  $\mu$ -strongly convex and  $L$ -smooth and  $\mathcal{X}$  is a closed convex set. In line with the previous sections, let  $\bar{x}$  be the minimizer of the problem and let  $\kappa = \frac{L}{\mu}$  denote the condition number. Suppose that we have available an algorithm  $\mathcal{M}(\epsilon)$  that generates a point  $x_\epsilon$  satisfying

$$\mathbb{P} \left( g(x_\epsilon) - \min_{x \in \mathcal{X}} g \leq \epsilon \right) \geq \frac{2}{3}. \quad (29)$$

Our immediate goal is to explain how to efficiently boost this low-probability guarantee to a high confidence outcome, albeit with the degraded accuracy  $\kappa \epsilon$ .


 Figure 2: Geometry of the region  $\Lambda$ .

We begin as in the unconstrained setting with the two-sided bound:

$$\langle \nabla g(\bar{x}), x - \bar{x} \rangle + \frac{\mu}{2} \|x - \bar{x}\|^2 \leq g(x) - g(\bar{x}) \leq \langle \nabla g(\bar{x}), x - \bar{x} \rangle + \frac{L}{2} \|x - \bar{x}\|^2 \quad (30)$$

for all  $x \in \mathcal{X}$ . In particular, if the minimizer  $\bar{x}$  lies in the interior of  $\mathcal{X}$  the gradient  $\nabla g(\bar{x})$  vanishes and the estimate (30) reduces to (7). In the more general constrained setting, however, the additive term  $\langle \nabla g(\bar{x}), x - \bar{x} \rangle$  plays an important role. Note that optimality conditions at  $\bar{x}$  immediately imply that this term is nonnegative

$$\langle \nabla g(\bar{x}), x - \bar{x} \rangle \geq 0 \quad \text{for all } x \in \mathcal{X}.$$

Moreover, we see from the estimate (30) that the point  $x_\epsilon$  returned by  $\mathcal{M}(\epsilon)$ , with probability  $2/3$ , lies in the region

$$\Lambda := \left\{ x \in \mathcal{X} : \|x - \bar{x}\| \leq \sqrt{\frac{2\epsilon}{\mu}} \quad \text{and} \quad 0 \leq \langle \nabla g(\bar{x}), x - \bar{x} \rangle \leq \epsilon \right\}. \quad (31)$$

Thus  $x_\epsilon$  simultaneously lies in the ball around  $\bar{x}$  of radius  $\sqrt{2\epsilon/\mu}$  and is sandwiched between two parallel hyperplanes with normal  $\nabla g(\bar{x})$ . See Figure 2 for an illustration.

Naturally, our goal is to generate a point  $x$  that lies in  $\Lambda$ , or a slight perturbation thereof, with probability  $1 - p$ . As the first attempt, suppose for the moment that we know the value of the gradient  $\nabla g(\bar{x})$ . Then we can define the metric

$$\rho(x, x') = \max \left\{ \sqrt{\frac{\epsilon\mu}{2}} \cdot \|x - x'\|, \quad |\langle \nabla g(\bar{x}), x - x' \rangle| \right\}$$

and form the robust distance estimator (Algorithm 1) with  $\rho(\cdot, \cdot)$  replacing the Euclidean norm  $\|\cdot\|$ . In particular, the confidence bound (29) and the left inequality in (30) imply that  $\mathcal{M}(\epsilon)$  is a weak distance oracle, that is,  $\mathbb{P}(\rho(x, \bar{x}) \leq \epsilon) \geq \frac{2}{3}$ . A direct extension of Lemma 1 shows that with  $m$  calls to the oracle  $\mathcal{M}(\epsilon)$ , the robust distance estimator returns a point  $x \in \mathcal{X}$  satisfying

$$\mathbb{P}(\rho(x, \bar{x}) \leq 3\epsilon) \geq 1 - \exp(-m/18).$$

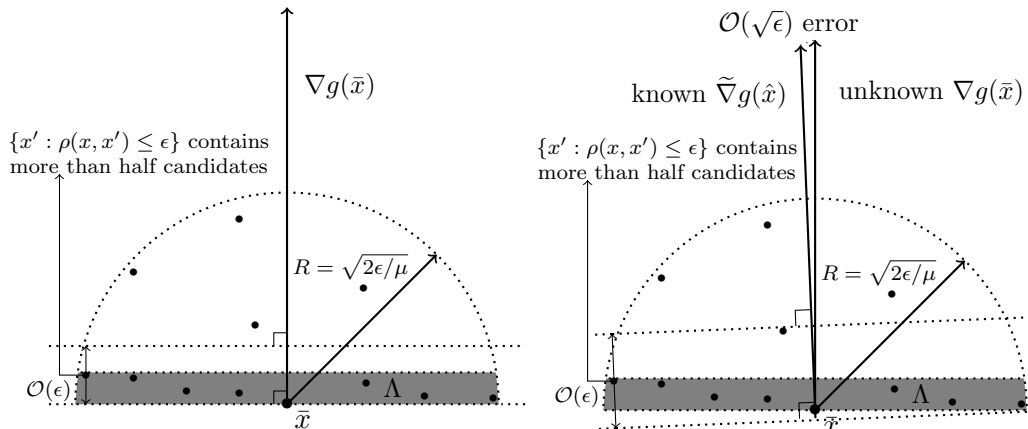


Figure 3: The left side illustrates the region  $\Lambda$ ; the right side depicts the perturbation of  $\Lambda$  obtained by replacing the exact gradient  $\nabla g(\bar{x})$  with an estimator  $\tilde{\nabla} g(\hat{x}) \approx \nabla g(\hat{x})$ .

Consequently, appealing to the right-hand-side of (30), we obtain the desired guarantee

$$\mathbb{P}(g(x) - g(\bar{x}) \leq 3(1 + \kappa)\epsilon) \geq 1 - p.$$

The assumption that we know the gradient  $\nabla g(\bar{x})$  is of course unrealistic. Therefore, the strategy we propose will instead replace the gradient  $\nabla g(\bar{x})$  with some estimate of the gradient  $\nabla g(\hat{x})$  at a nearby point  $\hat{x}$ , which we denote by  $\tilde{\nabla} g(\hat{x})$ . See Figure 3 for an illustration. Indeed, a natural candidate for  $\hat{x}$  is the robust distance estimator of  $\bar{x}$  in the Euclidean norm. We will see that in order for the proposed procedure to work, it suffices for the gradient estimator  $\tilde{\nabla} g(\hat{x})$  to approximate  $\nabla g(\hat{x})$  only up to the very loose accuracy  $\kappa\sqrt{\mu\epsilon}$ . In particular, if we have access to a stochastic gradient estimator of  $\nabla g(\hat{x})$  with variance  $\sigma^2$ , then  $\tilde{\nabla} g$  can be formed using only  $\frac{1}{\kappa^2} \cdot \frac{\sigma^2}{\mu\epsilon}$  samples. This overhead in sample complexity is negligible compared to the cost of executing typical algorithms  $\mathcal{M}(\epsilon)$ , e.g., as given in (25) and (26).

## 6.2 Convex composite setting

In this section, we formally develop the procedure that turns low probability guarantees on the function gap for composite problems to high probability outcomes.

