

# Adaptive Sequential Optimization with Applications to Machine Learning

Craig Wilson and Venugopal V. Veeravalli\*  
 Coordinated Science Lab and Electrical and Computer Engineering  
 University of Illinois at Urbana-Champaign  
 Urbana, IL 61801, USA  
 {wilson60, vvv}@illinois.edu

October 25, 2018

## Abstract

A framework is introduced for solving a sequence of slowly changing optimization problems, including those arising in regression and classification applications, using optimization algorithms such as stochastic gradient descent (SGD). The optimization problems change slowly in the sense that the minimizers change at either a fixed or bounded rate. A method based on estimates of the change in the minimizers and properties of the optimization algorithm is introduced for adaptively selecting the number of samples needed from the distributions underlying each problem in order to ensure that the excess risk, i.e., the expected gap between the loss achieved by the approximate minimizer produced by the optimization algorithm and the exact minimizer, does not exceed a target level. Experiments with synthetic and real data are used to confirm that this approach performs well.

## 1 Introduction

Consider solving a sequence of machine learning problems such as regression or classification by minimizing the expected value of a fixed loss function  $\ell(\mathbf{x}, \mathbf{z})$  at each time  $ns$ :

$$\min_{\mathbf{x} \in \mathcal{X}} \left\{ f_n(\mathbf{x}) \triangleq \mathbb{E}_{\mathbf{z}_n \sim p_n} [\ell(\mathbf{x}, \mathbf{z}_n)] \right\} \quad \forall n \geq 1 \quad (1)$$

For regression,  $\mathbf{z}_n$  corresponds to the predictors and response pair at time  $n$  and  $\mathbf{x}$  parameterizes the regression model. For classification  $\mathbf{z}_n$  corresponds to the feature and label pair at time  $n$  and  $\mathbf{x}$  parameterizes the classifier. Although, motivated by regression and classification, our framework works for any loss function  $\ell(\mathbf{x}, \mathbf{z})$  that satisfies certain properties discussed later. In the learning context, a *task* consists of the loss function  $\ell(\mathbf{x}, \mathbf{z})$  and the distribution  $p_n$ , and so our problem can be viewed as learning a sequence of tasks.

The problems change slowly at a constant but unknown rate in the sense that

$$\|\mathbf{x}_n^* - \mathbf{x}_{n-1}^*\| = \rho \quad \forall n \geq 2 \quad (2)$$

with  $\mathbf{x}_n^*$  the minimizer of  $f_n(\mathbf{x})$ . In an extended version of this paper [?], we also consider slow changes at a bounded but unknown rate

$$\|\mathbf{x}_n^* - \mathbf{x}_{n-1}^*\| \leq \rho \quad \forall n \geq 2 \quad (3)$$

Under this model, we find approximate minimizers  $\mathbf{x}_n$  of each function  $f_n(\mathbf{x})$  using  $K_n$  samples from distribution  $p_n$  by applying an optimization algorithm. We evaluate the quality of our approximate minimizers  $\mathbf{x}_n$  through an excess risk criterion  $\varepsilon_n$ , i.e.,

$$\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*) \leq \varepsilon_n$$

---

\*This work was supported by the NSF under award CCF 11-11342 through the University of Illinois at Urbana-Champaign.

which is a standard criterion for optimization and learning problems [1]. Our goal is to determine adaptively the number of samples  $K_n$  required to achieve a desired excess risk  $\varepsilon$  for each  $n$  with  $\rho$  unknown. As  $\rho$  is unknown, we will construct estimates of  $\rho$ . Given an estimate of  $\rho$ , we determine selection rules for the number of samples  $K_n$  to achieve a target excess risk  $\varepsilon$ .

## 1.1 Related Work

Our problem has connections with *multi-task learning* (MTL) and *transfer learning*. In multi-task learning, one tries to learn several tasks simultaneously as in [2],[3], and [4] by exploiting the relationships between the tasks. In transfer learning, knowledge from one source task is transferred to another target task either with or without additional training data for the target task [5]. Multi-task learning could be applied to our problem by running a MTL algorithm each time a new task arrives, while remembering all prior tasks. However, this approach incurs a memory and computational burden. Transfer learning lacks the sequential nature of our problem. For multi-task and transfer learning, there are theoretical guarantees on regret for some algorithms [6].

We can also consider the *concept drift* problem in which we observe a stream of incoming data that potentially changes over time, and the goal is to predict some property of each piece of data as it arrives. After prediction, we incur a loss that is revealed to us. For example, we could observe a feature  $w_n$  and predict the label  $y_n$  as in [7]. Some approaches for concept drift use iterative algorithms such as SGD, but without specific models on how the data changes. As a result, only simulation results showing good performance are available. There are also some bandit approaches in which one of a finite number of predictors must be applied to the data as in [8]. For this approach, there are regret guarantees using techniques for analyzing bandit problems.

Another relevant model is *sequential supervised learning* (see [9]) in which we observe a stream of data consisting of feature/label pairs  $(w_n, y_n)$  at time  $n$ , with  $w_n$  being the feature vector and  $y_n$  being the label. At time  $n$ , we want to predict  $y_n$  given  $x_n$ . One approach to this problem, studied in [10] and [11], is to look at  $L$  consecutive pairs  $\{(w_{n-i}, y_{n-i})\}_{i=1}^L$  and develop a predictor at time  $n$  by applying a supervised learning algorithm to this training data. Another approach is to assume that there is an underlying hidden Markov model (HMM) [12]. The label  $y_n$  represents the hidden state and the pair  $(w_n, \bar{y}_n)$  represents the observation with  $\bar{y}_n$  being a noisy version of  $y_n$ . HMM inference techniques are used to estimate  $y_n$ .

## 2 Adaptive Sequential Optimization With $\rho$ Known

For analysis, we need the following assumptions on our functions  $f_n(x)$  and the optimization algorithm:

**A.1** For the optimization algorithm under consideration, there is a function  $b(d_0, K_n)$  such that

$$\mathbb{E}[f_n(x_n)] - f_n(x_n^*) \leq b(d_0, K_n)$$

with  $K_n$  the number of samples from  $p_n$  and  $\mathbb{E}\|x_n(0) - x_n^*\|^2 \leq d_0$ , where  $x_n(0)$  is the initial point of the optimization algorithm at time  $n$ . Finally,  $b(d_0, K_n)$  is non-decreasing in  $d_0$ .

**A.2** Each loss function  $\ell(x, z)$  is differentiable in  $x$ . Each  $f_n(x)$  is strongly convex with parameter  $m$ , i.e.,

$$f_n(y) \geq f_n(x) + \langle \nabla_x f_n(x), y - x \rangle + \frac{1}{2}m\|y - x\|^2$$

**A.3**  $\text{diam}(\mathcal{X}) < +\infty$

**A.4** We can find initial points  $x_1$  and  $x_2$  that satisfy the excess risk criterion with  $\varepsilon_1$  and  $\varepsilon_2$  known, i.e.,

$$\mathbb{E}[f_i(x_i)] - f_i(x_i^*) \leq \varepsilon_i \quad i = 1, 2$$

*Remarks:* For assumption A.1, we assume that the bound  $b(d_0, K_n)$  depends on the number of samples  $K_n$  and not the number of iterations. For SGD, generally the number of iterations equals  $K_n$  as each sample is used to produce a noisy gradient. In addition, we often set  $\mathbf{x}_n(0) = \mathbf{x}_{n-1}$ . See Appendix A for a discussion of useful  $b(d_0, K_n)$  bounds. For assumption A.4, we can fix  $K_i$  and set  $\varepsilon_i = b(\text{diam}(\mathcal{X})^2, K_i)$  for  $i = 1, 2$ .

Now, we examine the case when the change in minimizers,  $\rho$  in (2) or (3), is known. For the analysis of the section, whether (2) or (3) holds does not affect the analysis. Later we will estimate  $\rho$  and in this case whether (2) or (3) holds matters substantially.

We want to find a bound  $\varepsilon_n$  on the excess risk at time  $n$  in terms of  $K_n$  and  $\rho$ , i.e.,  $\varepsilon_n$  such that  $\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*) \leq \varepsilon_n$ . The idea is to start with the bounds from assumption A.4 and proceed inductively using the previous  $\varepsilon_{n-1}$  and  $\rho$  from (2). Suppose that  $\varepsilon_{n-1}$  bounds the excess risk at time  $n-1$ . Using the triangle inequality, strong convexity, and (2) we have

$$\begin{aligned} \mathbb{E}\|\mathbf{x}_{n-1} - \mathbf{x}_n^*\|^2 &\leq (\|\mathbf{x}_{n-1} - \mathbf{x}_{n-1}^*\| + \|\mathbf{x}_{n-1}^* - \mathbf{x}_n^*\|)^2 \\ &\leq \left( \sqrt{\frac{2}{m} \mathbb{E}[f_{n-1}(\mathbf{x}_{n-1})] - f_{n-1}(\mathbf{x}_{n-1}^*)} + \|\mathbf{x}_{n-1}^* - \mathbf{x}_n^*\| \right)^2 \\ &\leq \left( \sqrt{\frac{2\varepsilon_{n-1}}{m}} + \rho \right)^2 \end{aligned} \quad (4)$$

In comparison, we could use the estimate  $\text{diam}^2(\mathcal{X})$  to bound  $\mathbb{E}\|\mathbf{x}_{n-1} - \mathbf{x}_n^*\|^2$  and select  $K_n$ . If the bound in (4) is much smaller than  $\text{diam}(\mathcal{X})^2$ , then we need significantly fewer samples  $K_n$  to guarantee a desired excess risk. Now, by using the bound  $b(d_0, K_n)$  from assumption A.1, we can set

$$\varepsilon_n = b\left(\left(\sqrt{\frac{2\varepsilon_{n-1}}{m}} + \rho\right)^2, K_n\right) \quad \forall n \geq 3$$

which yields a sequence of bounds on the excess risk. Note that this recursion only relies on the immediate past at time  $n-1$  through  $\varepsilon_{n-1}$ . To achieve  $\varepsilon_n \leq \varepsilon$  for all  $n$ , we set

$$K_1 = \min\{K \geq 1 \mid b(\text{diam}(\mathcal{X})^2, K) \leq \varepsilon\}$$

and  $K_n = K^*$  for  $n \geq 2$  with

$$K^* = \min\left\{K \geq 1 \mid b\left(\left(\sqrt{\frac{2\varepsilon}{m}} + \rho\right)^2, K\right) \leq \varepsilon\right\} \quad (5)$$

### 3 Estimating $\rho$

In practice, we do not know  $\rho$ , so we must construct an estimate  $\hat{\rho}_n$  using the samples from each distribution  $p_n$ . We introduce two approaches to estimate  $\rho$  at one time step,  $\|\mathbf{x}_i^* - \mathbf{x}_{i-1}^*\|$ , and methods to combine these estimates under assumptions (2) and (3). We show that for our estimate  $\hat{\rho}_n$  and appropriately chosen sequences  $\{t_n\}$  for all  $n$  large enough  $\hat{\rho}_n + t_n \geq \rho$  almost surely. With this property, analysis similar to that in Section 2 holds.

#### 3.1 Allowed Ways to Choose $K_n$

One of the sources of difficulty in estimating  $\rho$  is that we will allow  $K_n$  to be selected in a data dependent way, so  $K_n$  is itself a random variable. We make the assumption that  $K_n$  is selected using only information available at the end of time  $n-1$ . To make this precise we define a filtration of sigma algebras to describe the available information. First, we define the sigma algebra  $\mathcal{K}_0$  containing all the information on the initial conditions of our algorithm. For example, we may start at a random point  $\mathbf{x}_0$  and then

$$\mathcal{K}_0 = \sigma(\mathbf{x}_0)$$

The sigma algebra  $\mathcal{K}_0$  may also contain information about  $K_1$  and  $K_2$ . Next, we define the filtration

$$\mathcal{K}_n = \sigma\left(\{z_n(k)\}_{k=1}^{K_n}\right) \vee \mathcal{K}_{n-1} \quad \forall n \geq 1 \quad (6)$$

where

$$\mathcal{F} \vee \mathcal{G} = \sigma(\mathcal{F} \cup \mathcal{G})$$

is the merge operator for sigma algebras. The sigma algebra  $\mathcal{K}_n$  contains all the information available to us at the end of time  $n$ . We assume that  $K_n$  is  $\mathcal{K}_{n-1}$ -measurable to capture the idea that  $K_n$  is chosen only using information available at the end of time  $n-1$ .

### 3.2 Estimating One Step Change

First, we estimate the one step changes  $\|\mathbf{x}_i^* - \mathbf{x}_{i-1}^*\|$  denoted by  $\tilde{\rho}_i$ . Implicitly, we assume that all one step estimates are capped by  $\text{diam}(\mathcal{X})$ , since trivially  $\|\mathbf{x}_n^* - \mathbf{x}_{n-1}^*\| \leq \text{diam}(\mathcal{X})$ .

#### 3.2.1 Direct Estimate

First, we construct an estimate  $\tilde{\rho}_i$  of the one step changes  $\|\mathbf{x}_i^* - \mathbf{x}_{i-1}^*\|$ . Using the triangle inequality and variational inequalities from [13] yields

$$\begin{aligned} \|\mathbf{x}_i^* - \mathbf{x}_{i-1}^*\| &\leq \|\mathbf{x}_i - \mathbf{x}_{i-1}\| + \|\mathbf{x}_i - \mathbf{x}_i^*\| + \|\mathbf{x}_{i-1} - \mathbf{x}_{i-1}^*\| \\ &\leq \|\mathbf{x}_i - \mathbf{x}_{i-1}\| + \frac{1}{m} \|\nabla_{\mathbf{x}} f_i(\mathbf{x}_i)\| + \frac{1}{m} \|\nabla_{\mathbf{x}} f_i(\mathbf{x}_{i-1})\| \end{aligned}$$

We then approximate  $\|\nabla_{\mathbf{x}} f_i(\mathbf{x}_i)\| = \|\mathbb{E}_{\mathbf{z}_i \sim p_i} [\nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i)]\|$  by

$$\left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) \right\|$$

to yield the following estimate that we call the *direct estimate*:

$$\tilde{\rho}_i \triangleq \|\mathbf{x}_i - \mathbf{x}_{i-1}\| + \frac{1}{m} \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) \right\| + \frac{1}{m} \left\| \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} \nabla_{\mathbf{x}} \ell(\mathbf{x}_{i-1}, \mathbf{z}_{i-1}(k)) \right\|$$

#### 3.2.2 Vector Integral Probability Metric Estimate

Given a class of functions  $\mathcal{F}$  where each  $f \in \mathcal{F}$  maps  $\mathcal{Z} \rightarrow \mathbb{R}$ , an integral probability metric (IPM) [14] between two distributions  $p$  and  $q$  is defined to be

$$\gamma_{\mathcal{F}}(p, q) \triangleq \sup_{f \in \mathcal{F}} |\mathbb{E}_{\mathbf{z} \sim p}[f(\mathbf{z})] - \mathbb{E}_{\tilde{\mathbf{z}} \sim q}[f(\tilde{\mathbf{z}})]|$$

We consider an extension of this idea, which we call a *vector IPM*, in which the class of functions  $\mathcal{F}$  maps  $\mathcal{Z} \rightarrow \mathcal{X}$ :

$$\gamma_{\mathcal{F}}^V(p, q) \triangleq \sup_{f \in \mathcal{F}} \|\mathbb{E}_{\mathbf{z} \sim p}[f(\mathbf{z})] - \mathbb{E}_{\tilde{\mathbf{z}} \sim q}[f(\tilde{\mathbf{z}})]\| \quad (7)$$

Lemma 1 shows that a vector IPM can be used to bound the change in minimizer at time  $i$  and follows from variational inequalities in [13] and the assumption that  $\{\nabla_{\mathbf{x}} \ell(\mathbf{x}, \cdot) : \mathbf{x} \in \mathcal{X}\} \subset \mathcal{F}$ .

**Lemma 1.** Assume that  $\{\nabla_{\mathbf{x}} \ell(\mathbf{x}, \cdot) : \mathbf{x} \in \mathcal{X}\} \subset \mathcal{F}$ . Then  $\|\mathbf{x}_i^* - \mathbf{x}_{i-1}^*\| \leq \frac{1}{m} \gamma_{\mathcal{F}}^V(p_i, p_{i-1})$ .

*Proof.* By exploiting variational inequalities from [13], we can show that

$$\begin{aligned}\|\mathbf{x}_i^* - \mathbf{x}_{i-1}^*\| &\leq \frac{1}{m} \|\nabla_{\mathbf{x}} f_i(\mathbf{x}_{i-1}^*) - \nabla_{\mathbf{x}} f_{i-1}(\mathbf{x}_{i-1}^*)\| \\ &= \frac{1}{m} \|\mathbb{E}_{\mathbf{z}_i \sim p_i} [\nabla_{\mathbf{x}} \ell(\mathbf{x}_{i-1}^*, \mathbf{z}_i)] - \mathbb{E}_{\mathbf{z}_{i-1} \sim p_{i-1}} [\nabla_{\mathbf{x}} \ell(\mathbf{x}_{i-1}^*, \mathbf{z}_{i-1})]\| \end{aligned}$$

By assumption  $\{\nabla_{\mathbf{x}} \ell(\mathbf{x}_{i-1}^*, \cdot) : \mathbf{x} \in \mathcal{X}\} \subset \mathcal{F}$ , so

$$\begin{aligned}\|\nabla_{\mathbf{x}} f_i(\mathbf{x}_{i-1}^*) - \nabla_{\mathbf{x}} f_{i-1}(\mathbf{x}_{i-1}^*)\| &= \|\mathbb{E}_{\mathbf{z}_i \sim p_i} [\ell(\mathbf{x}_{i-1}^*, \mathbf{z}_i)] - \mathbb{E}_{\mathbf{z}_{i-1} \sim p_{i-1}} [\ell(\mathbf{x}_{i-1}^*, \mathbf{z}_{i-1})]\| \\ &\leq \sup_{f \in \mathcal{F}} \|\mathbb{E}_{\mathbf{z}_i \sim p_i} [f(\mathbf{z}_i)] - \mathbb{E}_{\mathbf{z}_{i-1} \sim p_{i-1}} [f(\mathbf{z}_{i-1})]\| \\ &= \gamma_{\mathcal{F}}^{\mathbf{V}}(p_i, p_{i-1}) \end{aligned}$$

□

We cannot compute this vector IPM, since we do not know the distributions  $p_i$  and  $p_{i-1}$ . Instead, we plug in the empiricals  $\hat{p}_i$  and  $\hat{p}_{i-1}$  to yield the estimate  $\frac{1}{m} \gamma_{\mathcal{F}}^{\mathbf{V}}(\hat{p}_i, \hat{p}_{i-1})$ . This estimate is biased upward, which ensures that  $\|\mathbf{x}_i^* - \mathbf{x}_{i-1}^*\| \leq \mathbb{E} \left[ \frac{1}{m} \gamma_{\mathcal{F}}^{\mathbf{V}}(\hat{p}_i, \hat{p}_{i-1}) \right]$ .

Our estimate is still not in a closed form since there is a supremum over  $\mathcal{F}$  in the computation of  $\gamma_{\mathcal{F}}^{\mathbf{V}}(\hat{p}_i, \hat{p}_{i-1})$ . For the class of functions

$$\mathcal{F} = \{f \mid \|f(\mathbf{z}) - f(\tilde{\mathbf{z}})\| \leq r(\mathbf{z}, \tilde{\mathbf{z}})\}. \quad (8)$$

we can compute an upper bound  $\Gamma_i$  on  $\gamma_{\mathcal{F}}^{\mathbf{V}}(\hat{p}_i, \hat{p}_{i-1})$  yielding a computable estimate  $\tilde{\rho}_i = \frac{1}{m} \Gamma_i$ . Set  $\tilde{\mathbf{z}}_i(k) = \mathbf{z}_i(k)$  if  $1 \leq k \leq K_i$  and  $\tilde{\mathbf{z}}_i(k) = \mathbf{z}_{i-1}(k)$  if  $K_i + 1 \leq k \leq K_i + K_{i-1}$ . From (7), we have

$$\gamma_{\mathcal{F}}^{\mathbf{V}}(\hat{p}_i, \hat{p}_{i-1}) = \sup_{f \in \mathcal{F}} \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} f(\tilde{\mathbf{z}}_i(k)) - \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} f(\tilde{\mathbf{z}}_i(K_i + k)) \right\|$$

We can relax this supremum by maximizing over the function value  $f(\tilde{\mathbf{z}}_i(k))$  denoted by  $\alpha_k$  in the following non-convex quadratically constrained quadratic program (QCQP):

$$\begin{aligned} &\text{maximize} \quad \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \alpha_k - \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} \alpha_{K_i+k} \right\| \\ &\text{subject to} \quad \|\alpha_k - \alpha_j\| \leq r(\tilde{\mathbf{z}}_i(k), \tilde{\mathbf{z}}_i(j)) \quad \forall k < j \end{aligned}$$

The constraints are imposed to ensure that the function values  $\alpha_k$  can correspond to a function in  $\mathcal{F}$  from (8). The value of this QCQP exactly may not equal the vector IPM but at least provides an upper bound. Finally, we note that this QCQP can be converted to its dual form to yield an SDP, which is often easier to solve.

### 3.2.3 Comparison of Estimates

The direct estimate is easier to compute but may be loose if  $\|\mathbf{x}_n - \mathbf{x}_n^*\|$  is large. If  $\|\mathbf{x}_n - \mathbf{x}_n^*\|$  is large, then the vector IPM approach is in general tighter. However, the vector IPM is more difficult to compute due to need to solve a QCQP or SDP and check the inclusion conditions in Lemma 1. Also, the number of constraints in the QCQP or SDP grows quadratically in the number of samples.

## 3.3 Combining One Step Estimates For Constant Change

Assuming that  $\|\mathbf{x}_i^* - \mathbf{x}_{i-1}^*\| = \rho$  from (2), we average the one step estimates  $\tilde{\rho}_i$  to yield a better estimate

$$\hat{\rho}_n = \frac{1}{n-1} \sum_{i=2}^n \tilde{\rho}_i$$

of  $\rho$  at each time  $n$  under (2). To analyze the behavior of our combined estimates, we use sub-Gaussian concentration inequalities detailed in Appendix B. Lemma 22 is of particular importance to our analysis.

