

Probabilistic Guarantees of Stochastic Recursive Gradient in Non-Convex Finite Sum Problems

Yanjie Zhong¹, Jiaqi Li^{*1}, and Soumendhra Lahiri¹

¹*Department of Statistics and Data Science, Washington University in St. Louis*

January 31, 2024

Abstract

This paper develops a new dimension-free Azuma-Hoeffding type bound on summation norm of a martingale difference sequence with random individual bounds. With this novel result, we provide high-probability bounds for the gradient norm estimator in the proposed algorithm Prob-SARAH, which is a modified version of the Stochastic Recursive Gradient algorithm (SARAH), a state-of-art variance reduced algorithm that achieves optimal computational complexity in expectation for the finite sum problem. The in-probability complexity by Prob-SARAH matches the best in-expectation result up to logarithmic factors. Empirical experiments demonstrate the superior probabilistic performance of Prob-SARAH on real datasets compared to other popular algorithms.

Keywords: machine learning, variance-reduced method, stochastic gradient descent, non-convex optimization

*Corresponding author: Jiaqi Li, lijiaqi@wustl.edu

1 Introduction

We consider the popular non-convex finite sum optimization problem in this work, that is, estimating $\mathbf{x}^* \in \mathcal{D} \subseteq \mathbb{R}^d$ minimizing the following loss function

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}), \quad \mathbf{x} \in \mathcal{D}, \quad (1)$$

where $f_i : \mathbb{R}^d \mapsto \mathbb{R}$ is a potentially non-convex function on some compact set \mathcal{D} . Such non-convex problems lie at the heart of many applications of statistical learning James et al. (2013) and machine learning Goodfellow et al. (2016).

Unlike convex optimization problems, in general, non-convex problems are intractable and the best we can expect is to find a stationary point. Given a target error ε , since $\nabla f(\mathbf{x}^*) = 0$, we aim to find an estimator $\hat{\mathbf{x}}$ such that roughly $\|\nabla f(\hat{\mathbf{x}})\| \leq \varepsilon$, where $\nabla f(\cdot)$ denotes the gradient vector the loss function f and $\|\cdot\|$ is the operator norm. With a non-deterministic algorithm, the output $\hat{\mathbf{x}}$ is always stochastic, and the most frequently considered measure of error bound is in expectation, i.e.,

$$\mathbb{E}\|\nabla f(\hat{\mathbf{x}})\|^2 \leq \varepsilon^2. \quad (2)$$

There has been a substantial amount of work providing upper bounds on computational complexity needed to achieve the in-expectation bound. However, in practice, we only run a stochastic algorithm for once and an in-expectation bound cannot provide a convincing bound in this situation. Instead, a high-probability bound is more appropriate by nature. Given a pair of target errors (ε, δ) , we want to obtain an estimator $\hat{\mathbf{x}}$ such that with probability at least $1 - \delta$, $\|\nabla f(\hat{\mathbf{x}})\| \leq \varepsilon$, that is

$$\mathbb{P}(\|\nabla f(\hat{\mathbf{x}})\| \leq \varepsilon) \geq 1 - \delta. \quad (3)$$

Though the Markov inequality might help, in general, an in-expectation bound cannot be simply converted to an in-probability bound with a desirable dependency on δ . It would be important to prove upper bounds on high-probability complexity, which ideally should be polylogarithmic in δ and with polynomial terms comparable to the in-expectation complexity bound.

Gradient-based methods are favored by practitioners due to simplicity and efficiency and have been widely studied by researchers in the non-convex setting (Nesterov 2003; Ghadimi and Lan 2013; Allen-Zhu and Hazan 2016; Reddi et al. 2016; Fang et al. 2018; Wang et al. 2019). Among numerous gradient-based methods, the Stochastic Recursive Gradient algorithm (SARAH) (Nguyen et al. 2017a; Nguyen et al. 2017b;

Wang et al. (2019) is the one with the best first-order guarantee as given an in-expectation error target, in both of convex and non-convex finite sum problems. It is worth noticing that Li (2019) attempted to show that a modified version of SARAH is able to approximate the second-order stationary point with a high probability. However, we believe that their application of the martingale Azuma-Hoeffding inequality is unjustifiable because the bounds are potentially random and uncontrollable. In this paper, we shall provide a correct dimension-free martingale Azuma-Hoeffding inequality with rigorous proofs and leverage it to show in-probability properties for SARAH-based algorithms in the non-convex setting.

1.1 Related Works

- High-Probability Bounds:** While most works in the literature of optimization provide in-expectation bounds, there is only a small fraction of works discussing bounds in the high probability sense. Kakade and Tewari (2009) provide a high-probability bound on the excess risk given a bound on the regret. Jain et al. (2019), Harvey et al. (2019a) and Harvey et al. (2019b) derive some high-probability bounds for SGD in convex online optimization problems. Zhou et al. (2018) and Li and Orabona (2020) prove high-probability bounds for several adaptive methods, including AMSGrad, RMSProp and Delayed AdaGrad with momentum. All these works rely on (generalized) Freedman’s inequality or the concentration inequality given in Lemma 6 in Jin et al. (2019). Different from them, our high-probability results are built on a novel Azuma-Hoeffding type inequality proved in this work and Corollary 8 from Jin et al. (2019). In addition, we notice that Li (2019) provide some probabilistic bounds on a SARAH-based algorithm. However, we believe their use of the plain martingale Azuma-Hoeffding inequality is not justifiable. Fang et al. (2018) show in-probability upper bound for SPIDER. Nevertheless, SPIDER’s practical performance is inferior due to its accuracy-dependent small step size Tran-Dinh et al. (2019) and Wang et al. (2019).
- Variance-Reduced Methods in Non-Convex Finite Sum Problems:** Since the invention of the variance-reduction technique in Le Roux et al. (2012), Johnson and Zhang (2013), and Defazio et al. (2014), there has been a large amount of work incorporating this efficient technique to methods targeting the non-convex finite-sum problem. Subsequent methods, including SVRG (Allen-Zhu and Hazan 2016; Reddi et al. 2016; Li and Li 2018), SARAH (Nguyen et al. 2017a; Nguyen et al. 2017b), SCSG (Lei and Jordan 2017; Lei et al. 2017; Horváth et al. 2020), SNVRG (Zhou et al. 2018), SPIDER (Fang et al. 2018), SpiderBoost (Wang et al. 2019) and PAGE (Li et al. 2021), have greatly reduced computational complexity in non-convex problems.

1.2 Our Contributions

- Dimension-Free Martingale Azuma-Hoeffding inequality:** To facilitate our probabilistic analysis, we provide a novel Azuma-Hoeffding type bound on the summation norm of a martingale difference sequence. The novelty is two-fold. Firstly, same as the plain martingale Azuma-Hoeffding inequality, it provides a dimension-free bound. In a recent paper, a sub-Gaussian type bound has been developed by Jin et al. (2019). However, their results are not dimension-free. Our technique in the proof is built on a classic paper by Pinelis (1992) and is completely different from the random matrix technique used in Jin et al. (2019). Secondly, our concentration inequality allows random bounds on each element of the martingale difference sequence, which is much tighter than a large deterministic bound. It should be highlighted that our novel concentration result perfectly suits the nature of SARAH-style methods where the increment can be characterized as a martingale difference sequence and it can be further used to analyze other algorithms beyond the current paper.
- In-probability error bounds of stochastic recursive gradient:** We design a SARAH-based algorithm, named Prob-SARAH, adapted to the high-probability target and provably show its good in-probability properties. Under appropriate parameter setting, the first order complexity needed to achieve the in-probability target is $\tilde{\mathcal{O}}\left(\frac{1}{\varepsilon^3} \wedge \frac{\sqrt{n}}{\varepsilon^2}\right)$, which matches the best known in-expectation upper bound up to some logarithmic factors (Zhou et al. 2018; Wang et al. 2019; Horváth et al. 2020). We would like to point out that the parameter setting used to achieve such complexity is semi-adaptive to ε . That is, only the final stopping rule relies on ε while other key parameters are independent of ε , including step size, mini-batch sizes, and lengths of loops.
- Probabilistic analysis of SARAH for non-convex finite sum:** Existing literature on the bounds of SARAH is mostly focusing on the strongly convex or general convex settings. We extend the case to the non-convex scenarios, which can be considered as a complimentary study to the stochastic recursive gradient in probability.

1.3 Notation

For a sequence of sets $\mathcal{A}_1, \mathcal{A}_2, \dots$, we denote the smallest sigma algebra containing $\mathcal{A}_i, i \geq 1$, by $\sigma(\bigcup_{i=1}^{\infty} \mathcal{A}_i)$. By abuse of notation, for a random variable \mathbf{X} , we denote the sigma algebra generated by \mathbf{X} by $\sigma(\mathbf{X})$. We define constant $C_e = \sum_{i=0}^{\infty} i^{-2}$. For two scalars $a, b \in \mathbb{R}$, we denote $a \wedge b = \min\{a, b\}$ and $a \vee b = \max\{a, b\}$. When we say a quantity T is $\mathcal{O}_{\theta_1, \theta_2}(\theta_3)$ for some $\theta_1, \theta_2, \theta_3 \in \mathbb{R}$, there exists a $g \in \mathbb{R}$ polylogarithmic in θ_1 and θ_2 such that $T \leq g \cdot \theta_3$, and similarly $\tilde{\mathcal{O}}_{(\cdot)}(\cdot)$ is defined the same but up to a logarithm factor.

2 Prob-SARAH Algorithm

The algorithm Prob-SARAH proposed in our work is a modified version of SpiderBoost (Wang et al. 2019) and SARAH (Nguyen et al. 2017a; Nguyen et al. 2017b). Since the key update structure is originated from (Nguyen et al. 2017a), we call our modified algorithm Prob-SARAH. In fact, it can also be viewed as a generalization of the SPIDER algorithm introduced in (Fang et al. 2018).

We present the Prob-SARAH in Algorithm 1, and here, we provide some explanation of the key steps. Following other SARAH-based algorithms, we adopt a similar gradient approximation design with nested loops, specifically with a checkpoint gradient estimator $\nu_0^{(j)}$ using a large mini-batch size B_j in Line 4 and a recursive gradient estimator $\nu_k^{(j)}$ updated in Line 9. When the mini-batch size B_j is large, we can regard the checkpoint gradient estimator $\nu_0^{(j)}$ as a solid approximation to the true gradient at $\tilde{\mathbf{x}}_{j-1}$. With this checkpoint, we can update the gradient estimator $\nu_k^{(j)}$ with a small mini-batch size b_j while maintaining a desirable estimation accuracy.

To emphasize, our stopping rules in Line 11 of Algorithm 1 is newly proposed, which ensures a critical enhancement of the performance compared to previous literature. In particular, with this new design, we can control the gradient norm of the output with high probability. For a more intuitive understanding of these stopping rules, we will see in our proof sketch section that the gradient norm of iterates in the j -th outer iteration, $\|\nabla f\|$, can be bounded by a linear combination of $\{\nu_k^{(j)}\}_{k=1}^{K_j}$ with a small remainder. The first stopping rule, therefore, strives to control the magnitude of the linear combination of $\{\nu_k^{(j)}\}_{k=1}^{K_j}$, while the second stopping rule is specifically designed to control the size of remainder terms. For this purpose, ε_j should be set as a credible controller of the remainder term, with an example given in Theorems 3.1. In this way, with small preset constants $\tilde{\varepsilon}$ and ε , we guarantee that the output has a desirably small gradient norm, dependent on $\tilde{\varepsilon}$ and ε , when the designed stopping rules are activated. Indeed, Proposition B.1 in Appendix B offers a guarantee that the stopping rule will be definitively satisfied at some point. More refined quantitative results regarding the number of steps required for stopping will follow in Theorems 3.1 and Appendix D.3.

3 Theoretical Results

This section is devoted to the main theoretical result of our proposed algorithm Prob-SARAH. We provide the stop guarantee of the algorithm along with the upper bound of the steps. The high-probability error bound of the estimated gradient is also established. The discussion of the dependence of our algorithm on the parameters is available after we introduce our main theorems.

Algorithm 1 Probabilistic Stochastic Recursive Gradient (Prob-SARAH)

```
1: Input: sample size  $n$ , constraint area  $\mathcal{D}$ , initial point  $\tilde{\mathbf{x}}_0 \in \mathcal{D}$ , large batch size  $\{B_j\}_{j \geq 1}$ , mini batch size  $\{b_j\}_{j \geq 1}$ , inner loop length  $\{K_j\}_{j \geq 1}$ , auxiliary error estimator  $\{\varepsilon_j\}_{j \geq 1}$ , errors  $\tilde{\varepsilon}^2, \varepsilon^2$ 
2: for  $j = 1, 2, \dots$  do
3:   Uniformly sample a batch  $\mathcal{I}_j \subseteq \{1, \dots, n\}$  without replacement,  $|\mathcal{I}_j| = B_j$ ;
4:    $\boldsymbol{\nu}_0^{(j)} \leftarrow \frac{1}{B_j} \sum_{i \in \mathcal{I}_j} \nabla f_i(\tilde{\mathbf{x}}_{j-1})$ ;
5:    $\mathbf{x}_0^{(j)} \leftarrow \tilde{\mathbf{x}}_{j-1}$ ;
6:   for  $k = 1, 2, \dots, K_j$  do
7:      $\mathbf{x}_k^{(j)} \leftarrow \text{Proj}(\mathbf{x}_{k-1}^{(j)} - \eta_j \boldsymbol{\nu}_{k-1}^{(j)}, \mathcal{D})$ , project the update back to  $\mathcal{D}$ ;
8:     Uniformly sample a mini-batch  $\mathcal{I}_k^{(j)} \subseteq \{1, \dots, n\}$  with replacement and  $|\mathcal{I}_k^{(j)}| = b_j$ ;
9:      $\boldsymbol{\nu}_k^{(j)} \leftarrow \frac{1}{b_j} \sum_{i \in \mathcal{I}_k^{(j)}} \nabla f_i(\mathbf{x}_k^{(j)}) - \frac{1}{b_j} \sum_{i \in \mathcal{I}_k^{(j)}} \nabla f_i(\mathbf{x}_{k-1}^{(j)}) + \boldsymbol{\nu}_{k-1}^{(j)}$ ;
10:    end for
11:    if  $\frac{1}{K_j} \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2 \leq \tilde{\varepsilon}^2$  and  $\varepsilon_j \leq \frac{1}{2} \varepsilon^2$  then
12:       $\hat{k} \leftarrow \arg \min_{0 \leq k \leq K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2$ ;
13:      Return  $\hat{\mathbf{x}} \leftarrow \mathbf{x}_{\hat{k}}^{(j)}$ ;
14:    end if
15:     $\tilde{\mathbf{x}}_j \leftarrow \mathbf{x}_{K_j}^{(j)}$ ;
16: end for
```

3.1 Technical Assumptions

We shall introduce some necessary regularized assumptions. Most assumptions are commonly used in the optimization literature. We have further clarifications in Appendix A.

Assumption 3.1 (Existence of achievable minimum). *Assume that for each $i = 1, 2, \dots, n$, f_i has continuous gradient on \mathcal{D} and \mathcal{D} is a compact subset of \mathbb{R}^d . Then, there exists a constant $\alpha_M < \infty$ such that*

$$\max_{1 \leq i \leq n} \sup_{\mathbf{x} \in \mathcal{D}} \|\nabla f_i(\mathbf{x})\| \leq \alpha_M. \quad (4)$$

Also, assume that there exists an interior point \mathbf{x}^* of the set \mathcal{D} such that

$$f(\mathbf{x}^*) = \inf_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}).$$

Assumption 3.2 (L -smoothness). *For each $i = 1, 2, \dots, n$, $f_i : \mathcal{D} \rightarrow \mathbb{R}$ is L -smooth for some constant $L > 0$, i.e.,*

$$\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}')\| \leq L \|\mathbf{x} - \mathbf{x}'\|, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathcal{D}.$$

Assumption 3.3 (L -smoothness extension). *There exists a L -smooth function $\tilde{f} : \mathcal{D} \rightarrow \mathbb{R}$ such that*

$$\tilde{f}(\mathbf{x}) = f(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{D}, \quad \text{and} \quad \tilde{f}(\text{Proj}(\mathbf{x}, \mathcal{D})) \leq \tilde{f}(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^d,$$

where $\text{Proj}(\mathbf{x}, \mathcal{D})$ is the Euclidean projection of \mathbf{x} on some compact set \mathcal{D} .