**Assumption 4 (Convex composite problem)** We consider the optimization problem

$$\min_{x \in \mathbf{R}^d} f(x) := g(x) + h(x) \quad (32)$$

where the function  $g: \mathbf{R}^d \rightarrow \mathbf{R}$  is  $L$ -smooth and  $\mu$ -strongly convex and  $h: \mathbf{R}^d \rightarrow \mathbf{R} \cup \{+\infty\}$  is closed and convex. We denote the minimizer of  $f$  by  $\bar{x}$ , its minimal value by  $f^* := \min f$ , and its condition number by  $\kappa := L/\mu$ .

In particular, we may model minimization of a smooth and strongly convex function  $g$  over a closed convex set  $\mathcal{X}$  by declaring  $h$  to take value zero on  $\mathcal{X}$  and  $+\infty$  off it.

Before stating the proposed algorithm, we require three elementary ingredients: a two-sided bound akin to (7), robust distance estimation with a “pseudometric,” and a robust gradient estimator.

### 6.2.1 THE TWO-SIDED BOUND

Observe that optimality conditions at  $\bar{x}$  imply the inclusion  $-\nabla g(\bar{x}) \in \partial h(\bar{x})$ , where  $\partial h(\bar{x})$  denotes the *subdifferential* of  $h$  at  $\bar{x}$ , and therefore the nonnegativity of the quantity

$$D_h(x, \bar{x}) := h(x) - h(\bar{x}) + \langle \nabla g(\bar{x}), x - \bar{x} \rangle.$$

Indeed, optimization specialists may recognize  $D_h(x, \bar{x})$  as a Bregman divergence induced by  $h$ —hence the notation. The term  $D_h(x, \bar{x})$  appears naturally in a two-sided bound similar to (7). Specifically, adding  $h(x) - h(\bar{x})$  to the two-sided bound (30) throughout, we obtain the key two-sided estimate

$$\boxed{D_h(x, \bar{x}) + \frac{\mu}{2} \|x - \bar{x}\|^2 \leq f(x) - f^* \leq D_h(x, \bar{x}) + \frac{L}{2} \|x - \bar{x}\|^2.} \quad (33)$$

### 6.2.2 ROBUST DISTANCE ESTIMATION WITH A PSEUDOMETRIC

As the second ingredient, we will require a slight modification of the robust distance estimation technique of (Nemirovsky and Yudin, 1983, p. 243) and Hsu and Sabato (2016). In particular, it will be convenient to replace the Euclidean norm  $\|\cdot\|$  with a more general distance measure.

**Definition 10 (Pseudometric)** A mapping  $\rho: \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  is a *pseudometric* on a set  $\mathcal{X}$  if for all  $x, y, z \in \mathcal{X}$  it satisfies:

1. (nonnegative)  $\rho(x, y) \geq 0$  and  $\rho(x, x) = 0$ ,
2. (symmetry)  $\rho(x, y) = \rho(y, x)$ ,
3. (triangle inequality)  $\rho(x, y) \leq \rho(x, z) + \rho(z, y)$ .

The symbol  $B_r^\rho(x) = \{y \in \mathcal{X} : \rho(x, y) \leq r\}$  will denote the  $r$ -radius ball around  $x$  in the pseudometric  $\rho$ .

With this notation, we record Algorithm 8, which is in the same spirit as the robust distance estimator, Algorithm 1. The differences are that the Euclidean norm is replaced with a pseudometric  $\rho$ , an index set is returned instead of a single point, and we leave the origin of the vectors  $y_i$  unspecified for the moment.

We will need the following elementary lemma, akin to Lemma 1. The main difference is that the lemma provides at least  $m/2$  points, instead of a single point, that are close to the target with high probability. The proof is identical to that of Nemirovsky and Yudin (1983, p. 243) and Hsu and Sabato (2016, Propositions 8 and 9); we provide details for the sake of completeness.



<p><b>Algorithm 8:</b> <math>\text{Extract}(\{y_i\}_{i=1}^m, \rho)</math></p> <p><b>Input:</b> A set of <math>m</math> points <math>Y = \{y_1, \dots, y_m\} \subset \mathcal{X}</math>, a pseudometric <math>\rho</math> on <math>\mathcal{X}</math>.</p> <p><b>Step</b> <math>i = 1, \dots, m</math>:              Compute <math>r_i = \min\{r \geq 0 :  B_r^\rho(y_i) \cap Y  &gt; \frac{m}{2}\}</math>.</p> <p>    Compute the median <math>\hat{r} = \text{median}(r_1, \dots, r_m)</math>.</p> <p><b>Return</b> <math>\mathcal{I} = \{i \in [1, m] : r_i \leq \hat{r}\}</math>.</p>
--

**Lemma 11 (Robust distance estimation)** *Let  $\rho$  be a pseudometric on a set  $\mathcal{X}$ . Consider a set of points  $Y = \{y_1, \dots, y_m\} \subset \mathcal{X}$  and a point  $\bar{y} \in \mathcal{X}$  satisfying  $|B_\varepsilon^\rho(\bar{y}) \cap Y| > \frac{m}{2}$  for some  $\varepsilon > 0$ . Then the index set  $\mathcal{I} = \text{Extract}(\{y_i\}_{i=1}^m, \rho)$  satisfies the guarantee*

$$\rho(y_i, \bar{y}) \leq 3\varepsilon \quad \text{for all } i \in \mathcal{I}.$$

**Proof** Note that for any points  $y_i, y_j \in B_\varepsilon^\rho(\bar{y})$ , the triangle inequality implies the estimate

$$\rho(y_i, y_j) \leq \rho(y_i, \bar{y}) + \rho(\bar{y}, y_j) \leq 2\varepsilon.$$

This means that any point  $y_i \in B_\varepsilon^\rho(\bar{y})$ , at least  $\frac{m}{2}$  of them, satisfies  $|B_{2\varepsilon}^\rho(y_i) \cap Y| > \frac{m}{2}$  and consequently  $r_i \leq 2\varepsilon$ . Therefore the inequality  $\hat{r} = \text{median}(r_1, \dots, r_m) \leq 2\varepsilon$  holds.

Fix an index  $i \in \mathcal{I}$ . Since both  $B_{r_i}^\rho(y_i)$  and  $B_\varepsilon^\rho(\bar{y})$  contain a strict majority of the points in  $Y$ , there must exist some point in the intersection  $y \in B_{r_i}^\rho(y_i) \cap B_\varepsilon^\rho(\bar{y})$ . Using the triangle inequality, we conclude  $\rho(y_i, \bar{y}) \leq \rho(y_i, y) + \rho(y, \bar{y}) \leq 3\varepsilon$ , thereby completing the proof. ■

### 6.2.3 ROBUST GRADIENT ESTIMATOR

The need for this last ingredient is explained at the end of Section 6.1. Namely, we will need to estimate the gradient  $\nabla g(\bar{x})$ , thereby perturbing the term  $D_h(x, \bar{x})$  in the two-sided bound (30). For this purpose, we make the following mild assumption that is standard in applications. Indeed, we already encountered this assumption when paring `proxBoost` with stochastic gradient methods for unconstrained optimization.

**Assumption 5 (Stochastic first-order oracle)** Fix a probability space  $(\Omega, \mathcal{F}, \mathcal{P})$  and let  $G: \mathbf{R}^d \times \Omega \rightarrow \mathbf{R}$  be a measurable map satisfying

$$\mathbb{E}_z G(x, z) = \nabla g(x) \quad \text{and} \quad \mathbb{E}_z \|G(x, z) - \nabla g(x)\|^2 \leq \sigma^2.$$

We suppose that for any point  $x$ , we may sample  $z \in \Omega$  and compute the vector  $G(x, z)$ , which serves as an unbiased estimator of the gradient  $\nabla g(x)$ .