### 3.3.1 Direct Estimate

The difficulty in analyzing the direct estimate comes because in approximating  $\frac{1}{m}\|\nabla f_i(\mathbf{x}_i)\|$  by

$$\frac{1}{m}\left\|\frac{1}{K_i}\sum_{k=1}^{K_i}\nabla_{\mathbf{x}}\ell(\mathbf{x}_i, \mathbf{z}_i(k))\right\|$$

$\mathbf{x}_i$  is dependent on all the samples  $\{\mathbf{z}_i(k)\}_{k=1}^{K_i}$ . To illustrate the problem further, consider drawing two independent copies  $\{\mathbf{z}_i(k)\}_{k=1}^{K_i} \stackrel{\text{iid}}{\sim} p_i$  and  $\{\tilde{\mathbf{z}}_i(k)\}_{k=1}^{K_i} \stackrel{\text{iid}}{\sim} p_i$  of the samples. Suppose that we use the second copy  $\{\tilde{\mathbf{z}}_i(k)\}_{k=1}^{K_i}$  to compute  $\mathbf{x}_i$  using our optimization algorithm of choice starting from  $\mathbf{x}_{i-1}$ . Then we approximate  $\frac{1}{m}\|\nabla f_i(\mathbf{x}_i)\|$  by

$$\frac{1}{m}\left\|\frac{1}{K_i}\sum_{k=1}^{K_i}\nabla_{\mathbf{x}}\ell(\mathbf{x}_i, \mathbf{z}_i(k))\right\|$$

Now, since  $\mathbf{x}_i$  is independent of  $\{\mathbf{z}_i(k)\}_{k=1}^{K_i}$  the quantity

$$\frac{1}{m}\left\|\frac{1}{K_i}\sum_{k=1}^{K_i}\nabla_{\mathbf{x}}\ell(\mathbf{x}_i, \mathbf{z}_i(k))\right\|$$

is the norm of an average of independent random variables conditioned on  $\mathbf{x}_i$ . This allows us to apply standard concentration inequalities for norms of random variables as in [15]. In this section, we argue that re-using the samples  $\{\mathbf{z}_i(k)\}_{k=1}^{K_i}$  to compute  $\mathbf{x}_i$  is not too far from using a second independent draw  $\{\tilde{\mathbf{z}}_i(k)\}_{k=1}^{K_i}$ .

For analysis, we need the following additional assumptions:

**B.1** The loss function  $\ell(\mathbf{x}, \mathbf{z})$  has uniform Lipschitz continuous gradients in  $\mathbf{x}$  with modulus  $L$ , i.e.

$$\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}}\ell(\tilde{\mathbf{x}}, \mathbf{z})\| \leq L\|\mathbf{x} - \tilde{\mathbf{x}}\| \quad \forall \mathbf{z} \in \mathcal{Z}$$

**B.2** Assuming  $\mathcal{Z}$  is  $d$ -dimensional, each component  $j$  of the gradient error  $\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z}_n) - f_n(\mathbf{x})$  satisfies

$$\mathbb{E}\left[\exp\left\{s(\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z}_n) - \nabla f_n(\mathbf{x}))_j\right\} \middle| \mathbf{x}\right] \leq \exp\left\{\frac{1}{2}\frac{C_g}{d^2}s^2\right\}$$

Assumption B.1 is reasonable if the space  $\mathcal{Z}$  containing  $\mathbf{z}$  is compact. Although in practice, the distribution of gradient error could depend on  $\mathbf{x}$ , we assume that the bound  $C_g$  does not depend on  $\mathbf{x}$ . We can view this as a pessimistic assumption corresponding to choosing the worst case bound as a function of  $\mathbf{x}$  and the resulting  $C_g$ . This is a common assumption for in high probability analysis of optimization algorithms as in [16] for example.

To proceed, we first define two other useful estimates for  $\rho$ . As discussed before, suppose that we make a second independent draw of samples  $\{\tilde{\mathbf{z}}_i(k)\}_{k=1}^{K_i}$  from  $p_i$ . We use these samples to compute  $\tilde{\mathbf{x}}_i$  in the same manner as  $\mathbf{x}_i$  starting from  $\mathbf{x}_{i-1}$  except with  $\{\tilde{\mathbf{z}}_i(k)\}_{k=1}^{K_i}$  used in place of  $\{\mathbf{z}_i(k)\}_{k=1}^{K_i}$ . Then define

$$\tilde{\rho}_i^{(2)} \triangleq \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_{i-1}\| + \frac{1}{m}\left\|\frac{1}{K_i}\sum_{k=1}^{K_i}\nabla_{\mathbf{x}}\ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k))\right\| + \frac{1}{m}\left\|\frac{1}{K_{i-1}}\sum_{k=1}^{K_{i-1}}\nabla_{\mathbf{x}}\ell(\tilde{\mathbf{x}}_{i-1}, \mathbf{z}_{i-1}(k))\right\|$$

This is the same form as the direct estimate with  $\tilde{\mathbf{x}}_i$  in place of  $\mathbf{x}_i$ . Next, define

$$\tilde{\rho}_i^{(3)} \triangleq \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_{i-1}\| + \frac{1}{m}\|\nabla f_i(\mathbf{x}_i)\| + \frac{1}{m}\|\nabla f_{i-1}(\mathbf{x}_{i-1})\|$$

This is in fact the bound that inspired the direct estimate. We also define the averaged estimates

$$\hat{\rho}_n^{(2)} \triangleq \frac{1}{n-1}\sum_{i=2}^n \tilde{\rho}_i^{(2)}$$

and

$$\hat{\rho}_n^{(3)} \triangleq \frac{1}{n-1} \sum_{i=2}^n \tilde{\rho}_i^{(3)}$$

We know that  $\hat{\rho}_n^{(3)} \geq \rho$ . Thus, if we can control the gap between the pair  $\hat{\rho}_n$  and  $\hat{\rho}_n^{(2)}$  and the pair  $\hat{\rho}_n^{(2)}$  and  $\hat{\rho}_n^{(3)}$ , then we can ensure that  $\hat{\rho}_n$  plus an appropriate constant upper bounds  $\rho$  for all  $n$  large enough as desired.

First, we show that  $\hat{\rho}_n^{(2)}$  upper bounds  $\rho$  eventually.

**Lemma 2.** *Suppose that the following conditions hold:*

1. *B.1 -B.2 hold*
2. *The sequence  $\{t_n\}$  satisfies*

$$\sum_{n=2}^{\infty} \exp \left\{ -\frac{(n-1)m^2 t_n^2}{72C_g} \right\} < \infty$$

*Then for all  $n$  large enough it holds that  $\hat{\rho}_n^{(2)} + \hat{C}_n^{(2)} + t_n \geq \rho$  almost surely with*

$$\hat{C}_n^{(2)} \triangleq \frac{1}{dm(n-1)} \left( \sqrt{\frac{C_g}{K_1}} + 2 \sum_{i=1}^n \sqrt{\frac{C_g}{K_i}} + \sqrt{\frac{C_g}{K_n}} \right)$$

*Proof.* First, we have by the triangle equality and reverse triangle inequality

$$\begin{aligned} & m|\tilde{\rho}_i^{(2)} - \tilde{\rho}_i^{(3)}| \\ &= \left| \left( \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) \right\| - \|\nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)\| \right) + \left( \left\| \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_{i-1}, \mathbf{z}_{i-1}(k)) \right\| - \|\nabla_{\mathbf{x}} f_{i-1}(\tilde{\mathbf{x}}_{i-1})\| \right) \right| \\ &\leq \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) \right\| - \|\nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)\| + \left\| \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_{i-1}, \mathbf{z}_{i-1}(k)) \right\| - \|\nabla_{\mathbf{x}} f_{i-1}(\tilde{\mathbf{x}}_{i-1})\| \\ &\leq \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| + \left\| \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_{i-1}, \mathbf{z}_{i-1}(k)) - \nabla_{\mathbf{x}} f_{i-1}(\tilde{\mathbf{x}}_{i-1})) \right\| \end{aligned}$$

Then by the triangle inequality, we have

$$\begin{aligned} |\hat{\rho}_n^{(2)} - \hat{\rho}_n^{(3)}| &\leq \frac{1}{m(n-1)} \sum_{i=2}^n \left( \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| \right. \\ &\quad \left. + \left\| \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_{i-1}, \mathbf{z}_{i-1}(k)) - \nabla_{\mathbf{x}} f_{i-1}(\tilde{\mathbf{x}}_{i-1})) \right\| \right) \\ &\leq \frac{1}{m(n-1)} \left( \left\| \frac{1}{K_1} \sum_{k=1}^{K_1} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_1, \mathbf{z}_1(k)) - \nabla_{\mathbf{x}} f_1(\tilde{\mathbf{x}}_1)) \right\| \right. \\ &\quad + 2 \sum_{i=2}^{n-1} \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| \\ &\quad \left. + \left\| \frac{1}{K_n} \sum_{k=1}^{K_n} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_n, \mathbf{z}_n(k)) - \nabla_{\mathbf{x}} f_n(\tilde{\mathbf{x}}_n)) \right\| \right) \end{aligned} \quad (9)$$

We will analyze the behavior of this bound on  $|\hat{\rho}_i^{(2)} - \hat{\rho}_i^{(3)}|$  using Lemma 22 in Appendix B. Define the filtration

$$\mathcal{F}_i = \sigma \left( \bigcup_{j=1}^i \{\mathbf{z}_j(k)\}_{k=1}^{K_j} \cup \bigcup_{j=1}^{i+1} \{\tilde{\mathbf{z}}_j(k)\}_{k=1}^{K_j} \right) \vee \mathcal{K}_0 \quad i = 0, \dots, n \quad (10)$$

with  $\mathcal{H}_0$  from (6). Note that  $\mathcal{H}_{i-1} \subset \mathcal{F}_{i-1}$ , so  $K_i$  is  $\mathcal{F}_{i-1}$ -measurable. In addition,  $\tilde{\mathbf{x}}_i$  but not  $\mathbf{x}_i$  is  $\mathcal{F}_{i-1}$ -measurable. Define the random variables

$$V_i = \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| - \mathbb{E} \left[ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| \middle| \mathcal{F}_{i-1} \right] \quad i = 1, \dots, n$$

Clearly,  $V_i$  is  $\mathcal{F}_i$ -measurable, since  $V_i$  is a function of  $\tilde{\mathbf{x}}_i$ ,  $K_i$ , and  $\{\mathbf{z}_i(k)\}_{k=1}^{K_i}$  all of which are  $\mathcal{F}_i$ -measurable. Conditioned on  $\mathcal{F}_{i-1}$ , the sum

$$\frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \quad (11)$$

is a sum of iid random variables. We now work with the conditional measure  $\mathbb{P}\{\cdot \mid \mathcal{F}_{i-1}\}$  to compute sub-Gaussian norms of (11) define in (24) and (25) of Appendix B. By assumption B.2, we have

$$\tau^2 \left( (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i))_j \right) \leq \frac{C_g}{d^2}$$

Therefore, applying Lemma 24 yields

$$B \left( \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right) \leq \sqrt{\frac{C_g}{K_i}}$$

due to the independence conditioned on  $\mathcal{F}_{i-1}$ . By applying Lemma 25 from [17] to the conditional distribution  $\mathbb{P}\{\cdot \mid \mathcal{F}_{i-1}\}$ , we have

$$\begin{aligned} \mathbb{P} \left\{ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| > t \middle| \mathcal{F}_{i-1} \right\} &\leq 2 \exp \left\{ -\frac{t^2}{2(\sqrt{C_g}/K_i)^2} \right\} \\ &= 2 \exp \left\{ -\frac{K_i t^2}{2C_g} \right\} \end{aligned}$$

Since

$$\mathbb{E} \left[ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| \middle| \mathcal{F}_{i-1} \right] \geq 0,$$

we have

$$\begin{aligned} &\mathbb{P} \left\{ V_i > t \middle| \mathcal{F}_{i-1} \right\} \\ &= \mathbb{P} \left\{ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| \right. \\ &\quad \left. - \mathbb{E} \left[ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| \middle| \mathcal{F}_{i-1} \right] > t \middle| \mathcal{F}_{i-1} \right\} \\ &\leq \mathbb{P} \left\{ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| > t \middle| \mathcal{F}_{i-1} \right\} \\ &\leq 2 \exp \left\{ -\frac{K_i t^2}{2C_g} \right\} \\ &\leq 2 \exp \left\{ -\frac{t^2}{2C_g} \right\} \end{aligned}$$

Since  $\mathbb{E}[V_i | \mathcal{F}_{i-1}] = 0$ , we can apply Lemma 26 with  $c = 1/(2C_g)$  to yield

$$\mathbb{E} [e^{sV_i} | \mathcal{F}_{i-1}] \leq \exp \left\{ \frac{1}{2} (18C_g) s^2 \right\}$$

This shows that the collection of random variables  $\{V_i\}_{i=1}^n$  and the filtration  $\{\mathcal{F}_i\}_{i=0}^n$  satisfies the conditions of Lemma 22. Before applying Lemma 22, we bound the conditional expectations

$$\mathbb{E} \left[ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\|^2 \middle| \mathcal{F}_{i-1} \right]$$

By a straightforward calculation conditioned on  $\mathcal{F}_{i-1}$ , we have

$$\begin{aligned} & \mathbb{E} \left[ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\|^2 \middle| \mathcal{F}_{i-1} \right] \\ &= \frac{1}{K_i^2} \sum_{k=1}^{K_i} \sum_{j=1}^{K_i} \mathbb{E} [\langle \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i), \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(j)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i) \rangle | \mathcal{F}_{i-1}] \\ &= \frac{1}{K_i^2} \sum_{k=1}^{K_i} \mathbb{E} [\| \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i) \|^2 | \mathcal{F}_{i-1}] \\ &\stackrel{(a)}{=} \frac{1}{K_i^2} \sum_{k=1}^{K_i} \sum_{q=1}^d \mathbb{E} [(\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i))_q^2 | \mathcal{F}_{i-1}] \\ &\stackrel{(b)}{\leq} \frac{1}{K_i^2} \sum_{k=1}^{K_i} d \frac{C_g}{d^2} \\ &\leq \frac{C_g}{dK_i} \end{aligned}$$

where (a) is a decomposition into each component of the vector and (b) follows since a centered sub-Gaussian random variable with parameter  $C_g/d^2$  satisfies

$$\mathbb{E} [(\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i))_q^2 | \mathcal{F}_{i-1}] \leq \frac{C_g}{d^2}$$

Then by Jensen's inequality

$$\mathbb{E} \left[ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\|^2 \middle| \mathcal{F}_{i-1} \right] \leq \frac{C_g}{dK_i}$$

Define the constants

$$\begin{aligned} a_1 &= a_n = \frac{1}{m(n-1)} \\ a_2 &= \dots = a_{n-1} = \frac{2}{m(n-1)} \end{aligned}$$

resulting in

$$\|\mathbf{a}\|_2^2 = \frac{2}{m^2(n-1)}$$

Using the bound in (9) and Lemma 22 from Appendix B with this choice of  $\mathbf{a}$ , it holds that

$$\begin{aligned}
& \mathbb{P} \left\{ |\hat{\rho}_n^{(2)} - \hat{\rho}_n^{(3)}| > \sum_{i=1}^n a_i \sqrt{\frac{C_g}{dK_i}} + t \right\} \\
& \leq \mathbb{P} \left\{ \sum_{i=1}^n a_i \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| \right. \\
& \quad \left. > \sum_{i=1}^n a_i \mathbb{E} \left[ \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} f_i(\tilde{\mathbf{x}}_i)) \right\| \middle| \mathcal{F}_{i-1} \right] + t \right\} \\
& = \mathbb{P} \left\{ \sum_{i=1}^n a_i V_i > t \right\} \\
& \leq \exp \left\{ -\frac{m^2(n-1)t^2}{72C_g} \right\}
\end{aligned}$$

Combining this bound with  $\hat{\rho}_n^{(3)} \geq \rho$  yields

$$\begin{aligned}
\sum_{n=2}^{\infty} \mathbb{P} \left\{ \hat{\rho}_n^{(2)} < \rho - \sum_{i=1}^n a_i \sqrt{\frac{C_g}{dK_i}} - t_n \right\} & \leq \sum_{n=2}^{\infty} \mathbb{P} \left\{ \hat{\rho}_n^{(2)} < \hat{\rho}_n^{(3)} - \sum_{i=1}^n a_i \sqrt{\frac{C_g}{dK_i}} - t_n \right\} \\
& \leq \sum_{n=2}^{\infty} \mathbb{P} \left\{ |\hat{\rho}_n^{(2)} - \hat{\rho}_n^{(3)}| > \sum_{i=1}^n a_i \sqrt{\frac{C_g}{dK_i}} + t_n \right\} \\
& \leq \sum_{n=2}^{\infty} \exp \left\{ -\frac{m^2(n-1)t_n^2}{72C_g} \right\} < \infty
\end{aligned}$$

The result follows from the Borel-Cantelli lemma. Note that as claimed

$$\hat{C}_n^{(2)} = \frac{1}{dm(n-1)} \left( \sqrt{\frac{C_g}{K_1}} + 2 \sum_{i=2}^{n-1} \sqrt{\frac{C_g}{K_i}} + \sqrt{\frac{C_g}{K_n}} \right)$$

□

Next, we show that  $\hat{\rho}_n$  upper bounds  $\hat{\rho}_n^{(2)}$  eventually with a general assumption on the optimization algorithm. When the conditions of Lemmas 2 and 3 are satisfied, it holds that  $\hat{\rho}_n$  plus a constant upper bounds  $\rho$ .

**Lemma 3.** *Suppose the following conditions hold:*

1. *B.1-B.2 hold*
2. *There exist bounds*

$$\mathbb{E} [\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \mid \mathcal{F}_{i-1}] \leq C(K_i) \quad i = 1, \dots, n$$

3. *The sequence  $\{t_n\}$  satisfies*

$$\sum_{n=2}^{\infty} \exp \left\{ -\frac{(n-1)^2 t_n^2}{2n \left(1 + \frac{L}{m}\right)^2 \text{diam}^2(\mathcal{X}^c)} \right\} < +\infty$$

Then for all  $n$  large enough it holds that  $\hat{\rho}_n + \hat{C}_n + t_n \geq \hat{\rho}_n^{(2)}$  almost surely with

$$\hat{C}_n \triangleq \frac{\left(1 + \frac{L}{m}\right)}{n-1} \left( C(K_1) + 2 \sum_{i=2}^{n-1} C(K_i) + C(K_n) \right)$$

*Proof.* We have by the triangle inequality, reverse triangle inequality, and the Lipschitz continuity of  $\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z})$  in  $\mathbf{x}$  from assumption B.1

$$\begin{aligned}
|\tilde{\rho}_i - \tilde{\rho}_i^{(2)}| &\leq \|\mathbf{x}_i - \mathbf{x}_{i-1}\| - \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_{i-1}\| \\
&\quad + \left\| \frac{1}{m} \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) \right\| - \frac{1}{m} \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) \right\| \right\| \\
&\quad + \left\| \frac{1}{m} \left\| \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} \nabla_{\mathbf{x}} \ell(\mathbf{x}_{i-1}, \mathbf{z}_{i-1}(k)) \right\| - \frac{1}{m} \left\| \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_{i-1}, \mathbf{z}_{i-1}(k)) \right\| \right\| \\
&\leq \|(\mathbf{x}_i - \tilde{\mathbf{x}}_i) - (\mathbf{x}_{i-1} - \tilde{\mathbf{x}}_{i-1})\| \\
&\quad + \frac{1}{m} \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k))) \right\| \\
&\quad + \frac{1}{m} \left\| \frac{1}{K_{i-1}} \sum_{k=1}^{K_{i-1}} (\nabla_{\mathbf{x}} \ell(\mathbf{x}_{i-1}, \mathbf{z}_{i-1}(k)) - \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_{i-1}, \mathbf{z}_{i-1}(k))) \right\| \\
&\leq \left(1 + \frac{L}{m}\right) (\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| + \|\mathbf{x}_{i-1} - \tilde{\mathbf{x}}_{i-1}\|)
\end{aligned}$$

so

$$\begin{aligned}
|\hat{\rho}_n - \hat{\rho}_n^{(2)}| &\leq \frac{1}{n-1} \sum_{i=2}^n |\tilde{\rho}_i - \tilde{\rho}_i^{(2)}| \\
&\leq \frac{(1 + \frac{L}{m})}{n-1} \sum_{i=2}^n (\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| + \|\mathbf{x}_{i-1} - \tilde{\mathbf{x}}_{i-1}\|) \\
&= \frac{(1 + \frac{L}{m})}{n-1} \left( \|\mathbf{x}_1 - \tilde{\mathbf{x}}_1\| + 2 \sum_{i=2}^{n-1} \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| + \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\| \right)
\end{aligned}$$

We will again apply Lemma 22 of Appendix B to analyze this upper bound using the sigma algebra

$$\mathcal{F}_i = \sigma \left( \bigcup_{j=1}^i \{z_j(k)\}_{k=1}^{K_j} \cup \bigcup_{j=1}^i \{\tilde{z}_j(k)\}_{k=1}^{K_j} \right) \vee \mathcal{K}_0 \quad i = 0, \dots, n \quad (12)$$

Define the random variable

$$V_i = \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| - \mathbb{E}[\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \mid \mathcal{F}_{i-1}]$$

Clearly,  $V_i$  is  $\mathcal{F}_i$ -measurable. Since

$$-\text{diam}(\mathcal{X}) \leq V_i \leq \text{diam}(\mathcal{X}),$$

and  $\mathbb{E}[V_i \mid \mathcal{F}_{i-1}] = 0$ , we can apply the conditional version Hoeffding's Lemma from Lemma 23 to yield

$$\mathbb{E}[e^{sV_i} \mid \mathcal{F}_{i-1}] \leq \exp \left\{ \frac{1}{2} \text{diam}^2(\mathcal{X}) s^2 \right\}$$

The collection of random variables  $\{V_i\}_{i=1}^n$  and the filtration  $\{\mathcal{F}_i\}_{i=0}^n$  satisfy the conditions of Lemma 22. Before applying Lemma 22, we bound the conditional expectations

$$\mathbb{E}[\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \mid \mathcal{F}_{i-1}]$$

By assumption, we have

$$\mathbb{E}[\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \mid \mathcal{F}_{i-1}] \leq C(K_i) \quad i = 1, \dots, n$$

and so

$$\begin{aligned} & \frac{(1 + \frac{L}{m})}{n-1} \left( \mathbb{E} [\|\mathbf{x}_1 - \tilde{\mathbf{x}}_1\| \mid \mathcal{F}_0] + 2 \sum_{i=2}^{n-1} \mathbb{E} [\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \mid \mathcal{F}_{i-1}] + \mathbb{E} [\|\mathbf{x}_n - \tilde{\mathbf{x}}_n\| \mid \mathcal{F}_{n-1}] \right) \\ & \leq \frac{(1 + \frac{L}{m})}{n-1} \left( C(K_1) + 2 \sum_{i=2}^{n-1} C(K_i) + C(K_n) \right) \triangleq \hat{C}_n \end{aligned}$$

Set

$$a_1 = a_n = \frac{(1 + \frac{L}{m})}{n-1}$$

and

$$a_2 = \dots = a_{n-1} = \frac{(1 + \frac{L}{m})}{n-1}$$

resulting in

$$\|\mathbf{a}\|_2^2 = \frac{n(1 + \frac{L}{m})^2}{(n-1)^2}$$

Applying our bound in (12) and Lemma 22 with this choice of  $\mathbf{a}$  yields

$$\begin{aligned} & \mathbb{P} \left\{ |\hat{\rho}_n - \hat{\rho}_n^{(2)}| > \hat{C}_n + t \right\} \\ & \leq \mathbb{P} \left\{ \frac{(1 + \frac{L}{m})}{n-1} \left( \|\mathbf{x}_1 - \tilde{\mathbf{x}}_1\| + 2 \sum_{i=2}^{n-1} \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| + \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\| \right) \right. \\ & \quad \left. > \frac{(1 + \frac{L}{m})}{n-1} \left( \mathbb{E} [\|\mathbf{x}_1 - \tilde{\mathbf{x}}_1\| \mid \mathcal{F}_0] + 2 \sum_{i=2}^{n-1} \mathbb{E} [\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \mid \mathcal{F}_{i-1}] + \mathbb{E} [\|\mathbf{x}_n - \tilde{\mathbf{x}}_n\| \mid \mathcal{F}_{n-1}] \right) + t \right\} \\ & = \mathbb{P} \left\{ \frac{(1 + \frac{L}{m})}{n-1} \left( V_1 + 2 \sum_{i=2}^{n-1} V_i + V_n \right) > t \right\} \\ & = \mathbb{P} \left\{ \sum_{i=1}^n a_i V_i > t \right\} \\ & \leq \exp \left\{ - \frac{(n-1)^2 t^2}{2n(1 + \frac{L}{m})^2 \text{diam}^2(\mathcal{X})} \right\} \end{aligned}$$

Finally, we have

$$\begin{aligned} \sum_{n=2}^{\infty} \mathbb{P} \left\{ \hat{\rho}_n < \hat{\rho}_n^{(2)} - \hat{C}_n - t_n \right\} & \leq \sum_{n=2}^{\infty} \mathbb{P} \left\{ |\hat{\rho}_n - \hat{\rho}_n^{(2)}| > \hat{C}_n + t_n \right\} \\ & \leq \sum_{n=2}^{\infty} \exp \left\{ - \frac{(n-1)^2 t_n^2}{2n(1 + \frac{L}{m})^2 \text{diam}^2(\mathcal{X})} \right\} < +\infty \end{aligned}$$

The claim follows from the Borel-Cantelli Lemma. □

If Lemmas 2 and 3 hold for the sequence  $\{t_n/2\}$ , then for all  $n$  large enough it holds that

$$\hat{\rho}_n + \hat{C}_n + \hat{C}_n^{(2)} + t_n \geq \rho$$

almost surely.