Assumption 3.4. Assume that the following conditions hold.

1. $\varepsilon \leq \frac{1}{e}$ and $\alpha_M^2 \geq \frac{1}{10240}$, where ε is the target error bound in (3) and α_M is defined in (4).
2. The diameter of \mathcal{D} is at least 1, i.e. $d_1 \triangleq \max\{\|\mathbf{x} - \mathbf{x}'\| : \mathbf{x}, \mathbf{x}' \in \mathcal{D}\} \geq 1$.

Assumption 3.1 also indicates that there exists a positive number Δ_f such that $\sup_{\mathbf{x} \in \mathcal{D}} [f(\mathbf{x}) - f(\mathbf{x}^*)] \leq \Delta_f$. Assumptions 3.1–3.3 are commonly used in the optimization literature, and Assumption 3.4 can be easily satisfied in practical use as long as the initial points are not too far from the optimum. See more comments on assumptions in Appendix A.

3.2 Main Results on Complexity

According to the definition given in Lei and Jordan (2020), an algorithm is called ε -independent if it can guarantee convergence at all target accuracies ε in expectation without explicitly using ε in the algorithm. This is a very favorable property because it means that we no longer need to set the target error beforehand. Here, we introduce a similar property regarding the dependency on ε .

Definition 3.1 (ε -semi-independence). *An algorithm is ε -semi-independent, given δ , if it can guarantee convergence at all target accuracies ε with probability at least δ and the knowledge of ε is only needed in the post-processing. That is, the algorithm can iterate without knowing ε and we can select an appropriate iterate out afterwards.*

The newly introduced property can be perceived as the probabilistic equivalent of ε -independence. As stated in the succeeding theorem, under the given conditions, Prob-SARAH can achieve ε -semi-independence, given δ .

Theorem 3.1. *Suppose that Assumptions 3.1, 3.2, 3.3 and 3.4 are valid. Given a pair of errors (ε, δ) , in Algorithm 1 (Prob-SARAH), set hyperparameters*

$$\eta_j = \frac{1}{4L}, \quad K_j = \sqrt{B_j} = \sqrt{j^2 \wedge n}, \quad b_j = l_j K_j, \quad \varepsilon_j = 8L^2 \tau_j + 2q_j, \quad \tilde{\varepsilon}^2 = \frac{1}{5} \varepsilon^2, \quad (5)$$

for $j \geq 1$, where

$$\tau_j = \frac{1}{j^3}, \delta'_j = \frac{\delta}{4C_e j^4}, \quad l_j = 18 \left(\log\left(\frac{2}{\delta'_j}\right) + \log \log\left(\frac{2d_1}{\tau_j}\right) \right), \quad q_j = \frac{128\alpha_M^2}{B_j} \log\left(\frac{3}{\delta'_j}\right) \mathbf{1}\{B_j < n\}.$$

Then,

$$\text{Comp}(\varepsilon, \delta) = \tilde{\mathcal{O}}_{L, \Delta_f, \alpha_M} \left(\frac{1}{\varepsilon^3} \wedge \frac{\sqrt{n}}{\varepsilon^2} \right),$$

where $\text{Comp}(\varepsilon, \delta)$ represents the number of computations needed to get an output $\hat{\mathbf{x}}$ satisfying $\|\nabla f(\hat{\mathbf{x}})\|^2 \leq \varepsilon^2$ with probability at least $1 - \delta$.

More detailed results can be found in Appendix C. In appendix C, we also introduce another hyper-parameter setting that can lead to a complexity with better dependency on α_M^2 , which could be implicitly affected by the choice of constraint region \mathcal{D} .

3.3 Proof Sketch

In this part, we explain the idea of the proof of Theorem 3.1. Same proofing strategy can be applied to other hyper-parameter settings. First, we bound the difference between $\boldsymbol{\nu}_k^{(j)}$ and $\nabla f(\mathbf{x}_k^{(j)})$ by a linear combination of $\{\|\boldsymbol{\nu}_m^{(j)}\|\}_{m=0}^{k-1}$ and small remainders, with which we can have a good control on $\|\nabla f(\mathbf{x}_k^{(j)})\|$ when the stopping rules are met. Second, we bound the number of steps we need to meet the stopping rules. Combining these 2 key components, we can smoothly get the final conclusions.

Let us firstly introduce a novel Azuma-Hoeffding type inequality, which is key to our analysis.

Theorem 3.2 (Martingale Azuma-Hoeffding Inequality with Random Bounds). *Suppose $\mathbf{z}_1, \dots, \mathbf{z}_K \in \mathbb{R}^d$ is a martingale difference sequence adapted to $\mathcal{F}_0, \dots, \mathcal{F}_K$. Suppose $\{r_k\}_{k=1}^K$ is a sequence of random variables such that $\|\mathbf{z}_k\| \leq r_k$ and r_k is measurable with respect to \mathcal{F}_k , $k = 1, \dots, K$. Then, for any fixed $\delta > 0$, and $B > b > 0$, with probability at least $1 - \delta$, for $1 \leq t \leq K$, either*

$$\exists 1 \leq t \leq K, \sum_{k=1}^t r_k^2 \geq B \text{ or } \left\| \sum_{k=1}^t \mathbf{z}_k \right\|^2 \leq 9 \max \left\{ \sum_{k=1}^t r_k^2, b \right\} \left(\log \left(\frac{2}{\delta} \right) + \log \log \left(\frac{B}{b} \right) \right).$$

Remark 3.1. *It is noteworthy that this probabilistic bound on large-deviation is dimension-free, which is a nontrivial extension of Theorem 3.5 in Pinelis (1994). If r_1, r_2, \dots, r_K are not random, we can let $B = \sum_{k=1}^K r_k^2 + \zeta_1$ and $b = \zeta_2 B$ with $\zeta_1 > 0$, $0 < \zeta_2 < 1$. Since ζ_1 can be arbitrarily close to 0 and ζ_2 can be arbitrarily close to 1, we can recover Theorem 3.5 in Pinelis (1994). Compared with Corollary 8 in Jin et al. (2019), which can be viewed as a sub-Gaussian counterpart of our result, a key feature of our Theorem 3.2 is its dimension-independence. We are also working towards improving the bound in Corollary 8 from Jin et al. (2019) to a dimension-free one.*

The success of Algorithm 1 is largely because $\nabla f(\mathbf{x}_k^{(j)})$ is well-approximated by $\boldsymbol{\nu}_k^{(j)}$, and meanwhile $\boldsymbol{\nu}_k^{(j)}$ can be easily updated. We can observe that $\boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)})$ is actually sum of a sequence of martingale

difference as

$$\begin{aligned}
\boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) &= \left[\frac{1}{b_j} \sum_{i \in \mathcal{I}_k^{(j)}} \nabla f_i(\mathbf{x}_k^{(j)}) - \frac{1}{b_j} \sum_{i \in \mathcal{I}_k^{(j)}} \nabla f_i(\mathbf{x}_{k-1}^{(j)}) + \nabla f(\mathbf{x}_{k-1}^{(j)}) - \nabla f(\mathbf{x}_k^{(j)}) \right] \\
&+ \left[\boldsymbol{\nu}_{k-1}^{(j)} - \nabla f(\mathbf{x}_{k-1}^{(j)}) \right] = \sum_{m=1}^k \left[\frac{1}{b_j} \sum_{i \in \mathcal{I}_m^{(j)}} \nabla f_i(\mathbf{x}_m^{(j)}) - \frac{1}{b_j} \sum_{i \in \mathcal{I}_m^{(j)}} \nabla f_i(\mathbf{x}_{m-1}^{(j)}) \right. \\
&\left. + \nabla f(\mathbf{x}_{m-1}^{(j)}) - \nabla f(\mathbf{x}_m^{(j)}) \right] + \left[\boldsymbol{\nu}_0^{(j)} - \nabla f(\mathbf{x}_0^{(j)}) \right]. \tag{6}
\end{aligned}$$

To be more specific, let $\mathcal{F}_0 = \{\emptyset, \Omega\}$, and iteratively define $\mathcal{F}_{j,-1} = \mathcal{F}_{j-1}$, $\mathcal{F}_{j,0} = \sigma(\mathcal{F}_{j-1} \cup \sigma(\mathcal{I}_j))$, $\mathcal{F}_{j,k} = \sigma(\mathcal{F}_{j,0} \cup \sigma(\mathcal{I}_k^{(j)}))$, $\mathcal{F}_j = \sigma(\bigcup_{k=1}^{\infty} \mathcal{F}_{j,k})$, $j \geq 1, k \geq 1$. We also denote $\boldsymbol{\epsilon}_0^{(j)} \triangleq \boldsymbol{\nu}_0^{(j)} - \nabla f(\mathbf{x}_0^{(j)})$, $\boldsymbol{\epsilon}_m^{(j)} \triangleq \frac{1}{b_j} \sum_{i \in \mathcal{I}_m^{(j)}} \nabla f_i(\mathbf{x}_m^{(j)}) - \nabla f(\mathbf{x}_m^{(j)}) + \nabla f(\mathbf{x}_{m-1}^{(j)}) - \frac{1}{b_j} \sum_{i \in \mathcal{I}_m^{(j)}} \nabla f_i(\mathbf{x}_{m-1}^{(j)})$, $m \geq 1$. Then, we can see that $\{\boldsymbol{\epsilon}_m^{(j)}\}_{m=0}^k$ is a martingale difference sequence adapted to $\{\mathcal{F}_{j,m}\}_{m=-1}^k$. With the help of our new Martingale Azuma-Hoeffding inequality, we can control the difference between $\boldsymbol{\nu}_k^{(j)}$ and $\nabla f(\mathbf{x}_k^{(j)})$ by a linear combination of $\{\|\boldsymbol{\nu}_m^{(j)}\|\}_{m=0}^{k-1}$ and small remainders, with details given in Appendix D.1. Then, given the stopping rules in line 11 and selection method specified in line 12 of Algorithm 1, it would be not hard for us to obtain $\|\nabla f(\hat{\mathbf{x}})\|^2 \leq \varepsilon^2$ with a high probability. More details can be found in Appendix D.2.

Another key question needed to be resolved is, when the algorithm can stop? The following analysis can build some intuitions for us. Given a $T \in \mathbb{Z}_+$, with the bound given in Proposition F.1 in Appendix F, with a high probability,

$$-\Delta_f \leq f(\tilde{\mathbf{x}}_{2T}) - f(\tilde{\mathbf{x}}_T) \leq A_T - \frac{1}{16L} \sum_{j=T+1}^{2T} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2, \tag{7}$$

where A_T is upper bounded by a value polylogarithmic in T . As for the second summation, if $\varepsilon_j \leq \frac{1}{2}\varepsilon^2$ for $j = T, T+1, \dots, 2T$ (which is obviously true when T is moderately large) and our algorithm doesn't stop in $2T$ outer iterations,

$$\frac{1}{16L} \sum_{j=T+1}^{2T} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 \geq \frac{\tilde{\varepsilon}^2}{16L} \sum_{j=T+1}^{2T} K_j \geq \frac{\tilde{\varepsilon}^2}{16L} \sum_{j=T+1}^{2T} (T \wedge \sqrt{n}) = \frac{\tilde{\varepsilon}^2}{16L} T^2 \wedge (\sqrt{n}T),$$

which grows at least linear in T . Consequently, when T is sufficiently large, the RHS of (7) can be smaller than $-\Delta_f$, which leads to a contradiction. Roughly, we can see that the stopping time T cannot exceed the order of $\tilde{\mathcal{O}}\left(\frac{1}{\varepsilon} \vee \frac{1}{\sqrt{n\varepsilon^2}}\right)$. More details can be found in Appendix D.3.

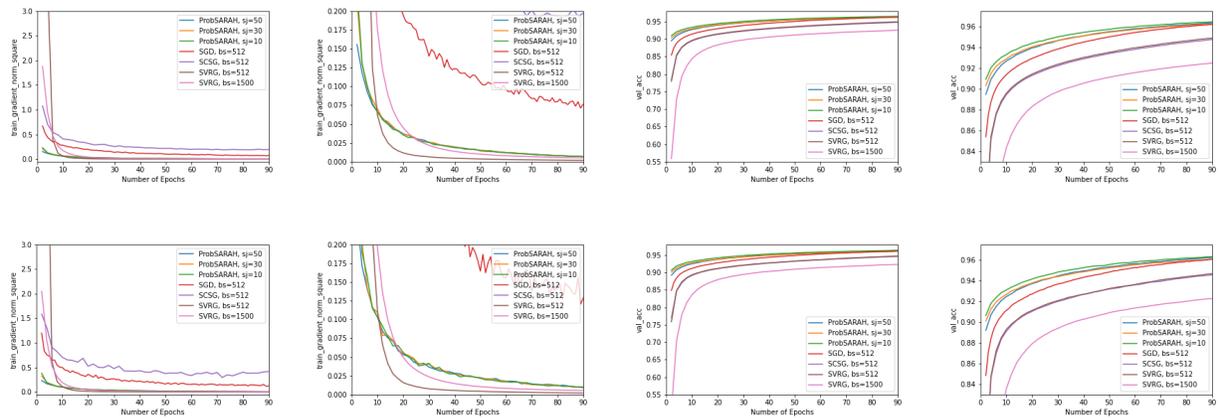


Figure 1: Comparison of convergence with respect to $(1 - \delta)$ -quantile of square of gradient norm ($\|\nabla f\|^2$) and δ -quantile of validation accuracy on the **MNIST** dataset for $\delta = 0.1$ and $\delta = 0.01$. The second (fourth) column presents zoom-in figures of those in the first (third) column. Top: $\delta = 0.1$. Bottom: $\delta = 0.01$. 'bs' stands for batch size. 'sj=x' means that the smallest batch size $\approx x \log x$.

4 Numerical Experiments

In order to validate our theoretical results and show good probabilistic property for the newly-introduced Prob-SARAH, we conduct some numerical experiments where the objectives are possibly non-convex.

4.1 Logistic Regression with Non-Convex Regularization

In this part, we consider to add a non-convex regularization term to the commonly-used logistic regression. Specifically, given a sequence of observations $(\mathbf{w}_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$, $i = 1, 2, \dots, n$ and a regularized parameter $\lambda > 0$, the objective is

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{w}_i^T \mathbf{x}} \right) + \frac{\lambda}{2} \sum_{j=1}^d \frac{x_j^2}{1 + x_j^2}.$$

Such an objective has also been considered in other works like Horváth et al. (2020) and Ji et al. (2020). Same as other works, we set the regularized parameter $\lambda = 0.1$ across all experiments. We compare the newly-introduced Prob-SARAH against three popular methods including SGD (Ghadimi and Lan 2013), SVRG (Reddi et al. 2016) and SCSG (Lei et al. 2017). Based on results given in Theorem 3.1, we let the length of the inner loop $K_j \sim j \wedge \sqrt{n}$, the inner loop batch size $b_j \sim \log j$ ($j \wedge \sqrt{n}$), the outer loop batch size $B_j \sim j^2 \wedge n$. For fair comparison, we determine the batch size (inner loop batch size) for SGD (SCSG and SVRG) based on the sample size n and the number of epochs needed to have sufficient decrease in gradient

norm. For example, for the w7a dataset, the sample size is 24692 and we run 60 epochs in total. In the 20th epoch, the inner loop batch size of Prob-SARAH is approximately $67 \log 67 \approx 281$. Thus, we set batch size 256 for SGD, SCSG and SVRG so that they can be roughly matched. In addition, based on the theoretical results from Reddi et al. (2016), we also consider a large inner loop batch size comparable to $n^{2/3}$ for SVRG. In addition, we set step size $\eta = 0.01$ for all algorithms across all experiments for simplicity.

Results are displayed in Figure 2, from which we can see that Prob-SARAH has superior probabilistic guarantee in controlling the gradient norm in all experiments. It is significantly better than SCSG and SVRG under our current setting. Prob-SARAH can achieve a lower gradient norm than SGD at the early stage while SGD has a slight advantage when the number of epochs is large.