Under this assumption, we can define a *weak gradient oracle*  $\mathcal{G}_\sigma(\cdot, \varepsilon)$  as the average of a finite sample of stochastic gradients, i.e., for any  $\hat{x} \in \mathcal{X}$ ,

$$\mathcal{G}_\sigma(\hat{x}, \varepsilon) := \frac{1}{s} \sum_{i=1}^s G(\hat{x}, z_i) \quad \text{where} \quad s = \left\lceil \frac{3\sigma^2}{\varepsilon^2} \right\rceil.$$

Taking into account the variance reduction by a factor of  $s$  and using Markov's inequality, we have

$$\mathbb{P}\left(\left\|\frac{1}{s}\sum_{i=1}^s G(\hat{x}, z_i) - \nabla g(\hat{x})\right\|^2 \geq \varepsilon^2\right) \leq \frac{\sigma^2/s}{\varepsilon^2} \leq \frac{1}{3}.$$

That is,  $\mathbb{P}(\|\mathcal{G}_\sigma(\hat{x}, \varepsilon) - \nabla g(\hat{x})\| < \varepsilon) \geq \frac{2}{3}$ , confirming that  $\mathcal{G}_\sigma(\cdot, \varepsilon)$  is indeed a weak distance oracle in the sense of (9). Based on this oracle, we can use Algorithm 1 to construct a *robust gradient estimator*  $\mathcal{G}_\sigma(\cdot, \varepsilon, m)$ , which returns an estimate  $\tilde{\nabla}g(\cdot)$ . By Lemma 1,

$$\mathbb{P}(\|\tilde{\nabla}g(\hat{x}) - \nabla g(\hat{x})\| \leq 3\varepsilon) \geq 1 - \exp(-m/18). \quad (34)$$

#### 6.2.4 ROBUST FUNCTION GAP ESTIMATION

Equipped with the three ingredients described above, we present in Algorithm 9 a procedure to robustly estimate the gap  $f(x) - f^*$ .

<p><b>Algorithm 9:</b> <math>\text{RobustGap}(\mathcal{M}_f(\cdot), m, \epsilon)</math>.</p> <p><b>Input:</b> Minimization oracle <math>\mathcal{M}_f(\cdot)</math>, integer <math>m \in \mathbb{N}</math>, accuracy <math>\epsilon &gt; 0</math>.</p> <p><b>Step 1:</b> Independently generate <math>x_1, \dots, x_m</math> by the oracle <math>\mathcal{M}_f(\epsilon)</math> so that</p> $\mathbb{P}(f(x_i) - f^* \leq \epsilon) \geq \frac{2}{3}, \quad \text{for all } i \in [1, m]. \quad (35)$ <p><b>Step 2:</b> Set <math>\rho_1 = \ \cdot\ </math> to be the usual Euclidean norm and compute</p> $\mathcal{I}_1 := \text{Extract}(\{x_i\}_{i=1}^m, \rho_1). \quad (36)$ <p><b>Step 3:</b> Fix arbitrary <math>i \in \mathcal{I}_1</math> and set <math>\hat{x} := x_i</math>. Use the robust gradient estimator to generate</p> $\tilde{\nabla}g(\hat{x}) := \mathcal{G}_\sigma(\hat{x}, \kappa\sqrt{\mu\epsilon}, m).$ <p><b>Step 4:</b> Define the pseudometric <math>\rho_2(x, x') :=  h(x) - h(x') + \langle \tilde{\nabla}g(\hat{x}), x - x' \rangle </math> on <math>\text{dom } h</math> and compute</p> $\mathcal{I}_2 = \text{Extract}(\{x_t\}_{t=1}^m, \rho_2).$ <p><b>Return:</b> <math>x_i</math> for an arbitrary <math>i \in \mathcal{I}_1 \cap \mathcal{I}_2</math>.</p>
---

Thus, the first step of  $\text{RobustGap}(\mathcal{M}_f(\cdot), m, \epsilon)$  generates  $m$  statistically independent points  $x_1, \dots, x_m$  satisfying (35). The second step determines a set of points  $\{x_i\}_{i \in \mathcal{I}}$  that are all close to  $\bar{x}$  with high probability. We then choose a distinguished point  $\hat{x} := x_i$  for an arbitrary  $i \in \mathcal{I}_1$ , and estimate the gradient  $\nabla g(\hat{x})$  with  $\tilde{\nabla}g(\hat{x})$ . The next step approximates  $D_h(\cdot, \bar{x})$  with a pseudometric  $\rho_2$  by replacing  $\nabla g(\bar{x})$  with  $\tilde{\nabla}g(\hat{x})$ , and then performs robust distance estimation to find a set of points  $\{x_i\}_{i \in \mathcal{I}_2}$  with low value of  $D_h(x_i, \bar{x})$ . Finally a point  $x_i$  is returned, for any  $i \in \mathcal{I}_1 \cap \mathcal{I}_2$ . The intuition is that this  $x_i$  simultaneously achieves low values of  $\|x_i - \bar{x}\|$  and  $D_h(x_i, \bar{x})$ , thus allowing us to use (33) for robust gap estimation.

The following theorem summarizes the guarantees of the  $\text{RobustGap}$  procedure.

**Theorem 12 (Robust function gap estimation)** *With probability at least  $1 - 2 \exp(-\frac{m}{18})$ , the point  $x = \text{RobustGap}(\mathcal{M}_f(\cdot), m, \epsilon)$  satisfies the guarantee*

$$\|x - \bar{x}\| \leq 3\sqrt{\frac{2\epsilon}{\mu}}, \quad D_h(x, \bar{x}) \leq 65\kappa\epsilon, \quad f(x) - f^* \leq 74\kappa\epsilon.$$

*In total, the procedure queries  $m$  times the oracle  $\mathcal{M}_f(\epsilon)$  and evaluates  $m \cdot \left\lceil \frac{3\sigma^2}{\kappa^2\mu\epsilon} \right\rceil$  times the stochastic gradient oracle  $G(\hat{x}, \cdot)$ .*

**Proof** Define the index set  $\mathcal{J} = \{i \in [1, m] : f(x_i) - f^* \leq \epsilon\}$  and define the event

$$E := \left\{ |\mathcal{J}| > \frac{m}{2} \right\}.$$

Hoeffding's inequality for Bernoulli random variables guarantees

$$\mathbb{P}(E) \geq 1 - \exp(-m/18).$$

Moreover, using the left inequality in (33), we deduce

$$\|x_i - \bar{x}\| \leq \sqrt{\frac{2\epsilon}{\mu}} \quad \text{and} \quad D_h(x_i, \bar{x}) \leq \epsilon \quad \text{for all } i \in \mathcal{J}. \quad (37)$$

Henceforth, suppose that the event  $E$  occurs. Then Lemma 11 implies

$$\|x_i - \bar{x}\| \leq 3\sqrt{\frac{2\epsilon}{\mu}} \quad \text{for all } i \in \mathcal{I}_1. \quad (38)$$

As discussed in Section 6.2.3, specifically (34), the estimate  $\tilde{\nabla}g(\hat{x})$  generated by the robust gradient estimator  $\mathcal{G}_\sigma(\hat{x}, \kappa\sqrt{\mu\epsilon}, m)$  satisfies

$$\mathbb{P}(\|\tilde{\nabla}g - \nabla g(\hat{x})\| \leq 3\kappa\sqrt{\mu\epsilon} \mid E) \geq 1 - \exp(-m/18). \quad (39)$$