**Lemma 4.** *It always holds that*

$$\mathbb{E} [\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \mid \mathcal{F}_{i-1}] \leq 2\sqrt{\frac{1}{m}b(\text{diam}^2(\mathcal{X}), K_i)}$$

Therefore, the choice

$$C(K_i) \triangleq 2\sqrt{\frac{2}{m}b(\text{diam}^2(\mathcal{X}), K_i)}$$

satisfies the conditions of Lemma 3.

*Proof.* Using the sigma algebras defined in (12) yields

$$\begin{aligned} \mathbb{E} [\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \mid \mathcal{F}_{i-1}] &\leq \mathbb{E} [\|\mathbf{x}_i - \mathbf{x}_i^*\| \mid \mathcal{F}_{i-1}] + \mathbb{E} [\|\tilde{\mathbf{x}}_i - \mathbf{x}_i^*\| \mid \mathcal{F}_{i-1}] \\ &\leq \mathbb{E} \left[ \sqrt{\frac{2}{m}(f_i(\mathbf{x}_i) - f_i(\mathbf{x}_i^*))} \mid \mathcal{F}_{i-1} \right] + \mathbb{E} \left[ \sqrt{\frac{2}{m}(f_i(\tilde{\mathbf{x}}_i) - f_i(\mathbf{x}_i^*))} \mid \mathcal{F}_{i-1} \right] \\ &\leq \sqrt{\frac{2}{m}\mathbb{E}[(f_i(\mathbf{x}_i) - f_i(\mathbf{x}_i^*)) \mid \mathcal{F}_{i-1}]} + \sqrt{\frac{2}{m}\mathbb{E}[(f_i(\tilde{\mathbf{x}}_i) - f_i(\mathbf{x}_i^*)) \mid \mathcal{F}_{i-1}]} \\ &\leq 2\sqrt{\frac{2}{m}b(\text{diam}^2(\mathcal{X}), K_i)} \end{aligned}$$

where the third inequality follows from Jensen's inequality.  $\square$

This choice of  $C(K_n)$  works for any algorithm with the associated  $b(d_0, K)$ . For any particular algorithm, we believe that we can produce tighter bounds independent of  $\text{diam}(\mathcal{X})$  by copying the Lyapunov analysis used to analyze SGD as in Appendix A. The analysis becomes algorithm dependent in this case and is omitted.

Finally, we state an overall theorem for the direct estimate that gives general combined conditions under which  $\hat{\rho}_n$  upper bounds  $\rho$ .

**Theorem 1.** *If B.1-B.2 hold and the sequence  $\{t_n\}$  satisfies  $\sum_{n=2}^{\infty} e^{-Cm_n^2} < \infty$  for all  $C > 0$ , then for a sequence of constants  $\{C_n\}$  and for all  $n$  large enough it holds that  $\hat{\rho}_n + C_n + t_n \geq \rho$  almost surely.*

*Proof.* Combine Lemmas 2 and 3 to yield the result with

$$C_n = \hat{C}_n + \hat{C}_n^{(2)}$$

$\square$

### 3.3.2 Vector IPM Estimate

We first derive a version of Hoeffding's inequality that allows for some dependence among the random variables. We use this concentration inequality to analyze  $\hat{\rho}_n$  for the IPM estimate. Given an integer  $W$ , we construct a cover of  $\{1, 2, \dots, n\}$  by dividing the set into  $W$  groups of integers spaced by  $W$ , i.e.,

$$\mathcal{A}_j = \left\{ j, j+W, j+2W, \dots, j + \left\lfloor \frac{n-j}{W} \right\rfloor W \right\} \quad j = 1, \dots, W \quad (13)$$

Note that

$$\{1, 2, \dots, n\} = \bigcup_{j=1}^W \mathcal{A}_j$$

and  $\mathcal{A}_i \cap \mathcal{A}_j = \emptyset$  for  $i \neq j$ . The proof of Lemma 5 is nearly identical to the proof of the extension of Hoeffding's inequality from [18] with Lemma 22 used instead. We assume that if we refer to a filtration  $\mathcal{F}_i$  with  $i < 0$ , then we implicitly refer to  $\mathcal{F}_0$ .

**Lemma 5** (Dependent Hoeffding's Inequality). *Suppose we are given a collection of random variable  $\{V_i\}_{i=1}^n$  and a filtration  $\{\mathcal{F}_i\}_{i=0}^n$  such that*

1.  $a_i \leq V_i \leq b_i$  for constants  $a_i$  and  $b_i$   $i = 1, \dots, n$
2.  $V_i$  is  $\mathcal{F}_i$ -measurable  $i = 1, \dots, n$
3. Given an integer  $W$  and a cover  $\{\mathcal{A}_j\}_{j=1}^W$  as in (13) for each  $j$  it holds that

$$\mathbb{E} \left[ V_{j+iW} \mid \mathcal{F}_{j+(i-1)W} \right] = 0 \quad i = 1, \dots, \left\lfloor \frac{n-j}{W} \right\rfloor$$

and

$$\mathbb{E} \left[ V_j \mid \mathcal{F}_0 \right] = 0$$

Then it holds that

$$\mathbb{P} \left\{ \sum_{i=1}^n V_i > t \right\} \leq \exp \left\{ -\frac{2t^2}{W \sum_{i=1}^n (b_i - a_i)^2} \right\}$$

and

$$\mathbb{P} \left\{ \sum_{i=1}^n V_i < -t \right\} \leq \exp \left\{ -\frac{2t^2}{W \sum_{i=1}^n (b_i - a_i)^2} \right\}$$

*Proof.* Define

$$U_j \triangleq \sum_{i=0}^{\left\lfloor \frac{n-j}{W} \right\rfloor} V_{j+iW}$$

for  $j = 1, \dots, W$ . Let  $\{p_j\}_{j=1}^W$  be a probability distribution on  $\{1, \dots, W\}$  to be specified later. By Jensen's inequality, we have

$$\begin{aligned} \exp \left\{ s \sum_{i=1}^n V_i \right\} &= \exp \left\{ \sum_{j=1}^W p_j \frac{s}{p_j} U_j \right\} \\ &\leq \sum_{j=1}^W p_j \exp \left\{ \frac{s}{p_j} U_j \right\} \end{aligned}$$

Then it holds that

$$\mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^n V_i \right\} \right] \leq \sum_{j=1}^W p_j \mathbb{E} \left[ \exp \left\{ \frac{s}{p_j} U_j \right\} \right]$$

Now consider one term

$$\mathbb{E} \left[ \exp \left\{ \frac{s}{p_j} U_j \right\} \right] = \mathbb{E} \left[ \exp \left\{ \frac{s}{p_j} \sum_{i=0}^{\left\lfloor \frac{n-j}{W} \right\rfloor} V_{j+iW} \right\} \right]$$

Since  $a_{j+iW} \leq V_{j+iW} \leq b_{j+iW}$  and

$$\mathbb{E} \left[ V_{j+iW} \mid \mathcal{F}_{j+(i-1)W} \right] = 0,$$

we can apply the conditional version Hoeffding's Lemma from Lemma 23 to yield

$$\mathbb{E} \left[ e^{sV_{j+iW}} \mid \mathcal{F}_{j+(i-1)W} \right] \leq \exp \left\{ \frac{1}{8} (b_{j+iW} - a_{j+iW})^2 s^2 \right\}$$

Then we can apply Lemma 22 to  $\{V_{j+iW}\}_{i=0}^{\lfloor \frac{n-j}{W} \rfloor}$  and  $\{\mathcal{F}_{j+iW}\}_{i=0}^{\lfloor \frac{n-j}{W} \rfloor}$  to yield

$$\begin{aligned} \mathbb{E} \left[ \exp \left\{ \frac{s}{p_j} U_j \right\} \right] &\leq \exp \left\{ \frac{s^2}{8p_j^2} \sum_{i=0}^{\lfloor \frac{n-j}{W} \rfloor} (b_{j+iW} - a_{j+iW})^2 \right\} \\ &= \prod_{i=0}^{\lfloor \frac{n-j}{W} \rfloor} \exp \left\{ \frac{s^2}{8p_j^2} (b_\alpha - a_\alpha)^2 \right\} \end{aligned}$$

Then we have

$$\begin{aligned} \mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^n V_i \right\} \right] &\leq \sum_{j=1}^W p_j \prod_{i=0}^{\lfloor \frac{n-j}{W} \rfloor} \exp \left\{ \frac{s^2}{8p_j^2} (b_\alpha - a_\alpha)^2 \right\} \\ &= \sum_{j=1}^W p_j \exp \left\{ \frac{s^2 c_j}{8p_j^2} \right\} \end{aligned}$$

with

$$c_j = \sum_{i=0}^{\lfloor \frac{n-j}{W} \rfloor} (b_{j+iW} - a_{j+iW})^2$$

Let  $p_j = \sqrt{c_j}/T$  and

$$T = \sum_{j=1}^W \sqrt{c_j}.$$

Therefore, we have

$$\mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^n V_i \right\} \right] \leq \exp \left\{ \frac{1}{8} T^2 s^2 \right\}$$

Applying the Chernoff bound [19] and optimizing yields

$$\mathbb{P} \left\{ \sum_{i=1}^n V_i > t \right\} \leq \exp \left\{ -2t^2/T^2 \right\}$$

Bounding  $T$  with Cauchy-Schwarz yields

$$T^2 \leq \left( \sum_{j=1}^W 1 \right) \left( \sum_{j=1}^W c_j \right) = W \sum_{i=1}^n (b_i - a_i)^2$$

and the results follows. The proof for the other tail is nearly identical.  $\square$

If we do not have the condition 3 of Lemma 5, then it holds that

$$\mathbb{P} \left\{ \sum_{i=1}^n V_i > \sum_{j=1}^W \sum_{i=0}^{\lfloor \frac{n-j}{W} \rfloor} \mathbb{E} [V_{j+iW} \mid \mathcal{F}_{j+(i-1)W}] + t \right\} \leq \exp \left\{ -\frac{2t^2}{W \sum_{i=1}^n (b_i - a_i)^2} \right\}$$

If we can bound the conditional expectation

$$\mathbb{E} [V_{j+iW} \mid \mathcal{F}_{j+(i-1)W}] \leq C_{j+iW},$$

by a  $\mathcal{F}_{j+(i-1)W}$ -measurable random variable, then we have

$$\begin{aligned}
\mathbb{P} \left\{ \sum_{i=1}^n V_i > \sum_{i=1}^n C_i + t \right\} &= \mathbb{P} \left\{ \sum_{i=1}^n V_i > \sum_{j=1}^W \sum_{i=0}^{\lfloor \frac{n-j}{W} \rfloor} C_{j+iW} + t \right\} \\
&\leq \mathbb{P} \left\{ \sum_{i=1}^n V_i > \sum_{j=1}^W \sum_{i=0}^{\lfloor \frac{n-j}{W} \rfloor} \mathbb{E} [V_{j+iW} \mid \mathcal{F}_{j+(i-1)W}] + t \right\} \\
&\leq \mathbb{P} \left\{ \sum_{j=1}^W \sum_{i=0}^{\lfloor \frac{n-j}{W} \rfloor} (V_{j+iW} - \mathbb{E} [V_{j+iW} \mid \mathcal{F}_{j+(i-1)W}]) > t \right\} \\
&\leq \exp \left\{ -\frac{2t^2}{W \sum_{i=1}^n (b_i - a_i)^2} \right\}
\end{aligned}$$

We have the following lemma characterizing the performance of the IPM estimate.

**Lemma 6.** *For the IPM estimate and any sequence  $\{t_n\}$  such that*

$$\sum_{n=2}^{\infty} \exp \left\{ -\frac{nt_n^2}{4\text{diam}(\mathcal{X})^2} \right\} < \infty$$

for all  $n$  large enough it holds that  $\hat{\rho}_n + t_n \geq \rho$  almost surely.

*Proof.* Define the random variables

$$V_i = \tilde{\rho}_i - \mathbb{E} [\tilde{\rho}_i \mid \mathcal{H}_{i-2}]$$

with  $\{\mathcal{H}_i\}_{i=1}^n$  defined in (6). We have

$$-\text{diam}(\mathcal{X}) \leq V_i \leq \text{diam}(\mathcal{X})$$

Clearly,  $V_i$  is  $\mathcal{H}_i$ -measurable and  $\mathbb{E}[V_i \mid \mathcal{H}_{i-2}] = 0$ . Now, we can apply Lemma 5 with  $W = 2$  to yield

$$\begin{aligned}
\mathbb{P} \left\{ \sum_{i=1}^n V_i < -nt \right\} &\leq \exp \left\{ -\frac{2(nt)^2}{(2)(4\text{diam}^2(\mathcal{X}))} \right\} \\
&= \exp \left\{ -\frac{nt^2}{4\text{diam}^2(\mathcal{X})} \right\}
\end{aligned}$$

None of the random variables  $\{z_i(k)\}_{k=1}^{K_i}$  and  $\{z_{i-1}(k)\}_{k=1}^{K_{i-1}}$  are  $\mathcal{H}_{i-2}$  measurable. Also, regardless of how many samples  $K_i$  and  $K_{i-1}$  are taken, the IPM estimate is biased upward. Thus, it holds that

$$\mathbb{E} [\tilde{\rho}_i \mid \mathcal{H}_{i-2}] \geq \rho$$

Therefore, it follows that

$$\begin{aligned}
\mathbb{P} \{ \hat{\rho}_n < \rho - t \} &\leq \mathbb{P} \left\{ \sum_{i=1}^n \tilde{\rho}_i < \sum_{i=1}^n \mathbb{E} [\tilde{\rho}_i \mid \mathcal{H}_{i-2}] - nt \right\} \\
&= \mathbb{P} \left\{ \sum_{i=1}^n V_i < -nt \right\} \\
&\leq \exp \left\{ -\frac{nt^2}{4\text{diam}^2(\mathcal{X})} \right\}
\end{aligned}$$

Note that we pay a price of two in the exponent due to  $\tilde{\rho}_i$  and  $\tilde{\rho}_{i-1}$  both depending on the samples from  $p_{i-1}$ . Since

$$\sum_{n=2}^{\infty} \exp \left\{ -\frac{nt_n^2}{4\text{diam}(\mathcal{X})^2} \right\} < \infty$$

it follows that

$$\sum_{n=2}^{\infty} \mathbb{P} \{ \hat{\rho}_n + t_n < \rho \} < +\infty,$$

This in turn guarantees by way of the Borel-Cantelli Lemma that for  $n$  large enough

$$\hat{\rho}_n + t_n \geq \rho$$

almost surely. □

### 3.4 Combining One Step Estimates For Bounded Change

We now look at estimating  $\rho$  in the case that

$$\|\mathbf{x}_n^* - \mathbf{x}_{n-1}^*\| \leq \rho.$$

We set

$$\rho_i \triangleq \|\mathbf{x}_i^* - \mathbf{x}_{i-1}^*\|$$

**B.3** Assume that we have estimators  $\hat{h}_W : \mathbb{R}^W \rightarrow \mathbb{R}$  such that

1.  $\mathbb{E}[\hat{h}_W(\rho_j, \dots, \rho_{j-W+1})] \geq \rho$  for all  $j \geq 1$  and  $W \geq 1$
2. For any random variables  $\{\tilde{\rho}_i\}$  such that  $\mathbb{E}[\tilde{\rho}_i] \geq \mathbb{E}[\rho_i]$ , we have

$$\mathbb{E}[\hat{h}_W(\tilde{\rho}_j, \dots, \tilde{\rho}_{j-W+1})] \geq \mathbb{E}[\hat{h}_W(\rho_j, \dots, \rho_{j-W+1})]$$

For example, if  $\rho_i \stackrel{\text{iid}}{\sim} \text{Unif}[0, \rho]$ , then

$$\hat{h}_W(\rho_i, \rho_{i+1}, \dots, \rho_{i+W-1}) = \frac{W+1}{W} \max\{\rho_i, \rho_{i+1}, \dots, \rho_{i+W-1}\}$$

is an estimator of  $\rho$  with the required properties. Also, note that the two conditions on the estimator in B.3 imply that

$$\mathbb{E}[\hat{h}_W(\tilde{\rho}_j, \dots, \tilde{\rho}_{j-W+1})] \geq \mathbb{E}[\hat{h}_W(\rho_j, \dots, \rho_{j-W+1})] \geq \rho$$

Given an estimator satisfying assumption B.3, we compute

$$\tilde{\rho}^{(i)} = \hat{h}_W(\tilde{\rho}_i, \tilde{\rho}_{i-1}, \dots, \tilde{\rho}_{i-W+1})$$

and set

$$\hat{\rho}_n = \frac{1}{n-1} \sum_{i=2}^n \tilde{\rho}^{(i)} = \frac{1}{n-1} \sum_{i=2}^n \hat{h}_{\min\{W, i-1\}}(\tilde{\rho}_i, \tilde{\rho}_{i-1}, \dots, \tilde{\rho}_{\max\{i-W+1, 2\}}) \quad (14)$$

We have

$$\mathbb{E}[\hat{\rho}_n] = \frac{1}{n-1} \sum_{i=2}^n \mathbb{E}[\tilde{\rho}^{(i)}] \geq \rho$$

**Lemma 7** (IPM Single Step Estimates). *For the estimator in (14) computed using the IPM estimate for  $\tilde{\rho}_i$  and any sequence  $\{t_n\}$  such that*

$$\sum_{n=2}^{\infty} \exp \left\{ -\frac{2(n-1)t_n^2}{(W+1)\text{diam}(\mathcal{X})^2} \right\} < \infty$$

*it holds that for all  $n$  large enough  $\hat{\rho}_n + t_n \geq \rho$  almost surely.*

*Proof.* We copy the proof of Lemma 6 with  $W + 1$  in place of 2 and note that  $\tilde{\rho}^{(i)}$  and  $\tilde{\rho}^{(j)}$  with  $|i - j| > W + 1$  do not depend on the same samples. Lemma 5 and some simple algebra yields

$$\mathbb{P}\{\hat{\rho}_n < \rho - t\} \leq \exp\left\{-\frac{2(n-1)t^2}{(W+1)\text{diam}(\mathcal{X})^2}\right\}$$

We pay a price of  $W + 1$  in the denominator of the exponent due to the dependence of the  $\tilde{\rho}^{(i)}$ . By the Borel-Cantelli Lemma, for all  $n$  large enough it holds that  $\hat{\rho}_n + t_n \geq \rho$  almost surely as long as

$$\sum_{n=2}^{\infty} \exp\left\{-\frac{2(n-1)t_n^2}{(W+1)\text{diam}(\mathcal{X})^2}\right\} < \infty$$

□

To analyze the direct estimate, we need the following assumption

**B.4** Suppose that there exists absolute constants  $\{b_i\}_{i=1}^W$  for any fixed  $W$  such that

$$|\hat{h}_W(p_1, \dots, p_W) - \hat{h}_W(q_1, \dots, q_W)| \leq \sum_{i=1}^W b_i |p_i - q_i| \quad \forall \mathbf{p}, \mathbf{q} \in \mathbb{R}_{\geq 0}^W$$

For the uniform case, we have

$$\begin{aligned} \left| \frac{W+1}{W} \max\{p_1, \dots, p_W\} - \frac{W+1}{W} \max\{q_1, \dots, q_W\} \right| &\leq \frac{W+1}{W} \max\{|p_1 - q_1|, \dots, |p_W - q_W|\} \\ &\leq \frac{W+1}{W} \sum_{i=1}^W |p_i - q_i| \end{aligned}$$

so

$$b_1 = \dots = b_W = \frac{W+1}{W}$$

Under assumption B.4, we can then show that

$$\hat{\rho}_n = \frac{1}{n-W} \sum_{i=W+1}^n \tilde{\rho}^{(i)}$$

eventually upper bounds  $\rho$  by copying the proofs of the lemmas behind Theorem 1.