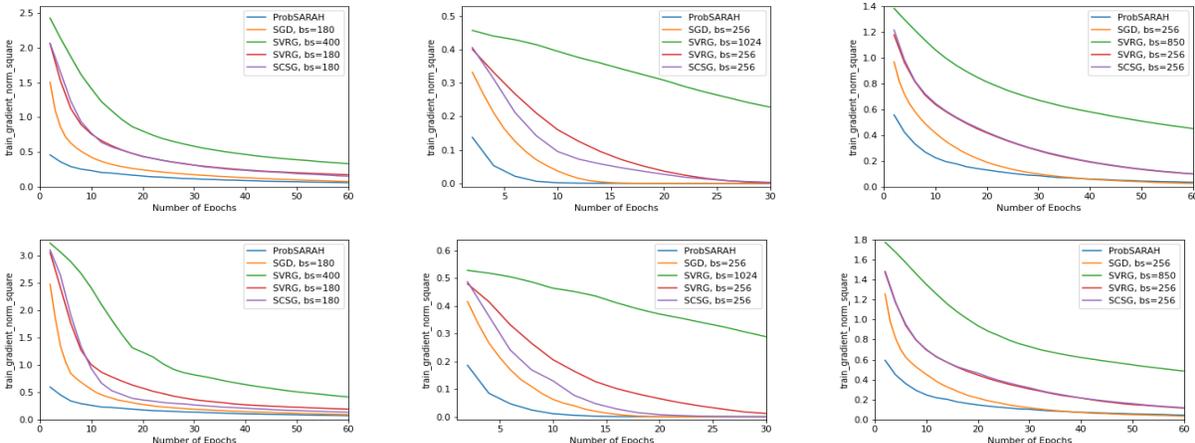


Figure 2: Comparison of convergence with respect to $(1 - \delta)$ -quantile of square of gradient norm ($\|\nabla f\|^2$) over 3 datasets for $\delta = 0.1$ and $\delta = 0.01$. Top: $\delta = 0.1$. Bottom: $\delta = 0.01$. Datasets: **mushrooms**, **ijcnn1**, **w7a** (from left to right). 'bs' stands for batch size.

4.2 Two-Layer Neural Network

We also evaluate the performance of Prob-SARAH, SGD, SVRG and SCSG on the MNIST dataset with a simple 2-layer neural network. The two hidden layers respectively have 128 and 64 neurons. We include a GELU activation layer following each hidden layer. We use the negative log likelihood as our loss function. Under this setting, the objective is possibly non-convex and smooth on any given compact set. The step size is fixed to be 0.01 for all algorithms. For Prob-SARAH, we still have the length of the inner loop $K_j \sim j \wedge \sqrt{n}$, the inner loop batch size $b_j \sim \log j (j \wedge \sqrt{n})$, the outer loop batch size $B_j \sim j^2 \wedge n$. But to reduce computational time, we let j start from 10, 30 and 50 respectively. Based on the same rule described in the previous subsection, we let the batch size (or inner loop batch size) for SGD, SVRG and SCSG be 512.

Results are given in Figure 1. In terms of gradient norm, Prob-SARAH has the best performance among algorithms considered here when the number of epochs is relatively small. With increasing number of epochs, SVRG tends to be better in finding first-order stationary points. However, based on the 3rd and 4th columns in Figure 1, SVRG apparently has an inferior performance on the validation set, which indicates that it could be trapped at local minima. In brief, Prob-SARAH achieves the best tradeoff between finding a first-order stationary point and generalization.

We also consider another set of experiments by replacing the GELU activation function with ReLU, resulting in a non-smooth objective. The results are shown in Appendix G, which resemble those in Figure 1 and the similar conclusions can be drawn.

5 Conclusion

In this paper, we propose a SARAH-based variance reduction algorithm called Prob-SARAH and provide high-probability bounds on gradient norm for estimator resulted from Prob-SARAH. Under appropriate assumptions, the high-probability first order complexity nearly match the one in the in-expectation sense. The main tool used in the theoretical analysis is a novel Azuma-Hoeffding type inequality. We believe that similar probabilistic analysis can be applied to SARAH-based algorithms in other settings.

A Remarks and Examples for Assumptions

A.1 More comments on Assumptions 3.1–3.4

Remark A.1 (Convexity and smoothness). *It is worth noticing that Assumption 3.1 is widely used in many non-convex optimization works and can be met for most applications in practice. Assumption 3.2 is also needed in deriving in-expectation bound for many non-convex variance-reduced methods, including state-of-art ones like SPIDER and SpiderBoost. As for Assumption 3.3, it is a byproduct of the compact constraint and can be satisfied with some commonly-seen f and usual choices of \mathcal{D} . For more discussions on Assumption 3.3, please see Appendix A.2.*

Remark A.2 (Compact set \mathcal{D}). *Compared with other works in the literature of non-convex optimization, the compact constraint region $\mathcal{D} \in \mathbb{R}^d$ imposed in the finite sum problem (1) may seem somewhat restrictive. In fact, such constraint is largely due to technical convenience and it can be removed with additional condition on gradients. We will elaborate on this point in subsection C.1. Besides, in many practical applications, it is reasonable to restrict estimators to a compact set when certain prior knowledge is available.*

A.2 An Example of Assumption 3.3

Let us consider the logistic regression with non-convex regularization where the object function can be characterized as

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \langle \mathbf{w}_i, \mathbf{x} \rangle)) + \frac{\lambda}{2} \Phi(\mathbf{x}),$$

where $\Phi(\mathbf{x}) = \sum_{j=1}^d (x_j^2)^{\frac{1}{4}}$, x_j is the j th element of \mathbf{x} , $\lambda > 0$ is the regularization parameter, $\{y_i\}_{i=1}^n$ are labels and $\{\mathbf{w}_i\}_{i=1}^n$ are normalized covariates with norm 1. In fact, for any fixed $\lambda > 0$, Assumption 3.3 holds with $\tilde{f} = f$ and $\mathcal{D} = \{\mathbf{x} : \|\mathbf{x}\| \leq R\}$ when R is sufficiently large. Since smoothness is easy to show, we focus on the second part of Assumption 3.3. To show that

$$f(\text{Proj}(\mathbf{x}, \mathcal{D})) \leq f(\mathbf{x})$$

holds for any $\mathbf{x} \in \mathbb{R}^d$, since the projection direction is pointed towards the origin, it suffices to show that for any $\boldsymbol{\nu} \in \mathbb{R}^d$ with $\|\boldsymbol{\nu}\| = 1$,

$$\frac{d}{dt} f_i(t\boldsymbol{\nu}) = \frac{d}{dt} \left(\log(1 + \exp(-ty_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle)) + \frac{\lambda}{2} \sum_{j=1}^d \sqrt{t} (\nu_j^2)^{\frac{1}{4}} \right) \geq 0,$$

when $t \geq R$ for $i = 1, 2, \dots, n$, where ν_j is the j th element of $\boldsymbol{\nu}$. To see this,

$$\begin{aligned} & \frac{d}{dt} f_i(t\boldsymbol{\nu}) \\ &= \frac{-y_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle \exp(-ty_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle)}{1 + \exp(-ty_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle)} + \frac{\lambda}{2} \sum_{j=1}^d \frac{(\nu_j^2)^{\frac{1}{4}}}{2\sqrt{t}} \\ &= \frac{-y_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle}{1 + \exp(ty_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle)} + \frac{\lambda}{2} \sum_{j=1}^d \frac{(\nu_j^2)^{\frac{1}{4}}}{2\sqrt{t}} \\ &\geq \frac{-y_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle}{1 + \exp(ty_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle)} + \frac{\lambda}{2} \sum_{j=1}^d \frac{\nu_j^2}{2\sqrt{t}} \\ &= \frac{-y_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle}{1 + \exp(ty_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle)} + \frac{\lambda}{4\sqrt{t}}. \end{aligned}$$

If $y_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle \leq 0$, we can immediately know that $\frac{d}{dt} f_i(t\boldsymbol{\nu}) \geq 0$ for any $t > 0$.

If $y_i \langle \mathbf{w}_i, \boldsymbol{\nu} \rangle > 0$, let us consider an auxiliary function

$$g(b) = \frac{-b}{1 + e^{tb}}.$$

Then,

$$g'(b) \propto -(1 + e^{tb}) + bte^{bt},$$

from where we can know the minimum of $g(b)$ is achieved for some $b^* \in [\frac{1}{t}, \frac{2}{t}]$. Thus,

$$g(b) \geq g(b^*) \geq \frac{-2}{(1 + e^{tb})t} \geq \frac{-2}{(1 + e)t}.$$

Therefore,

$$\frac{d}{dt} f_i(t\boldsymbol{\nu}) \geq \frac{-2}{(1 + e)t} + \frac{\lambda}{4\sqrt{t}},$$

which is positive when $t \geq \left(\frac{8}{(1+e)\lambda}\right)^2$.

If we consider other non-convex regularization terms in logistic regression, such as $\Phi(\mathbf{x}) = \sum_{j=1}^d \frac{x_j^2}{1+x_j^2}$, we may no longer enjoy Assumption 3.3 because monotony may not hold for a few projection directions even when the constraint region is large. Nevertheless, such theoretical flaw can be easily remedied by adding an extra regularization term like $\frac{\lambda_e}{2} \|\mathbf{x}\|^2$ with appropriate $\lambda_e > 0$.

B Stop Guarantee

We would like to point out that, under appropriate parameter setting, Prob-SARAH is guaranteed to stop. Actually, we can have the stopping guarantee under more general conditions than those stated in the following proposition. But for simplicity, we only present conditions naturally matched parameter settings given in the next two subsections.

Proposition B.1 (Stop guarantee of Prob-SARAH). *Suppose that Assumptions 3.1, 3.2, 3.3 and 3.4 are satisfied. Let step size $\eta_j \equiv 1/(4L)$ and suppose that $b_j \geq K_j$, $j \geq 1$. The large batch size $\{B_j\}_{j \geq 1}$ is set appropriately such that $B_j = n$ when j is sufficiently large. If the limit of $\{\varepsilon_j\}_{j \geq 1}$ is 0, then, for any fixed $\tilde{\varepsilon}$ and ε , with probability 1, Prob-SARAH (Algorithm 1) stops. In settings where we always have $\varepsilon_j \leq \frac{1}{2}\varepsilon^2$, we also have the result that Prob-SARAH (Algorithm 1) stops with probability 1.*

C Detailed Results on Complexity

Theorem C.1. *Suppose that Assumptions 3.1, 3.2, 3.3 and 3.4 are valid. Given a pair of errors (ε, δ) , in Algorithm 1 (Prob-SARAH), set hyperparameters*

$$\eta_j = \frac{1}{4L}, \quad K_j = \sqrt{B_j} = \sqrt{j^2 \wedge n}, \quad b_j = l_j K_j, \quad \varepsilon_j = 8L^2 \tau_j + 2q_j, \quad \tilde{\varepsilon}^2 = \frac{1}{5}\varepsilon^2, \quad (8)$$

for $j \geq 1$, where

$$\tau_j = \frac{1}{j^3}, \delta'_j = \frac{\delta}{4C_e j^4}, \quad l_j = 18 \left(\log\left(\frac{2}{\delta'_j}\right) + \log \log\left(\frac{2d_1}{\tau_j}\right) \right), \quad q_j = \frac{128\alpha_M^2}{B_j} \log\left(\frac{3}{\delta'_j}\right) \mathbf{1}\{B_j < n\}.$$

Then, with probability at least $1 - \delta$, Prob-SARAH stops in at most

$$2(T_1 \vee T_2 \vee T_3 \vee T_4) = \tilde{\mathcal{O}}_{L, \Delta_f, \alpha_M} \left(\frac{1}{\varepsilon} + \frac{1}{\sqrt{n\varepsilon^2}} \right)$$

outer iterations and the output satisfies $\|\nabla f(\hat{\mathbf{x}})\|^2 \leq \varepsilon^2$. Detailed definitions of T_1, T_2, T_3 and T_4 can be found in Propositions [D.3](#) and [D.4](#).

Corollary C.1. Under parameter settings in [Theorem C.1](#),

$$\text{Comp}(\varepsilon, \delta) = \tilde{\mathcal{O}}_{L, \Delta_f, \alpha_M} \left(\frac{1}{\varepsilon^3} \wedge \frac{\sqrt{n}}{\varepsilon^2} \right).$$

We introduce another setting that can help to reduce the dependence on α_M^2 , which could be implicitly affected by the choice of constraint region \mathcal{D} . We should also notice that, under such setting, the algorithm is no longer ε -semi-independent.

Theorem C.2. Suppose that [Assumptions 3.1, 3.2, 3.3](#) and [3.4](#) are valid. We denote $\Delta_f^0 \triangleq f(\tilde{\mathbf{x}}_0) - f(\mathbf{x}^*)$. Given a pair of errors (ε, δ) , in [Algorithm 1](#) (Prob-SARAH), set parameters

$$\eta_j = \frac{1}{4L}, \quad K_j = \sqrt{B_j} = \sqrt{n}, \quad b_j = l_j K_j, \quad \varepsilon_j = \frac{1}{2} \tilde{\varepsilon}^2 = \frac{1}{10} \varepsilon^2, \quad (9)$$

for $j \geq 1$, where

$$\tau_j = \frac{1}{40L^2} \varepsilon^2, \delta'_j = \frac{\delta}{4C_e j^4}, \quad l_j = 18 \left(\log\left(\frac{2}{\delta'_j}\right) + \log \log\left(\frac{2d_1}{\tau_j}\right) \right).$$

Then, with probability at least $1 - \delta$, Prob-SARAH stops in at most

$$T_5 = \frac{160L(\Delta_f^0 + 1)}{\sqrt{n\varepsilon^2}} = \mathcal{O}_{L, \Delta_f^0} \left(\frac{1}{\sqrt{n\varepsilon^2}} \right)$$

outer iterations and the output satisfies $\|\nabla f(\hat{\mathbf{x}})\|^2 \leq \varepsilon^2$.

Corollary C.2. Under parameter settings in [Theorem C.2](#),

$$\text{Comp}(\varepsilon, \delta) = \tilde{\mathcal{O}}_{L, \Delta_f^0} \left(\frac{\sqrt{n}}{\varepsilon^2} \right).$$

Comparing the complexities in Corollary C.1 and Corollary C.2, we can notice that under the second setting, we can get rid of the dependence on α_M at the expense of losing some adaptivity to ε . We would also like to point out that, in the low precision region, i.e. when $\frac{1}{\varepsilon} = o(\sqrt{n})$, the second complexity is inferior.

C.1 Dependency on Parameters

To apply our newly-introduced Azuma-Hoeffding type inequality (see Theorem 3.2), it is necessary to impose a compact constraint region \mathcal{D} . Therefore, let us provide a delicate analysis on how \mathcal{D} can affect the convergence guarantee.

Dependency on d_1 : d_1 , the diameter of \mathcal{D} , is a parameter directly related to the choice of \mathcal{D} . Shown in theoretical results presented above, the in-probability first-order complexities always have a polylogarithmic dependency on d_1 , which implies that as long as d_1 is polynomial in n or $\frac{1}{\varepsilon}$, it should only have a minor effect on the complexity. With certain prior knowledge, we should be able to control d_1 at a reasonable scale.

Dependency on Δ_f and Δ_f^0 : Under the setting given in Theorem 3.1, the first-order complexity is polynomial in Δ_f . Such dependency implicates that the complexity would not deteriorate much if Δ_f is of a small order, which is definitely true when the loss function is bounded. As for the setting given in Theorem C.2, the first-order complexity is polynomial in Δ_f^0 , which is conventionally assumed to be $\mathcal{O}(1)$ and will not be affected by \mathcal{D} .

D Postponed Proofs for the Results in Section 3

D.1 Bounding the Difference between $\boldsymbol{\nu}_k^{(j)}$ and $\nabla f(\mathbf{x}_k^{(j)})$

Proposition D.1. For $k \geq 0, j \geq 1$, denote

$$(\tilde{\sigma}_k^{(j)})^2 \triangleq \frac{4L^2\eta_j^2}{b_j} \sum_{m=1}^k \|\boldsymbol{\nu}_{m-1}^{(j)}\|^2.$$

Under Assumptions 3.2 and 3.4, for any prescribed constant $\delta' \in (0, 1), \tau \in (0, 1), k \geq 0, j \geq 1$,

$$\begin{aligned} & \|\boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)})\|^2 \\ & \leq 18 \left((\tilde{\sigma}_k^{(j)})^2 + \frac{4L^2\tau k}{b_j} \right) \left(\log \frac{2}{\delta'} + \log \log \frac{2d_1^2}{\tau} \right) \\ & \quad + \frac{128\alpha_M^2}{B_j} \log \frac{3}{\delta'} \mathbf{1}\{B_j < n\} \end{aligned} \tag{10}$$

with probability at least $1 - 2\delta'$.