Define the event  $\hat{E} := \{\|\tilde{\nabla}g(\hat{x}) - \nabla g(\hat{x})\| \leq 3\kappa\sqrt{\mu\epsilon}\}$  and suppose that  $E \cap \hat{E}$  occurs. Then, we compute

$$\begin{aligned} \|\tilde{\nabla}g(\hat{x}) - \nabla g(\bar{x})\| &\leq \|\tilde{\nabla}g(\hat{x}) - \nabla g(\hat{x})\| + \|\nabla g(\hat{x}) - \nabla g(\bar{x})\| \\ &\leq 3\kappa\sqrt{\mu\epsilon} + L\|\hat{x} - \bar{x}\| \end{aligned} \quad (40)$$

$$\leq 3\kappa\sqrt{\mu\epsilon} + 3L\sqrt{2\epsilon/\mu} = 3(1 + \sqrt{2})\kappa\sqrt{\mu\epsilon}, \quad (41)$$

where (40) follows from (39) and Lipschitz continuity of  $\nabla g$ , while (41) follows from (37). Consequently, for each index  $i \in \mathcal{J}$ , we successively deduce

$$\begin{aligned} \rho_2(x_i, \bar{x}) &= |h(x_i) - h(\bar{x}) + \langle \tilde{\nabla}g(\hat{x}), x_i - \bar{x} \rangle| \\ &\leq D_h(x_i, \bar{x}) + |\langle \tilde{\nabla}g(\hat{x}) - \nabla g(\bar{x}), x_i - \bar{x} \rangle| \\ &\leq \epsilon + 3(1 + \sqrt{2})\kappa\sqrt{\mu\epsilon} \cdot \sqrt{2\epsilon/\mu} \\ &= (1 + (3\sqrt{2} + 6)\kappa)\epsilon, \end{aligned} \quad (42)$$

where (42) follows from (37) and (41). Therefore, appealing to Lemma 1 in the event  $E \cap \hat{E}$ , we conclude

$$\rho_2(x_i, \bar{x}) \leq 3(1 + (3\sqrt{2} + 6)\kappa)\epsilon \quad \text{for all } i \in \mathcal{I}_2. \quad (43)$$

Finally, fix an arbitrary index  $i \in \mathcal{I}_1 \cap \mathcal{I}_2$ . We therefore deduce

$$\begin{aligned} D_h(x_i, \bar{x}) &\leq \rho_2(x_i, \bar{x}) + |\langle \nabla g(\bar{x}) - \tilde{\nabla} g, x_i - \bar{x} \rangle| \\ &\leq 3(1 + (3\sqrt{2} + 6)\kappa)\epsilon + 3(1 + \sqrt{2})\kappa\sqrt{\mu\epsilon} \cdot 3\sqrt{\frac{2\epsilon}{\mu}} \end{aligned} \quad (44)$$

$$= 3(1 + (6\sqrt{2} + 12)\kappa)\epsilon \leq 65\kappa\epsilon, \quad (45)$$

where (44) follows from the estimates (38), (41), and (43). Using the right side of (33), we therefore conclude

$$f(x_i) - f(\bar{x}) \leq D_h(x_i, \bar{x}) + \frac{L}{2}\|x_i - \bar{x}\|^2 \leq 65\kappa\epsilon + 9\kappa\epsilon = 74\kappa\epsilon,$$

where the last inequality follows from the estimates (38) and (45). Noting

$$\mathbb{P}(E \cap \hat{E}) = \mathbb{P}(\hat{E} \mid E)\mathbb{P}(E) \geq (1 - \exp(-\frac{m}{18})) (1 - \exp(-\frac{m}{18})) \geq 1 - 2\exp(-\frac{m}{18}),$$

completes the proof. ■

With Theorem 12 at hand, we can now replace robust distance estimation with `RobustGap` within the `proxBoost` framework, thereby making `proxBoost` applicable to convex composite problems. The following two sections illustrate the consequences of the resulting method for regularized empirical risk minimization and (proximal) stochastic approximation algorithms.

### 6.3 Consequences for empirical risk minimization

In this section, we explore the consequences of `RobustGap` and `proxBoost` for regularized empirical risk minimization. In particular, we will boost the low-probability guarantees developed in the seminal work of Shalev-Shwartz et al. (2009) for strongly convex problems to high confidence outcomes. The following assumption summarizes the setting of this section.

**Assumption 6** Fix a probability space  $(\Omega, \mathcal{F}, \mathcal{P})$  and equip  $\mathbf{R}^d$  with the Borel  $\sigma$ -algebra. Throughout, we consider the optimization problem

$$\min_x f(x) := g(x) + h(x) \quad \text{where} \quad g(x) = \mathbb{E}_{z \sim \mathcal{P}}[g(x, z)],$$

under the following assumptions.

1. **(Measurability)** The function  $g: \mathbf{R}^d \times \Omega \rightarrow \mathbf{R}$  is measurable.
2. **(Strong convexity)** The function  $h: \mathbf{R}^d \rightarrow \mathbf{R} \cup \{+\infty\}$  is convex and there exists  $\mu > 0$  such that the function  $g(x, z) + h(x)$  is  $\mu$ -strongly convex for a.e.  $z \sim \mathcal{P}$ .

3. (**Lipschitz continuity**) There exists a measurable map  $\ell: \Omega \rightarrow \mathbf{R}$  and a real  $\bar{\ell} > 0$  satisfying the moment bound  $\sqrt{\mathbb{E}_z \ell(z)^2} \leq \bar{\ell}$  and the Lipschitz condition

$$|g(x, z) - g(y, z)| \leq \ell(z) \|x - y\| \quad \forall x \in U, z \in \Omega,$$

where  $U$  is some open neighborhood of  $\text{dom } h$ .

4. (**Smoothness**) The function  $g: \mathbf{R}^d \rightarrow \mathbf{R}$  is  $L$ -smooth.

The first three assumptions are slight modifications of those used by Shalev-Shwartz et al. (2009), while the additional smoothness assumption on  $g$  will be necessary in the sequel to obtain high-confidence guarantees. That being said, the sample efficiency will depend only polylogarithmically on  $L$ , and therefore we will be able to treat nonsmooth loss function  $g(x, z)$  using standard smoothing techniques. Under the first three assumptions, Shalev-Shwartz et al. (2009) obtained the following guarantee for the accuracy of empirical risk minimization.

**Lemma 13 (Shalev-Shwartz et al. (2009, Theorem 6))** *Let the set  $S \subset \Omega$  consist of  $n$  i.i.d. samples drawn from  $\mathcal{P}$ . Then the minimizer of the regularized empirical risk*

$$x_S := \operatorname{argmin}_x \frac{1}{n} \sum_{z \in S} g(x, z) + h(x)$$

*satisfies the generalization bound*

$$\mathbb{E}_S [f(x_S) - f^*] \leq \frac{2\bar{\ell}^2}{\mu n}.$$

We will see now how to equip this guarantee with a high confidence bound using `proxBoost`. Recall that to apply the `RobustGap` algorithm, we require an unbiased gradient estimator  $G: \mathbf{R}^d \times \Omega \rightarrow \mathbf{R}^d$  for  $g$ . Let us therefore simply declare

$$G(x, z) := \nabla g(x, z).$$

Then we can upper-bound the variance by the second moment

$$\mathbb{E}_z \|G(x, z) - \nabla g(x)\|^2 \leq 2(\mathbb{E}_z \|\nabla g(x, z)\|^2 + \mathbb{E}_z \|\nabla g(x)\|^2) \leq 4\bar{\ell}^2.$$

We are now ready to present Algorithm 10 as an instantiation of the `proxBoost` procedure for *regularized* empirical risk minimization. In particular, we can still use `ERM()` in Algorithm 3 for the proximal subproblem, with the definition  $f(y, z_i) := g(y, z_i) + h(y)$ . The robust distance estimator `ERM-R()` in Algorithm 4 can be used without any change.