**Lemma 8** (Direct Single Step Estimates). *Suppose that the following conditions hold:*

1. B.1 -B.4 hold
2. The sequence  $\{t_n\}$  satisfies

$$\sum_{n=W+1}^{\infty} \exp\left\{-\frac{(n-W)^2 t_n^2}{32n \left(1 + \frac{L}{m}\right)^2 \left(\sum_{j=1}^W b_j\right)^2 \text{diam}^2(\mathcal{X})}\right\} < +\infty$$

and

$$\sum_{n=W+1}^{\infty} \exp\left\{-\frac{(n-W)^2 m^2 t_n^2}{144nC_g \left(\sum_{j=1}^W b_j\right)^2}\right\} < +\infty$$

3. There are bounds  $C(K)$  such that

$$\mathbb{E}[\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \mid \mathcal{F}_{i-1}] \leq C(K_i)$$

Then for all  $n$  large enough it holds that  $\hat{\rho}_n + \hat{U}_n + \hat{V}_n + t_n \geq \rho$  almost surely with

$$\hat{U}_n = \frac{2(1 + \frac{L}{m}) \sum_{j=1}^W b_j}{n - W} \sum_{i=1}^n C(K_i)$$

and

$$\hat{V}_n = \frac{2 \sum_{j=1}^W b_j}{m(n - W)} \sum_{i=1}^n \sqrt{\frac{C_g}{dK_i}}$$

*Proof.* Define  $\tilde{\rho}_i^{(2)}$ ,  $\tilde{\rho}_i^{(3)}$ ,  $\hat{\rho}_i^{(2)}$ , and  $\hat{\rho}_i^{(3)}$  as in Lemmas 2 and 3. First, we have

$$\begin{aligned} |\hat{\rho}_n - \hat{\rho}_n^{(3)}| &\leq \frac{1}{n - W} \sum_{i=W+1}^n |\tilde{\rho}^{(i)} - \tilde{\rho}_3^{(i)}| \\ &\leq \frac{1}{n - W} \sum_{i=W+1}^n \sum_{j=i-W+1}^i b_j |\tilde{\rho}_j - \tilde{\rho}_j^{(3)}| \\ &\leq \frac{1}{n - W} \sum_{i=W+1}^n \sum_{j=i-W+1}^i b_j \left( |\tilde{\rho}_j - \tilde{\rho}_j^{(2)}| + |\tilde{\rho}_j^{(2)} - \tilde{\rho}_j^{(3)}| \right) \\ &\leq \frac{\sum_{j=1}^W b_j}{n - W} \sum_{i=2}^n \left( |\tilde{\rho}_i - \tilde{\rho}_i^{(2)}| + |\tilde{\rho}_i^{(2)} - \tilde{\rho}_i^{(3)}| \right) \end{aligned}$$

Second, define

$$U_i \triangleq \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|$$

and

$$V_i \triangleq \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla f_i(\tilde{\mathbf{x}}_i)) \right\|$$

Then we have

$$\begin{aligned} |\tilde{\rho}_i - \tilde{\rho}_i^{(2)}| &\leq \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| + \frac{1}{m} \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k))) \right\| \\ &\leq \left( 1 + \frac{L}{m} \right) (U_i + V_{i-1}) \end{aligned}$$

and

$$|\tilde{\rho}_i^{(2)} - \tilde{\rho}_i^{(3)}| \leq \frac{1}{m} (V_i + V_{i-1})$$

Then it follows that

$$\begin{aligned} |\hat{\rho}_n - \hat{\rho}_n^{(3)}| &\leq \frac{\sum_{j=1}^W b_j}{n - W} \sum_{i=2}^n \left( |\tilde{\rho}_i - \tilde{\rho}_i^{(2)}| + |\tilde{\rho}_i^{(2)} - \tilde{\rho}_i^{(3)}| \right) \\ &\leq \frac{2(1 + \frac{L}{m}) \sum_{j=1}^W b_j}{n - W} \sum_{i=1}^n U_i + \frac{2 \sum_{j=1}^W b_j}{m(n - W)} \sum_{i=1}^n V_i \end{aligned}$$

Suppose that

$$\frac{2(1 + \frac{L}{m}) \sum_{j=1}^W b_j}{n - W} \sum_{i=1}^n \mathbb{E}[U_i | \mathcal{F}_{i-1}] \leq \hat{U}_n$$

and

$$\frac{2 \sum_{j=1}^W b_j}{m(n - W)} \sum_{i=1}^n \mathbb{E}[V_i | \mathcal{F}_{i-1}] \leq \hat{V}_n$$

Then it holds that

$$\begin{aligned}
& \mathbb{P} \left\{ |\hat{\rho}_n - \hat{\rho}_n^{(3)}| > \hat{U}_n + \hat{V}_n + t \right\} \\
& \leq \mathbb{P} \left\{ \frac{2(1 + \frac{L}{m}) \sum_{j=1}^W b_j}{n - W} \sum_{i=1}^n U_i + \frac{2 \sum_{j=1}^W b_j}{m(n - W)} \sum_{i=1}^n V_i > \hat{U}_n + \hat{V}_n + t \right\} \\
& \leq \mathbb{P} \left\{ \frac{2(1 + \frac{L}{m}) \sum_{j=1}^W b_j}{n - W} \sum_{i=1}^n U_i > \hat{U}_n + \frac{t}{2} \right\} + \mathbb{P} \left\{ \frac{2 \sum_{j=1}^W b_j}{m(n - W)} \sum_{i=1}^n V_i > \hat{V}_n + \frac{t}{2} \right\}
\end{aligned}$$

We can apply Lemma 22 to each term to yield

$$\mathbb{P} \left\{ \frac{2(1 + \frac{L}{m}) \sum_{j=1}^W b_j}{n - W} \sum_{i=1}^n U_i > \hat{U}_n + \frac{t}{2} \right\} \leq \exp \left\{ - \frac{(n - W)^2 t^2}{32n(1 + \frac{L}{m})^2 \left( \sum_{j=1}^W b_j \right)^2 \text{diam}^2(\mathcal{X})} \right\}$$

and

$$\mathbb{P} \left\{ \frac{2 \sum_{j=1}^W b_j}{m(n - W)} \sum_{i=1}^n V_i > \hat{V}_n + \frac{t}{2} \right\} \leq \exp \left\{ - \frac{(n - W)^2 m^2 t^2}{144n C_g \left( \sum_{j=1}^W b_j \right)^2} \right\}$$

Then it holds that

$$\begin{aligned}
& \mathbb{P} \left\{ |\hat{\rho}_n - \hat{\rho}_n^{(3)}| > \hat{U}_n + \hat{V}_n + t \right\} \\
& \leq \exp \left\{ - \frac{(n - W)^2 t^2}{32n(1 + \frac{L}{m})^2 \left( \sum_{j=1}^W b_j \right)^2 \text{diam}^2(\mathcal{X})} \right\} + \exp \left\{ - \frac{(n - W)^2 m^2 t^2}{144n C_g \left( \sum_{j=1}^W b_j \right)^2} \right\}
\end{aligned}$$

We have by straightforward computation

$$\hat{U}_n = \frac{2(1 + \frac{L}{m}) \sum_{j=1}^W b_j}{n - W} \sum_{i=1}^n C(K_i)$$

and

$$\hat{V}_n = \frac{2 \sum_{j=1}^W b_j}{m(n - W)} \sum_{i=1}^n \sqrt{\frac{C_g}{dK_i}}$$

Then it holds that

$$\begin{aligned}
& \sum_{n=W+1}^{\infty} \mathbb{P} \left\{ \hat{\rho}_n < \rho - \hat{U}_n - \hat{V}_n - t_n \right\} \\
& \leq \sum_{n=W+1}^{\infty} \mathbb{P} \left\{ \hat{\rho}_n < \hat{\rho}_n^{(3)} - \hat{U}_n - \hat{V}_n - t_n \right\} \\
& \leq \sum_{n=W+1}^{\infty} \mathbb{P} \left\{ |\hat{\rho}_n - \hat{\rho}_n^{(3)}| > \hat{U}_n + \hat{V}_n + t_n \right\} \\
& \leq \sum_{n=W+1}^{\infty} \exp \left\{ - \frac{(n - W)^2 t_n^2}{32n(1 + \frac{L}{m})^2 \left( \sum_{j=1}^W b_j \right)^2 \text{diam}^2(\mathcal{X})} \right\} + \sum_{n=W+1}^{\infty} \exp \left\{ - \frac{(n - W)^2 m^2 t_n^2}{144n C_g \left( \sum_{j=1}^W b_j \right)^2} \right\} \\
& < \infty
\end{aligned}$$

By the Borel-Cantelli lemma, it follows that for all  $n$  large enough

$$\hat{\rho}_n + \hat{U}_n + \hat{V}_n + t_n \leq \rho$$

almost surely. □

### 3.5 Parameter Estimation

We may need to estimate parameters of the functions  $\{f_n\}$  such as the strong convexity parameter  $m$  to compute  $b(d_0, K)$ . We need the following assumption on our bound:

**D.1** Suppose that our bound  $b(d_0, K, \psi)$  is parameterized by  $\psi$ , which depends on properties of the function  $\ell(\mathbf{x}, \mathbf{z})$  and the distributions  $\{p_n\}_{n=1}^\infty$ . Suppose that

$$\psi_1 \leq \psi_2 \Leftrightarrow b(d_0, K, \psi_1) \leq b(d_0, K, \psi_2)$$

**D.2** There exists a true set of parameters  $\psi^*$  such that

$$\psi_n = \psi^* \quad \forall n \geq 1$$

**D.3** The spaces  $\mathcal{X}$  and  $\mathcal{Z}$  are compact

**D.4** There exists a constant  $L$  such that

$$\|\nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}, \mathbf{z})\| \leq L \|\mathbf{x} - \tilde{\mathbf{x}}\|$$

**D.5** Suppose that we know that the parameters  $\psi \in \mathcal{P}$  with  $\mathcal{P}$  compact

**D.6** Suppose that  $\nabla f_n(\mathbf{x}_n)$  has Lipschitz continuous gradients with modulus  $M$

As a consequence of Assumption D.4, it follows that there exists a constant  $G$  such that there exists a constant  $G$  such that

$$\|\nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z})\| \leq G \quad \forall \mathbf{x} \in \mathcal{X}, \mathbf{z} \in \mathcal{Z}$$

Satisfying Assumption D.5 is usually easy due to the compactness assumptions in Assumption D.4.

In most cases, we have

$$\psi = \begin{bmatrix} -m \\ M \\ A \\ B \end{bmatrix}$$

where  $m$  is the parameter of strong convexity,  $M$  is the Lipschitz gradient modulus, and the pair  $(A, B)$  controls gradient growth, i.e.,

$$\mathbb{E} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z})\|^2 \leq A + B \|\mathbf{x} - \mathbf{x}^*\|^2$$

We parameterize using  $-m$ , since smaller  $m$  increase the bound  $b(d_0, K)$ . We present several general methods for estimating these parameters, although in practice, problem specific estimators based on the form of the function may offer better performance. As an example, we present problem specific estimates for

$$\ell(\mathbf{x}, \mathbf{z}) = \frac{1}{2} (y - \mathbf{w}^\top \mathbf{x})^2 + \frac{1}{2} \lambda \|\mathbf{x}\|^2$$

As in estimating  $\rho$ , we produce one time instant estimates  $\tilde{m}_i$ ,  $\tilde{M}_i$ ,  $\tilde{A}_i$ , and  $\tilde{B}_i$  at time  $i$  and combine them. We only examine the case under Assumption D.4, although we could examine an inequality constraints as with estimating  $\rho$ . We combine estimates by averaging to yield

1.  $\hat{m}_n = \frac{1}{n} \sum_{i=1}^n \tilde{m}_i$
2.  $\hat{M}_n = \frac{1}{n} \sum_{i=1}^n \tilde{M}_i$
3.  $\hat{A}_n = \frac{1}{n} \sum_{i=1}^n \tilde{A}_i$
4.  $\hat{B}_n = \frac{1}{n} \sum_{i=1}^n \tilde{B}_i$

### 3.5.1 Estimating Strong Convexity Parameter and Lipschitz Gradient Modulus

We seek one step estimators  $\tilde{m}_n$  and  $\tilde{M}_n$  such that

$$\mathbb{E}[\tilde{m}_n \mid \mathcal{K}_{n-1}] \leq m$$

and

$$\mathbb{E}[\tilde{M}_n \mid \mathcal{K}_{n-1}] \geq M$$

with  $\{\mathcal{K}_n\}$  defined in (6).

*Hessian Method:* We exploit the fact that

$$\nabla_{\mathbf{x}\mathbf{x}}^2 f_n(\mathbf{x}) \succeq m\mathbf{I} \quad \forall \mathbf{x} \in \mathcal{X}$$

This in turn implies that

$$\lambda_{\min}(\nabla_{\mathbf{x}\mathbf{x}}^2 f_n(\mathbf{x})) \geq m \quad \forall \mathbf{x} \in \mathcal{X}$$

This suggests that given  $\{\mathbf{z}_n(k)\}_{k=1}^{K_n}$  we set

$$\tilde{m}_n \triangleq \min_{\mathbf{x} \in \mathcal{X}} \lambda_{\min} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}, \mathbf{z}_n(k)) \right)$$

Since

$$\lambda_{\min}(A) = \min_{\mathbf{v}: \|\mathbf{v}\|=1} \langle A\mathbf{v}, \mathbf{v} \rangle,$$

$\lambda_{\min}(A)$  is a concave function of  $A$ . Then by Jensen's inequality, we have

$$\begin{aligned} \mathbb{E}[\tilde{m}_n] &= \mathbb{E} \left[ \min_{\mathbf{x} \in \mathcal{X}} \lambda_{\min} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}, \mathbf{z}_n(k)) \right) \mid \mathcal{K}_{n-1} \right] \\ &\leq \min_{\mathbf{x} \in \mathcal{X}} \mathbb{E} \left[ \lambda_{\min} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}, \mathbf{z}_n(k)) \right) \mid \mathcal{K}_{n-1} \right] \\ &\leq \min_{\mathbf{x} \in \mathcal{X}} \lambda_{\min} \left( \mathbb{E} \left[ \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}, \mathbf{z}_n(k)) \mid \mathcal{K}_{n-1} \right] \right) \\ &= \min_{\mathbf{x} \in \mathcal{X}} \lambda_{\min}(\nabla_{\mathbf{x}\mathbf{x}}^2 f_n(\mathbf{x})) \\ &= m \end{aligned}$$

Similarly, we can set

$$\tilde{M}_n \triangleq \max_{\mathbf{x} \in \mathcal{X}} \lambda_{\max} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}, \mathbf{z}_n(k)) \right)$$

Since

$$\lambda_{\max}(A) = \max_{\mathbf{v}: \|\mathbf{v}\|=1} \langle A\mathbf{v}, \mathbf{v} \rangle,$$

$\lambda_{\max}(A)$  is a convex function of  $A$ . By Jensen's inequality, it holds that

$$\mathbb{E}[\tilde{M}_n \mid \mathcal{K}_{n-1}] \geq M$$

*Gradient Method To Compute  $\tilde{m}_n$ :* To actually minimize over  $\mathbf{x}$ , we can use gradient descent. To apply gradient descent, we use eigenvalue perturbation results [20]. Suppose that we have a base matrix  $T_0$  with eigenvectors  $\mathbf{v}_{0i}$  and eigenvalues  $\lambda_{0i}$ . We want to find the eigenvectors  $\mathbf{v}_i$  and eigenvalues  $\lambda_i$  of a perturbed matrix  $T$ :

$$\begin{aligned} T_0 \mathbf{v}_{0i} &= \lambda_{0i} \mathbf{v}_{0i} \\ T \mathbf{v}_i &= \lambda_i \mathbf{v}_i \end{aligned}$$

In particular, we want to relate  $\lambda_{0i}$  to  $\lambda_i$ . With

$$\delta \mathbf{T} \triangleq \mathbf{T} - \mathbf{T}_0,$$

we have

$$\delta \lambda_i = \mathbf{v}_{0i}^\top (\delta \mathbf{T}) \mathbf{v}_{0i}$$

and

$$\frac{\partial \lambda_i}{\partial \mathbf{T}_{ij}} = \mathbf{v}_{0i}(i) \mathbf{v}_{0j}(2 - \delta_{ij})$$

Suppose we are given a matrix-valued function  $\mathbf{T}(\mathbf{x})$  with

$$\mathbf{T}(\mathbf{x}) \mathbf{v}(\mathbf{x}) = \lambda_{\min}(\mathbf{x}) \mathbf{v}(\mathbf{x})$$

Then it holds that

$$\begin{aligned} \nabla_{\mathbf{x}} \lambda_{\min}(\mathbf{T}(\mathbf{x})) &= \sum_{i,j} \frac{\partial \lambda_{\min}}{\partial \mathbf{T}_{ij}} \nabla_{\mathbf{x}} \mathbf{T}_{ij}(\mathbf{x}) \\ &= \sum_{i,j} \mathbf{v}_i(\mathbf{x}) \mathbf{v}_j(\mathbf{x}) (2 - \delta_{ij}) \nabla_{\mathbf{x}} \mathbf{T}_{ij}(\mathbf{x}) \end{aligned}$$

Then we can use gradient descent to solve

$$\min_{\mathbf{x} \in \mathcal{X}} \lambda_{\min} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}_n(k)) \right)$$

Starting from any  $\mathbf{x}(0)$ , we can compute

$$\mathbf{x}(p) = \Pi_{\mathcal{X}} \left[ \mathbf{x}(p-1) - \mu \nabla_{\mathbf{x}} \lambda_{\min} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}, \mathbf{z}_n(k)) \right) \right] \quad p = 1, \dots, P$$

and set

$$\hat{m}_n \triangleq \lambda_{\min} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}(P), \mathbf{z}_n(k)) \right) \quad (15)$$

*Heuristic Method:* For any two points  $\mathbf{x}$  and  $\mathbf{y}$ , we have by strong convexity

$$f_n(\mathbf{y}) \geq f_n(\mathbf{x}) + \langle \nabla f_n(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} m \|\mathbf{y} - \mathbf{x}\|^2$$

Suppose that we have  $N$  points  $\mathbf{x}(1), \dots, \mathbf{x}(N)$ . Then we know that for any two distinct points  $\mathbf{x}_i$  and  $\mathbf{x}_j$

$$m \leq \frac{f_n(\mathbf{x}(i)) - f_n(\mathbf{x}(j)) - \langle \nabla f_n(\mathbf{x}(j)), \mathbf{x}(i) - \mathbf{x}(j) \rangle}{\frac{1}{2} \|\mathbf{x}(i) - \mathbf{x}(j)\|^2}$$

This suggests the estimator

$$\hat{m}_n \triangleq \min_{i \neq j} \frac{\frac{1}{K_n} \sum_{k=1}^{K_n} \ell(\mathbf{x}(i), \mathbf{z}_n(k)) - \frac{1}{K_n} \sum_{k=1}^{K_n} \ell(\mathbf{x}(j), \mathbf{z}_n(k)) - \left\langle \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}} \ell(\mathbf{x}(j), \mathbf{z}_n(k)), \mathbf{x}(i) - \mathbf{x}(j) \right\rangle}{\frac{1}{2} \|\mathbf{x}(i) - \mathbf{x}(j)\|^2} \quad (16)$$

for the strong convexity parameter. Then we have

$$\begin{aligned}
& \mathbb{E}[\hat{m}_n] \\
&= \mathbb{E} \left[ \min_{i \neq j} \frac{\frac{1}{K_n} \sum_{k=1}^{K_n} \ell(\mathbf{x}(i), \mathbf{z}_n(k)) - \frac{1}{K_n} \sum_{k=1}^{K_n} \ell(\mathbf{x}(j), \mathbf{z}_n(k)) - \left\langle \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}} \ell(\mathbf{x}(j), \mathbf{z}_n(k)), \mathbf{x}(i) - \mathbf{x}(j) \right\rangle}{\frac{1}{2} \|\mathbf{x}(i) - \mathbf{x}(j)\|^2} \right] \\
&\leq \min_{i \neq j} \mathbb{E} \left[ \frac{\frac{1}{K_n} \sum_{k=1}^{K_n} \ell(\mathbf{x}(i), \mathbf{z}_n(k)) - \frac{1}{K_n} \sum_{k=1}^{K_n} \ell(\mathbf{x}(j), \mathbf{z}_n(k)) - \left\langle \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}} \ell(\mathbf{x}(j), \mathbf{z}_n(k)), \mathbf{x}(i) - \mathbf{x}(j) \right\rangle}{\frac{1}{2} \|\mathbf{x}(i) - \mathbf{x}(j)\|^2} \right] \\
&\leq \min_{i \neq j} \frac{f_n(\mathbf{x}(i)) - f_n(\mathbf{x}(j)) - \langle \nabla f_n(\mathbf{x}(j)), \mathbf{x}(i) - \mathbf{x}(j) \rangle}{\frac{1}{2} \|\mathbf{x}(i) - \mathbf{x}(j)\|^2}
\end{aligned}$$

It is difficult to compare this estimator to  $m$  exactly. All we can say is that

$$m \leq \min_{i \neq j} \frac{f_n(\mathbf{x}(i)) - f_n(\mathbf{x}(j)) - \langle \nabla f_n(\mathbf{x}(j)), \mathbf{x}(i) - \mathbf{x}(j) \rangle}{\frac{1}{2} \|\mathbf{x}(i) - \mathbf{x}(j)\|^2}$$

as well. In practice, this method produces estimates close to  $m$ .

Similarly, we can set

$$\hat{M}_n \triangleq \max_{i \neq j} \frac{\frac{1}{K_n} \sum_{k=1}^{K_n} \ell(\mathbf{x}(i), \mathbf{z}_n(k)) - \frac{1}{K_n} \sum_{k=1}^{K_n} \ell(\mathbf{x}(j), \mathbf{z}_n(k)) - \left\langle \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}} \ell(\mathbf{x}(j), \mathbf{z}_n(k)), \mathbf{x}(i) - \mathbf{x}(j) \right\rangle}{\frac{1}{2} \|\mathbf{x}(i) - \mathbf{x}(j)\|^2} \quad (17)$$

*Problem Specific:* For the penalized quadratic, we have

$$\nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}, \mathbf{z}) = \lambda \mathbf{I} + \mathbf{w}\mathbf{w}^\top$$

so

$$\nabla_{\mathbf{x}\mathbf{x}}^2 f_n(\mathbf{x}) = \lambda \mathbf{I} + \mathbb{E}[\mathbf{w}_n \mathbf{w}_n^\top]$$

This suggests the simple closed-form estimates

$$\tilde{m}_n = \lambda + \lambda_{\min} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \mathbf{w}_n(k) \mathbf{w}_n(k)^\top \right)$$

and

$$\tilde{M}_n = \lambda + \lambda_{\max} \left( \frac{1}{K_n} \sum_{k=1}^{K_n} \mathbf{w}_n(k) \mathbf{w}_n(k)^\top \right)$$

Again, by Jensen's inequality, it holds that

$$\mathbb{E}[\tilde{m}_n \mid \mathcal{K}_{n-1}] \leq m$$

and

$$\mathbb{E}[\tilde{M}_n \mid \mathcal{K}_{n-1}] \geq M$$

*Combining Estimates:* We now look at combining the single time instant estimates of the strong convexity parameter and the Lipschitz gradient modulus.

**Lemma 9.** Choose  $t_n$  such that for all  $C > 0$  it holds that

$$\sum_{n=1}^{\infty} e^{-Cm_n^2} < +\infty$$

Then for all  $n$  large enough it holds that

1.  $\hat{m}_n - t_n \leq m$
2.  $\hat{M}_n + t_n \geq M$

almost surely.

*Proof.* By the compactness of the space  $\mathcal{P}$  containing  $\psi$ , we can apply the dependent version of Hoeffding's lemma (Lemma 23) to yield

$$\mathbb{E} \left[ e^{s\tilde{m}_i} \mid \mathcal{H}_{i-1} \right] \leq \exp \left\{ \frac{1}{2} \sigma_m^2 s^2 \right\}$$

and

$$\mathbb{E} \left[ e^{s\tilde{M}_i} \mid \mathcal{H}_{i-1} \right] \leq \exp \left\{ \frac{1}{2} \sigma_M^2 s^2 \right\}$$

for some constants  $\sigma_m^2$  and  $\sigma_M^2$  derived from Hoeffding's lemma. Then applying Lemma 22, it follows that

$$\mathbb{P} \left\{ \hat{m}_n > \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\tilde{m}_i \mid \mathcal{H}_{i-1}] + t_n \right\} \leq \exp \left\{ -\frac{nt_n^2}{2\sigma_m^2} \right\}$$

We know that

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\tilde{m}_i \mid \mathcal{H}_{i-1}] > m$$

so it follows that

$$\mathbb{P} \{ \hat{m}_n > m + t_n \} \leq \exp \left\{ -\frac{nt_n^2}{2\sigma_m^2} \right\}$$

Similarly, for the Lipschitz gradient modulus, it holds that

$$\mathbb{P} \{ \hat{M}_n < M - t_n \} \leq \exp \left\{ -\frac{nt_n^2}{2\sigma_M^2} \right\}$$

As before, we have

$$\sum_{n=1}^{\infty} \mathbb{P} \{ \hat{m}_n > m + t_n \} \leq \sum_{n=1}^{\infty} \exp \left\{ -\frac{nt_n^2}{2\sigma_m^2} \right\} < +\infty$$

and

$$\sum_{n=1}^{\infty} \mathbb{P} \{ \hat{M}_n < M - t_n \} \leq \sum_{n=1}^{\infty} \exp \left\{ -\frac{nt_n^2}{2\sigma_M^2} \right\} < +\infty$$

to ensure that almost surely for all  $n$  large enough it holds that

$$\hat{m}_n - t_n \leq m$$

and

$$\hat{M}_n + t_n \geq m$$

□

For Lemma 9, we need  $t_n$  to decay no faster than  $\mathcal{O}(n^{-1/2})$ .