Remark D.1. Let us briefly explain this high-probability bound on $\|\boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)})\|^2$. When $k = o(b_j)$ and $L = \mathcal{O}(1)$, by letting τ be of appropriate n^{-1} -polynomial order, $4L^2\tau k/b_j$ will be roughly $o(1)$. If further we have d_1 be of n -polynomial order and let δ' be of n^{-1} -polynomial order, $\log(2/\delta') + \log \log(2d_1^2/\tau)$ will be $\tilde{\mathcal{O}}(1)$. As a result, the upper bound is roughly $(\tilde{\sigma}_k^{(j)})^2 = (4L^2\eta_j^2/b_j) \sum_{m=1}^k \|\boldsymbol{\nu}_{m-1}^{(j)}\|^2$ when B_j is sufficiently large so that the last term in the bound (10) is negligible. Bounding $\|\boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)})\|^2$ by linear combination of $\{\|\boldsymbol{\nu}_m^{(j)}\|\}_{m=0}^\infty$ is the key to our analysis.

D.2 Analysis on the Output $\hat{\mathbf{x}}$

Under parameter setting specified in Theorem 3.1, if we suppose that the algorithm stops at the j -th outer iteration, i.e.

$$\frac{1}{K_j} \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2 \leq \tilde{\varepsilon}^2, \quad \varepsilon_j \leq \frac{1}{2}\varepsilon, \quad (11)$$

there must exist a $0 \leq k' \leq K_j - 1$, such that $\|\boldsymbol{\nu}_{k'}^{(j)}\|^2 \leq \tilde{\varepsilon}^2$.

Then, on the event

$$\Omega_j \triangleq \left\{ \omega : \left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 \leq l_j \left((\tilde{\sigma}_k^{(j)})^2 + \frac{4L^2\tau_j k}{b_j} \right) + q_j, 0 \leq k \leq K_j \right\}, \quad (12)$$

where $l_j = 18 \left(\log \frac{2}{\delta'_j} + \log \log \frac{2d_1^2}{\tau_j} \right)$ and $q_j = \frac{128\alpha_M^2}{B_j} \log \frac{3}{\delta'_j} \mathbf{1}\{B_j < n\}$, we can easily derive an upper bound on $\left\| \nabla f(\mathbf{x}_{k'}^{(j)}) \right\|^2$,

$$\begin{aligned} & \left\| \nabla f(\mathbf{x}_{k'}^{(j)}) \right\|^2 \\ & \leq 2\|\boldsymbol{\nu}_{k'}^{(j)}\|^2 + 2\left\| \boldsymbol{\nu}_{k'}^{(j)} - \nabla f(\mathbf{x}_{k'}^{(j)}) \right\|^2 \\ & \leq 2\tilde{\varepsilon}^2 + 2l_j \left((\tilde{\sigma}_{k'}^{(j)})^2 + \frac{4L^2\tau_j k'}{b_j} \right) + 2q_j \\ & = 2\tilde{\varepsilon}^2 + 2l_j \left(\frac{4L^2\eta_j^2}{b_j} \sum_{m=1}^{k'} \|\boldsymbol{\nu}_{m-1}^{(j)}\|^2 + \frac{4L^2\tau_j k'}{b_j} \right) + 2q_j \\ & \leq 2\tilde{\varepsilon}^2 + 2l_j \left(\frac{4L^2\eta_j^2}{b_j} \sum_{m=1}^{K_j} \|\boldsymbol{\nu}_{m-1}^{(j)}\|^2 + \frac{4L^2\tau_j K_j}{b_j} \right) + 2q_j \\ & \leq 2\tilde{\varepsilon}^2 + 2l_j \left(\frac{4L^2\eta_j^2 K_j}{b_j} \tilde{\varepsilon}^2 + \frac{4L^2\tau_j K_j}{b_j} \right) + 2q_j \\ & = 2\tilde{\varepsilon}^2 + 2l_j \left(\frac{K_j}{4b_j} \tilde{\varepsilon}^2 + \frac{4L^2\tau_j K_j}{b_j} \right) + 2q_j \\ & = 2.5\tilde{\varepsilon}^2 + 8L^2\tau_j + 2q_j = 2.5\tilde{\varepsilon}^2 + \varepsilon_j \leq \varepsilon^2, \end{aligned}$$

where the 2nd step is based on (10) with definitions of l_j and q_j given in Theorem 3.1, the 5th step is based on (11), the 6th step is based on the choice of $\eta_j = \frac{1}{4L}$ and the 7th step is based on the choice of $b_j = l_j K_j$. In addition, based on Proposition D.1, the union event $\bigcup_{j=1}^{\infty} \Omega_j$ occurs with probability at least

$$1 - 2 \sum_{j=1}^{\infty} \sum_{k=0}^{K_j} \delta'_j \geq 1 - \sum_{j=1}^{\infty} \frac{\delta K_j}{2C_e j^4} \geq 1 - \sum_{j=1}^{\infty} \frac{\delta}{C_e j^2} = 1 - \delta.$$

In one word, it is highly likely to control the norm of gradient at our desired level when the algorithm stops.

The above results can sufficiently explain our choice of stopping rule imposed in Algorithm 1. We can summarize them as the following proposition.

Proposition D.2. *Suppose that Assumptions 3.2 and 3.4 are true. Under the parameter setting given in Theorem 3.1, the output of Algorithm 1 satisfies*

$$\|\nabla f(\hat{\mathbf{x}})\|^2 \leq \varepsilon^2,$$

with probability at least $1 - \delta$.

D.3 Upper-bounding the Stopping Time

Proposition D.3 (First Stopping Rule). *Suppose that Assumptions 3.2, 3.3 and 3.4 are valid. Let*

$$\begin{aligned} T_1 &= \left\lceil \frac{\sqrt{320L(c_1 + \Delta_f)}}{\varepsilon} + \frac{320L(c_1 + \Delta_f)}{\sqrt{n}\varepsilon^2} \right\rceil, \\ T_2 &= \left\lceil 3 \left(\frac{\sqrt{320Lc_2}}{\varepsilon} \log \frac{\sqrt{320Lc_2}}{\varepsilon} + \frac{640Lc_2}{\sqrt{n}\varepsilon^2} \log \frac{320Lc_2}{\varepsilon^2} + 1 \right) \right\rceil, \end{aligned} \quad (13)$$

where

$$c_1 = \frac{C_e L}{4} + \frac{16\alpha_M^2}{L} \log \frac{192C_e}{\delta}, \quad c_2 = \frac{64\alpha_M^2}{L}.$$

Under the parameter setting given in Theorem 3.1, on Ω , when $T \geq T_1 \vee T_2$, there exists a $T + 1 \leq j \leq 2T$ such that

$$\frac{1}{K_j} \sum_{k=0}^{K_j-1} \|\mathbf{v}_k^{(j)}\|^2 \leq \tilde{\varepsilon}^2.$$

Proposition D.4 (Second Stopping Rule). *Let*

$$T_3 = \left\lceil \frac{2\sqrt{c_3}}{\varepsilon} \right\rceil, \quad T_4 = \left\lceil \frac{6\sqrt{c_4}}{\varepsilon} \log \frac{2\sqrt{c_4}}{\varepsilon} \right\rceil, \quad (14)$$

where

$$c_3 = 8L^2 + 256\alpha_M^2 \log \frac{12C_e}{\delta}, \quad c_4 = 1024\alpha_M^2.$$

Under the parameter setting given in Theorem 3.1, on Ω , when $T \geq T_3 \vee T_4$,

$$\varepsilon_T \leq \frac{1}{2}\varepsilon^2.$$

Proposition D.5 (Stop Guarantee). *Under the parameter setting and assumptions given in Theorem 3.1, on Ω , when $T \geq T_1 \vee T_2 \vee T_3 \vee T_4$, Algorithm 1 stops in at most $2T$ outer iterations.*

Proof. If Algorithm 1 stops in T outer iterations, our conclusion is obviously true. If not, according to Proposition D.3, there must exist a $j \in [T + 1, 2T]$ such that the first stopping rule is met, i.e.

$$\frac{1}{K_j} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 \leq \tilde{\varepsilon}^2.$$

According to Proposition D.4, the second stopping rule is also met, i.e. $\varepsilon_j \leq \frac{1}{2}\varepsilon^2$.

Consequently, the algorithm stops at the j -th outer iteration. □

E Technical Lemmas

Lemma E.1 (Theorem 4 in Hoeffding (1963)). *Let $\{\boldsymbol{\epsilon}_1, \boldsymbol{\epsilon}_2, \dots, \boldsymbol{\epsilon}_n\}$ be a set of fixed vectors in \mathbb{R}^d . $\mathcal{I}, \mathcal{J} \subseteq \{1, 2, \dots, n\}$ are 2 random index sets sampled respectively with replacement and without replacement, with size $|\mathcal{I}| = |\mathcal{J}| = k$. For any continuous and convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$,*

$$\mathbb{E}f\left(\sum_{j \in \mathcal{J}} \boldsymbol{\epsilon}_j\right) \leq \mathbb{E}f\left(\sum_{i \in \mathcal{I}} \boldsymbol{\epsilon}_i\right).$$

Lemma E.2 (Proposition 1.2 in Boucheron et al. (2013)¹). *Let X be real random variable such that $\mathbb{E}X = 0$ and $a \leq X \leq b$ for some $a, b \in \mathbb{R}$. Then, for all $t \in \mathbb{R}$,*

$$\log \mathbb{E}e^{tX} \leq \frac{t^2(b-a)^2}{8}.$$

¹See also Lemma 1.3 in Bardenet and Maillard (2015).

Lemma E.3 (Theorem 3.5 in Pinelis (1994)²). Let $\{\epsilon_k\}_{k=1}^K \subseteq \mathbb{R}^d$ be a vector-valued martingale difference sequence with respect to \mathcal{F}_k , $k = 0, 1, \dots, K$, i.e. for $k = 1, \dots, K$, $\mathbb{E}[\epsilon_k | \mathcal{F}_{k-1}] = \mathbf{0}$. Assume $\|\epsilon_k\|^2 \leq B_k^2$, $k = 1, 2, \dots, K$. Then,

$$\mathbb{P}\left(\left\|\sum_{k=1}^K \epsilon_k\right\| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2 \sum_{k=1}^K B_k^2}\right),$$

$\forall t \in \mathbb{R}$.

Proposition E.1 (Norm-Hoeffding, Sampling without Replacement). Let $\{\epsilon_1, \epsilon_2, \dots, \epsilon_n\}$ be a set of n fixed vectors in \mathbb{R}^d such that $\|\epsilon_i\|^2 \leq \sigma^2$, $\forall 1 \leq i \leq n$, for some $\sigma^2 > 0$. Let $\mathcal{J} \subseteq \{1, 2, \dots, n\}$ be a random index sets sampled without replacement from $\{1, 2, \dots, n\}$, with size $|\mathcal{J}| = k$. Then,

$$\mathbb{P}\left(\left\|\frac{1}{k} \sum_{j \in \mathcal{J}} \epsilon_j - \frac{1}{n} \sum_{j=1}^n \epsilon_j\right\| \geq t\right) \leq 3 \exp\left(-\frac{kt^2}{64\sigma^2}\right),$$

$\forall t \in \mathbb{R}$. In addition,

$$\mathbb{E}\left\|\frac{1}{k} \sum_{j \in \mathcal{I}} \epsilon_j - \frac{1}{n} \sum_{j=1}^n \epsilon_j\right\|^2 \leq \frac{16\sigma^2}{k}.$$

Proof. Firstly, we start with developing moment bounds. Let \mathcal{I} be a random index sets sampled with replacement from $\{1, 2, \dots, n\}$, independent of \mathcal{J} , with size $|\mathcal{I}| = k$. For any $p \in \mathbb{Z}_+$,

$$\begin{aligned} & \mathbb{E}\left\|\frac{1}{k} \sum_{j \in \mathcal{J}} \epsilon_j - \frac{1}{n} \sum_{j=1}^n \epsilon_j\right\|^p \\ & \leq \mathbb{E}\left\|\frac{1}{k} \sum_{j \in \mathcal{I}} \epsilon_j - \frac{1}{n} \sum_{j=1}^n \epsilon_j\right\|^p \\ & = \int_0^\infty \mathbb{P}\left(\left\|\frac{1}{k} \sum_{j \in \mathcal{I}} \epsilon_j - \frac{1}{n} \sum_{j=1}^n \epsilon_j\right\|^p \geq r\right) dr \\ & = \int_0^\infty \mathbb{P}\left(\left\|\frac{1}{k} \sum_{j \in \mathcal{I}} \epsilon_j - \frac{1}{n} \sum_{j=1}^n \epsilon_j\right\| \geq r^{1/p}\right) dr \\ & \leq \int_0^\infty 2 \exp\left(-\frac{kr^{2/p}}{8\sigma^2}\right) dr \\ & = p \cdot \left(\frac{8\sigma^2}{k}\right)^{p/2} \cdot \Gamma\left(\frac{p}{2}\right), \end{aligned}$$

²See also Theorem 3 in Pinelis (1992) and Proposition 2 in Fang et al. (2018).

where the 1st step is based on Lemma E.1 and the 4th step is based on the fact that $\|\epsilon_j - \frac{1}{n} \sum_{i=1}^n \epsilon_i\| \leq 2\sigma, \forall j$ and Lemma E.3.

Then, for any $s > 0$,

$$\begin{aligned}
& \mathbb{E} \exp \left(s \left\| \frac{1}{k} \sum_{j \in \mathcal{J}} \epsilon_j - \frac{1}{n} \sum_{j=1}^n \epsilon_j \right\| \right) \\
& \leq 1 + \sum_{p=1}^{\infty} \frac{s^p \mathbb{E} \left\| \frac{1}{k} \sum_{j \in \mathcal{J}} \epsilon_j - \frac{1}{n} \sum_{j=1}^n \epsilon_j \right\|^p}{p!} \\
& \leq 1 + \sum_{p=2}^{\infty} \frac{s^p \mathbb{E} \left\| \frac{1}{k} \sum_{j \in \mathcal{J}} \epsilon_j - \frac{1}{n} \sum_{j=1}^n \epsilon_j \right\|^p}{p!} + s \sqrt{\frac{8\pi\sigma^2}{k}} \\
& = 1 + s \sqrt{\frac{8\pi\sigma^2}{k}} + \sum_{p=1}^{\infty} \frac{\left(\frac{8\sigma^2 s^2}{k}\right)^p \cdot (2p) \cdot \Gamma(p)}{(2p)!} \\
& \quad + \sum_{p=1}^{\infty} \frac{\left(\frac{8\sigma^2 s^2}{k}\right)^{\frac{2p+1}{2}} \cdot (2p+1) \cdot \Gamma\left(p + \frac{1}{2}\right)}{(2p+1)!} \\
& = 1 + s \sqrt{\frac{8\pi\sigma^2}{k}} + 2 \sum_{p=1}^{\infty} \frac{\left(\frac{8\sigma^2 s^2}{k}\right)^p \cdot (p!)}{(2p)!} \\
& \quad + \sqrt{\frac{8s^2\sigma^2}{k}} \sum_{p=1}^{\infty} \frac{\left(\frac{8\sigma^2 s^2}{k}\right)^p \cdot \Gamma\left(p + \frac{1}{2}\right)}{(2p)!} \\
& \leq 1 + s \sqrt{\frac{8\pi\sigma^2}{k}} + \left(2 + \sqrt{\frac{8\pi s^2 \sigma^2}{k}}\right) \sum_{p=1}^{\infty} \frac{\left(\frac{8\sigma^2 s^2}{k}\right)^p \cdot (p!)}{(2p)!} \\
& \leq 1 + s \sqrt{\frac{8\pi\sigma^2}{k}} + \left(1 + \sqrt{\frac{2\pi s^2 \sigma^2}{k}}\right) \sum_{p=1}^{\infty} \frac{\left(\frac{8\sigma^2 s^2}{k}\right)^p}{p!} \\
& = 1 + s \sqrt{\frac{8\pi\sigma^2}{k}} + \left(1 + \sqrt{\frac{2\pi s^2 \sigma^2}{k}}\right) \left[\exp\left(\frac{8s^2\sigma^2}{k}\right) - 1\right] \\
& \leq \sqrt{\frac{8\pi s^2 \sigma^2}{k}} + \left(1 + \sqrt{\frac{2\pi s^2 \sigma^2}{k}}\right) \exp\left(\frac{8s^2\sigma^2}{k}\right) \\
& \leq \exp\left(\frac{8s^2\sigma^2}{k}\right) + 2\exp\left(\frac{16s^2\sigma^2}{k}\right) \\
& \leq 3\exp\left(\frac{16s^2\sigma^2}{k}\right),
\end{aligned}$$

where the 1st step is based on Taylor's expansion and the second to the last step is based on the fact that $x \leq e^{\frac{x^2}{\pi}}, \sqrt{\frac{\pi}{2}} x e^{x^2} \leq e^{2x^2}, \forall x \geq 0$.