<b>Algorithm 10:</b> <code>BoostERM</code> ( $\delta, T, m$ )
<b>Input:</b> accuracy $\delta > 0$ , iterations $m, T \in \mathbb{N}$
Set $\lambda_{-1} = 0, x_{-1} = 0$
<b>Step</b> $j = 0, \dots, T$ :
$x_j = \text{ERM-R} \left( \frac{54\bar{\ell}^2}{(\mu + \lambda_{j-1})\delta}, m, \lambda_{j-1}, x_{j-1} \right)$
Define the minimization oracle $\mathcal{M}_T(\epsilon) := \text{ERM} \left( \frac{6\bar{\ell}^2}{(\mu + \lambda_T)\delta}, \lambda_T, x_T \right)$ .
<b>Return</b> $x_{T+1} = \text{RobustGap} \left( \mathcal{M}_T \left( \frac{\delta(\mu + \lambda_T)}{222(L + \lambda_T)} \right), m, \frac{\delta(\mu + \lambda_T)}{222(L + \lambda_T)} \right)$

Notice that in Algorithm 10, we only need to call `RobustGap` in the last cleanup stage, since the intermediate iterations of `proxBoost` only rely on distance estimates to the optimal solutions and not on the function gap. The following theorem and its corollary are immediate consequences of Theorems 3 and 12.

**Theorem 14 (Efficiency of BoostERM)** *Fix  $\delta > 0$  and integers  $T, m \in \mathbb{N}$ . Then with probability at least  $1 - (T + 3) \exp(-\frac{m}{18})$ , the point  $x_{T+1} = \text{BoostERM}(\delta, T, m)$  satisfies*

$$f(x_{T+1}) - f^* \leq \left(1 + \sum_{i=0}^T \frac{\lambda_i}{\mu + \lambda_{i-1}}\right) \delta.$$

**Corollary 15 (Efficiency of BoostERM with geometric decay)** *Fix a target accuracy  $\epsilon > 0$  and a probability of failure  $p \in (0, 1)$ . Define the algorithm parameters:*

$$T = \lceil \log_2(\kappa) \rceil, \quad m = \left\lceil 18 \ln \left( \frac{T+3}{p} \right) \right\rceil, \quad \delta = \frac{\epsilon}{4+2T}, \quad \lambda_i = \mu 2^i.$$

*Then the point  $x_{T+1} = \text{BoostERM}(\delta, T, m)$  satisfies*

$$\mathbb{P}(f(x_{T+1}) - f(x^*) \leq \epsilon) \geq 1 - p.$$

*Moreover, the total number of samples used by the algorithm is*

$$\mathcal{O} \left( \ln^2(\kappa) \ln \left( \frac{\ln(\kappa)}{p} \right) \cdot \frac{\bar{\ell}^2}{\epsilon \mu} \right). \tag{46}$$

Thus, `proxBoost` endows regularized empirical risk minimization with high confidence guarantees at an overhead cost that is only polylogarithmic in  $\kappa$  and logarithmic in  $1/p$ . In particular, observe that the sample complexity (46) established in Corollary 15 depends on the smoothness parameter  $L$  (through  $\kappa = L/\mu$ ) only polylogarithmically. Consequently, it appears plausible that if the losses  $g(\cdot, z)$  are nonsmooth, we may simply replace them by a smooth approximation and apply `BoostERM`. The price to pay should then only be polylogarithmic in the target accuracy  $\epsilon$ . Let us formally see how this can be done. To this end, we will assume that the optimization problem in question is to minimize a sum of an expectation of convex functions, a deterministic smooth and strongly convex function (e.g. squared  $\ell_2$  norm), and a nonsmooth regularizer.

**Assumption 7** Consider the optimization problem

$$\min_x f(x) := \mathbb{E}_{z \sim \mathcal{P}}[g(x, z)] + \varphi(x) + h(x) \tag{47}$$

under the following assumptions.

1. **(Measurability)** The function  $g: \mathbf{R}^d \times \Omega \rightarrow \mathbf{R}$  is measurable and the assignment  $x \mapsto g(x, z)$  is convex for a.e.  $z \in \Omega$ .
2. **(Strong convexity)** The function  $h: \mathbf{R}^d \rightarrow \mathbf{R} \cup \{+\infty\}$  is convex and there exist parameters  $\mu, \beta > 0$  such that the function  $\varphi: \mathbf{R}^d \rightarrow \mathbf{R}$  is  $\mu$ -strongly convex and  $\beta$ -smooth for a.e.  $z \in \Omega$ .

3. **(Lipschitz continuity)** There exists a measurable map  $\ell: \Omega \rightarrow \mathbf{R}$  and a real  $\bar{\ell} > 0$  satisfying the moment bound  $\sqrt{\mathbb{E}_z \ell(z)^2} \leq \bar{\ell}$  and the Lipschitz condition

$$|g(x, z) - g(y, z)| \leq \ell(z) \|x - y\| \quad \forall x \in \mathbf{R}^d, z \in \Omega.$$

The strategy we follow is to simply replace  $g(\cdot, z)$  by a smooth approximation and then apply `BoostERMCMC`. We now make precise what we mean by a smooth approximation. Assumptions of this type are classical in convex optimization; see for example Nesterov (2005) and Beck and Teboulle (2012).

**Assumption 8 (Smoothing)** Suppose that for any parameter  $\epsilon > 0$ , there exist measurable functions  $g_\epsilon: \mathbf{R}^d \times \Omega \rightarrow \mathbf{R}$  and  $\ell_\epsilon, L_\epsilon: \Omega \rightarrow \mathbf{R}_+$  such that  $g_\epsilon(\cdot, z)$  is Lipschitz continuous with constant  $\ell_\epsilon(z)$  and its gradient is Lipschitz continuous with constant  $L_\epsilon(z)$ , and the estimate holds:

$$|g(x, z) - g_\epsilon(x, z)| \leq \epsilon \quad \text{for all } x \in \mathbf{R}^d, z \in \Omega. \quad (48)$$

We suppose moreover that the moment conditions,  $\sqrt{\mathbb{E}_z \ell_\epsilon^2(z)} \leq \bar{\ell}_\epsilon$  and  $\mathbb{E}_z L_\epsilon(z) \leq \bar{L}_\epsilon$ , hold for some constants  $\bar{\ell}_\epsilon, \bar{L}_\epsilon > 0$ .

Let us look at two standard examples of smoothings of convex functions.