### 3.5.2 Estimating Gradient Parameters

From Assumption D.6 , it holds that

$$\begin{aligned}\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z})\|^2 &= \mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}^*, \mathbf{z}) + (\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}}\ell(\mathbf{x}^*, \mathbf{z}))\|^2 \\ &\leq 2\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}^*, \mathbf{z})\|^2 + 2\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}}\ell(\mathbf{x}^*, \mathbf{z})\|^2 \\ &\leq 2\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}^*, \mathbf{z})\|^2 + 2M^2\|\mathbf{x} - \mathbf{x}^*\|^2\end{aligned}$$

Thus, we can set

$$B = 2M^2$$

and

$$A = 2\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}^*, \mathbf{z})\|^2$$

This suggests that given an estimate  $\tilde{M}_n$  for  $M$ , we set

$$\tilde{B}_n = 2\tilde{M}_n^2$$

Then by Jensen's inequality, we have

$$\begin{aligned}\mathbb{E}[\tilde{B}_n \mid \mathcal{K}_{n-1}] &= 2\mathbb{E}[\tilde{M}_n^2 \mid \mathcal{K}_{n-1}] \\ &\geq 2(\mathbb{E}[\tilde{B}_n \mid \mathcal{K}_{n-1}])^2 \\ &\geq 2M^2 \\ &= B\end{aligned}$$

**Lemma 10.** Choose  $t_n$  such that for all  $C > 0$  it holds that

$$\sum_{n=1}^{\infty} e^{-Cn t_n^2} < +\infty$$

Then for all  $n$  large enough it holds that

$$\hat{B}_n + t_n \geq B$$

almost surely.

*Proof.* By identical reasoning for the strong convexity and Lipschitz continuous gradients, it holds that

$$\mathbb{P}\{\hat{B}_n < B - t_n\} \leq \exp\left\{-\frac{nt_n^2}{2\sigma_B^2}\right\}$$

Since we have

$$\sum_{n=1}^{\infty} \exp\left\{-\frac{nt_n^2}{2\sigma_B^2}\right\} < +\infty$$

for all  $n$  large enough it holds that

$$\hat{B}_n + t_n \geq B$$

almost surely. □

To estimate  $A$ , consider using a point  $\mathbf{x}$  to approximate  $\mathbf{x}^*$ . It holds that

$$\begin{aligned}\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}^*, \mathbf{z})\|^2 &= \mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z}) + (\nabla_{\mathbf{x}}\ell(\mathbf{x}^*, \mathbf{z}) - \nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z}))\|^2 \\ &\leq 2\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z})\|^2 + 2\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}^*, \mathbf{z}) - \nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z})\|^2 \\ &\leq 2\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z})\|^2 + 2M^2\mathbb{E}\|\mathbf{x} - \mathbf{x}^*\|^2 \\ &\leq 2\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z})\|^2 + 2\left(\frac{M}{m}\right)^2 \|\nabla f(\mathbf{x})\|^2 \\ &\leq 2\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z})\|^2 + 2\left(\frac{M}{m}\right)^2 \|\nabla f(\mathbf{x})\|^2\end{aligned}$$

This suggests the estimate

$$\tilde{A}_n(\mathbf{x}) = \frac{2}{K_n} \sum_{k=1}^{K_n} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}_n(k))\|^2 + 4 \left( \frac{\tilde{M}_{n-1} + t_{n-1}}{\tilde{m}_{n-1} - t_{n-1}} \right)^2 \left\| \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}_n(k)) \right\|^2$$

**Lemma 11.** For any  $\mathbf{x}$  possibly random but not a function of  $\{\mathbf{z}_n(k)\}_{k=1}^{K_n}$  and all  $n$  large enough, it holds that

$$\mathbb{E}[\tilde{A}_n \mid \mathcal{K}_{n-1}] \geq A$$

*Proof.* For any  $\mathbf{x}$  possibly random but not a function of  $\{\mathbf{z}_n(k)\}_{k=1}^{K_n}$ , it holds that

$$\begin{aligned} & \mathbb{E}[\tilde{A}_n \mid \mathcal{K}_{n-1}] \\ &= \mathbb{E} \left[ \frac{2}{K_n} \sum_{k=1}^{K_n} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}_n(k))\|^2 + 4 \left( \frac{\tilde{M}_{n-1} + t_{n-1}}{\tilde{m}_{n-1} - t_{n-1}} \right)^2 \left\| \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}_n(k)) \right\|^2 \mid \mathcal{K}_{n-1} \right] \\ &= \mathbb{E} \left[ \frac{2}{K_n} \sum_{k=1}^{K_n} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}_n(k))\|^2 \mid \mathcal{K}_{n-1} \right] + 4 \left( \frac{\tilde{M}_{n-1} + t_{n-1}}{\tilde{m}_{n-1} - t_{n-1}} \right)^2 \mathbb{E} \left[ \left\| \frac{1}{K_n} \sum_{k=1}^{K_n} \nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}_n(k)) \right\|^2 \mid \mathcal{K}_{n-1} \right] \\ &\geq 2\mathbb{E} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}_n)\|^2 + 4 \left( \frac{\tilde{M}_{n-1} + t_{n-1}}{\tilde{m}_{n-1} - t_{n-1}} \right)^2 \|\nabla f_n(\mathbf{x})\|^2 \end{aligned}$$

The last inequality uses Jensen's inequality. Then by our prior analysis, almost surely for all  $n$  sufficiently large it holds that

$$\frac{\tilde{M}_{n-1} + t_{n-1}}{\tilde{m}_{n-1} - t_{n-1}} \geq \frac{M}{m}$$

and so for all  $n$  sufficiently large

$$\begin{aligned} \mathbb{E}[\tilde{A}_n \mid \mathcal{K}_{n-1}] &\geq 2\mathbb{E} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z}_n)\|^2 + 4 \left( \frac{M}{m} \right)^2 \|\nabla f_n(\mathbf{x})\|^2 \\ &= 2\mathbb{E} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}_n^*, \mathbf{z}_n)\|^2 \\ &= A \end{aligned}$$

Therefore, for all  $n$  sufficiently large (dependent on estimation of  $m$  and  $M$ ), it holds that

$$\mathbb{E}[\tilde{A}_n \mid \mathcal{K}_{n-1}] \geq A$$

□

*Combining Estimates for A:* In practice, we use  $\tilde{A}_n(\mathbf{x}_n)$ , which complicates the analysis due to the fact that  $\mathbf{x}_n$  is computed using the same samples  $\{\mathbf{z}_n(k)\}_{k=1}^{K_n}$ .

**Lemma 12.** Choose  $t_n$  such that for all  $C > 0$  it holds that

$$\sum_{n=1}^{\infty} e^{-Cn t_n^2} < +\infty$$

Then for all  $n$  large enough it holds that

$$\hat{A}_n + t_n \geq A$$

almost surely.

*Proof.* Consider the following three estimates of  $A$  all computed with knowledge of  $m$  and  $M$  and  $\tilde{\mathbf{x}}_n$  as in Lemma 2:

$$\begin{aligned}\tilde{A}_i^{(2)} &= \frac{2}{K_i} \sum_{k=1}^{K_i} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k))\|^2 + 4 \left(\frac{M}{m}\right)^2 \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) \right\|^2 \\ \tilde{A}_i^{(3)} &= \frac{2}{K_i} \sum_{k=1}^{K_i} \|\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k))\|^2 + 4 \left(\frac{M}{m}\right)^2 \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) \right\|^2 \\ \tilde{A}_i^{(4)} &= 2\mathbb{E} \|\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i)\|^2 + 4 \left(\frac{M}{m}\right)^2 \|\nabla f_i(\tilde{\mathbf{x}}_i)\|^2\end{aligned}$$

Define the averaged estimates

$$\begin{aligned}\hat{A}_n^{(2)} &= \frac{1}{n} \sum_{i=1}^n \tilde{A}_i^{(2)} \\ \hat{A}_n^{(3)} &= \frac{1}{n} \sum_{i=1}^n \tilde{A}_i^{(3)} \\ \hat{A}_n^{(4)} &= \frac{1}{n} \sum_{i=1}^n \tilde{A}_i^{(4)}\end{aligned}$$

We always have

$$\tilde{A}_i^{(4)} \geq A$$

so

$$\hat{A}_n^{(4)} \geq A$$

First, we show that  $\hat{A}_n^{(2)}$  is close to  $\hat{A}_n^{(3)}$ . We have

$$\begin{aligned}& |\tilde{A}_i^{(2)} - \tilde{A}_i^{(3)}| \\ & \leq 2 \left| \frac{1}{K_i} \sum_{k=1}^{K_i} (\|\nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k))\|^2 - \|\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k))\|^2) \right| \\ & \quad + 4 \left(\frac{M}{m}\right)^2 \left| \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) \right\|^2 - \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) \right\|^2 \right| \\ & \leq 4G \frac{1}{K_i} \sum_{k=1}^{K_i} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k))\| + 8G \left(\frac{M}{m}\right)^2 \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) - \nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k))) \right\|^2 \\ & \leq \left(4 + 8 \left(\frac{M}{m}\right)^2\right) GM \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|\end{aligned}$$

yielding

$$|\hat{A}_n^{(2)} - \hat{A}_n^{(3)}| \leq \left(4 + 8 \left(\frac{M}{m}\right)^2\right) GM \left(\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|\right)$$

Second, we have

$$\begin{aligned}& |\hat{A}_n^{(3)} - \hat{A}_n^{(4)}| \\ & \leq \left| \frac{1}{n} \sum_{i=1}^n \left( \frac{2}{K_i} \sum_{k=1}^{K_i} (\|\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k))\|^2 - \mathbb{E} [\|\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i)\|^2 | \mathcal{F}_{n-1}]) \right) \right| \\ & \quad + 8 \left(\frac{M}{m}\right)^2 G \frac{1}{n} \sum_{i=1}^n \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla f_i(\tilde{\mathbf{x}}_i)) \right\|\end{aligned}$$

Combining both inequalities, we know that

$$\begin{aligned}
& |\hat{A}_n^{(2)} - \hat{A}_n^{(4)}| \\
& \leq \left( 4 + 8 \left( \frac{M}{m} \right)^2 \right) GM \left( \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| \right) \\
& \quad + \left| \frac{1}{n} \sum_{i=1}^n \left( \frac{2}{K_i} \sum_{k=1}^{K_i} (\|\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k))\|^2 - \mathbb{E} [\|\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i)\|^2 | \mathcal{F}_{n-1}]) \right) \right| \\
& \quad + 8 \left( \frac{M}{m} \right)^2 G \frac{1}{n} \sum_{i=1}^n \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} (\nabla_{\mathbf{x}} \ell(\tilde{\mathbf{x}}_i, \mathbf{z}_i(k)) - \nabla f_i(\tilde{\mathbf{x}}_i)) \right\|
\end{aligned}$$

The first and third terms in this bound can be controlled by the analysis of the direct estimate and the second term by Lemma (22). This shows that

$$\begin{aligned}
& \mathbb{P} \left\{ \hat{A}_n^{(2)} < A - \frac{1}{n} \sum_{i=1}^n \frac{C_i}{\sqrt{K_i}} - t_n \right\} \\
& \leq \mathbb{P} \left\{ \hat{A}_n^{(2)} < \hat{A}_n^{(4)} - \frac{1}{n} \sum_{i=1}^n \frac{C_i}{\sqrt{K_i}} - t_n \right\} \\
& \leq \mathbb{P} \left\{ |\hat{A}_n^{(2)} - \hat{A}_n^{(4)}| > \frac{1}{n} \sum_{i=1}^n \frac{C_i}{\sqrt{K_i}} t_n \right\} \\
& \leq 2 \exp \left\{ -\frac{nt_n^2}{2\sigma_{A2}^2} \right\}
\end{aligned}$$

Since

$$\sum_{n=1}^{\infty} \mathbb{P} \left\{ \hat{A}_n^{(2)} < A - \frac{1}{n} \sum_{i=1}^n \frac{C_i}{\sqrt{K_i}} - t_n \right\} \leq \sum_{n=1}^{\infty} C \exp \left\{ -\frac{nt_n^2}{2\sigma_{A2}^2} \right\} < +\infty$$

almost surely for all  $n$  large enough, it holds that

$$\hat{A}_n^{(2)} + \frac{1}{n} \sum_{i=1}^n \frac{C_i}{\sqrt{K_i}} + t_n \geq A$$

In addition, we have

$$\hat{A}_n^{(2)} + \frac{1}{n} \sum_{i=1}^n \frac{C_i}{\sqrt{K_i}} + 2t_n \geq A$$

There exists a random variable  $\tilde{N}$  such that

$$n \geq \tilde{N} \Rightarrow \frac{M_n + t_n}{m_n - t_n} \geq \frac{M}{m}$$

Then for  $n \geq \tilde{N}$ , it holds that

$$\begin{aligned}
& \hat{A}_n - \hat{A}_n^{(2)} \\
& = \frac{4}{n} \sum_{i=1}^n \left[ \left( \frac{\hat{M}_{i-1} + t_{i-1}}{\hat{m}_{i-1} - t_{i-1}} \right)^2 - \left( \frac{M}{m} \right)^2 \right] \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) \right\|^2 \\
& \geq \frac{4}{n} \sum_{i=1}^{\tilde{N}-1} \left[ \left( \frac{\hat{M}_{i-1} + t_{i-1}}{\hat{m}_{i-1} - t_{i-1}} \right)^2 - \left( \frac{M}{m} \right)^2 \right] \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) \right\|^2
\end{aligned}$$

Since our choice of  $t_n$  can decay only as fast as  $C/\sqrt{n}$ , it follows that

$$\frac{4}{n} \sum_{i=1}^{\tilde{N}-1} \left[ \left( \frac{\hat{M}_{i-1} + t_{i-1}}{\hat{m}_{i-1} - t_{i-1}} \right)^2 - \left( \frac{M}{m} \right)^2 \right] \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) \right\|^2 - t_n < 0$$

for all  $n$  large enough. This implies that

$$\begin{aligned} \hat{A}_n + \frac{1}{n} \sum_{i=1}^n \frac{C_i}{\sqrt{K_i}} + t_n &\geq \hat{A}_n - \left( \frac{4}{n} \sum_{i=1}^{\tilde{N}-1} \left[ \left( \frac{M}{m} \right)^2 - \left( \frac{\hat{M}_{i-1} + t_{i-1}}{\hat{m}_{i-1} + t_{i-1}} \right)^2 \right] \left\| \frac{1}{K_i} \sum_{k=1}^{K_i} \nabla_{\mathbf{x}} \ell(\mathbf{x}_i, \mathbf{z}_i(k)) \right\|^2 - t_n \right) + \frac{1}{n} \sum_{i=1}^n \frac{C_i}{\sqrt{K_i}} + t_n \\ &\geq \hat{A}_n^{(2)} + \frac{1}{n} \sum_{i=1}^n \frac{C_i}{\sqrt{K_i}} + 2t_n \\ &\geq A \end{aligned}$$

for  $n$  large enough. □

Using these estimates, we have constructed estimates  $\hat{\psi}_n$  such that for all  $n$  large enough it holds that

$$\hat{\psi}_n + C_n + t_n \mathbf{1} \geq \psi^*$$

for appropriate constants  $C_n$  almost surely. Therefore, by assumption for all  $n$  large enough it holds that

$$b(d_0, K, \psi^*) \leq b(d_0, K, \hat{\psi}_n + t_n)$$

### 3.5.3 Effect on $\rho$ Estimation

Our analysis of estimating  $\rho$  assumes that we know the parameters of the function and in particular the strong convexity parameter  $m$ . We now argue that the effect of using estimated parameters instead is minimal. This happens because we know that for all  $n$  large enough it holds that

$$\hat{\psi}_n \geq \psi^*$$

almost surely.

**Lemma 13.** *We want to estimate a non-negative parameter  $\phi^*$  by producing a sequence of estimates  $\phi_i$  for all  $i \geq 1$  and averaging to produce*

$$\hat{\phi}_n = \frac{1}{n} \sum_{i=1}^n \phi_i$$

where the estimates  $\phi_i$  are dependent on an auxiliary sequence  $\psi_i$  in the sense that  $\phi_i(\psi_i)$ . Suppose that the following conditions hold:

1. Suppose that there exists a random variable  $\tilde{N}$  such that  $n \geq \tilde{N}$  implies that  $\hat{\psi}_n \geq \psi^*$
2.  $\mathbb{E}[\phi_i(\psi^*)] \geq \phi^*$

Then it follows that

$$\liminf_{n \rightarrow \infty} \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \phi_i \right] \geq \phi^*$$

*Proof.* It holds that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \phi_i &= \frac{1}{n} \sum_{i=1}^{\tilde{N}-1} \phi_i(\psi_i) + \frac{1}{n} \sum_{i=\tilde{N}}^n \phi_i(\psi_i) \\ &\geq \frac{1}{n} \sum_{i=1}^{\tilde{N}-1} \phi_i(\psi_i) + \frac{1}{n} \sum_{i=\tilde{N}}^n \phi_i(\psi_i^*) \end{aligned} \quad (18)$$

Therefore, it follows that

$$\begin{aligned} \liminf_{n \rightarrow \infty} \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \phi_i \right] &\geq \liminf_{n \rightarrow \infty} \mathbb{E} \left[ \frac{1}{n} \sum_{i=\tilde{N}}^n \phi_i(\psi_i^*) \right] \\ &\geq \phi^* \end{aligned}$$

□

We can extend all the concentration inequalities for estimating  $\rho$  as well by extending the inequality in (18) to yield

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \phi_i &= \frac{1}{n} \sum_{i=1}^{\tilde{N}-1} \phi_i(\psi_i) + \frac{1}{n} \sum_{i=\tilde{N}}^n \phi_i(\psi_i) \\ &\geq \frac{1}{n} \sum_{i=1}^{\tilde{N}-1} \phi_i(\psi_i) + \frac{1}{n} \sum_{i=\tilde{N}}^n \phi_i(\psi_i^*) \\ &\geq \frac{1}{n} \sum_{i=1}^{\tilde{N}-1} (\phi_i(\psi_i) - \phi_i(\psi_i^*)) + \frac{1}{n} \sum_{i=\tilde{N}}^n \phi_i(\psi_i^*) \\ &= \frac{1}{n} \sum_{i=1}^n \phi_i(\psi_i^*) + o(1) \end{aligned}$$

Before, we have analyzed

$$\frac{1}{n} \sum_{i=1}^n \phi_i(\psi_i^*)$$

so for large enough  $n$ , we recover previous results, since the  $o(1)$  term goes to 0.

## 4 Adaptive Sequential Optimization With $\rho$ Unknown

We now examine the case with  $\rho$  unknown. We extend the work of Section 2 using the estimates of  $\rho$  in Section 3. Our analysis depends on the following crucial assumption:

**C.1** For appropriate sequences  $\{t_n\}$ , for all  $n$  sufficiently large it holds that  $\hat{\rho}_n + t_n \geq \rho$  almost surely.

**C.2**  $b(d_0, K_n)$  factors as  $b(d_0, K_n) = \alpha(K_n)d_0 + \beta(K_n)$

We have demonstrated that assumption C.1 that holds for the direct and IPM estimates of  $\rho$  under (2) and (3). Note that whether we assume (2) or (3) does not matter for analysis.

### 4.1 General Condition on $K_n$

We start with a general result showing that for any choice of  $K_n$  such that  $K_n \geq K^*$  for all  $n$  large enough the excess risk is controlled in the sense that

$$\limsup_{n \rightarrow \infty} (\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*)) \leq \varepsilon$$

We then apply this result to two different selection rules for Kn.

Consider the function

$$\phi_K(v) = \alpha(K) \left( \sqrt{\frac{2}{m}} v + \rho \right)^2 + \beta(K)$$

derived from assumption C.2. Note that as a function of  $v$ ,  $\phi_K(v)$  is clearly increasing and strictly concave. First, suppose that we select  $K^*$  defined in (5). Then by definition it holds that

$$\phi_{K^*}(\varepsilon) \leq \varepsilon$$

We study fixed points of the function  $\phi_{K^*}(v)$ :

**Lemma 14.** *The function  $\phi_{K^*}(v)$  has a unique positive fixed point  $\bar{v}$  with*

1.  $\bar{v} = \phi_{K^*}(\bar{v}) \leq \varepsilon$
2.  $\phi'_{K^*}(\bar{v}) < 1$

*Proof.* We have

$$\phi_{K^*}(0) = \alpha(K^*)\rho^2 + \beta(K^*) > 0$$

Since

$$\lim_{v \rightarrow 0} \phi_{K^*}(v) = \phi_{K^*}(0)$$

and  $\phi_{K^*}(0) > 0$ , there exists a positive  $a$  sufficiently small that

$$\phi_{K^*}(a) > a$$

Next, expanding  $\phi_K(v)$  yields

$$\phi_K(v) = \frac{2}{m} \alpha(K)v + 2\alpha(K)\rho \sqrt{\frac{2}{m}} \sqrt{v} + \alpha(K)\rho^2 + \beta(K)$$

Since  $\phi_{K^*}(\varepsilon) \leq \varepsilon$ , we obviously must have  $\frac{2}{m} \alpha(K^*) \leq 1$ . Suppose that

$$\frac{2}{m} \alpha(K^*) = 1$$

Then it holds that

$$\phi_{K^*}(\varepsilon) = \varepsilon + \sqrt{2m}\rho\sqrt{\varepsilon} + \frac{m}{2}\rho^2 + \beta(K) > \varepsilon$$

This is a contradiction, so it holds that

$$\frac{2}{m} \alpha(K^*) < 1$$

It is thus readily apparent that

$$v - \phi_{K^*}(v) \rightarrow \infty$$

as  $v \rightarrow \infty$ . Therefore, there exists a point  $b > a$  such that

$$\phi_{K^*}(b) < b$$

It is easy to check that  $\phi_{K^*}(v)$  is increasing and strictly concave. Therefore, we can apply Theorem 3.3 from [21] to conclude that there exists a unique, positive fixed point  $\bar{v}$  of  $\phi_{K^*}(v)$ .