For any $s > 0$,

$$\begin{aligned}
& \mathbb{P}\left(\left\|\frac{1}{k}\sum_{j\in\mathcal{J}}\epsilon_j - \frac{1}{n}\sum_{j=1}^n\epsilon_j\right\|\geq t\right) \\
& \leq \mathbb{E}\exp\left(s\left\|\frac{1}{k}\sum_{j\in\mathcal{J}}\epsilon_j - \frac{1}{n}\sum_{j=1}^n\epsilon_j\right\| - ts\right) \\
& \leq 3\exp\left(\frac{16s^2\sigma^2}{k} - ts\right).
\end{aligned} \tag{15}$$

By letting $s = \frac{kt}{32\sigma^2}$ in (15),

$$\mathbb{P}\left(\left\|\frac{1}{k}\sum_{j\in\mathcal{J}}\epsilon_j - \frac{1}{n}\sum_{j=1}^n\epsilon_j\right\|\geq t\right) \leq 3\exp\left(-\frac{kt^2}{64\sigma^2}\right).$$

□

Definition E.1. A random vector $\epsilon \in \mathbb{R}^d$ is (a, σ^2) -norm-subGaussian (or $nSG(a, \sigma^2)$), if $\exists a, \sigma^2 > 0$ such that

$$\mathbb{P}(\|\epsilon - \mathbb{E}\epsilon\| \geq t) \leq a \cdot \exp\left(-\frac{t^2}{2\sigma^2}\right),$$

$\forall t \in \mathbb{R}$.

Definition E.2. A sequence of random vectors $\epsilon_1, \dots, \epsilon_K \in \mathbb{R}^d$ is $(a, \{\sigma_k^2\}_{k=1}^K)$ -norm-subGaussian martingale difference sequence adapted to $\mathcal{F}_0, \mathcal{F}_1, \dots, \mathcal{F}_K$, if $\exists a, \sigma_1^2, \dots, \sigma_K^2 > 0$ such that for $k = 1, 2, \dots, K$,

$$\mathbb{E}[\epsilon | \mathcal{F}_{k-1}] = \mathbf{0}, \quad \sigma_k \in \mathcal{F}_{k-1}, \quad \epsilon_k \in \mathcal{F}_k,$$

and $\epsilon_k | \mathcal{F}_{k-1}$ is (a, σ_k^2) -norm-subGaussian.

Lemma E.4 (Corollary 8 in Jin et al. (2019)). Suppose $\epsilon_1, \dots, \epsilon_K \in \mathbb{R}^d$ is $(a, \{\sigma_k^2\}_{k=1}^K)$ -norm-subGaussian martingale difference sequence adapted to $\mathcal{F}_0, \mathcal{F}_1, \dots, \mathcal{F}_K$. Then for any fixed $\delta > 0$, and $B > b > 0$, with probability at least $1 - \delta$, either

$$\sum_{k=1}^K \sigma_k^2 \geq B$$

or,

$$\left\|\sum_{k=1}^K \epsilon_k\right\| \leq \frac{a}{2} e^{1/e} \sqrt{\max\left\{\sum_{k=1}^K \sigma_k^2, b\right\} \left(\log \frac{2d}{\delta} + \log \log \frac{B}{b}\right)}.$$

Lemma E.5. For any $\varepsilon > 0, n \in \mathbb{Z}_+$,

$$\frac{1}{T^2 \wedge (\sqrt{n}T)} \leq \varepsilon^2,$$

when $T \geq \lceil \frac{1}{\varepsilon} + \frac{1}{\sqrt{n\varepsilon^2}} \rceil$.

Proof.

$$\begin{aligned}
\frac{1}{T^2 \wedge (\sqrt{n}T)} &= \max \left\{ \frac{1}{T^2}, \frac{1}{\sqrt{n}T} \right\} \\
&\leq \max \left\{ \frac{1}{T'^2} \Big|_{T'=\frac{1}{\varepsilon}}, \frac{1}{\sqrt{n}T'} \Big|_{T'=\frac{1}{\sqrt{n\varepsilon^2}}} \right\} \\
&= \varepsilon^2.
\end{aligned}$$

□

Lemma E.6. For any $\varepsilon \in (0, e^{-1}]$, $n \in \mathbb{Z}_+$,

$$\frac{\log T}{T^2 \wedge (\sqrt{n}T)} \leq \varepsilon^2,$$

when $T \geq \left\lceil 3 \left(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon} + \frac{2}{\sqrt{n\varepsilon^2}} \log \frac{1}{\varepsilon^2} + \mathbf{1} \left\{ \frac{1}{\sqrt{n\varepsilon^2}} \log \frac{1}{\varepsilon^2} \leq \frac{\varepsilon}{6} \right\} \right) \right\rceil$.

Proof. Function $h(T) = \frac{\log T}{T^2}$ is monotonically decreasing when $T \geq \sqrt{e}$. Since $T \geq \frac{3}{\varepsilon} \log \frac{1}{\varepsilon} \geq \sqrt{e}$,

$$\begin{aligned}
\frac{\log T}{T^2} &\leq \frac{\log \left(\frac{3}{\varepsilon} \log \frac{1}{\varepsilon} \right)}{\left(\frac{3}{\varepsilon} \log \frac{1}{\varepsilon} \right)^2} = \frac{\log 3 + \log \frac{1}{\varepsilon} + \log \log \frac{1}{\varepsilon}}{9 \left(\log \frac{1}{\varepsilon} \right)^2} \varepsilon^2 \\
&\leq \frac{1 + \log 3 + 2 \log \frac{1}{\varepsilon}}{9 \left(\log \frac{1}{\varepsilon} \right)^2} \varepsilon^2 \leq \frac{3 + \log 3}{9} \varepsilon^2 \leq \varepsilon^2.
\end{aligned}$$

Define a function $\tilde{h}(T) = \frac{\log T}{\sqrt{n}T}$. It is monotonically decreasing when $T \geq e$. Thus, if $\frac{6}{\sqrt{n\varepsilon^2}} \log \frac{1}{\varepsilon^2} \geq e$, we know $T \geq e$ and consequently,

$$\begin{aligned}
\frac{\log T}{\sqrt{n}T} &\leq \frac{\log \left(\frac{6}{\sqrt{n\varepsilon^2}} \log \frac{1}{\varepsilon^2} \right)}{\sqrt{n} \left(\frac{6}{\sqrt{n\varepsilon^2}} \log \frac{1}{\varepsilon^2} \right)} \\
&= \frac{\log \frac{1}{\sqrt{n\varepsilon^2}} + \log \log \frac{1}{\varepsilon^2} + \log 6}{6 \log \frac{1}{\varepsilon^2}} \varepsilon^2 \\
&\leq \frac{\log \frac{1}{\varepsilon^2} + \log \log \frac{1}{\varepsilon^2} + \log 6}{6 \log \frac{1}{\varepsilon^2}} \varepsilon^2 \\
&\leq \frac{2 \log \frac{1}{\varepsilon^2} + 1 + \log 6}{6 \log \frac{1}{\varepsilon^2}} \varepsilon^2 \\
&\leq \frac{3 + \log 6}{6} \varepsilon^2 \\
&\leq \varepsilon^2.
\end{aligned}$$

If $\frac{6}{\sqrt{n\varepsilon^2}} \log \frac{1}{\varepsilon^2} \leq e$, $T \geq 3$. Hence,

$$\begin{aligned} \frac{\log T}{\sqrt{nT}} &\leq \frac{\log 3}{\sqrt{n3}} \leq \frac{6}{\sqrt{ne}} \log \frac{1}{\varepsilon^2} \left(\frac{\log 3}{3} \cdot \frac{e}{6} \right) \\ &\leq \frac{6}{\sqrt{ne}} \log \frac{1}{\varepsilon^2} \leq \varepsilon^2. \end{aligned}$$

Based on the above results,

$$\frac{\log T}{T^2 \wedge (\sqrt{nT})} \leq \max \left\{ \frac{\log T}{T^2}, \frac{\log T}{\sqrt{nT}} \right\} \leq \varepsilon^2.$$

□

F Proofs of Main Theorems

Proof of Proposition B.1. It is not hard to conclude that we only need to show

$$\mathbb{P} \left(\exists j \geq 1, \frac{1}{K_j} \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2 \leq \tilde{\varepsilon}^2 \right) = 1. \quad (16)$$

For simplicity, we denote $V_j \triangleq \frac{1}{K_j} \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2$. To show (16), we firstly derive the in-expectation bound on V_j , which has been covered in works like Wang et al. (2019).

With our basic assumptions, we have

$$\begin{aligned} &f(\mathbf{x}_{k+1}^{(j)}) \\ &= \tilde{f}(\text{Proj}(\mathbf{x}_k^{(j)} - \eta_j \boldsymbol{\nu}_k^{(j)}, \mathcal{D})) \\ &\leq \tilde{f}(\mathbf{x}_k^{(j)} - \eta_j \boldsymbol{\nu}_k^{(j)}) \\ &\leq \tilde{f}(\mathbf{x}_k^{(j)}) - \langle \nabla \tilde{f}(\mathbf{x}_k^{(j)}), \eta_j \boldsymbol{\nu}_k^{(j)} \rangle + \frac{L}{2} \eta_j^2 \|\boldsymbol{\nu}_k^{(j)}\|^2 \\ &= f(\mathbf{x}_k^{(j)}) - \langle \nabla f(\mathbf{x}_k^{(j)}), \eta_j \boldsymbol{\nu}_k^{(j)} \rangle + \frac{L}{2} \eta_j^2 \|\boldsymbol{\nu}_k^{(j)}\|^2 \\ &= f(\mathbf{x}_k^{(j)}) + \frac{\eta_j}{2} \|\boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)})\|^2 - \frac{\eta_j}{2} \|\nabla f(\mathbf{x}_k^{(j)})\|^2 \\ &\quad - \frac{\eta_j}{2} (1 - L\eta_j) \|\boldsymbol{\nu}_k^{(j)}\|^2, \end{aligned}$$

where the 2nd and 3rd step is based on Assumption 3.3. Then, summing the above inequality from $k = 0$ to $K_j - 1$,

$$f(\tilde{\mathbf{x}}_j) - f(\tilde{\mathbf{x}}_{j-1})$$

$$\begin{aligned}
&= f(\mathbf{x}_{K_j}^{(j)}) - f(\mathbf{x}_0^{(j)}) \\
&\leq \frac{\eta_j}{2} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 - \frac{\eta_j}{2} (1 - L\eta_j) K_j V_j.
\end{aligned} \tag{17}$$

Then,

$$\begin{aligned}
&\mathbb{E}(f(\tilde{\mathbf{x}}_j) - f(\tilde{\mathbf{x}}_{j-1}) | \mathcal{F}_{j-1}) \\
&\leq \mathbb{E} \left(\frac{\eta_j}{2} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 | \mathcal{F}_{j-1} \right) \\
&\quad - \frac{\eta_j}{2} (1 - L\eta_j) K_j \mathbb{E}(V_j | \mathcal{F}_{j-1}).
\end{aligned} \tag{18}$$

For convenience, we abbreviate $\mathbb{E}(\cdot | \mathcal{F}_{j-1})$ as $\mathbb{E}_{j-1}(\cdot)$. For $k = 1, 2, \dots, K_j - 1$,

$$\begin{aligned}
&\mathbb{E}_{j-1} \left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 \\
&= \mathbb{E}_{j-1} \mathbb{E} \left(\left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 | \mathcal{F}_{j,k-1} \right) \\
&= \mathbb{E}_{j-1} \mathbb{E} \left(\left\| \frac{1}{b_j} \sum_{i \in \mathcal{I}_k^{(j)}} \nabla f_i(\mathbf{x}_k^{(j)}) - \frac{1}{b_j} \sum_{i \in \mathcal{I}_k^{(j)}} \nabla f_i(\mathbf{x}_{k-1}^{(j)}) \right. \right. \\
&\quad \left. \left. + \nabla f(\mathbf{x}_{k-1}^{(j)}) - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 | \mathcal{F}_{j,k-1} \right) \\
&\quad + \mathbb{E}_{j-1} \left\| \boldsymbol{\nu}_{k-1}^{(j)} - \nabla f(\mathbf{x}_{k-1}^{(j)}) \right\|^2 \\
&= \mathbb{E}_{j-1} \frac{1}{b_j^2} \sum_{i \in \mathcal{I}_k^{(j)}} \mathbb{E} \left(\left\| \nabla f_i(\mathbf{x}_k^{(j)}) - \frac{1}{b_j} \sum_{i \in \mathcal{I}_k^{(j)}} \nabla f_i(\mathbf{x}_{k-1}^{(j)}) \right. \right. \\
&\quad \left. \left. + \nabla f(\mathbf{x}_{k-1}^{(j)}) - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 | \mathcal{F}_{j,k-1} \right) \\
&\quad + \mathbb{E}_{j-1} \left\| \boldsymbol{\nu}_{k-1}^{(j)} - \nabla f(\mathbf{x}_{k-1}^{(j)}) \right\|^2 \\
&\leq \frac{4L^2}{b_j} \mathbb{E}_{j-1} \left\| \mathbf{x}_k^{(j)} - \mathbf{x}_{k-1}^{(j)} \right\|^2 + \mathbb{E}_{j-1} \left\| \boldsymbol{\nu}_{k-1}^{(j)} - \nabla f(\mathbf{x}_{k-1}^{(j)}) \right\|^2 \\
&\leq \frac{4L^2 \eta_j^2}{b_j} \mathbb{E}_{j-1} \left\| \boldsymbol{\nu}_{k-1}^{(j)} \right\|^2 + \mathbb{E}_{j-1} \left\| \boldsymbol{\nu}_{k-1}^{(j)} - \nabla f(\mathbf{x}_{k-1}^{(j)}) \right\|^2 \\
&\leq \frac{4L^2 \eta_j^2}{b_j} \sum_{t=0}^{k-1} \mathbb{E}_{j-1} \left\| \boldsymbol{\nu}_t^{(j)} \right\|^2 + \mathbb{E}_{j-1} \left\| \boldsymbol{\nu}_0^{(j)} - \nabla f(\mathbf{x}_0^{(j)}) \right\|^2 \\
&\leq \frac{4L^2 \eta_j^2 K_j}{b_j} \mathbb{E}_{j-1} V_j + \mathbb{E}_{j-1} \left\| \boldsymbol{\nu}_0^{(j)} - \nabla f(\mathbf{x}_0^{(j)}) \right\|^2.
\end{aligned} \tag{19}$$

Based on (18) and (19),

$$\begin{aligned}
& \mathbb{E}(f(\tilde{\mathbf{x}}_j) - f(\tilde{\mathbf{x}}_{j-1})) \\
& \leq -\frac{\eta_j K_j}{2} \left(1 - L\eta_j - \frac{4L^2 \eta_j^2 K_j}{b_j} \right) \mathbb{E}V_j \\
& \quad + \frac{\eta_j K_j}{2} \mathbb{E} \left\| \boldsymbol{\nu}_0^{(j)} - \nabla f(\mathbf{x}_0^{(j)}) \right\|^2 \\
& \leq -\frac{K_j}{16L} \mathbb{E}V_j + \frac{K_j}{8L} \mathbb{E} \left\| \boldsymbol{\nu}_0^{(j)} - \nabla f(\mathbf{x}_0^{(j)}) \right\|^2, \tag{20}
\end{aligned}$$

where the second step is based on the choice of $\eta_j \equiv \frac{1}{4L}$ and $b_j \geq K_j$, $j \geq 1$. Let us define $J_0 = \min\{j : B_j = n\}$. Then, for $j \geq J_0$, based on (20),

$$\frac{1}{16L} \mathbb{E}V_j \leq \frac{K_j}{16L} \mathbb{E}V_j \leq \mathbb{E}(f(\tilde{\mathbf{x}}_{j-1}) - f(\tilde{\mathbf{x}}_j)).$$

Then for any $m \in \mathbb{Z}_+$,

$$\begin{aligned}
& \mathbb{P}(V_j > \tilde{\epsilon}^2, j \geq 1) \\
& \leq \mathbb{P}(V_{J_0} + V_{J_0+1} + \dots + V_{J_0+m} > (m+1)\tilde{\epsilon}^2) \\
& \leq \frac{\mathbb{E}(V_{J_0} + V_{J_0+1} + \dots + V_{J_0+m})}{(m+1)\tilde{\epsilon}^2} \\
& \leq \frac{16L}{(m+1)\tilde{\epsilon}^2} \sum_{j=J_0}^{J_0+m} \mathbb{E}(f(\tilde{\mathbf{x}}_{j-1}) - f(\tilde{\mathbf{x}}_j)) \\
& \leq \frac{16L\Delta_f}{(m+1)\tilde{\epsilon}^2}.
\end{aligned}$$

Since m can be arbitrarily large, we know

$$\mathbb{P}(V_j > \tilde{\epsilon}^2, j \geq 1) = 0,$$

which can directly lead to (16). □

Proof of Theorem 3.2. In this proof, for simplicity, we denote $\mathbb{E}[\cdot | \mathcal{F}_k]$ by $\mathbb{E}_k[\cdot]$. Let $\mathbf{s}_k = \sum_{i=1}^k \mathbf{z}_i$, $k \geq 1$.