**Example 1 (Moreau envelope)** A classical approach to smoothing a convex function is based on the Moreau envelope (Moreau, 1965). Namely, fix a convex function  $\psi: \mathbf{R}^d \rightarrow \mathbf{R}$ . The Moreau envelope of  $\psi$  with parameter  $\nu > 0$  is defined to be

$$M_\nu^\psi(x) = \min_y \psi(y) + \frac{1}{2\nu} \|y - x\|^2.$$

It is well-known that  $M_\nu^\psi$  is  $\frac{1}{\nu}$ -smooth. Moreover if  $\psi$  is Lipschitz continuous with constant  $\text{lip}(\psi)$ , then  $M_\nu^\psi$  is also Lipschitz continuous with the same constant and the bound holds:

$$0 \leq \psi(x) - \psi_\nu(x) \leq \nu(\text{lip}(\psi))^2.$$

Coming back to our target problem (47), we may define  $g_\epsilon(\cdot, z)$  to be the Moreau envelope of  $g(\cdot, z)$  with parameter  $\nu(z) := \frac{\epsilon}{\bar{l}(z)^2}$ , or more explicitly

$$g_\epsilon(x, z) = \min_y g(y, z) + \frac{l(z)^2}{2\epsilon} \|y - x\|^2.$$

Then the parameters from Assumption 7 become  $\ell_\epsilon(z) := \ell(z)$  and  $L_\epsilon(z) := \frac{l(z)^2}{\epsilon}$ .

**Example 2 (Compositional smoothing)** Often, the Moreau envelope of  $g(\cdot, z)$  may be difficult to compute explicitly. In typical circumstances, however, the function  $g(\cdot, z)$  may be written as a composition of a simple nonsmooth convex function with a linear map. It then suffices to replace only the outer function with its Moreau envelope—a technique famously explored by Nesterov (2005).

To illustrate on a concrete example, suppose that the population data consists of tuples  $z = (a, b) \sim \mathcal{P}$  and the loss takes the form  $g(x, z) = h(\langle a, x \rangle, b)$  for some measurable function

$h(\cdot, \cdot)$  that is convex and 1-Lipschitz in its first argument. In order to control the Lipschitz constant, suppose also the moment bound  $\sqrt{\mathbb{E}_a \|a\|^2} \leq A$  for some constant  $A > 0$ . Let us now define the smoothing

$$g_\epsilon(x, z) = h_\epsilon(\langle a, x \rangle, b),$$

where  $h_\epsilon(\cdot, b)$  is the Moreau envelope of  $h(\cdot, b)$  with parameter  $\nu = \epsilon$ . It is straightforward to verify that the estimate (48) holds and that we may set  $\ell_\epsilon(z) = A$  and  $L_\epsilon(z) = \frac{A^2}{\epsilon}$ .

With Assumptions 7 and 8 at hand, we may now simply apply **BoostERM**C to the smoothed problem

$$\min_x f_\epsilon(x) := g(x) + h(x) \quad \text{where} \quad g(x) = \mathbb{E}_{z \sim \mathcal{P}} [g_\epsilon(x, z)] + \varphi(x).$$

Using Corollary 15, we deduce that the procedure will find a point  $x$  satisfying

$$\mathbb{P}(f_\epsilon(x) - f_\epsilon^* \leq \epsilon) \geq 1 - p,$$

using

$$\mathcal{O} \left( \ln^2 \left( \frac{\bar{L}_\epsilon + \beta}{\mu} \right) \ln \left( \frac{\ln \left( \frac{\bar{L}_\epsilon + \beta}{\mu} \right)}{p} \right) \cdot \frac{\bar{\ell}_\epsilon^2}{\epsilon \mu} \right)$$

samples. Observe that with probability  $1 - p$  the returned point  $x$  satisfies:

$$f(x) - f^* \leq f_\epsilon(x) - f_\epsilon^* + (f(x) - f_\epsilon(x)) + (f_\epsilon^* - f^*) \leq 3\epsilon.$$

In particular, in the setup of Examples 1 and 2, the sample complexities become:

$$\mathcal{O} \left( \ln^2 \left( \frac{\bar{\ell}^2 / \epsilon + \beta}{\mu} \right) \ln \left( \frac{\ln \left( \frac{\bar{\ell}^2 / \epsilon + \beta}{\mu} \right)}{p} \right) \cdot \frac{\bar{\ell}^2}{\epsilon \mu} \right) \quad \text{and} \quad \mathcal{O} \left( \ln^2 \left( \frac{A^2 / \epsilon + \beta}{\mu} \right) \ln \left( \frac{\ln \left( \frac{A^2 / \epsilon + \beta}{\mu} \right)}{p} \right) \cdot \frac{A^2}{\epsilon \mu} \right),$$

respectively. Hence, the price to pay for nonsmoothness is only polylogarithmic in  $1/\epsilon$ .

#### 6.4 Consequences for stochastic approximation

We now extend the results of Section 5 to the convex composite setting. In addition to Assumptions 4 and 5, in this section we will use the following composite analogue of Assumption 3. At the end of the section, we will let  $\mathbf{Alg}(\cdot)$  be the (accelerated) proximal stochastic gradient method.

**Assumption 9** Consider the proximal minimization problem

$$\min_y \varphi_x(y) := g(y) + \frac{\lambda}{2} \|y - x\|^2 + h(y),$$

Let  $\Delta > 0$  be a real number satisfying  $\varphi_x(x) - \min \varphi_x \leq \Delta$ . We will let  $\mathbf{Alg}(\delta, \lambda, \Delta, x)$  be a procedure that returns a point  $y$  satisfying

$$\mathbb{P}[\varphi_x(y) - \min \varphi_x \leq \delta] \geq \frac{2}{3}.$$



The following algorithm is a direct extension of **BoostAlg** (Algorithm 6) to the convex composite setting; the only difference is that **BoostAlgC** (Algorithm 9) replaces the distance estimator **Alg-R**( $\cdot$ ) with **RobustGap**( $\cdot$ ).

<p><b>Algorithm 11:</b> <b>BoostAlgC</b>(<math>\delta, \Delta_{\text{in}}, x_{\text{in}}, T, m</math>)</p> <p><b>Input:</b> accuracy <math>\delta &gt; 0</math>, upper bound <math>\Delta_{\text{in}} &gt; 0</math>, initial <math>x_{\text{in}} \in \mathbf{R}^d</math>, numbers <math>m, T \in \mathbb{N}</math>                  Set <math>\lambda_{-1} = 0, \Delta_{-1} = \Delta_{\text{in}}, x_{-1} = x_{\text{in}}</math>  <b>Step</b> <math>j = 0, \dots, T</math>:                      Define the minimization oracle for the proximal subproblem</p> $\mathcal{M}_{j-1}(\cdot) := \text{Alg}(\cdot, \lambda_{j-1}, \Delta_{j-1}, x_{j-1}).$ <p>    Set <math>x_j = \text{RobustGap}(\mathcal{M}_{j-1}(\delta/9), m, \delta/9)</math>                      Set <math>\Delta_j = \delta \left( 9 \cdot \frac{L + \lambda_{j-1}}{\mu + \lambda_{j-1}} + \sum_{i=0}^{j-1} \frac{\lambda_i}{\mu + \lambda_{i-1}} \right)</math>  <b>Return</b> <math>x_{T+1} = \text{RobustGap} \left( \mathcal{M}_T \left( \frac{\delta(\mu + \lambda_T)}{74(L + \lambda_T)} \right), m, \frac{\delta(\mu + \lambda_T)}{74(L + \lambda_T)} \right)</math></p>
---

In **BoostAlgC**, we need to use **RobustGap** in every proximal iteration as well as the cleanup step, because the stochastic proximal gradient method encoded as **Alg** typically requires robust gap estimation on the initialization gap  $\Delta_j$ . The proof of the following theorem is almost identical to that of Theorem 8, with Theorem 12 playing the role of Lemma 1.