Next, suppose that  $\phi'_{K^*}(\bar{v}) > 1$ . Then by Taylor's Theorem for  $v > \bar{v}$  sufficiently close to  $\bar{v}$ , we have

$$\phi_{K^*}(v) > v$$

However, we know that as  $v \rightarrow \infty$ , it holds that  $v - \phi_{K^*}(v) \rightarrow \infty$ . By the Intermediate Value Theorem, this implies that there is another fixed point on  $[v, \infty)$ . This is a contradiction, since  $\bar{v}$  is the unique, positive fixed point. Therefore, it holds that  $\phi'_{K^*}(\bar{v}) \leq 1$ . Now, suppose that  $\phi'_{K^*}(\bar{v}) = 1$ . Since  $\phi_{K^*}(v)$  is strictly concave, its derivative is decreasing [22]. Therefore, on  $[0, \bar{v})$ , it holds that

$$\phi'_{K^*}(v) > 1$$

This implies that

$$\begin{aligned} \phi_{K^*}(\bar{v}) &= \phi_{K^*}(0) + \int_0^{\bar{v}} \phi'_{K^*}(v) dx \\ &\geq \phi_{K^*}(0) + \bar{v} \\ &> \bar{v} \end{aligned}$$

This is a contradiction, so it must be that  $\phi'_{K^*}(\bar{v}) < 1$ . □

As a simple consequence of the concavity of  $\phi_{K^*}(v)$ , we can study a fixed point iteration involving  $\phi_K(v)$ . Define the  $n$ -fold composition mapping

$$\phi_K^{(n)}(v) \triangleq (\phi_K \circ \dots \circ \phi_K)(v)$$

**Lemma 15.** *For any  $v > 0$ , it holds that*

$$\lim_{n \rightarrow \infty} \phi_{K^*}^{(n)}(v) = \bar{v}$$

*Proof.* Following [23], for any fixed point  $\bar{v}$ , it holds that

$$|\phi_{K^*}(v) - \bar{v}| \leq \phi'_{K^*}(\bar{v})|v - \bar{v}|$$

Therefore, applying the fixed point property repeatedly yields

$$|\phi_{K^*}^{(n)}(v) - \bar{v}| \leq (\phi'_{K^*}(\bar{v}))^n |v - \bar{v}|$$

By Lemma 14, it holds that

$$\phi'_{K^*}(\bar{v}) < 1$$

and so the result follows. □

Now, we show that we appropriately control the excess risk when we estimate  $\rho$ . The extension of this argument to the case when we also estimate function parameters  $\psi$  is straightforward. If we have

$$p(\{z_n(k)\}_{k=1}^{K_n} | \mathbf{x}_{n-1}, K_n) = \prod_{k=1}^{K_n} p_n(z_n(k))$$

then

$$\mathbb{E}[f_n(\mathbf{x}_n) | \mathbf{x}_{n-1}, K_n] - f_n(\mathbf{x}_n^*) \leq b \left( \left( \sqrt{\frac{2}{m}} (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)) + \rho \right)^2, K_n \right)$$

Therefore, it holds that

$$\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*) \leq \mathbb{E} \left[ b \left( \left( \sqrt{\frac{2}{m}} (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)) + \rho \right)^2, K_n \right) \right]$$

Suppose that we set

$$\mathcal{H}_\infty = \sigma(\{K_n\}_{n=1}^\infty \cup \{\hat{\rho}_n\}_{n=2}^\infty)$$

This sigma algebra contains all the information about  $\{\hat{\rho}_n\}$  and thus  $\{K_n\}$ . Then, we do not have

$$p(\{z_n(k)\}_{k=1}^{K_n} | \mathcal{H}_\infty) = \prod_{k=1}^{K_n} p_n(z_n(k))$$

since  $K_{n+1}, K_{n+2}, \dots$  are a function of  $\{K_n\}_{k=1}^{K_n}$ . We do not even have

$$\mathbb{E}[f_n(\mathbf{x}_n) | \mathcal{H}_\infty] - f_n(\mathbf{x}_n^*) \leq b \left( \left( \sqrt{\frac{2}{m}} (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)) + \rho \right)^2, K_n \right)$$

However, we would expect that this is not too far from true. Conceptually, we consider running our approach twice on independent samples. The first run determines the required number of samples  $\{K_n\}_{n=1}^\infty$ . We then run our process for a second run with these fixed choices of  $\{K_n\}_{n=1}^\infty$  and independent samples as in Figure 1. For the second run, it is true that

$$p(\{z_n^{(2)}(k)\}_{k=1}^{K_n} | \mathcal{H}_\infty) = \prod_{k=1}^{K_n} p_n(z_n^{(2)}(k))$$

and

$$\mathbb{E}[f_n(\mathbf{x}_n^{(2)}) | \mathcal{H}_\infty] - f_n(\mathbf{x}_n^*) \leq b \left( \left( \sqrt{\frac{2}{m}} (f_{n-1}(\mathbf{x}_{n-1}^{(2)}) - f_{n-1}(\mathbf{x}_{n-1}^*)) + \rho \right)^2, K_n \right)$$

In practice, we do not need to run our process twice. This is only a proof technique. Now, for the second run the recursion

$$\varepsilon_n^{(2)} = b \left( \left( \sqrt{\frac{2}{m}} \varepsilon_{n-1}^{(2)} + \rho \right)^2, K_n \right) \quad \forall n \geq 3 \quad (19)$$

with  $\varepsilon_1$  and  $\varepsilon_2$  from Assumption A.4 bounds the excess risk of the second run

$$\mathbb{E}[f_n(\mathbf{x}_n^{(2)}) | \mathcal{H}_\infty] - f_n(\mathbf{x}_n^*) \leq \varepsilon_n^{(2)}$$

Then it follows that

$$\mathbb{E}[f_n(\mathbf{x}_n^{(2)})] - f_n(\mathbf{x}_n^*) \leq \mathbb{E}[\varepsilon_n^{(2)}]$$

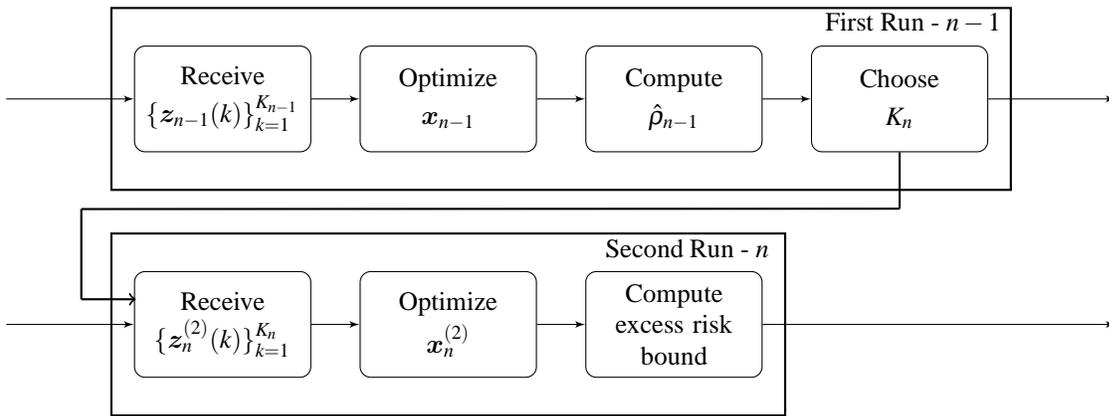


Figure 1: Two Run Process

We now argue that  $\mathbb{E}[\varepsilon_n^{(2)}]$  also bounds the excess risk of the first run.

**Lemma 16.** *For the first run, it holds that*

$$\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*) \leq \mathbb{E}[\varepsilon_n^{(2)}]$$

*Proof.* We proceed by induction. For  $n = 1, 2$ , we know that

$$\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*) \leq \mathbb{E}[\varepsilon_n^{(2)}]$$

by definition. Next, suppose that

$$\mathbb{E}[f_{n-1}(\mathbf{x}_{n-1})] - f_{n-1}(\mathbf{x}_{n-1}^*) \leq \mathbb{E}[\varepsilon_{n-1}^{(2)}]$$

We have

$$\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*) \leq \mathbb{E} \left[ \alpha(K_n) \left( \sqrt{f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)} + \rho \right)^2 + \beta(K_n) \right]$$

so it holds that

$$\begin{aligned} & \mathbb{E}[\varepsilon_n^{(2)}] - (\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*)) \\ & \geq \mathbb{E} \left[ \alpha(K_n) \left( \sqrt{\varepsilon_{n-1}^{(2)}} + \rho \right)^2 - \alpha(K_n) \left( \sqrt{f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)} + \rho \right)^2 \right] \\ & = \mathbb{E} \left[ \alpha(K_n) \left( \varepsilon_{n-1}^{(2)} - (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)) \right) \right] \\ & \quad + \mathbb{E} \left[ 2\rho\alpha(K_n) \left( \sqrt{\varepsilon_{n-1}^{(2)}} - \sqrt{f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)} \right) \right] \end{aligned}$$

By the Monotone Convergence Theorem, it holds that

$$\begin{aligned} & \mathbb{E} \left[ \alpha(K_n) \left( \varepsilon_{n-1}^{(2)} - (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)) \right) \right] \\ & = \lim_{q \rightarrow \infty} \mathbb{E} \left[ \max\{\alpha(K_n), 1/q\} \left( \varepsilon_{n-1}^{(2)} - (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)) \right) \right] \\ & \geq \liminf_{q \rightarrow \infty} \frac{1}{q} \mathbb{E} \left[ \varepsilon_{n-1}^{(2)} - (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)) \right] \\ & \geq 0 \end{aligned}$$

where the last line follows, since by hypothesis

$$\mathbb{E}[f_{n-1}(\mathbf{x}_{n-1})] - f_{n-1}(\mathbf{x}_{n-1}^*) \leq \mathbb{E}[\varepsilon_{n-1}^{(2)}]$$

Similarly, it holds that

$$\begin{aligned}
& \mathbb{E} \left[ 2\rho \alpha(K_n) \left( \sqrt{\varepsilon_{n-1}^{(2)}} - \sqrt{f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)} \right) \right] \\
&= \mathbb{E} \left[ 2\rho \alpha(K_n) \frac{\varepsilon_{n-1}^{(2)} - (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*))}{\sqrt{\varepsilon_{n-1}^{(2)} + \sqrt{f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)}}} \right] \\
&= \lim_{q \rightarrow \infty} \mathbb{E} \left[ 2\rho \max\{\alpha(K_n), 1/q\} \frac{\varepsilon_{n-1}^{(2)} - (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*))}{\sqrt{\varepsilon_{n-1}^{(2)} + \sqrt{f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)}}} \right] \\
&\geq \limsup_{q \rightarrow \infty} \frac{2\rho}{q} \mathbb{E} \left[ \frac{\varepsilon_{n-1}^{(2)} - (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*))}{\sqrt{\varepsilon_{n-1}^{(2)} + \sqrt{f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)}}} \right] \\
&\geq \limsup_{q \rightarrow \infty} \frac{2\rho}{q} \lim_{\tau \rightarrow \infty} \mathbb{E} \left[ \frac{\varepsilon_{n-1}^{(2)} - (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*))}{\sqrt{\varepsilon_{n-1}^{(2)} + \sqrt{f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)}}} \mathbb{1}_{\{\sqrt{\varepsilon_{n-1}^{(2)}} + \sqrt{f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)} \leq \tau\}} \right] \\
&\geq \limsup_{q \rightarrow \infty} \frac{2\rho}{q} \limsup_{\tau \rightarrow \infty} \frac{1}{\tau} \mathbb{E} \left[ \varepsilon_{n-1}^{(2)} - (f_{n-1}(\mathbf{x}_{n-1}) - f_{n-1}(\mathbf{x}_{n-1}^*)) \right] \\
&\geq 0
\end{aligned}$$

Therefore, we conclude that

$$\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*) \leq \mathbb{E}[\varepsilon_n^{(2)}]$$

□

**Theorem 2.** Under assumptions C.1 - C.2 and with  $K_n \geq K^*$  for all  $n$  large enough almost surely with  $K^*$  from (20), we have

$$\limsup_{n \rightarrow \infty} (\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*)) \leq \varepsilon$$

*Proof.* Let  $\bar{v}$  be the fixed point associated with  $\phi_{K^*}(v)$  from Lemma 14. We know that

$$\bar{v} = \phi_{K^*}(\bar{v}) \leq \varepsilon$$

and

$$\phi_{K^*}^{(n)}(v) \rightarrow \bar{v} \leq \varepsilon$$

with  $\bar{v} \leq \varepsilon$ . Since we have  $K_n \geq K^*$  for all  $n$  large enough almost surely, there exists a random variable  $\tilde{N}$  such that

$$n \geq \tilde{N} \Rightarrow K_n \geq K^*$$

Then we have almost surely

$$\begin{aligned}
\limsup_{n \rightarrow \infty} \varepsilon_n^{(2)} &\leq \limsup_{n \rightarrow \infty} (\phi_{K_n} \circ \dots \circ \phi_{K_{\tilde{N}}})(\varepsilon_{\tilde{N}-1}) \\
&\leq \limsup_{n \rightarrow \infty} \phi_{K^*}^{(n-\tilde{N}+1)}(\varepsilon_{\tilde{N}-1}) \\
&= \bar{v} \\
&\leq \varepsilon
\end{aligned}$$

Finally, applying Lemma 19 and Fatou's lemma yields

$$\begin{aligned} \limsup_{n \rightarrow \infty} (\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*)) &\leq \limsup_{n \rightarrow \infty} \mathbb{E} \left[ \varepsilon_n^{(2)} \right] \\ &\leq \mathbb{E} \left[ \limsup_{n \rightarrow \infty} \varepsilon_n^{(2)} \right] \\ &\leq \varepsilon \end{aligned}$$

□

## 4.2 Update Past Excess Risk Bounds

We first consider updating all past excess risk bounds as we go. At time  $n$ , we plug-in  $\hat{\rho}_{n-1} + t_{n-1}$  in place of  $\rho$  and follow the analysis of Section 2. Define for  $i = 1, \dots, n$

$$\hat{\varepsilon}_i^{(n)} = b \left( \left( \sqrt{\frac{2}{m} \hat{\varepsilon}_{i-1}^{(n)}} + (\hat{\rho}_{n-1} + t_{n-1}) \right)^2, K_i \right)$$

If it holds that  $\hat{\rho}_{n-1} + t_{n-1} \geq \rho$ , then  $\mathbb{E}[f_n(\mathbf{x}_n)] - f_n(\mathbf{x}_n^*) \leq \hat{\varepsilon}_n^{(i)}$  for  $i = 1, \dots, n$ . Assumption C.1 guarantees that this holds for all  $n$  large enough almost surely. We can thus set  $K_n$  equal to the smallest  $K$  such that

$$b \left( \left( \sqrt{\frac{2}{m} \max\{\hat{\varepsilon}_{n-1}^{(n-1)}, \varepsilon\}} + (\hat{\rho}_{n-1} + t_{n-1}) \right)^2, K \right) \leq \varepsilon$$

for all  $n \geq 3$  to achieve excess risk  $\varepsilon$ . The maximum in this definition ensures that when  $\hat{\rho}_{n-1} + t_{n-1} \geq \rho$ ,  $K_n \geq K^*$  with  $K^*$  from (5). We can therefore apply Theorem 2.

## 4.3 Do Not Update Past Excess Risk Bounds

Updating all past estimates of the excess risk bounds from time 1 up to  $n$  imposes a computational and memory burden. Suppose that for all  $n \geq 3$  we set

$$K_n = \min \left\{ K \geq 1 \mid b \left( \left( \sqrt{\frac{2\varepsilon}{m}} + (\hat{\rho}_{n-1} + t_{n-1}) \right)^2, K \right) \leq \varepsilon \right\} \quad (20)$$

This is the same form as the choice in (5) with  $\hat{\rho}_{n-1} + t_{n-1}$  in place of  $\rho$ . Due to assumption C.1, for all  $n$  large enough it holds that  $\hat{\rho}_n + t_n \geq \rho$  almost surely. Then by the monotonicity assumption in A.1, for all  $n$  large enough we pick  $K_n \geq K^*$  almost surely. We can therefore apply Theorem 2.

## 5 Experiments

We focus on two regression applications for synthetic and real data as well as two classification applications for synthetic and real data. For the synthetic regression problem, we can explicitly compute  $\rho$  and  $\mathbf{x}_n^*$  and exactly evaluate the performance of our method. It is straightforward to check that all requirements in A.1 -A.4 are satisfied for the problems considered in this section. We apply the do not update past excess risk choice of  $K_n$  here.

## 5.1 Synthetic Regression

Consider a regression problem with synthetic data using the penalized quadratic loss

$$\ell(\mathbf{x}, \mathbf{z}) = \frac{1}{2} (y - \mathbf{w}^\top \mathbf{x})^2 + \frac{1}{2} \lambda \|\mathbf{x}\|^2$$

with  $\mathbf{z} = (\mathbf{w}, y) \in \mathbb{R}^{d+1}$ . The distribution of  $\mathbf{z}_n$  is zero mean Gaussian with covariance matrix

$$\begin{bmatrix} \sigma_w^2 \mathbf{I} & r_{w_n, y_n} \\ r_{w_n, y_n}^\top & \sigma_{y_n}^2 \end{bmatrix}$$

Under these assumptions, we can analytically compute minimizers  $\mathbf{x}_n^*$  of  $f_n(\mathbf{x}) = \mathbb{E}_{\mathbf{z}_n \sim p_n} [\ell(\mathbf{x}, \mathbf{z}_n)]$ . We change only  $r_{w_n, y_n}$  and  $\sigma_{y_n}^2$  appropriately to ensure that  $\|\mathbf{x}_n^* - \mathbf{x}_{n-1}^*\| = \rho$  holds for all  $n$ . We find approximate minimizers using SGD with  $\lambda = 0.1$ . We estimate  $\rho$  using the direct estimate.

We let  $n$  range from 1 to 20 with  $\rho = 1$ , a target excess risk  $\varepsilon = 0.1$ , and  $K_n$  from (20). We average over twenty runs of our algorithm. Figure 2 shows  $\hat{\rho}_n$ , our estimate of  $\rho$ , which is above  $\rho$  in general. Figure 3 shows the number of samples  $K_n$ , which settles down. We can exactly compute  $f_n(\mathbf{x}_n) - f_n(\mathbf{x}_n^*)$ , and so by averaging over the twenty runs of our algorithm, we can estimate the excess risk (denoted ‘‘sample average estimate’’). Figure 4 shows this estimate of the excess risk, the target excess risk, and our bound on the excess risk from Section 4.3. We achieve at least our targeted excess risk

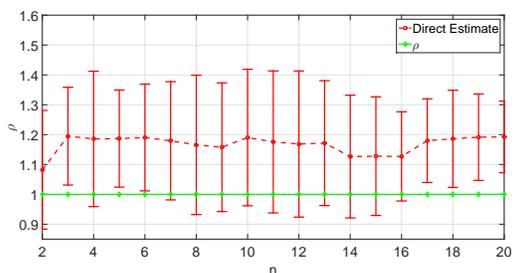


Figure 2:  $\rho$  Estimate

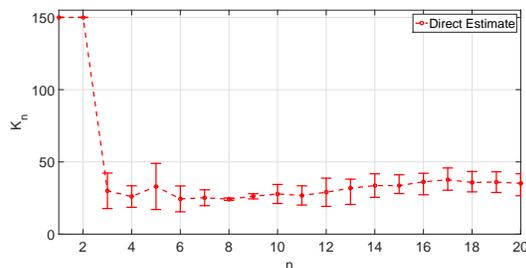


Figure 3:  $K_n$

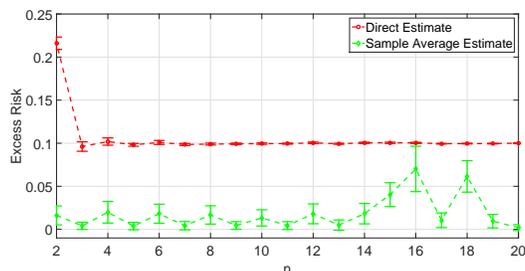


Figure 4: Excess Risk

## 5.2 Panel Study on Income Dynamics Income - Regression

The Panel Study of Income Dynamics (PSID) surveyed individuals every year to gather demographic and income data annually from 1981-1997 [24]. We want to predict an individual’s annual income ( $y$ ) from several demographic features ( $\mathbf{w}$ ) including age, education, work experience, etc. chosen based on previous economic studies in [25]. The

idea of this problem conceptually is to rerun the survey process and determine how many samples we would need if we wanted to solve this regression problem to within a desired excess risk criterion  $\varepsilon$ .

We use the same loss function, direct estimate for  $\rho$ , and minimization algorithm as the synthetic regression problem. The income is adjusted for inflation to 1997 dollars with mean \$20,294. We average over twenty runs of our algorithm by resampling without replacement [26]. We compare to taking an equivalent number of samples up front. Figure 5 shows the test losses over time evaluated over twenty percent of the available samples. The test loss for our approach is substantially less than taking the same number of samples up front. The square root of the average test loss over this time period for our approach and all samples up front are \$1153  $\pm$  352 and \$2805  $\pm$  424 respectively in 1997 dollars.

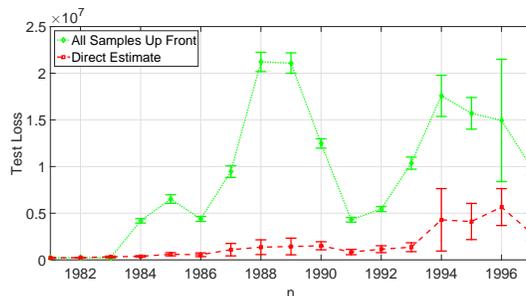


Figure 5: Test Loss

### 5.3 Synthetic Classification

Consider a binary classification problem using  $\ell(\mathbf{x}, \mathbf{z}) = \frac{1}{2}(1 - y(\mathbf{w}^\top \mathbf{x}))_+^2 + \frac{1}{2}\lambda \|\mathbf{x}\|^2$  with  $\mathbf{z} = (\mathbf{w}, y) \in \mathbb{R}^d \times \mathbb{R}$  and  $(y)_+ = \max\{y, 0\}$ . This is a smoothed version of the hinge loss used in support vector machines (SVM) [26]. We suppose that at time  $n$ , the two classes have features drawn from a Gaussian distribution with covariance matrix  $\sigma^2 \mathbf{I}$  but different means  $\mu_n^{(1)}$  and  $\mu_n^{(2)}$ , i.e.,  $\mathbf{w}_n | \{y_n = i\} \sim \mathcal{N}(\mu_n^{(i)}, \sigma^2 \mathbf{I})$ . The class means move slowly over uniformly spaced points on a unit sphere in  $\mathbb{R}^d$  as in Figure 6 to ensure that (2) holds. We find approximate minimizers using SGD with  $\lambda = 0.1$ . We estimate  $\rho$  using the direct estimate with  $t_n \propto 1/n^{3/8}$ .