For a $1 \leq k \leq K$, consider

$$f_k(t) = \mathbb{E}_{k-1} [\cosh(\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|)], \quad \lambda > 0, t > 0.$$

Then,

$$f'_k(t) = \frac{1}{2} \mathbb{E}_{k-1} \left[\frac{\lambda \langle \mathbf{z}_k, \mathbf{s}_{k-1} + t\mathbf{z}_k \rangle}{\|\mathbf{s}_{k-1} + t\mathbf{z}_k\|} \left(e^{\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|} - e^{-\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|} \right) \right],$$

and consequently,

$$\begin{aligned} f'_k(0) &= \frac{1}{2} \mathbb{E}_{k-1} \left[\frac{\lambda \langle \mathbf{z}_k, \mathbf{s}_{k-1} \rangle}{\|\mathbf{s}_{k-1}\|} \left(e^{\lambda \|\mathbf{s}_{k-1}\|} - e^{-\lambda \|\mathbf{s}_{k-1}\|} \right) \right] \\ &= 0. \end{aligned}$$

Next,

$$\begin{aligned} &f''_k(t) \\ &= \frac{1}{2} \mathbb{E}_{k-1} \left[\left(\frac{\lambda^2 \langle \mathbf{z}_k, \mathbf{s}_{k-1} + t\mathbf{z}_k \rangle^2}{\|\mathbf{s}_{k-1} + t\mathbf{z}_k\|^2} + \frac{\lambda \|\mathbf{z}_k\|^2}{\|\mathbf{s}_{k-1} + t\mathbf{z}_k\|} \right) e^{\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|} \right. \\ &\quad \left. + \left(\frac{\lambda^2 \langle \mathbf{z}_k, \mathbf{s}_{k-1} + t\mathbf{z}_k \rangle^2}{\|\mathbf{s}_{k-1} + t\mathbf{z}_k\|^2} - \frac{\lambda \|\mathbf{z}_k\|^2}{\|\mathbf{s}_{k-1} + t\mathbf{z}_k\|} \right) e^{-\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|} \right] \\ &= \mathbb{E}_{k-1} \left[\frac{\lambda^2 \langle \mathbf{z}_k, \mathbf{s}_{k-1} + t\mathbf{z}_k \rangle^2}{\|\mathbf{s}_{k-1} + t\mathbf{z}_k\|^2} \cosh(\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|) \right. \\ &\quad \left. + \frac{\lambda^2 \|\mathbf{z}_k\|^2}{\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|} \sinh(\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|) \right] \\ &\leq \mathbb{E}_{k-1} \left[\left(\frac{\lambda^2 \langle \mathbf{z}_k, \mathbf{s}_{k-1} + t\mathbf{z}_k \rangle^2}{\|\mathbf{s}_{k-1} + t\mathbf{z}_k\|^2} + \lambda^2 \|\mathbf{z}_k\|^2 \right) \cosh(\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|) \right] \\ &\leq 2\lambda^2 \mathbb{E}_{k-1} [\|\mathbf{z}_k\|^2 \cosh(\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|)] \\ &\leq 2\lambda^2 r_k^2 \mathbb{E}_{k-1} [\cosh(\lambda \|\mathbf{s}_{k-1} + t\mathbf{z}_k\|)] \\ &= 2\lambda^2 r_k^2 f_k(t), \end{aligned}$$

where the first inequality is based on the fact that if $y > 0$, $\frac{\sinh(y)}{y} \leq \cosh(y)$.

According to Lemma 3 in Pinelis (1992),

$$f_k(t) \leq f_k(0) \exp(\lambda^2 r_k^2 t^2) = \cosh(\lambda \|\mathbf{s}_{k-1}\|) \exp(\lambda^2 r_k^2 t^2).$$

Thus,

$$\begin{aligned} &\mathbb{E}_{k-1} [\cosh(\lambda \|\mathbf{s}_k\|)] \\ &= f_k(1) \leq \cosh(\lambda \|\mathbf{s}_{k-1}\|) \exp(\lambda^2 r_k^2 t^2). \end{aligned} \tag{21}$$

Now, let

$$G_k = \cosh(\lambda \|\mathbf{s}_k\|) \exp\left(-\lambda^2 \sum_{i=1}^k r_i^2\right), \quad k = 1, 2, \dots, K.$$

We can easily know that for $k = 1, 2, \dots, K$, G_k is measurable with respect to \mathcal{F}_k . According to (21),

$$\begin{aligned} \mathbb{E}_{k-1} G_k &= \exp\left(-\lambda^2 \sum_{i=1}^k r_i^2\right) \mathbb{E}_{k-1} [\cosh(\lambda \|\mathbf{s}_k\|)] \\ &\leq \cosh(\lambda \|\mathbf{s}_{k-1}\|) \exp\left(-\lambda^2 \sum_{i=1}^{k-1} r_i^2\right) \\ &= G_{k-1}, \end{aligned}$$

which implies that $\{G_k\}_{k=1}^K$ is a non-negative super-martingale adapted to $\mathcal{F}_0, \mathcal{F}_1, \dots, \mathcal{F}_K$.

For any constant $m > 0$, if we define stopping time $T_m = \inf\left\{t : \|\mathbf{s}_t\| \geq \lambda \sum_{i=1}^t r_i^2 + m\right\}$, we immediately know that $G_{T_m \wedge k}$, $k \geq 0$, is a supermartingale and

$$\begin{aligned} &\mathbb{P}\left(\exists 1 \leq t \leq k, \|\mathbf{s}_t\| \geq \lambda \sum_{i=1}^t r_i^2 + m\right) \\ &= \mathbb{P}\left(\|\mathbf{s}_{T_m}\| \geq \lambda \sum_{i=1}^{T_m} r_i^2 + m, 1 \leq T_m \leq k\right) \\ &= \mathbb{P}\left(\|\mathbf{s}_{T_m \wedge k}\| \geq \lambda \sum_{i=1}^{T_m \wedge k} r_i^2 + m, 1 \leq T_m \leq k\right) \\ &\leq \mathbb{P}\left(G_{T_m \wedge k} \geq \exp\left(-\lambda^2 \sum_{i=1}^{T_m \wedge k} r_i^2\right) \cosh\left(\lambda^2 \sum_{i=1}^{T_m \wedge k} r_i^2 + m\lambda\right)\right) \\ &\leq \mathbb{P}\left(G_{T_m \wedge k} \geq \frac{1}{2} \exp\left(-\lambda^2 \sum_{i=1}^{T_m \wedge k} r_i^2 + (\lambda^2 \sum_{i=1}^{T_m \wedge k} r_i^2 + m\lambda)\right)\right) \\ &= \mathbb{P}\left(2G_{T_m \wedge k} \geq e^{\lambda m}\right) \\ &\leq \frac{2\mathbb{E}G_{T_m \wedge k}}{e^{\lambda m}} \\ &\leq 2e^{-\lambda m} \mathbb{E}G_0 \\ &= 2e^{-\lambda m}, \end{aligned}$$

where the 2nd step is based on the fact that $\cosh(y) \geq \frac{1}{2}e^y, \forall y \in \mathbb{R}$, the 4th step is by Chebyshev's inequality and the 5th step is based on the supermartingale property.

Therefore, if we let $\lambda m = \log \frac{2}{\delta}$,

$$\mathbb{P}\left(\exists 1 \leq t \leq k, \|\mathbf{s}_t\| \geq \lambda \sum_{i=1}^t r_i^2 + \frac{1}{\lambda} \log \frac{2}{\delta}\right) \leq \delta.$$

Since k can be up to K ,

$$\mathbb{P}\left(\exists 1 \leq t \leq K, \|\mathbf{s}_t\| \geq \lambda \sum_{i=1}^t r_i^2 + \frac{1}{\lambda} \log \frac{2}{\delta}\right) \leq \delta.$$

The final conclusion can be obtained immediately by following similar steps given in the proof of Corollary 8 from Jin et al. (2019). \square

Proof of Proposition D.1. Recall that

$$\boldsymbol{\epsilon}_0^{(j)} = \boldsymbol{\nu}_0^{(j)} - \nabla f(\mathbf{x}_0^{(j)}) = \frac{1}{B_j} \sum_{i \in \mathcal{I}_j} \nabla f_i(\mathbf{x}_0^{(j)}) - \nabla f(\mathbf{x}_0^{(j)}),$$

where \mathcal{I}_j is sampled without replacement. Since $\|\nabla f_i(\mathbf{x}_0^{(j)})\| \leq \alpha_M$, $i = 1, 2, \dots, n$, based on Proposition E.1,

$$\mathbb{P}\left(\|\boldsymbol{\epsilon}_0^{(j)}\| \geq t | \mathcal{F}_{j,-1}\right) \leq 3 \exp\left(-\frac{B_j t^2}{64 \alpha_M^2}\right) \mathbf{1}\{B_j < n\}. \quad (22)$$

Next, if we suppose $\mathcal{I}_m^{(j)} = \{i_{m,1}^{(j)}, i_{m,2}^{(j)}, \dots, i_{m,b_j}^{(j)}\}$, where $i_{m,t_1}^{(j)} \neq i_{m,t_2}^{(j)}$ for any $1 \leq t_1 < t_2 \leq b_j$, we have

$$\begin{aligned} & \boldsymbol{\epsilon}_m^{(j)} \\ &= \frac{1}{b_j} \sum_{i \in \mathcal{I}_m^{(j)}} \left[\nabla f_i(\mathbf{x}_m^{(j)}) - \nabla f(\mathbf{x}_m^{(j)}) + \nabla f(\mathbf{x}_{m-1}^{(j)}) - \nabla f_i(\mathbf{x}_{m-1}^{(j)}) \right] \\ &= \sum_{r=1}^{b_j} \frac{1}{b_j} \left[\nabla f_{i_{m,r}^{(j)}}(\mathbf{x}_m^{(j)}) - \nabla f(\mathbf{x}_m^{(j)}) + \nabla f(\mathbf{x}_{m-1}^{(j)}) - \nabla f_{i_{m,r}^{(j)}}(\mathbf{x}_{m-1}^{(j)}) \right] \\ &\triangleq \sum_{r=1}^{b_j} \boldsymbol{\rho}_{(m-1)b_j+r}^{(j)}. \end{aligned}$$

Let

$$\tilde{\mathcal{F}}_0^{(j)} = \mathcal{F}_{j,0}$$

and

$$\tilde{\mathcal{F}}_{a_1 b_j + a_2}^{(j)} = \sigma\left(\tilde{\mathcal{F}}_{a_1 b_j + a_2 - 1}^{(j)} \cup \sigma(i_{a_1+1, a_2}^{(j)})\right)$$

for $a_1 = 0, 1, 2, \dots$ and $a_2 = 1, 2, \dots, b_j$. Then, we can see that $\{\boldsymbol{\rho}_s^{(j)}\}_{s=1}^{k b_j}$ is a martingale difference sequence adapted to $\{\tilde{\mathcal{F}}_s^{(j)}\}_{s=0}^{k b_j}$.

Notice that for $m = 1, 2, \dots, k$ and $r = 1, 2, \dots, b_j$,

$$\begin{aligned} & \|\boldsymbol{\rho}_{(m-1)b_j+r}^{(j)}\| \\ &= \left\| \frac{1}{b_j} \left[\nabla f_{i_{m,r}^{(j)}}(\mathbf{x}_m^{(j)}) - \nabla f(\mathbf{x}_m^{(j)}) + \nabla f(\mathbf{x}_{m-1}^{(j)}) - \nabla f_{i_{m,r}^{(j)}}(\mathbf{x}_{m-1}^{(j)}) \right] \right\| \end{aligned}$$

$$\leq \frac{2L}{b_j} \left\| \mathbf{x}_m^{(j)} - \mathbf{x}_{m-1}^{(j)} \right\|.$$

Therefore, based on Theorem 3.2, for any fixed $\delta' > 0$, $B > b > 0$, with probability at least $1 - \delta'$, either

$$\begin{aligned} (\sigma_k^{(j)})^2 &\triangleq \sum_{m=1}^k \sum_{r=1}^{b_j} \left(\frac{2L}{b_j} \left\| \mathbf{x}_m^{(j)} - \mathbf{x}_{m-1}^{(j)} \right\| \right)^2 \\ &= \frac{4L^2}{b_j} \sum_{m=1}^k \left\| \mathbf{x}_m^{(j)} - \mathbf{x}_{m-1}^{(j)} \right\|^2 \\ &\geq B, \end{aligned}$$

or

$$\begin{aligned} &\left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) - \boldsymbol{\epsilon}_0^{(j)} \right\|^2 \\ &= \left\| \sum_{s=1}^{kb_j} \boldsymbol{\rho}_s^{(j)} \right\|^2 \\ &\leq 9 \max \left\{ (\sigma_k^{(j)})^2, b \right\} \left(\log \frac{2}{\delta'} + \log \log \frac{B}{b} \right). \end{aligned}$$

Under the compact constraint,

$$(\sigma_k^{(j)})^2 \leq \frac{4L^2 d_1^2 k}{b_j}.$$

Thus, if we let $B = \frac{8L^2 d_1^2 k}{b_j}$ and $b = \frac{4L^2 \tau k}{b_j}$ for some $\tau \in (0, 1)$, it would be of probability 0 to have $(\sigma_k^{(j)})^2 \geq B$.

Thus, with probability at least $1 - \delta'$,

$$\begin{aligned} &\left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) - \boldsymbol{\epsilon}_0^{(j)} \right\|^2 \\ &\leq 9 \left((\sigma_k^{(j)})^2 + \frac{4L^2 \tau k}{b_j} \right) \left(\log \frac{2}{\delta'} + \log \log \frac{2d_1^2}{\tau} \right) \\ &\leq 9 \left((\tilde{\sigma}_k^{(j)})^2 + \frac{4L^2 \tau k}{b_j} \right) \left(\log \frac{2}{\delta'} + \log \log \frac{2d_1^2}{\tau} \right). \end{aligned} \tag{23}$$

According to (22), with probability at least $1 - \delta'$,

$$\left\| \boldsymbol{\epsilon}_0^{(j)} \right\|^2 \leq \frac{64\alpha_M^2}{B_j} \log \frac{3}{\delta'} \mathbf{1} \{B_j < n\}. \tag{24}$$

Thus, combining (23) and (24), with probability at least $1 - 2\delta'$,

$$\left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2$$

$$\begin{aligned} &\leq 18 \left((\tilde{\sigma}_k^{(j)})^2 + \frac{4L^2\tau k}{b_j} \right) \left(\log \frac{2}{\delta'} + \log \log \frac{2d_1^2}{\tau} \right) \\ &\quad + \frac{128\alpha_M^2}{B_j} \log \frac{3}{\delta'} \mathbf{1}\{B_j < n\}. \end{aligned}$$

□

Proposition F.1 (Inner Loop Analysis). *Given Assumptions 3.2, 3.3 and 3.4, under the parameter setting given in Theorem 3.1, let $\Omega = \bigcup_{j=1}^{\infty} \Omega_j$, where the definition of Ω_j is given in (12). On Ω ,*

$$\begin{aligned} &f(\mathbf{x}_{K_j}^{(j)}) - f(\mathbf{x}_0^{(j)}) \\ &\leq -\frac{1}{16L} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + \frac{L^2\eta_j\tau_j l_j K_j^2}{b_j} + \frac{\eta_j K_j q_j}{2}, \end{aligned}$$

for all $j \in \mathbb{Z}_+$. Such event Ω occurs with probability at least $1 - \delta$.