**Theorem 16 (Efficiency of BoostAlgC)** *Fix an arbitrary point  $x_{\text{in}} \in \mathbf{R}^d$  and let  $\Delta_{\text{in}}$  be any upper bound  $\Delta_{\text{in}} \geq f(x_{\text{in}}) - \min f$ . Fix natural numbers  $T, m \in \mathbb{N}$ . Then with probability at least  $1 - 2(T + 2) \exp(-\frac{m}{18})$ , the point  $x_{T+1} = \text{BoostAlgC}(\delta, \Delta_{\text{in}}, x_{\text{in}}, T, m)$  satisfies*

$$f(x_{T+1}) - \min f \leq \delta \left( 1 + \sum_{i=0}^T \frac{\lambda_i}{\mu + \lambda_{i-1}} \right).$$

When using the proximal parameters  $\lambda_i = \mu 2^i$ , we obtain the following guarantee, which generalizes Corollary 9 to the composite setting.

**Corollary 17 (Efficiency of BoostAlgC with geometric decay)** *Fix an arbitrary point  $x_{\text{in}} \in \mathbf{R}^d$  and let  $\Delta_{\text{in}}$  be any upper bound  $\Delta_{\text{in}} \geq f(x_{\text{in}}) - \min f$ . Fix a target accuracy  $\epsilon > 0$  and probability of failure  $p \in (0, 1)$ , and set the algorithm parameters*

$$T = \lceil \log_2(\kappa) \rceil, \quad m = \left\lceil 18 \ln \left( \frac{4 + 2T}{p} \right) \right\rceil, \quad \delta = \frac{\epsilon}{2 + 2T}, \quad \lambda_i = \mu 2^i.$$

*Then the point  $x_{T+1} = \text{BoostAlg}(\delta, \Delta_{\text{in}}, x_{\text{in}}, T, m)$  satisfies*

$$\mathbb{P}(f(x_{T+1}) - \min f \leq \epsilon) \geq 1 - p.$$

*Moreover, the total number of calls to **Alg**( $\cdot$ ) is*

$$\left\lceil 18 \ln \left( \frac{4 + 2 \lceil \log_2(\kappa) \rceil}{p} \right) \right\rceil \lceil 2 + \log_2(\kappa) \rceil,$$

the number of evaluations of the stochastic gradient oracle  $G(\cdot, \cdot)$  is at most<sup>1</sup>

$$\left\lceil 18 \ln \left( \frac{\lceil 4 + 2 \log_2(\kappa) \rceil}{p} \right) \right\rceil \cdot \left\lceil \frac{6\sigma^2}{L\epsilon} \cdot (2 + 2 \lceil \log_2(\kappa) \rceil) \right\rceil \cdot \lceil 2 + \log_2(\kappa) \rceil,$$

and the initialization errors satisfy

$$\max_{i=0, \dots, T+1} \Delta_i \leq \frac{9\kappa + 1 + 2 \lceil \log_2(\kappa) \rceil}{2 + 2 \lceil \log_2(\kappa) \rceil} \epsilon.$$

In particular, the stochastic gradient method and its accelerated variant (Ghadimi and Lan, 2013; Kulunchakov and Mairal, 2019) admit proximal extensions with exactly the same sample complexities as in the smooth case, (25) and (26), respectively. Clearly, we may use either of these two procedures as  $\text{Alg}(\cdot)$  within Algorithm 11. Corollary 17 then immediately shows that the two resulting algorithms will find a point  $x$  satisfying  $\mathbb{P}[f(x) - f^* \leq \epsilon] \geq 1 - p$  with the same sample complexities as in the smooth setting, (27) and (28), respectively.

## Acknowledgments

Research of Davis was supported by an Alfred P. Sloan Research Fellowship. Research of Drusvyatskiy was supported by the NSF DMS 1651851 and CCF 1740551 awards and a research visiting position at Microsoft Research, Redmond, WA. The authors would like to thank Sebastian Bubeck for insightful discussions.

## References

- Z. Allen-Zhu. How to make the gradients small stochastically: Even faster convex and nonconvex sgd. In *Advances in Neural Information Processing Systems*, pages 1157–1167, 2018.
- N. Alon, Y. Matias, and M. Szegedy. The space complexity of approximating the frequency moments. *J. Comput. System Sci.*, 58(1, part 2):137–147, 1999. ISSN 0022-0000. doi: 10.1006/jcss.1997.1545. URL <https://doi-org.offcampus.lib.washington.edu/10.1006/jcss.1997.1545>. Twenty-eighth Annual ACM Symposium on the Theory of Computing (Philadelphia, PA, 1996).
- H. Asi and J.C. Duchi. Stochastic (approximate) proximal point methods: Convergence, optimality, and adaptivity. *arXiv:1810.05633*, 2018.
- H. Asi and J.C. Duchi. The importance of better models in stochastic optimization. *arXiv:1903.08619*, 2019.
- Jean-Yves Audibert, Olivier Catoni, et al. Robust linear least squares regression. *The Annals of Statistics*, 39(5):2766–2794, 2011.

---

1. For the middle term, we use the observation,  $\min_{\lambda \geq 0} \frac{(L + \lambda)^2}{\mu + \lambda} \geq \frac{L}{2}$ , which is straightforward to verify.

- Necdet Serhat Aybat, Alireza Fallah, Mert Gurbuzbalaban, and Asuman Ozdaglar. A universally optimal multistage accelerated stochastic gradient method. *arXiv preprint arXiv:1901.08022*, 2019.
- P.L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- A. Beck. *First-order methods in optimization*, volume 25. SIAM, 2017.
- A. Beck and M. Teboulle. Smoothing and first order methods: a unified framework. *SIAM J. Optim.*, 22(2):557–580, 2012. ISSN 1052-6234. doi: 10.1137/100818327. URL <http://dx.doi.org/10.1137/100818327>.
- O. Bousquet and A. Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- S. Bubeck, N. Cesa-Bianchi, and G. Lugosi. Bandits with heavy tail. *IEEE Transactions on Information Theory*, 59(11):7711–7717, 2013.
- O. Catoni. Challenging the empirical mean and empirical variance: a deviation study. *Ann. Inst. Henri Poincaré Probab. Stat.*, 48(4):1148–1185, 2012. ISSN 0246-0203. doi: 10.1214/11-AIHP454. URL <https://doi.org/10.1214/11-AIHP454>.
- Y. Chen, L. Su, and J. Xu. Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 1(2):44, 2017.
- D. Davis and D. Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.
- D. Davis and B. Grimmer. Proximally guided stochastic method for nonsmooth, nonconvex problems. *arXiv:1707.03505*, 2017.
- Damek Davis and Dmitriy Drusvyatskiy. High probability guarantees for stochastic convex optimization. In *Conference on Learning Theory*, pages 1411–1427, 2020.
- J.C. Duchi and F. Ruan. Stochastic methods for composite and weakly convex optimization problems. *SIAM Journal on Optimization*, 28(4):3229–3259, 2018.
- R. Frostig, R. Ge, S. Kakade, and A. Sidford. Un-regularizing: approximate proximal point and faster stochastic algorithms for empirical risk minimization. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, 2015.
- R. Ge, S. Kakade, R. Kidambi, and P. Netrapalli. The step decay schedule: A near optimal, geometrically decaying learning rate procedure. *arXiv:1904.12838*, 2019.
- S. Ghadimi and G. Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization I: A generic algorithmic framework. *SIAM J. Optim.*, 22(4):1469–1492, 2012. ISSN 1052-6234. doi: 10.1137/110848864. URL <https://doi-org.offcampus.lib.washington.edu/10.1137/110848864>.