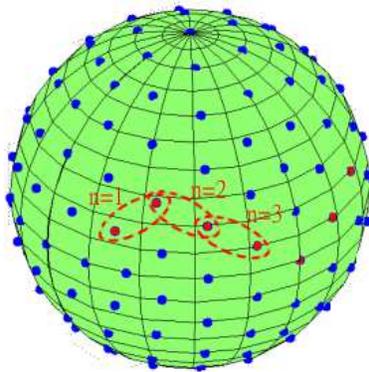


Figure 6: Evolution of Class Means

We let  $n$  range from 1 to 20 and target a excess risk  $\varepsilon = 0.1$ . We average over twenty runs of our algorithm. As a comparison, if our algorithm takes  $\{K_n\}_{n=1}^{20}$  samples, then we consider taking  $\sum_{n=1}^{20} K_n$  samples up front at  $n = 1$ . This is what we would do if we assumed that our problem is not time varying. Figure 7 shows  $\hat{\rho}_n$ , our estimate of  $\rho$ . Figure 8 shows the average test loss for both sampling strategies. To compute test loss we draw  $T_n$  additional samples

$\{z_n^{\text{test}}(k)\}_{k=1}^{T_n}$  from  $p_n$  and compute  $\frac{1}{T_n} \sum_{k=1}^{T_n} \ell(\mathbf{x}_n, z_n^{\text{test}}(k))$ . We see that our approach achieves substantially smaller test loss than taking all samples up front.

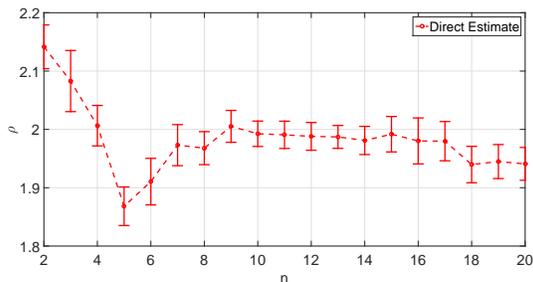


Figure 7:  $\rho$  Estimate

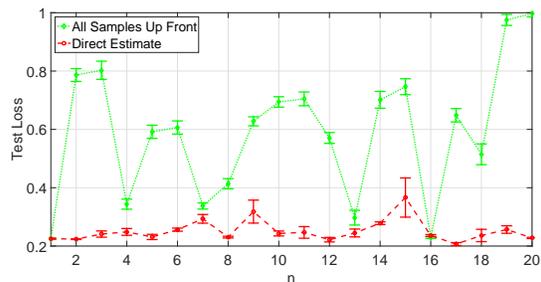


Figure 8: Test Loss

## 5.4 General Social Survey - Classification

The General Social Survey (GSS) surveyed individuals every year to gather socio-economic data annually from 1981-2013 [27]. We want to predict an individual’s marital status ( $y$ ) from several demographic features ( $\mathbf{w}$ ) including age, education, etc. We model this as a binary classification problem using loss

$$\ell(\mathbf{x}, \mathbf{z}) = \frac{1}{2}(1 - y(\mathbf{w}^\top \mathbf{x}))_+^2 + \frac{1}{2}\lambda \|\mathbf{x}\|^2$$

with  $\mathbf{z} = (\mathbf{w}, y) \in \mathbb{R}^d \times \mathbb{R}$  and  $(y)_+ = \max\{y, 0\}$ . This is a smoothed version of the hinge loss used in support vector machines [26]. We find approximate minimizers using SGD with  $\lambda = 0.1$ . Figure 9 shows the test loss. We see that our approach achieves smaller test loss than taking all samples up front. We also plot receiver operating characteristics (ROC) [26] to characterize the performance of our classifiers. In particular we plot the ROC for 1974 in Figure 10 and the ROC for 2012 in Figure 11. By examining the ROC, we see that taking all samples up front is much better in 1974 but much worse in 2012.

## 6 Conclusion

We introduced a framework for adaptively solving a sequence of optimization problems with applications to machine learning. We developed estimates of the change in the minimizers used to determine the number of samples  $K_n$  needed to achieve a target excess risk  $\epsilon$ . Experiments with synthetic and real data demonstrate that this approach is effective.

## References

- [1] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of Machine Learning*, The MIT Press, 2012.
- [2] A. Agarwal, H. Daumé, and S. Gerber, “Learning multiple tasks using manifold regularization.,” in *NIPS*, 2011, pp. 46–54.
- [3] T. Evgeniou and M. Pontil, “Regularized multi-task learning,” in *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, New York, NY, USA, 2004, KDD ’04, pp. 109–117, ACM.
- [4] Y. Zhang and D. Yeung, “A convex formulation for learning task relationships in multi-task learning,” *CoRR*, vol. abs/1203.3536, 2012.
- [5] S. Pan and Q. Yang, “A survey on transfer learning,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, Oct 2010.

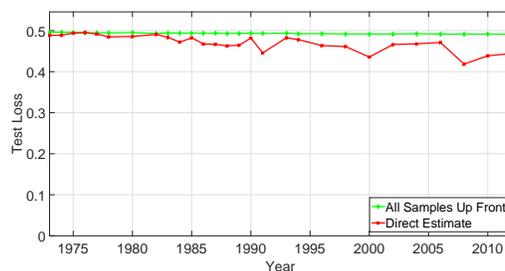


Figure 9: Test Loss

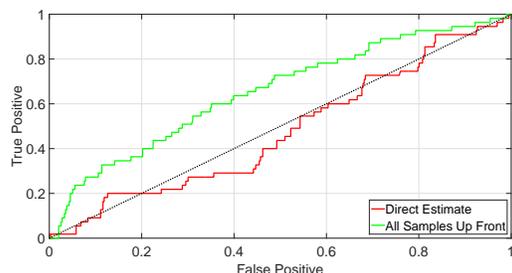


Figure 10: ROC for 1974

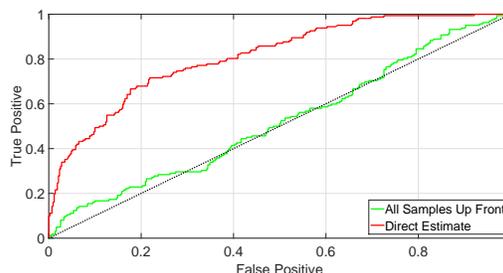


Figure 11: ROC for 2012

- [6] A. Agarwal, A. Rakhlin, and P. Bartlett, “Matrix regularization techniques for online multitask learning,” Tech. Rep. UCB/EECS-2008-138, EECS Department, University of California, Berkeley, Oct 2008.
- [7] Z. Towfic, J. Chu, and A. Sayed, “Online distributed online classification in the midst of concept drifts,” *Neurocomputing*, vol. 112, pp. 138–152, 2013.
- [8] C. Tekin, L. Canzian, and M. van der Schaar, “Context adaptive big data stream mining,” in *Allerton Conference*, 2014, pp. 46–54.
- [9] T. Dietterich, “Machine learning for sequential data: A review,” in *Structural, Syntactic, and Statistical Pattern Recognition*, 2002, pp. 15–30.
- [10] T. Fawcett and F. Provost, “Adaptive fraud detection,” *Data Min. Knowl. Discov.*, vol. 1, no. 3, pp. 291–316, 1997.
- [11] N. Qian and T. Sejnowski, “Predicting the secondary structure of globular proteins using neural network models,” *Journal of Molecular Biology*, vol. 202, pp. 865–884, Aug 1988.
- [12] Y. Bengio and P. Frasconi, “Input-output HMM’s for sequence processing,” *IEEE Transactions on Neural Networks*, vol. 7(5), pp. 1231–1249, 1996.
- [13] A. Dontchev and R. Rockafellar, *Implicit Functions and Solution Mappings: A View from Variational Analysis*, Springer, New York, New York, 2009.
- [14] B. Sriperumbudur, “On the empirical estimation of integral probability metrics,” *Electronic Journal of Statistics*, pp. 1550–1599, 2012.
- [15] R. Vershyn, “Introduction to non-asymptotic analysis of random matrices,” Tech. Rep., University of Michigan, 2012.
- [16] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, “Stochastic approximation approach to stochastic programming,” *SIAM Journal on Optimization*, vol. 19, pp. 1574–1609, 2009.

- [17] V.V Buldygin and E.D. Pechuk, “Inequalities for the distributions of functionals of sub-gaussian vectors,” *Theor. Probability and Math. Statist.*, pp. 25–36, 2010.
- [18] S. Janson, “Large deviations for sums of partly dependent random variables,” *Random Structures Algorithms*, vol. 24, pp. 234–248, 2004.
- [19] S. Boucheron, G. Lugosi, and P. Massart, *Concentration Inequalities: A Nonasymptotic Theory of Independence*, Oxford University Press, 2013.
- [20] L. Trefethen, *Numerical Linear Algebra*, SIAM, 1997.
- [21] J. Kennan, “Uniqueness of positive fixed points for increasing concave functions on  $m$ : An elementary result,” *Review of Economic Dynamics*, vol. 4, pp. 893–899, 2001.
- [22] Stephen Boyd and Lieven Vandenberghe, *Convex Optimization*, Cambridge University Press, New York, NY, USA, 2004.
- [23] A. Granas and J. Dugundji, *Fixed Point Theory*, Springer-Verlag, 2003.
- [24] “Panel study of income dynamics: public use dataset,” *Survey Research Center*, 2015.
- [25] S. Jenkins and P. Van Kerm, “Trends in income inequality, pro-poor income growth, and income mobility,” *Oxford Economic Papers*, vol. 58, no. 3, pp. 531–548, 2006.
- [26] T. Hastie, R. Tibshirani, and J.H. Friedman, *The elements of statistical learning: data mining, inference, and prediction: with 200 full-color illustrations*, New York: Springer-Verlag, 2001.
- [27] “General social survey,” *National Opinion Research Center*, 2015.
- [28] F. Bach and E. Moulines, “Non-Asymptotic Analysis of Stochastic Approximation Algorithms for Machine Learning,” in *Advances in Neural Information Processing Systems (NIPS)*, Spain, 2011.
- [29] D. Bertsekas, *Nonlinear Programming*, Athena Scientific, 1999.
- [30] Léon Bottou, “Online learning and stochastic approximations,” 1998.
- [31] A. Nedic and S. Lee, “Analysis of mirror descent for strongly convex functions,” *ArXiv*, 2013.
- [32] Yu. Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*, Kluwer Academic Publishers, Norwell, Massachusetts, USA, 2004.
- [33] R. Antonini and Y. Kozachenko, “A note on the asymptotic behavior of sequences of generalized subgaussian random vectors,” *Random Op. and Stoch. Equ.*, vol. 13, pp. 39–52, 2005.

## A Examples of $b(d_0, K)$ :

For this section, we drop the  $n$  index for convenience. The bounds of this form depend on the strong convexity parameter  $m$  and an assumption on how the gradients grow. In general, we assume that

$$\mathbb{E}_{\mathbf{z} \sim p} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{z})\|^2 \leq A + B \|\mathbf{x} - \mathbf{x}^*\|^2$$

The base algorithm we look at is SGD. First, we generate iterates  $\mathbf{x}(0), \dots, \mathbf{x}(K)$  through SGD as follows:

$$\mathbf{x}(\ell + 1) = \Pi_{\mathcal{X}} [\mathbf{x}(\ell) - \mu(\ell + 1) \nabla_{\mathbf{x}} \ell(\mathbf{x}(\ell), \mathbf{z}(\ell))] \quad \ell = 0, \dots, K - 1$$

with  $\mathbf{x}(0)$  fixed. We then combine the iterates to yield a final approximate minimizer

$$\bar{\mathbf{x}}(K) = \phi(\mathbf{x}(0), \dots, \mathbf{x}(K))$$

For our choice of  $\phi$ , we look at two cases:

1. No iterate averaging, i.e.,

$$\phi(\mathbf{x}(0), \dots, \mathbf{x}(K)) = \mathbf{x}(K)$$

2. Iterate averaging, i.e, for a convex combination  $\{\lambda(\ell)\}_{\ell=0}^K$

$$\phi(\mathbf{x}(0), \dots, \mathbf{x}(K)) = \sum_{\ell=0}^K \lambda(\ell) \mathbf{x}(\ell)$$

Define

$$d(\ell) \triangleq \|\mathbf{x}(\ell) - \mathbf{x}^*\|^2 \quad (21)$$

First we bound  $\mathbb{E}[d(\ell)]$  in Lemma 17.

**Lemma 17.** *Suppose that the function  $f(\mathbf{x})$  has Lipschitz continuous gradients. Then it holds that*

$$\mathbb{E}[d(\ell)] \leq \prod_{k=1}^{\ell} (1 - 2m\mu(k) + B\mu^2(k)) + \sum_{k=1}^{\ell} \prod_{i=k+1}^{\ell} (1 - 2m\mu(i) + B\mu^2(i)) \mu^2(k)$$

*Proof.* Following the standard SGD analysis (see [16]), it holds that

$$\begin{aligned} d(\ell) &\leq \|\mathbf{x}(\ell-1) - \mathbf{x}^* - \mu(\ell) \nabla_{\mathbf{x}} \ell(\mathbf{x}(\ell-1), \mathbf{z}(\ell))\|^2 \\ &\leq d(\ell-1) - 2\mu(\ell) \langle \mathbf{x}(\ell-1) - \mathbf{x}^*, \nabla_{\mathbf{x}} \ell(\mathbf{x}(\ell-1), \mathbf{z}(\ell)) \rangle + \mu^2(\ell) \|\nabla_{\mathbf{x}} \ell(\mathbf{x}(\ell-1), \mathbf{z}(\ell))\|^2 \end{aligned}$$

Then it follows that

$$\begin{aligned} &\mathbb{E}[d(\ell) \mid \mathbf{x}(\ell-1)] \\ &\leq d(\ell-1) - 2\mu(\ell) \langle \mathbf{x}(\ell-1) - \mathbf{x}^*, \nabla f(\mathbf{x}(\ell-1)) \rangle + \mu^2(\ell) \mathbb{E}[\|\nabla_{\mathbf{x}} \ell(\mathbf{x}(\ell-1), \mathbf{z}(\ell))\|^2 \mid \mathbf{x}(\ell-1)] \\ &\leq (1 - 2m\mu(\ell) + B\mu^2(\ell)) d(\ell-1) + \mu^2(\ell-1) A \end{aligned}$$

and

$$\mathbb{E}[d(\ell)] \leq (1 - 2m\mu(\ell) + B\mu^2(\ell)) \mathbb{E}[d(\ell-1)] + \mu^2(\ell-1) A$$

Since  $B > m$ , we have

$$2m\mu - B\mu^2 \leq 2\sqrt{\frac{B}{2}}\mu \left(1 - \sqrt{\frac{B}{2}}\mu\right) \leq 2\frac{1}{4} = \frac{1}{2}$$

and so

$$1 - 2m\mu(\ell) + B\mu^2(\ell) \geq 1 - \frac{1}{2} = \frac{1}{2}$$

Since this quantity is non-negative, we can unwind this recursion to yield

$$\mathbb{E}[d(\ell)] \leq \prod_{k=1}^{\ell} (1 - 2m\mu(k) + B\mu^2(k)) + \sum_{k=1}^{\ell} \prod_{i=k+1}^{\ell} (1 - 2m\mu(i) + B\mu^2(i)) \mu^2(k)$$

□

The bound in Lemma 17 can be further bounded into a closed form as follows from [28]: Define

$$\varphi_{\beta}(t) = \begin{cases} \frac{t^{\beta}-1}{\beta}, & \text{if } \beta \neq 0 \\ \log(t), & \text{if } \beta = 0 \end{cases}$$

Then with  $\mu(\ell) = C\ell^{-\alpha}$ , it holds that

$$\mathbb{E}[d(\ell)] \leq \begin{cases} 2 \exp\{2BC^2 \varphi_{1-2\alpha}(\ell)\} \exp\{-\frac{mC}{4} \ell^{1-\alpha}\} \left(\mathbb{E}[d(0)] + \frac{A}{B}\right) + \frac{2AC}{m\ell^{\alpha}}, & \text{if } 0 \leq \alpha < 1 \\ \frac{\exp\{BC^2\}}{\ell^{mC}} \left(\mathbb{E}[d(0)] + \frac{A}{B}\right) + AC^2 \frac{\varphi_{mC/2-1}(\ell)}{\ell^{mC/2}}, & \text{if } \alpha = 1 \end{cases}$$

Note that this bound is a closed form but is substantially looser than Lemma 17. In the case that the functions in question have Lipschitz continuous gradients, we introduce a bound on the excess risk using Lemma 17. This case corresponds to choosing

$$\phi(\mathbf{x}(0), \dots, \mathbf{x}(K)) = \mathbf{x}(K)$$

**Lemma 18.** *With arbitrary step sizes and assuming that  $f(\mathbf{x})$  has Lipschitz continuous gradients with modulus  $M$ , it holds that*

$$\mathbb{E}[f(\mathbf{x})] - f(\mathbf{x}^*) \leq \frac{1}{2}M\mathbb{E}[d(K)]$$

and therefore, we set

$$b(d_0, K) = \frac{1}{2}M \left( \prod_{\ell=1}^K (1 - 2m\mu(\ell) + B\mu^2(\ell)) + \sum_{\ell=1}^K \prod_{i=\ell+1}^K (1 - 2m\mu(i) + B\mu^2(i))\mu^2(\ell) \right)$$

*Proof.* Using the descent lemma from [29], it holds that

$$\mathbb{E}[f(\mathbf{x})] - f(\mathbf{x}^*) \leq \frac{1}{2}M\mathbb{E}[d(K)]$$

Plugging in the bound from Lemma 17 yields the bound  $b(d_0, K)$ . □

Next, we introduce a bound inspired by [30] for the case where  $\phi(\mathbf{x}(0), \dots, \mathbf{x}(K))$  corresponds to forming a convex combination of the iterates.

**Lemma 19.** *With a constant step size and averaging with*

$$\lambda(\ell) = \begin{cases} \frac{\gamma(\ell)}{\sum_{\tau=1}^K \gamma(\tau)}, & \text{if } \ell > 0 \\ 0, & \text{if } \ell = 0 \end{cases}$$

where

$$\gamma(\ell) = (1 - m\mu + B\mu^2)^{-\ell}$$

it holds that

$$b(d_0, K) = \frac{d_0}{2\mu \sum_{\ell=0}^K \gamma(\ell)} + \frac{1}{2}A\mu$$

*Proof.* By strong convexity, it holds that

$$-\langle \mathbf{x}(\ell-1) - \mathbf{x}^*, \nabla f(\mathbf{x}(\ell-1)) \rangle \leq -m\|\mathbf{x}(\ell-1) - \mathbf{x}^*\|^2 - (f(\mathbf{x}(\ell-1)) - f(\mathbf{x}^*))$$

Following the Lyapunov-style analysis of Lemma 17, it holds that

$$\mathbb{E}[d(\ell)] \leq (1 - m\mu + B\mu^2)\mathbb{E}[d(\ell-1)] - 2\mu(\mathbb{E}[f(\mathbf{x}(\ell-1))] - f(\mathbf{x}^*)) + A\mu^2$$

Rearranging, using the telescoping sum, and using convexity, it holds that

$$\mathbb{E}[f(\mathbf{x})] - f(\mathbf{x}^*) \leq \frac{d_0}{2\mu \sum_{\tau=0}^K \gamma(\tau)} + \frac{1}{2}A\mu$$

□

If we set  $\mu = \frac{1}{\sqrt{K}}$ , then it holds that

$$b(d_0, K) = \mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$$

for Lemma 19.

We consider an extension of the averaging scheme in [31]. The bound in this paper only works with  $B = 0$ , so we extend it slightly to handle  $B > 0$ .

**Lemma 20.** Consider the choice of step sizes given by

$$\mu(\ell) = \frac{1}{m\ell} \quad \forall \ell \geq 1$$

Then

$$b(d_0, K) = \frac{\frac{1}{2}d(0) + \frac{1}{2}(K+1)A + \frac{1}{2}B \sum_{\ell=0}^K \gamma(\ell)}{1 + \frac{1}{2}m(K+1)(K+2)}$$

where

$$\mathbb{E}[d(\ell)] \leq \gamma(\ell)$$

Note that we can use the bound in Lemma 17 here.

*Proof.* We have using Lyapunov style analysis

$$\mathbb{E}[d(\ell)] \leq (1 - 2m\mu(\ell) + B\mu^2(\ell))\mathbb{E}[d(\ell-1)] - 2\mu(\ell)(\mathbb{E}[f(\mathbf{x}(\ell))] - f(\mathbf{x}^*)) + A\mu^2(\ell)$$

Then we have

$$\frac{1}{\mu^2(\ell)}\mathbb{E}[d(\ell)] \leq \left( \frac{1 - 2m\mu(\ell)}{\mu^2(\ell)} + B \right) \mathbb{E}[d(\ell-1)] - \frac{2}{\mu(\ell)}(\mathbb{E}[f(\mathbf{x}(\ell))] - f(\mathbf{x}^*)) + A$$

It holds that

$$\begin{aligned} \frac{1 - 2m\mu(\ell)}{\mu^2(\ell)} - \frac{1}{\mu^2(\ell-1)} &= \frac{1}{\mu^2(\ell)} - 2m\frac{1}{\mu(\ell)} - \frac{1}{\mu^2(\ell-1)} \\ &= \frac{\ell^2}{C^2} - \frac{2m\ell}{C} - \frac{(\ell-1)^2}{C^2} \\ &= \frac{2(mC-1)L-1}{C^2} \end{aligned}$$

As long as we have

$$mC - 1 \leq 1 \Leftrightarrow C \leq \frac{2}{m}$$

then we get

$$\frac{1}{\mu^2(\ell)}\mathbb{E}[d(\ell)] - \frac{1}{\mu^2(\ell-1)}\mathbb{E}[d(\ell-1)] \leq B\mathbb{E}[d(\ell-1)] - \frac{2}{\mu(\ell)}(\mathbb{E}[f(\mathbf{x}(\ell))] - f(\mathbf{x}^*)) + A$$

Summing an rearranging yields

$$\sum_{\ell=0}^K \frac{1}{\mu(\ell)} (\mathbb{E}[f(\mathbf{x}(\ell))] - f(\mathbf{x}^*)) \leq \frac{1}{2}d(0) + \frac{1}{2}(K+1)A + \frac{1}{2}B \sum_{\ell=0}^K \mathbb{E}[d(\ell)]$$

with  $\mu(0) = 1$  by convention. With the weights

$$\gamma(\ell) = \frac{\frac{1}{\mu(\ell)}}{\sum_{j=0}^{\ell} \frac{1}{\mu(j)}}$$

we have

$$\mathbb{E}[f(\bar{\mathbf{x}}(K))] - f(\mathbf{x}^*) \leq \frac{\frac{1}{2}d(0) + \frac{1}{2}(K+1)A + \frac{1}{2}B \sum_{\ell=0}^K \mathbb{E}[d(\ell)]}{\sum_{\tau=0}^K \frac{1}{\mu(\tau)}}$$

Then it holds that

$$\sum_{\tau=0}^K = 1 + \sum_{\tau=1}^K m\tau = 1 + \frac{1}{2}m(K+1)(K+2)$$

so

$$\mathbb{E}[f(\bar{\mathbf{x}}(K))] - f(\mathbf{x}^*) \leq \frac{\frac{1}{2}d(0) + \frac{1}{2}(K+1)A + \frac{1}{2}B\sum_{\ell=0}^K \mathbb{E}[d(\ell)]}{1 + \frac{1}{2}m(K+1)(K+2)}$$

□

For the choice of step sizes in Lemma 20 from Lemma 17, it holds that

$$\mathbb{E}[d(\ell)] = \mathcal{O}\left(\frac{1}{\ell}\right)$$

Since

$$\sum_{\ell=1}^K \frac{1}{\ell} = \mathcal{O}(\log K)$$

it holds that

$$\mathbb{E}[f(\bar{\mathbf{x}}(K))] - f(\mathbf{x}^*) = \mathcal{O}\left(\frac{d(0)}{K^2} + \frac{\log(K)}{K^2} + \frac{1}{K}\right)$$

Note that a rate of  $\mathcal{O}(\frac{1}{K})$  is minimax optimal for stochastic minimization of a strongly convex function [32].