Proof of Proposition F.1. Firstly, as what we have shown in section 3, $\bigcup_{j=0}^{\infty} \Omega_j$ occurs with probability at least $1 - \delta$.

Based on (17), on $\bigcup_{j=0}^{\infty} \Omega_j$,

$$\begin{aligned} &f(\mathbf{x}_{K_j}^{(j)}) - f(\mathbf{x}_0^{(j)}) \\ &\leq \frac{\eta_j}{2} \sum_{k=0}^{K_j-1} \left[l_j \left((\tilde{\sigma}_k^{(j)})^2 + \frac{4L^2\tau_j k}{b_j} \right) + q_j \right] \\ &\quad - \frac{\eta_j}{2} (1 - L\eta_j) \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 \\ &\leq \frac{\eta_j}{2} \sum_{k=0}^{K_j-1} \left(\frac{4L^2\eta_j^2 l_j}{b_j} \sum_{m=1}^k \left\| \boldsymbol{\nu}_{m-1}^{(j)} \right\|^2 \right) + \frac{2L^2\eta_j\tau_j l_j K_j^2}{b_j} \\ &\quad + \frac{\eta_j K_j q_j}{2} - \frac{\eta_j}{2} (1 - L\eta_j) \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 \\ &\leq \frac{2L^2\eta_j^3 l_j K_j}{b_j} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + \frac{L^2\eta_j\tau_j l_j K_j^2}{b_j} + \frac{\eta_j K_j q_j}{2} \\ &\quad - \frac{\eta_j}{2} (1 - L\eta_j) \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 \\ &= -\frac{\eta_j}{2} \left(1 - L\eta_j - \frac{4L^2\eta_j^2 l_j K_j}{b_j} \right) \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 \\ &\quad + \frac{L^2\eta_j\tau_j l_j K_j^2}{b_j} + \frac{\eta_j K_j q_j}{2} \end{aligned}$$

$$\begin{aligned}
&= -\frac{1}{8L} \left(1 - \frac{1}{4} - \frac{1}{4}\right) \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + \frac{L^2 \eta_j \tau_j l_j K_j^2}{b_j} + \frac{\eta_j K_j q_j}{2} \\
&= -\frac{1}{16L} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + \frac{L^2 \eta_j \tau_j l_j K_j^2}{b_j} + \frac{\eta_j K_j q_j}{2},
\end{aligned}$$

where the 5th step is based on our choices of $\eta_j = \frac{1}{4L}$ and $b_j = l_j K_j$, $j = 1, 2, \dots$ □

Proof of Proposition D.3. Firstly,

$$\begin{aligned}
-\Delta_f &\leq f(\tilde{\mathbf{x}}_{2T}) - f(\tilde{\mathbf{x}}_T) = f(\mathbf{x}_{K_{2T}}^{(2T)}) - f(\mathbf{x}_0^{(T+1)}) \\
&\leq \sum_{j=T+1}^{2T} \left[\frac{-1}{16L} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + \frac{L^2 \eta_j \tau_j l_j K_j^2}{b_j} + \frac{\eta_j K_j q_j}{2} \right] \\
&= \sum_{j=T+1}^{2T} \left[\frac{L^2 \eta_j \tau_j l_j K_j^2}{b_j} + \frac{\eta_j K_j q_j}{2} \right] - \frac{1}{16L} \sum_{j=T+1}^{2T} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2,
\end{aligned} \tag{25}$$

where the 3rd step is based on Proposition F.1.

For simplifying notations, we denote

$$A_T \triangleq \sum_{j=T+1}^{2T} \left[\frac{L^2 \eta_j \tau_j l_j K_j^2}{b_j} + \frac{\eta_j K_j q_j}{2} \right].$$

Then,

$$\begin{aligned}
&A_T \\
&= \sum_{j=T+1}^{2T} \left(\frac{L \tau_j K_j}{4} + \frac{K_j q_j}{8L} \right) \\
&= \sum_{j=T+1}^{2T} \left(\frac{L \sqrt{j^2 \wedge n}}{4j^3} + \frac{16\alpha_M^2}{L \sqrt{j^2 \wedge n}} \log \frac{12C_e j^4}{\delta} \mathbf{1}\{j^2 < n\} \right) \\
&\leq \sum_{j=T+1}^{2T} \left(\frac{L}{4j^2} + \frac{16\alpha_M^2}{L \sqrt{j^2 \wedge n}} \log \frac{12C_e j^4}{\delta} \mathbf{1}\{j^2 < n\} \right) \\
&\leq \sum_{j=T+1}^{2T} \left(\frac{L}{4j^2} + \frac{16\alpha_M^2}{Lj} \log \frac{12C_e j^4}{\delta} \right) \\
&\leq \frac{C_e L}{4} + \sum_{j=T+1}^{2T} \frac{16\alpha_M^2}{Lj} \log \frac{12C_e j^4}{\delta} \\
&\leq \frac{C_e L}{4} + \sum_{j=T+1}^{2T} \frac{16\alpha_M^2}{LT} \log \frac{12C_e (2T)^4}{\delta} \\
&= \frac{C_e L}{4} + \frac{16\alpha_M^2}{L} \log \frac{192C_e T^4}{\delta}
\end{aligned} \tag{26}$$

$$= \frac{C_e L}{4} + \frac{16\alpha_M^2}{L} \log \frac{192C_e}{\delta} + \frac{64\alpha_M^2}{L} \log T,$$

where the 1st step is based on the choices that $\eta_j = \frac{1}{4L}$ and $b_j = l_j K_j$, the second step is based on the choices of $K_j = \sqrt{B_j} = \sqrt{j^2 \wedge n}$, $\delta'_j = \frac{\delta}{4C_e j^4}$. According to Lemma E.5, as $T \geq T_1$,

$$\frac{1}{T^2 \wedge (\sqrt{n}T)} \leq \frac{\varepsilon^2}{320L(c_1 + \Delta_f)}. \quad (27)$$

According to Lemma E.6, as $T \geq T_2$,

$$\frac{\log T}{T^2 \wedge (\sqrt{n}T)} \leq \frac{\varepsilon^2}{320Lc_2}. \quad (28)$$

If we suppose to the contrary that

$$\frac{1}{K_j} \sum_{k=0}^{K_j-1} \left\| \nu_k^{(j)} \right\|^2 > \tilde{\varepsilon}^2$$

holds for all $T+1 \leq j \leq 2T$, then we have

$$\frac{1}{16L} \sum_{j=T+1}^{2T} \sum_{k=0}^{K_j-1} \left\| \nu_k^{(j)} \right\|^2 \geq \frac{\tilde{\varepsilon}^2}{16L} \sum_{j=T+1}^{2T} K_j \geq \frac{\tilde{\varepsilon}^2}{16L} \sum_{j=T+1}^{2T} (T \wedge \sqrt{n}) = \frac{\tilde{\varepsilon}^2}{16L} T^2 \wedge (\sqrt{n}T).$$

By (26), (27), (28) and the above results,

$$\begin{aligned} & \frac{80L}{T^2 \wedge (\sqrt{n}T)} \left\{ \Delta_f + A_T - \frac{1}{16L} \sum_{j=T+1}^{2T} \sum_{k=0}^{K_j-1} \left\| \nu_k^{(j)} \right\|^2 \right\} \\ & \leq \frac{80L}{T^2 \wedge (\sqrt{n}T)} (\Delta_f + A_T) - 5\tilde{\varepsilon}^2 \\ & = \frac{80L}{T^2 \wedge (\sqrt{n}T)} (\Delta_f + A_T) - \varepsilon^2 \\ & \leq \frac{80L}{T^2 \wedge (\sqrt{n}T)} (\Delta_f + c_1) + \frac{80Lc_2 \log T}{T^2 \wedge (\sqrt{n}T)} - \varepsilon^2 \\ & \leq \frac{\varepsilon^2}{4} + \frac{\varepsilon^2}{4} - \varepsilon^2 \\ & = -\frac{\varepsilon^2}{2}, \end{aligned}$$

which contradicts (25). □

Proof of Proposition D.4.

$$\begin{aligned} & \varepsilon_T \\ & = 8L^2 \tau_T + 2q_T \end{aligned}$$

$$\begin{aligned}
&= \frac{8L^2}{T^3} + \frac{256\alpha_M^2}{B_T} \log \frac{3}{\delta'_T} \mathbf{1}\{B_T < n\} \\
&\leq \frac{8L^2}{T^3} + \frac{256\alpha_M^2}{T^2} \log \frac{3}{\delta'_T} \tag{29}
\end{aligned}$$

$$\begin{aligned}
&= \frac{8L^2}{T^3} + \frac{256\alpha_M^2}{T^2} \log \frac{12C_e T^4}{\delta} \\
&\leq \left(8L^2 + 256\alpha_M^2 \log \frac{12C_e}{\delta} \right) \frac{1}{T^2} + \frac{1024\alpha_M^2}{T^2} \log T \\
&= c_3 \frac{1}{T^2} + c_4 \frac{\log T}{T^2}, \tag{30}
\end{aligned}$$

where the 2nd step is based on our choice of $\tau_T = \frac{1}{T^3}$ and the 4th step is based on the choice of $\delta'_T = \frac{\delta}{4C_e T^4}$. According to Lemma E.5, where we can simply let $n = \infty$, as $T \geq T_3$,

$$\frac{1}{T^2} \leq \frac{\varepsilon^2}{4c_3}. \tag{31}$$

Similarly, according to Lemma E.6 and Assumption 3.4, as $T \geq T_4$,

$$\frac{\log T}{T^2} \leq \frac{\varepsilon^2}{4c_4}. \tag{32}$$

Combining (29), (31) and (32),

$$\varepsilon_T \leq \frac{\varepsilon^2}{2}.$$

□

Proof of Corollary C.1. This part follows a similar way as the complexity analysis in Horváth et al. (2020). It is easy to know that if Algorithm 1 stops in T outer iterations, the first order computational complexity is

$$\tilde{O}_{L, \Delta_f, \alpha_M} (T^3 \wedge (nT)).$$

Thus, it is sufficient to show

$$T_i^3 \wedge (nT_i) = \tilde{O}_{L, \Delta_f, \alpha_M} \left(\frac{1}{\varepsilon^3} \wedge \frac{\sqrt{n}}{\varepsilon^2} \right), \quad i = 1, 2, 3, 4.$$

- $T_1^3 \wedge (nT_1)$

For simplicity, we let $\tilde{c}_1 = \sqrt{320L(c_1 + \Delta_f)}$ and consequently $T_1 = \left\lceil \frac{\tilde{c}_1}{\varepsilon} + \frac{\tilde{c}_1^2}{\sqrt{n\varepsilon^2}} \right\rceil$.

When $\sqrt{n}\varepsilon \leq \tilde{c}_1$, $\frac{\tilde{c}_1}{\varepsilon} \leq \frac{\tilde{c}_1^2}{\sqrt{n}\varepsilon^2}$ and consequently $T = \mathcal{O}\left(\frac{\tilde{c}_1^2}{\sqrt{n}\varepsilon^2}\right)$. Hence,

$$T_1^3 \wedge (nT_1) = \mathcal{O}(nT_1) = \mathcal{O}\left(\frac{\sqrt{n}\tilde{c}_1^2}{\varepsilon^2}\right) = \mathcal{O}\left(\frac{\sqrt{n}\tilde{c}_1^2}{\varepsilon^2} \wedge \frac{\tilde{c}_1^3}{\varepsilon^3}\right),$$

where the last step is due to $\sqrt{n} \leq \frac{\tilde{c}_1}{\varepsilon}$.

When $\sqrt{n}\varepsilon \geq \tilde{c}_1$, $\frac{\tilde{c}_1}{\varepsilon} \geq \frac{\tilde{c}_1^2}{\sqrt{n}\varepsilon^2}$ and consequently $T = \mathcal{O}\left(\frac{\tilde{c}_1}{\varepsilon}\right)$. Hence,

$$T_1^3 \wedge (nT_1) = \mathcal{O}(T_1^3) = \mathcal{O}\left(\frac{\tilde{c}_1^3}{\varepsilon^3}\right) = \mathcal{O}\left(\frac{\sqrt{n}\tilde{c}_1^2}{\varepsilon^2} \wedge \frac{\tilde{c}_1^3}{\varepsilon^3}\right),$$

where the last step is due to $\tilde{c}_1 \leq \sqrt{n}\varepsilon$.

To sum up,

$$T_1^3 \wedge (nT_1) = \mathcal{O}\left(\frac{\sqrt{n}\tilde{c}_1^2}{\varepsilon^2} \wedge \frac{\tilde{c}_1^3}{\varepsilon^3}\right) = \tilde{\mathcal{O}}_{L, \Delta_f, \alpha_M}\left(\frac{1}{\varepsilon^3} \wedge \frac{\sqrt{n}}{\varepsilon^2}\right).$$

• $T_2^3 \wedge (nT_2)$

Secondly, if we let $\tilde{c}_2 = \sqrt{320Lc_2}$, we have $\tilde{c}_2 \geq 4$ based on Assumption 3.4. As a result, $T_2 = \Theta\left(\frac{3\tilde{c}_2}{\varepsilon} \log \frac{3\tilde{c}_2}{\varepsilon} + \frac{2\tilde{c}_2^2}{\sqrt{n}\varepsilon^2} \log \frac{\tilde{c}_2^2}{\varepsilon^2}\right)$. Therefore, it is equivalent to study $\bar{T}_2^3 \wedge (n\bar{T}_2)$ where $\bar{T}_2 = \frac{3\tilde{c}_2}{\varepsilon} \log \frac{3\tilde{c}_2}{\varepsilon} + \frac{2\tilde{c}_2^2}{\sqrt{n}\varepsilon^2} \log \frac{\tilde{c}_2^2}{\varepsilon^2}$.

When $\tilde{c}_2 \geq 1.5\sqrt{n}\varepsilon$,

$$\frac{3\tilde{c}_2}{\varepsilon} \log \frac{3\tilde{c}_2}{\varepsilon} \leq \frac{3\tilde{c}_2}{\varepsilon} \log \frac{\tilde{c}_2^2}{\varepsilon^2} \leq \frac{2\tilde{c}_2^2}{\sqrt{n}\varepsilon^2} \log \frac{\tilde{c}_2^2}{\varepsilon^2}.$$

Thus, $\bar{T}_2 = \mathcal{O}\left(\frac{2\tilde{c}_2^2}{\sqrt{n}\varepsilon^2} \log \frac{\tilde{c}_2^2}{\varepsilon^2}\right)$. Then,

$$\begin{aligned} \bar{T}_2^3 \wedge (n\bar{T}_2) &= \mathcal{O}(n\bar{T}_2) \\ &= \mathcal{O}\left(\frac{2\sqrt{n}\tilde{c}_2^2}{\varepsilon^2} \log \frac{\tilde{c}_2^2}{\varepsilon^2}\right) = \mathcal{O}\left(\frac{2\sqrt{n}\tilde{c}_2^2}{\varepsilon^2} \log \frac{\tilde{c}_2^2}{\varepsilon^2} \wedge \frac{4\tilde{c}_2^3}{3\varepsilon^3} \log \frac{\tilde{c}_2^2}{\varepsilon^2}\right). \end{aligned}$$

When $\tilde{c}_2 \leq 1.5\sqrt{n}\varepsilon$,

$$\begin{aligned} \frac{2\tilde{c}_2^2}{\sqrt{n}\varepsilon^2} \log \frac{\tilde{c}_2^2}{\varepsilon^2} &= \frac{4\tilde{c}_2^2}{\sqrt{n}\varepsilon^2} \log \frac{\tilde{c}_2}{\varepsilon} \leq \frac{4\tilde{c}_2^2}{\sqrt{n}\varepsilon^2} \log \frac{3\tilde{c}_2}{\varepsilon} \\ &\leq \frac{1.5\varepsilon}{\tilde{c}_2} \cdot \frac{4\tilde{c}_2^2}{\varepsilon^2} \log \frac{3\tilde{c}_2}{\varepsilon} = \frac{6\tilde{c}_2}{\varepsilon} \log \frac{3\tilde{c}_2}{\varepsilon}. \end{aligned}$$

Thus, $\bar{T}_2 = \mathcal{O}\left(\frac{\tilde{c}_2}{\varepsilon} \log \frac{3\tilde{c}_2}{\varepsilon}\right)$. Then

$$\bar{T}_2^3 \wedge (n\bar{T}_2) = \mathcal{O}(\bar{T}_2^3) = \mathcal{O}\left(\frac{\tilde{c}_2^3}{\varepsilon^3} \left(\log \frac{3\tilde{c}_2}{\varepsilon}\right)^3\right)$$

$$= \mathcal{O} \left(\frac{\tilde{c}_2^3}{\varepsilon^3} \left(\log \frac{3\tilde{c}_2}{\varepsilon} \right)^3 \wedge \frac{1.5\sqrt{n}\tilde{c}_2^2}{\varepsilon^2} \left(\log \frac{3\tilde{c}_2}{\varepsilon} \right)^3 \right).$$

To sum up,

$$T_2^3 \wedge (nT_2) = \tilde{\mathcal{O}}_{L,\Delta_f,\alpha_M} \left(\frac{1}{\varepsilon^3} \wedge \frac{\sqrt{n}}{\varepsilon^2} \right).$$

- $T_3^3 \wedge (nT_3)$

Since $T_3 = \tilde{\Theta}_{L,\Delta_f,\alpha_M} \left(\frac{1}{\varepsilon} \right)$, we can directly know that

$$T_3^3 \wedge (nT_3) = \tilde{\mathcal{O}}_{L,\Delta_f,\alpha_M} \left(\frac{1}{\varepsilon^3} \wedge \frac{n}{\varepsilon} \right) = \tilde{\mathcal{O}}_{L,\Delta_f,\alpha_M} \left(\frac{1}{\varepsilon^3} \wedge \frac{\sqrt{n}}{\varepsilon^2} \right).$$

- $T_4^3 \wedge (nT_4)$

Similar to the previous case,

$$T_4^3 \wedge (nT_4) = \tilde{\mathcal{O}}_{L,\Delta_f,\alpha_M} \left(\frac{1}{\varepsilon^3} \wedge \frac{\sqrt{n}}{\varepsilon^2} \right).$$

□

Proof of Theorem C.2. We can see that many results given under the setting of Theorem 3.1 can still apply under the current setting. If we still define Ω_j as (12), $\Omega = \bigcup_{j=1}^{\infty} \Omega_j$ occurs with probability at least $1 - \delta$.