- S. Ghadimi and G. Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization, II: shrinking procedures and optimal algorithms. *SIAM Journal on Optimization*, 23(4):2061–2089, 2013.
- J.L. Goffin. On convergence rates of subgradient optimization methods. *Math. Programming*, 13(3):329–347, 1977. ISSN 0025-5610. doi: 10.1007/BF01584346. URL <https://doi-org.offcampus.lib.washington.edu/10.1007/BF01584346>.
- E. Gorbunov, M. Danilova, and A. Gasnikov. Stochastic optimization with heavy-tailed noise via accelerated gradient clipping. *arXiv preprint arXiv:2005.10785*, 2020.
- Nicholas JA Harvey, Christopher Liaw, Yaniv Plan, and Sikander Randhawa. Tight analyses for non-smooth stochastic gradient descent. In *Conference on Learning Theory*, pages 1579–1613. PMLR, 2019a.
- Nicholas JA Harvey, Christopher Liaw, and Sikander Randhawa. Simple and optimal high-probability bounds for strongly-convex stochastic gradient descent. *arXiv preprint arXiv:1909.00843*, 2019b.
- Elad Hazan and Satyen Kale. Beyond the regret minimization barrier: optimal algorithms for stochastic strongly-convex optimization. *The Journal of Machine Learning Research*, 15(1):2489–2512, 2014.
- D. Hsu and S. Sabato. Loss minimization and parameter estimation with heavy tails. *The Journal of Machine Learning Research*, 17(1):543–582, 2016.
- M.R. Jerrum, L.G. Valiant, and V.V. Vazirani. Random generation of combinatorial structures from a uniform distribution. *Theoret. Comput. Sci.*, 43(2-3):169–188, 1986. ISSN 0304-3975. doi: 10.1016/0304-3975(86)90174-X. URL [https://doi-org.offcampus.lib.washington.edu/10.1016/0304-3975\(86\)90174-X](https://doi-org.offcampus.lib.washington.edu/10.1016/0304-3975(86)90174-X).
- E. Joly, G. Lugosi, and R.I. Oliveira. On the estimation of the mean of a random vector. *Electronic Journal of Statistics*, 11(1):440–451, 2017.
- A. Juditsky and Y. Nesterov. Deterministic and stochastic primal-dual subgradient algorithms for uniformly convex minimization. *Stoch. Syst.*, 4(1):44–80, 2014. ISSN 1946-5238. doi: 10.1214/10-SSY010. URL <https://doi-org.offcampus.lib.washington.edu/10.1214/10-SSY010>.
- A. Juditsky, A. Nazin, A. Nemirovsky, and A. Tsybakov. Algorithms of robust stochastic optimization based on mirror descent method. *arXiv:1907.02707*, 2019.
- A. Kulunchakov and J. Mairal. Estimate sequences for stochastic composite optimization: Variance reduction, acceleration, and robustness to noise. *arXiv:1901.08788*, 2019.
- G. Lan and Y. Zhou. An optimal randomized incremental gradient method. *Math. Program.*, 171(1-2, Ser. A):167–215, 2018. ISSN 0025-5610. doi: 10.1007/s10107-017-1173-0. URL <https://doi.org/10.1007/s10107-017-1173-0>.

- H. Lin, J. Mairal, and Z. Harchaoui. A universal catalyst for first-order optimization. In *Advances in Neural Information Processing Systems*, pages 3366–3374, 2015.
- G. Lugosi and S. Mendelson. Risk minimization by median-of-means tournaments. *arXiv:1608.00757*, 2016.
- G. Lugosi and S. Mendelson. Sub-gaussian estimators of the mean of a random vector. *The Annals of Statistics*, 47(2):783–794, 2019.
- B. Martinet. Régularisation d’inéquations variationnelles par approximations successives. *Rev. Française Informat. Rech. Opérationnelle*, 4(Sér. R-3):154–158, 1970.
- B. Martinet. Détermination approchée d’un point fixe d’une application pseudo-contractante. Cas de l’application prox. *C. R. Acad. Sci. Paris Sér. A-B*, 274:A163–A165, 1972.
- S. Minsker. Geometric median and robust estimation in Banach spaces. *Bernoulli*, 21(4):2308–2335, 2015. ISSN 1350-7265. doi: 10.3150/14-BEJ645. URL <https://doi-org.offcampus.lib.washington.edu/10.3150/14-BEJ645>.
- J.-J. Moreau. Proximité et dualité dans un espace hilbertien. *Bull. Soc. Math. France*, 93:273–299, 1965. ISSN 0037-9484. URL [http://www.numdam.org.offcampus.lib.washington.edu/item?id=BSMF\\_1965\\_\\_93\\_\\_273\\_0](http://www.numdam.org.offcampus.lib.washington.edu/item?id=BSMF_1965__93__273_0).
- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM J. Optim.*, 19(4):1574–1609, 2008. ISSN 1052-6234. doi: 10.1137/070704277. URL <https://doi-org.offcampus.lib.washington.edu/10.1137/070704277>.
- A. S. Nemirovsky and D. B. Yudin. *Problem complexity and method efficiency in optimization*. A Wiley-Interscience Publication. John Wiley & Sons, Inc., New York, 1983. ISBN 0-471-10345-4. Translated from the Russian and with a preface by E. R. Dawson, Wiley-Interscience Series in Discrete Mathematics.
- Yu. Nesterov. Smooth minimization of non-smooth functions. *Math. Program.*, 103(1, Ser. A):127–152, 2005. ISSN 0025-5610. doi: 10.1007/s10107-004-0552-5. URL <http://dx.doi.org/10.1007/s10107-004-0552-5>.
- Yu. Nesterov. *Lectures on convex optimization*, volume 137 of *Springer Optimization and Its Applications*. Springer, 2018. ISBN 978-3-319-91577-7; 978-3-319-91578-4. doi: 10.1007/978-3-319-91578-4. URL <https://doi.org/10.1007/978-3-319-91578-4>.
- Yu. Nesterov and J.-P. Vial. Confidence level solutions for stochastic programming. *Automatica*, 44:1559–1568, 2008.
- B.T. Polyak and A.B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM J. Control Optim.*, 30(4):838–855, 1992. ISSN 0363-0129. doi: 10.1137/0330046. URL <https://doi.org/10.1137/0330046>.

- R.T. Rockafellar. Monotone operators and the proximal point algorithm. *SIAM J. Control Optimization*, 14(5):877–898, 1976. ISSN 0363-0129. doi: 10.1137/0314056. URL <https://doi-org.offcampus.lib.washington.edu/10.1137/0314056>.
- S. Shalev-Shwartz and S. Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- S. Shalev-Shwartz and T. Zhang. Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization. *Math. Program.*, 155(1-2, Ser. A):105–145, 2016. ISSN 0025-5610. doi: 10.1007/s10107-014-0839-0. URL <https://doi.org/10.1007/s10107-014-0839-0>.
- S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Stochastic convex optimization. In *COLT*, 2009.
- A. Shapiro and A. Nemirovski. On complexity of stochastic programming problems. In V. Jeyakuman and A. M. Rubinov, editors, *Continuous optimization: Current trends and applications*, pages 111–144. Springer, 2005.
- Y. Xu, Q. Lin, and T. Yang. Accelerated stochastic subgradient methods under local error bound condition. *arXiv preprint arXiv:1607.01027*, 2016.
- T. Yang, Y. Yan, Z. Yuan, and R. Jin. Why does stagewise training accelerate convergence of testing error over SGD? *arXiv:1812.03934*, 2018.
- D. Yin, Y. Chen, R. Kannan, and P. Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5650–5659, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL <http://proceedings.mlr.press/v80/yin18a.html>.