Next, we look at a special case of averaging for functions such that

$$\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}}\ell(\tilde{\mathbf{x}}, \mathbf{z}) - \nabla_{\mathbf{xx}}^2\ell(\tilde{\mathbf{x}}, \mathbf{z})(\mathbf{x} - \tilde{\mathbf{x}})\|^2 = 0$$

from [28]. For example, quadratics satisfy this condition.

**Lemma 21.** *Assuming that*

$$\mathbb{E}\|\nabla_{\mathbf{x}}\ell(\mathbf{x}, \mathbf{z}) - \nabla_{\mathbf{x}}\ell(\tilde{\mathbf{x}}, \mathbf{z}) - \nabla_{\mathbf{xx}}^2\ell(\tilde{\mathbf{x}}, \mathbf{z})(\mathbf{x} - \tilde{\mathbf{x}})\|^2 = 0,$$

*we select step sizes*

$$\mu(\ell) = C\ell^{-\alpha}$$

*with  $\alpha > 1/2$ , and*

$$\lambda(\ell) = \begin{cases} \frac{1}{K}, & \text{if } \ell > 0 \\ 0, & \text{if } \ell = 0 \end{cases}$$

*it holds that*

$$\begin{aligned} & (\mathbb{E}[\bar{d}(K)])^{1/2} \\ & \leq \frac{1}{m^{1/2}} \sum_{k=1}^{K-1} \left| \frac{1}{\mu(k+1)} - \frac{1}{\mu(k)} \right| (\mathbb{E}[d(k)])^{1/2} + \frac{1}{m^{1/2}\mu(1)} (\mathbb{E}[d(0)])^{1/2} + \frac{1}{m^{1/2}\mu(K)} (\mathbb{E}[d(K)])^{1/2} \\ & \quad + \sqrt{\frac{A}{mK}} + \sqrt{\frac{2B}{mK^2} \sum_{k=1}^K \mathbb{E}[d(k-1)]} \end{aligned}$$

*with  $\bar{d}(K) = \|\bar{\mathbf{x}}(K) - \mathbf{x}^*\|^2$ . If in addition  $f$  has Lipschitz continuous gradients with modulus  $M$ , then it holds that*

$$\mathbb{E}[f(\bar{\mathbf{x}}(K))] - f(\mathbf{x}^*) \leq \frac{1}{2}M\mathbb{E}[\bar{d}(K)]$$

*Proof.* Suppose that we set

$$\bar{\mathbf{x}}(K) = \frac{1}{n} \sum_{k=1}^K \mathbf{x}(k)$$

Then it holds that

$$\begin{aligned}\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*)(\mathbf{x}(k) - \mathbf{x}^*) &= \nabla_{\mathbf{x}} \ell(\mathbf{x}(k-1), \mathbf{z}(k-1)) - \nabla_{\mathbf{x}} \ell(\mathbf{x}^*, \mathbf{z}(k-1)) \\ &\quad + [\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*) - \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}^*, \mathbf{z}(k-1))] (\mathbf{x}(k-1) - \mathbf{x}^*)\end{aligned}$$

yielding

$$\begin{aligned}\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*)(\bar{\mathbf{x}}(k) - \mathbf{x}^*) &= \frac{1}{K} \sum_{k=1}^K \nabla_{\mathbf{x}} \ell(\mathbf{x}(k-1), \mathbf{z}(k-1)) - \frac{1}{K} \sum_{k=1}^K \nabla_{\mathbf{x}} \ell(\mathbf{x}^*, \mathbf{z}(k-1)) \\ &\quad + \frac{1}{K} \sum_{k=1}^K [\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*) - \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}^*, \mathbf{z}(k-1))] (\mathbf{x}(k-1) - \mathbf{x}^*)\end{aligned}$$

First, we have

$$\begin{aligned}\frac{1}{K} \sum_{k=1}^K \nabla_{\mathbf{x}} \ell(\mathbf{x}(k-1), \mathbf{z}(k-1)) &= \frac{1}{K} \sum_{k=1}^K \nabla_{\mathbf{x}} \ell(\mathbf{x}(\ell-1), \mathbf{z}(\ell-1)) \\ &= \frac{1}{K} \sum_{k=1}^K \frac{1}{\mu(k)} (\mathbf{x}(\ell-1) - \mathbf{x}(\ell)) \\ &= \frac{1}{K} \sum_{k=1}^K \frac{1}{\mu(k)} (\mathbf{x}(\ell-1) - \mathbf{x}^*) - \frac{1}{K} \sum_{k=1}^K \frac{1}{\mu(k)} (\mathbf{x}(\ell) - \mathbf{x}^*) \\ &= \frac{1}{K} \sum_{k=1}^{K-1} \left( \frac{1}{\mu(k+1)} - \frac{1}{\mu(k)} \right) (\mathbf{x}(\ell) - \mathbf{x}^*) + \frac{1}{\mu(1)} (\mathbf{x}(0) - \mathbf{x}^*) \\ &\quad - \frac{1}{\mu(K)} (\mathbf{x}(K) - \mathbf{x}^*)\end{aligned}$$

Second, we have

$$\begin{aligned}\mathbb{E} \left\| \frac{1}{K} \sum_{k=1}^K \nabla_{\mathbf{x}} \ell(\mathbf{x}^*, \mathbf{z}(k-1)) \right\|^2 &= \frac{1}{K^2} \sum_{k=1}^K \mathbb{E} \|\nabla_{\mathbf{x}} \ell(\mathbf{x}^*, \mathbf{z}(k-1))\|^2 \\ &\leq \frac{A}{n^2}\end{aligned}$$

Third, we have

$$\mathbb{E} \left\| \frac{1}{K} \sum_{k=1}^K [\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*) - \nabla_{\mathbf{x}\mathbf{x}}^2 \ell(\mathbf{x}^*, \mathbf{z}(k-1))] (\mathbf{x}(k-1) - \mathbf{x}^*) \right\|^2 \leq \frac{2B}{K^2} \sum_{k=1}^K \mathbb{E}[d(k-1)]$$

Combining these bounds with Minkowski's inequality yields

$$\begin{aligned}&(m\mathbb{E}[\bar{d}(K)])^{1/2} \\ &\leq (\mathbb{E} \|\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*)(\bar{\mathbf{x}}(K) - \mathbf{x}^*)\|^2)^{1/2} \\ &\leq \sum_{k=1}^{K-1} \left| \frac{1}{\mu(k+1)} - \frac{1}{\mu(k)} \right| (\mathbb{E}[d(k)])^{1/2} + \frac{1}{\mu(1)} (\mathbb{E}[d(0)])^{1/2} + \frac{1}{\mu(K)} (\mathbb{E}[d(K)])^{1/2} \\ &\quad + \sqrt{\frac{A}{K}} + \sqrt{\frac{2B}{K^2} \sum_{k=1}^K \mathbb{E}[d(k-1)]}\end{aligned}$$

Then we have

$$\begin{aligned}
& (\mathbb{E}[\bar{d}(K)])^{1/2} \\
& \leq \frac{1}{m^{1/2}} \sum_{k=1}^{K-1} \left| \frac{1}{\mu(k+1)} - \frac{1}{\mu(k)} \right| (\mathbb{E}[d(k)])^{1/2} + \frac{1}{m^{1/2}\mu(1)} (\mathbb{E}[d(0)])^{1/2} + \frac{1}{m^{1/2}\mu(K)} (\mathbb{E}[d(K)])^{1/2} \\
& \quad + \sqrt{\frac{A}{mK}} + \sqrt{\frac{2B}{mK^2} \sum_{k=1}^K \mathbb{E}[d(k-1)]}
\end{aligned}$$

□

This decays at rate  $\mathcal{O}\left(\frac{1}{K}\right)$  as long as  $\mu(\ell) = C\ell^{-\alpha}$  with  $\frac{1}{2} \leq \alpha \leq 1$ .

## B Useful Concentration Inequalities

For our analysis of both the direct and IPM estimates, we need the following key technical lemma from [33]. This lemma controls the concentration of sums of random variables that are sub-Gaussian conditioned on a particular filtration  $\{\mathcal{F}_i\}_{i=0}^n$ . Such a collection of random variables is referred to as a *sub-Gaussian martingale sequence*. We include the proof for completeness.

**Lemma 22** (Theorem 7.5 of [33]). *Suppose we have a collection of random variables  $\{V_i\}_{i=1}^n$  and a filtration  $\{\mathcal{F}_i\}_{i=0}^n$  such that for each random variable  $V_i$  it holds that*

1.  $\mathbb{E}[e^{sV_i} \mid \mathcal{F}_{i-1}] \leq e^{\frac{1}{2}\sigma_i^2 s^2}$  with  $\sigma_i^2$  a constant
2.  $V_i$  is  $\mathcal{F}_i$ -measurable

Then for every  $\mathbf{a} \in \mathbb{R}^n$  it holds that

$$\mathbb{P}\left\{\sum_{i=1}^n a_i V_i > t\right\} \leq \exp\left\{-\frac{t^2}{2v}\right\} \quad \forall t > 0$$

and

$$\mathbb{P}\left\{\sum_{i=1}^n a_i V_i < -t\right\} \leq \exp\left\{-\frac{t^2}{2v}\right\} \quad \forall t > 0$$

with

$$v = \sum_{i=1}^n \sigma_i^2 a_i^2$$

*Proof.* We bound the moment generating function of  $\sum_{i=1}^n a_i V_i$  by induction. As a base case, we have

$$\begin{aligned}
\mathbb{E}[e^{sa_1 V_1}] &= \mathbb{E}\left[\mathbb{E}[e^{sa_1 V_1} \mid \mathcal{F}_0]\right] \\
&\leq e^{\frac{1}{2}\sigma_1^2 a_1^2 s^2}
\end{aligned}$$

Assume for induction that we have

$$\mathbb{E}\left[\exp\left\{s \sum_{i=1}^j a_i V_i\right\}\right] \leq \exp\left\{\frac{1}{2}\left(\sum_{i=1}^j \sigma_i^2 a_i^2\right) s^2\right\}$$

Then we have

$$\begin{aligned}
\mathbb{E} \left[ \exp \left\{ \sum_{i=1}^{j+1} a_i V_i \right\} \right] &= \mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^j a_i V_i \right\} e^{s a_{j+1} X_{j+1}} \right] \\
&= \mathbb{E} \left[ \mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^j a_i V_i \right\} e^{s a_{j+1} X_{j+1}} \mid \mathcal{F}_{j+1} \right] \right] \\
&\stackrel{(a)}{=} \mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^j a_i V_i \right\} \mathbb{E} \left[ e^{s a_{j+1} X_{j+1}} \mid \mathcal{F}_{j+1} \right] \right] \\
&\stackrel{(b)}{\leq} \mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^j a_i V_i \right\} \right] e^{\frac{1}{2} \sigma_{j+1}^2 a_{j+1}^2 s^2} \\
&\stackrel{(c)}{\leq} \exp \left\{ \frac{1}{2} \left( \sum_{i=1}^{j+1} \sigma_i^2 a_i^2 \right) s^2 \right\}
\end{aligned}$$

where (a) follows since  $\sum_{i=1}^j a_i V_i$  is  $\mathcal{F}_j$  measurable, (b) follows since

$$\mathbb{E} \left[ e^{s a_{j+1} X_{j+1}} \mid \mathcal{F}_{j+1} \right] \leq e^{\frac{1}{2} \sigma_{j+1}^2 a_{j+1}^2 s^2},$$

and (c) is the inductive assumption. This proves that

$$\mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^n a_i V_i \right\} \right] \leq \exp \left\{ \frac{1}{2} \left( \sum_{i=1}^n \sigma_i^2 a_i^2 \right) s^2 \right\} \leq \exp \left\{ \frac{1}{2} v s^2 \right\}$$

Using the Chernoff bound [19], we have

$$\mathbb{P} \left\{ \sum_{i=1}^n a_i V_i > t \right\} \leq e^{-st} \mathbb{E} \left[ \exp \left\{ s \sum_{i=1}^n a_i V_i \right\} \right] \leq \exp \left\{ -st + \frac{1}{2} v s^2 \right\}$$

Optimizing the bound over  $s$  yields

$$\mathbb{P} \left\{ \sum_{i=1}^n a_i V_i > t \right\} \leq \exp \left\{ -\frac{t^2}{2v} \right\}$$

The proof for the other tail is similar. □

If the random variables instead satisfy

1.  $\mathbb{E} \left[ \exp \left\{ s (V_i - \mathbb{E} [V_i \mid \mathcal{F}_{i-1}]) \right\} \mid \mathcal{F}_{i-1} \right] \leq e^{\frac{1}{2} \sigma_i^2 s^2}$  with  $\sigma_i^2$  a constant
2.  $V_i$  is  $\mathcal{F}_i$ -measurable

then Lemma 22 can be applied to  $\{V_i - \mathbb{E} [V_i \mid \mathcal{F}_{i-1}]\}_{i=1}^n$  to yield

$$\mathbb{P} \left\{ \sum_{i=1}^n a_i V_i > \sum_{i=1}^n a_i \mathbb{E} [V_i \mid \mathcal{F}_{i-1}] + t \right\} \leq \exp \left\{ -\frac{t^2}{2v} \right\}$$

If we can upper bound the conditional expectations

$$\mathbb{E} [V_i \mid \mathcal{F}_{i-1}] \leq C_i,$$

by  $\mathcal{F}_{i-1}$ -measurable random variables  $C_i$ , then we have

$$\mathbb{P} \left\{ \sum_{i=1}^n a_i V_i > \sum_{i=1}^n a_i C_i + t \right\} \leq \mathbb{P} \left\{ \sum_{i=1}^n a_i V_i > \sum_{i=1}^n a_i \mathbb{E} [V_i \mid \mathcal{F}_{i-1}] + t \right\} \leq \exp \left\{ -\frac{t^2}{2v} \right\}$$

For our analysis, we generally cannot compute  $\mathbb{E} [V_i \mid \mathcal{F}_{i-1}]$ , but we can find “nice”  $C_i$ .

To find  $\sigma_i^2$  for use in Lemma 22, we frequently use the following conditional version of Hoeffding’s Lemma.

**Lemma 23** (Conditional Hoeffding's Lemma). *If a random variable  $V$  and a sigma algebra  $\mathcal{F}$  satisfy  $a \leq V \leq b$  and  $\mathbb{E}[V | \mathcal{F}] = 0$ , then*

$$\mathbb{E}[e^{sV} | \mathcal{F}] \leq \exp\left\{\frac{1}{8}(b-a)^2 s^2\right\}$$

*Proof.* We follow standard proof of Hoeffding's Lemma from [19]. Since  $e^{sx}$  is convex, it follows that

$$e^{sx} \leq \frac{b-x}{b-a}e^{sa} + \frac{x-a}{b-a}e^{sb} \quad a \leq x \leq b$$

Therefore, taking the conditional expectation with respect to  $\mathcal{F}$  yields

$$\mathbb{E}[e^{sV} | \mathcal{F}] \leq \frac{b - \mathbb{E}[V | \mathcal{F}]}{b-a}e^{sa} + \frac{\mathbb{E}[V | \mathcal{F}] - a}{b-a}e^{sb} \quad (22)$$

Let  $h = s(b-a)$ ,  $p = -\frac{a}{b-a}$ , and  $L(h) = -hp + \log(1-p + pe^h)$ . Then we have

$$\begin{aligned} e^{L(h)} &= \frac{b}{b-a}e^{sa} + \frac{-a}{b-a}e^{sb} \\ &= \frac{b - \mathbb{E}[V | \mathcal{F}]}{b-a}e^{sa} + \frac{\mathbb{E}[V | \mathcal{F}] - a}{b-a}e^{sb} \end{aligned} \quad (23)$$

since  $\mathbb{E}[V | \mathcal{F}] = 0$ . Since  $L(h) = L'(h) = 0$  and  $L''(h) \leq \frac{1}{4}$ , it holds that  $L(h) \leq \frac{1}{8}(b-a)^2 s^2$ . Combining this bound on  $L(h)$  with (22) and (23) yields the result.  $\square$

Before proceeding with our analysis, we need to introduce a few useful concentration inequalities for sub-Gaussian vector-valued random variables. First, for a scalar random variable  $\xi$ , define the sub-Gaussian norm

$$\tau(\xi) = \inf\left\{a > 0 \mid \mathbb{E}[e^{s\xi}] \leq e^{\frac{1}{2}a^2 s^2} \quad \forall s \geq 0\right\} \quad (24)$$

Clearly, if  $\tau(\xi) < +\infty$ , then  $\xi$  is sub-Gaussian. Second, for a random vector  $\mathbf{v}$  in  $\mathbb{R}^d$ , define

$$B(\mathbf{v}) = \sum_{i=1}^d \tau((\mathbf{v})_i) \quad (25)$$

where  $(\mathbf{v})_i$  is the  $i^{\text{th}}$  component of  $\mathbf{v}$ . We define  $\mathbf{v}$  to be sub-Gaussian if  $B(\mathbf{v}) < +\infty$ .

Of crucial importance in our analysis is analyzing the norm of an average of vector-valued sub-Gaussian random variables. The following lemma describes how to control the sub-Gaussian norm in such a situation.

**Lemma 24.** *Suppose that  $\{\mathbf{v}_i\}_{i=1}^K$  is a collection of independent sub-Gaussian random variables in  $\mathbb{R}^d$ . Then it holds that*

$$B\left(\frac{1}{K} \sum_{i=1}^K \mathbf{v}_i\right) \leq \frac{1}{K} \sum_{j=1}^d \sqrt{\sum_{i=1}^K \tau^2((\mathbf{v}_i)_j)}$$

*If in addition the random variables  $\{\mathbf{v}_i\}_{i=1}^K$  satisfy*

$$\max_{i=1, \dots, K} \max_{j=1, \dots, d} \tau^2((\mathbf{v}_i)_j) \leq \tau^2$$

*then it holds that*

$$B\left(\frac{1}{K} \sum_{i=1}^K \mathbf{v}_i\right) \leq \frac{\tau d}{\sqrt{K}}$$

*Proof.* We analyze one component of the sum  $\frac{1}{K} \sum_{i=1}^K \mathbf{v}_i$ . It holds that

$$\begin{aligned} \mathbb{E} \left[ \exp \left\{ s \left( \frac{1}{K} \sum_{i=1}^K \mathbf{v}_i \right)_j \right\} \right] &= \mathbb{E} \left[ \exp \left\{ \frac{s}{K} \sum_{i=1}^K (\mathbf{v}_i)_j \right\} \right] \\ &= \prod_{i=1}^K \mathbb{E} \left[ \exp \left\{ \frac{s}{K} (\mathbf{v}_i)_j \right\} \right] \\ &\leq \prod_{i=1}^K \exp \left\{ \frac{1}{2} \frac{1}{K^2} \tau^2((\mathbf{v}_i)_j) s^2 \right\} \\ &= \exp \left\{ \frac{1}{2} \left( \frac{1}{K^2} \sum_{i=1}^K \tau^2((\mathbf{v}_i)_j) \right) s^2 \right\} \end{aligned}$$

This implies that

$$\tau \left( \left( \frac{1}{K} \sum_{i=1}^K \mathbf{v}_i \right)_j \right) \leq \frac{1}{K} \sqrt{\sum_{i=1}^K \tau^2((\mathbf{v}_i)_j)}$$

and so

$$B \left( \frac{1}{K} \sum_{i=1}^K \mathbf{v}_i \right) \leq \frac{1}{K} \sum_{j=1}^d \sqrt{\sum_{i=1}^K \tau^2((\mathbf{v}_i)_j)}$$

Finally, if  $\tau^2((\mathbf{v}_i)_j) \leq \tau^2$ , then we have

$$\begin{aligned} B \left( \frac{1}{K} \sum_{i=1}^K \mathbf{v}_i \right) &\leq \frac{1}{K} \sum_{j=1}^d \sqrt{\sum_{i=1}^K \tau^2((\mathbf{v}_i)_j)} \\ &\leq \frac{d}{K} \sqrt{\sum_{i=1}^K \tau^2} \\ &= \frac{\tau d}{\sqrt{K}} \end{aligned}$$

□

Example 3.2 from [17], a consequence of Theorem 3.1 in [17], is useful for the concentration of the norm of sub-Gaussian vector random variables.

**Lemma 25** (Example 3.2 of [17]). *If  $\mathbf{v}$  is a random vector in  $\mathbb{R}^d$  with  $B(\mathbf{v}) < +\infty$ , then*

$$\mathbb{P} \{ \|\mathbf{v}\| > t \} \leq 2 \exp \left\{ -\frac{t^2}{2B^2(\mathbf{v})} \right\}$$

Finally, we will also need to deal with dependent random variables that are sub-Gaussian with respect to a particular filtration.

**Lemma 26.** *Suppose that a random variable  $V$  and a sigma algebra  $\mathcal{F}$  satisfies*

1.  $\mathbb{E}[V \mid \mathcal{F}] = 0$
2.  $\mathbb{P} \{ |V| > t \mid \mathcal{F} \} \leq 2e^{-ct^2}$  with  $c$  a constant.

*Then it holds that*

$$\mathbb{E}[e^{sV} \mid \mathcal{F}] \leq \exp \left\{ \frac{1}{2} \left( \frac{9}{c} \right) s^2 \right\}$$

*for all  $s \geq 0$ .*

*Proof.* Adapted from the characterization of sub-Gaussian random variables in [15]. First, we have for any  $a < c$  that

$$\begin{aligned}\mathbb{E}\left[e^{aV^2} \mid \mathcal{F}\right] &\leq 1 + \int_0^\infty 2ate^{at^2} \mathbb{P}\{|V| > t \mid \mathcal{F}\} dt \\ &\leq 1 + \int_0^\infty 2ate^{-(c-a)t^2} dt \\ &= 1 + \frac{2a}{c-a}\end{aligned}$$

Setting  $a = \frac{c}{3}$  yields the bound

$$\mathbb{E}\left[e^{aV^2} \mid \mathcal{F}\right] \leq 2$$

Since  $\mathbb{E}[V \mid \mathcal{F}] = 0$ , by a Taylor expansion we have

$$\begin{aligned}\mathbb{E}\left[e^{sV} \mid \mathcal{F}\right] &= 1 + \int_0^\infty (1-y)\mathbb{E}\left[(sV)^2 e^{ysV} \mid \mathcal{F}\right] dy \\ &\leq \left(1 + \frac{s^2}{a}\right) e^{\frac{s^2}{2a}} \\ &\leq \exp\left\{\frac{5s^2}{2a}\right\} \\ &= \exp\left\{\frac{1}{2}\left(\frac{9}{c}\right)s^2\right\}\end{aligned}$$

□