Under the current setting, Proposition F.1 is still valid. Thus, on Ω , for any $j \in \mathbb{Z}_+$,

$$\begin{aligned} & f(\mathbf{x}_{K_j}^{(j)}) - f(\mathbf{x}_0^{(j)}) \\ & \leq -\frac{1}{16L} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + \frac{L^2 \eta_j \tau_j l_j K_j^2}{b_j} + \frac{\eta_j K_j q_j}{2} \\ & = -\frac{1}{16L} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + \frac{L^2 \eta_j \tau_j l_j K_j^2}{b_j} \\ & = -\frac{1}{16L} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + \frac{L \tau_j K_j}{4} \\ & = -\frac{1}{16L} \sum_{k=0}^{K_j-1} \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + \frac{\sqrt{n}L}{4} \tau_j, \end{aligned}$$

where the 2nd step is due to our choice of $B_j \equiv n$ and consequently $q_j \equiv 0$, the 3rd step is based on our choices of $\eta_j = \frac{1}{4L}$ and $b_j = l_j K_j$, the 4th step is based on our choice of $K_j = n$. Summing the above

inequality from $j = 1$ to T ,

$$\begin{aligned}
& -\Delta_f^0 \\
&= f(\mathbf{x}^*) - f(\mathbf{x}_0^{(1)}) \\
&\leq f(\mathbf{x}_{K_T}^{(T)}) - f(\mathbf{x}_0^{(1)}) \\
&= \sum_{j=1}^T \left(\frac{\sqrt{n}L}{4} \tau_j - \frac{1}{16L} \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2 \right) \\
&= \frac{\sqrt{n}T\tilde{\varepsilon}^2}{32L} - \frac{1}{16L} \sum_{j=1}^T \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2, \tag{33}
\end{aligned}$$

where the 2nd step is according to Assumption 3.1, the 4th step is based on our choice of $\tau_j \equiv \frac{\tilde{\varepsilon}^2}{8L^2}$. We assert that when $T \geq T_5$, there must exist a $1 \leq j \leq T$ such that

$$\frac{1}{K_j} \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2 \leq \tilde{\varepsilon}^2.$$

If not,

$$\begin{aligned}
& \frac{\sqrt{n}T\tilde{\varepsilon}^2}{32L} - \frac{1}{16L} \sum_{j=1}^T \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2 \\
&= \frac{\sqrt{n}T\tilde{\varepsilon}^2}{32L} - \frac{1}{16L} \sum_{j=1}^T \frac{K_j}{K_j} \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2 \\
&\leq \frac{\sqrt{n}T\tilde{\varepsilon}^2}{32L} - \frac{1}{16L} \sum_{j=1}^T \tilde{\varepsilon}^2 K_j \\
&= \frac{\sqrt{n}T\tilde{\varepsilon}^2}{32L} - \frac{\sqrt{n}\tilde{\varepsilon}^2 T}{16L} \\
&= -\frac{\sqrt{n}T\tilde{\varepsilon}^2}{32L} \\
&= -\frac{\sqrt{n}\tilde{\varepsilon}^2 T}{160L} \\
&\leq -(\Delta_f^0 + 1),
\end{aligned}$$

which is in conflict with (33). Thus, on Ω , the first stopping rule will be met in at most T outer iterations while the second stopping rule is always satisfied. When both stopping rules are met, we can show that the output is of desirable property. Let $1 \leq j \leq T$ and $0 \leq k \leq K_j$ such that

$$\frac{1}{K_j} \sum_{k=0}^{K_j-1} \|\boldsymbol{\nu}_k^{(j)}\|^2 \leq \tilde{\varepsilon}^2$$

and

$$\left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 \leq \tilde{\varepsilon}^2.$$

Then, on Ω ,

$$\begin{aligned}
& \|\nabla f(\hat{\mathbf{x}})\|^2 \\
&= \left\| \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 \\
&\leq 2 \left\| \boldsymbol{\nu}_k^{(j)} \right\|^2 + 2 \left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 \\
&\leq 2\tilde{\varepsilon}^2 + 2 \left\| \boldsymbol{\nu}_k^{(j)} - \nabla f(\mathbf{x}_k^{(j)}) \right\|^2 \\
&\leq 2\tilde{\varepsilon}^2 + 2l_j \left(\frac{4L^2\eta_j^2}{b_j} \sum_{m=1}^k \left\| \boldsymbol{\nu}_{m-1}^{(j)} \right\|^2 + \frac{4L^2\tau_j k}{b_j} \right) \\
&\leq 2\tilde{\varepsilon}^2 + 2l_j \left(\frac{4L^2\eta_j^2}{b_j} \sum_{m=1}^{K_j} \left\| \boldsymbol{\nu}_{m-1}^{(j)} \right\|^2 + \frac{4L^2\tau_j K_j}{b_j} \right) \\
&\leq 2\tilde{\varepsilon}^2 + 2l_j \left(\frac{4L^2\eta_j^2 K_j \tilde{\varepsilon}^2}{b_j} + \frac{4L^2\tau_j K_j}{b_j} \right) \\
&= 2\tilde{\varepsilon}^2 + 0.5\tilde{\varepsilon}^2 + 8L^2\tau_j \\
&= 3.5\tilde{\varepsilon}^2 \\
&\leq \varepsilon^2,
\end{aligned}$$

where the 4th step is based on Proposition [D.1](#), the 7th step is based on our choices of $\eta_j = \frac{1}{4L}$ and $b_j = l_j K_j$, the 8th step is based on our choice of $\tau_j \equiv \frac{\varepsilon^2}{8L^2}$. \square

G Supplementary Figures

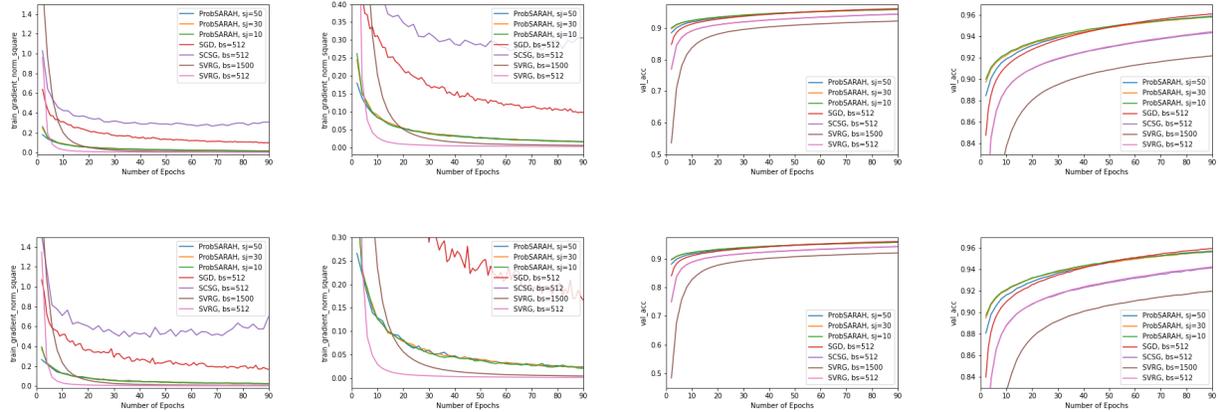


Figure 3: Comparison of convergence with respect to $(1 - \delta)$ -quantile of square of gradient norm ($\|\nabla f\|^2$) and δ -quantile of validation accuracy on the **MNIST** dataset for $\delta = 0.1$ and $\delta = 0.01$. The second (fourth) column presents zoom-in figures of those in the first (third) column. Top: $\delta = 0.1$. Bottom: $\delta = 0.01$. 'bs' stands for batch size. 'sj=x' means that the smallest batch size $\approx x \log x$.

References

- Allen-Zhu, Z. and E. Hazan (2016). “Variance reduction for faster non-convex optimization”. In: *International conference on machine learning*. PMLR, pp. 699–707.
- Bardenet, R. and O.-A. Maillard (2015). “Concentration inequalities for sampling without replacement”. In: *Bernoulli* 21.3, pp. 1361–1385.
- Boucheron, S., G. Lugosi, and P. Massart (2013). *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press.
- Defazio, A., F. Bach, and S. Lacoste-Julien (2014). “SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives”. In: *Advances in neural information processing systems*, pp. 1646–1654.
- Fang, C., C. J. Li, Z. Lin, and T. Zhang (2018). “SPIDER: near-optimal non-convex optimization via stochastic path integrated differential estimator”. In: *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 687–697.
- Ghadimi, S. and G. Lan (2013). “Stochastic first-and zeroth-order methods for nonconvex stochastic programming”. In: *SIAM Journal on Optimization* 23.4, pp. 2341–2368.
- Goodfellow, I., Y. Bengio, and A. Courville (2016). *Deep learning*. MIT press.

- Harvey, N. J., C. Liaw, Y. Plan, and S. Randhawa (2019a). “Tight analyses for non-smooth stochastic gradient descent”. In: *Conference on Learning Theory*. PMLR, pp. 1579–1613.
- Harvey, N. J., C. Liaw, and S. Randhawa (2019b). “Simple and optimal high-probability bounds for strongly-convex stochastic gradient descent”. In: *arXiv preprint arXiv:1909.00843*.
- Hoeffding, W. (1963). “Probability Inequalities for Sums of Bounded Random Variables”. In: *Journal of the American Statistical Association* 58.301, pp. 13–30.
- Horváth, S., L. Lei, P. Richtárik, and M. I. Jordan (2020). “Adaptivity of stochastic gradient methods for nonconvex optimization”. In: *arXiv preprint arXiv:2002.05359*.
- Jain, P., D. Nagaraj, and P. Netrapalli (2019). “Making the last iterate of sgd information theoretically optimal”. In: *Conference on Learning Theory*. PMLR, pp. 1752–1755.
- James, G., D. Witten, T. Hastie, and R. Tibshirani (2013). *An introduction to statistical learning*. Vol. 112. Springer.
- Ji, K., Z. Wang, B. Weng, Y. Zhou, W. Zhang, and Y. Liang (2020). “History-gradient aided batch size adaptation for variance reduced algorithms”. In: *International Conference on Machine Learning*. PMLR, pp. 4762–4772.
- Jin, C., P. Netrapalli, R. Ge, S. M. Kakade, and M. I. Jordan (2019). “A short note on concentration inequalities for random vectors with subgaussian norm”. In: *arXiv preprint arXiv:1902.03736*.
- Johnson, R. and T. Zhang (2013). “Accelerating stochastic gradient descent using predictive variance reduction”. In: *Advances in neural information processing systems* 26, pp. 315–323.
- Kakade, S. M. and A. Tewari (2009). “On the generalization ability of online strongly convex programming algorithms”. In: *Advances in Neural Information Processing Systems*, pp. 801–808.
- Le Roux, N., M. Schmidt, and F. Bach (2012). “A stochastic gradient method with an exponential convergence rate for finite training sets”. In: *Proceedings of the 25th International Conference on Neural Information Processing Systems-Volume 2*, pp. 2663–2671.
- Lei, L. and M. Jordan (2017). “Less than a single pass: Stochastically controlled stochastic gradient”. In: *Artificial Intelligence and Statistics*. PMLR, pp. 148–156.
- Lei, L. and M. I. Jordan (2020). “On the adaptivity of stochastic gradient-based optimization”. In: *SIAM Journal on Optimization* 30.2, pp. 1473–1500.
- Lei, L., C. Ju, J. Chen, and M. I. Jordan (2017). “Non-convex finite-sum optimization via SCSG methods”. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 2345–2355.
- Li, X. and F. Orabona (2020). “A high probability analysis of adaptive sgd with momentum”. In: *arXiv preprint arXiv:2007.14294*.

- Li, Z. (2019). “SSRGD: Simple Stochastic Recursive Gradient Descent for Escaping Saddle Points”. In: *Advances in Neural Information Processing Systems 32*, pp. 1523–1533.
- Li, Z., H. Bao, X. Zhang, and P. Richtárik (2021). “PAGE: A simple and optimal probabilistic gradient estimator for nonconvex optimization”. In: *International Conference on Machine Learning*. PMLR, pp. 6286–6295.
- Li, Z. and J. Li (2018). “A simple proximal stochastic gradient method for nonsmooth nonconvex optimization”. In: *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 5569–5579.
- Nesterov, Y. (2003). *Introductory lectures on convex optimization: A basic course*. Vol. 87. Springer Science & Business Media.
- Nguyen, L. M., J. Liu, K. Scheinberg, and M. Takáč (2017a). “SARAH: A novel method for machine learning problems using stochastic recursive gradient”. In: *International Conference on Machine Learning*. PMLR, pp. 2613–2621.
- Nguyen, L. M., J. Liu, K. Scheinberg, and M. Takáč (2017b). “Stochastic recursive gradient algorithm for nonconvex optimization”. In: *arXiv preprint arXiv:1705.07261*.
- Pinelis, I. (1992). “An approach to inequalities for the distributions of infinite-dimensional martingales”. In: *Probability in Banach Spaces, 8: Proceedings of the Eighth International Conference*. Springer, pp. 128–134.
- (1994). “Optimum bounds for the distributions of martingales in Banach spaces”. In: *The Annals of Probability*, pp. 1679–1706.
- Reddi, S. J., A. Hefny, S. Sra, B. Póczos, and A. Smola (2016). “Stochastic variance reduction for nonconvex optimization”. In: *International conference on machine learning*. PMLR, pp. 314–323.
- Tran-Dinh, Q., N. H. Pham, D. T. Phan, and L. M. Nguyen (2019). “Hybrid stochastic gradient descent algorithms for stochastic nonconvex optimization”. In: *arXiv preprint arXiv:1905.05920*.
- Wang, Z., K. Ji, Y. Zhou, Y. Liang, and V. Tarokh (2019). “Spiderboost and momentum: Faster variance reduction algorithms”. In: *Advances in Neural Information Processing Systems 32*, pp. 2406–2416.
- Zhou, D., J. Chen, Y. Cao, Y. Tang, Z. Yang, and Q. Gu (2018). “On the convergence of adaptive gradient methods for nonconvex optimization”. In: *arXiv preprint arXiv:1808.05671*.