

# On Biased Stochastic Gradient Estimation

**Derek Driggs**

D.DRIGGS@MATHS.CAM.AC.UK

*Department of Applied Mathematics and Theoretical Physics  
University of Cambridge  
Cambridge, CB3 0WA, UK*

**Jingwei Liang\***

JINGWEI.LIANG@SJTU.EDU.CN

*Institute of Natural Sciences and School of Mathematical Sciences  
Shanghai Jiao Tong University  
Shanghai, 200240, China*

**Carola-Bibiane Schönlieb**

CBS31@CAM.AC.UK

*Department of Applied Mathematics and Theoretical Physics  
University of Cambridge  
Cambridge, CB3 0WA, UK*

**Editor:** Julien Mairal

## Abstract

We present a uniform analysis of biased stochastic gradient methods for minimizing convex, strongly convex, and non-convex composite objectives, and identify settings where bias is useful in stochastic gradient estimation. The framework we present allows us to extend proximal support to biased algorithms, including SAG and SARAH, for the first time in the convex setting. We also use our framework to develop a new algorithm, Stochastic Average Recursive GradiEnt (SARGE), that achieves the oracle complexity lower-bound for non-convex, finite-sum objectives and requires strictly fewer calls to a stochastic gradient oracle per iteration than SVRG and SARAH. We support our theoretical results with numerical experiments that demonstrate the benefits of certain biased gradient estimators.

**Keywords:** stochastic gradient descent, variance reduction, biased gradient estimation

## 1. Introduction

In this paper, we focus on the following composite minimization problem:

$$\min_{x \in \mathbb{R}^p} \{F(x) \stackrel{\text{def}}{=} f(x) + g(x)\}. \quad (1)$$

Throughout, we assume:

- $g$  is proper and closed such that its proximity operator (see (3) in Section 2) is well-defined,
- $f$  admits a finite-sum structure  $f(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(x)$ , and for all  $i \in \{1, 2, \dots, n\}$ ,  $\nabla f_i$  is  $L$ -Lipschitz continuous for some  $L > 0$ .

We consider three settings: the convex setting, where all of  $\{f_i\}_{i=1}^n$  and  $g$  are convex; the strongly convex setting, where additionally  $g$  is strongly convex; and the non-convex setting, where  $\{f_i\}_{i=1}^n$  and  $g$  are not necessarily convex.

---

\*. Corresponding author

Problems of this form arise frequently in many areas of science and engineering, such as machine learning, statistics, operations research, and imaging. For instance, in machine learning, these problems often arise as empirical risk minimization problems from classification and regression tasks. Examples include ridge regression, logistic regression, Lasso, and  $\ell_1$ -regularized logistic regression (Bishop, 2006). Principal component analysis (PCA) can also be formulated as a problem with this structure, where the functions  $f_i$  are non-convex (Garber and Hazan, 2015; Allen-Zhu and Yuan, 2018). In imaging,  $\ell_1$  or total variation regularization is often combined with differentiable data discrepancy terms that appear in both convex and non-convex instances (Bredies and Lorenz, 2018).

### 1.1 Stochastic gradient methods

Two classical approaches to solve (1) are the proximal gradient descent method (PGD) (Lions and Mercier, 1979) and its accelerated variants, including inertial PGD (Liang et al., 2017) and FISTA (Beck and Teboulle, 2009). For these deterministic approaches, the full gradient of  $f$  must be evaluated at each iteration, which often requires huge computational resources when  $n$  is large. Such a drawback makes these schemes unsuitable for large-scale machine learning tasks.

By exploiting the finite sum structure of  $f$ , stochastic gradient methods enjoy low per-iteration complexity while achieving comparable convergence rates. These qualities make stochastic gradient methods the standard approach to solving many problems in machine learning, and are gaining popularity in other areas such as image processing (Chambolle et al., 2018). Stochastic gradient descent (SGD) was first proposed in the 1950’s (Robbins and Monro, 1951) and has experienced a renaissance in the past decade, with numerous variants of SGD proposed in the literature (see, for instance (Schmidt et al., 2017; Johnson and Zhang, 2013; Defazio et al., 2014a) and references therein). Most of these algorithms can be summarized into one general form, which is described below in Algorithm 1.

---

**Algorithm 1** Stochastic gradient descent framework

---

**Input:** starting point  $x_0 \in \mathbb{R}^p$ , gradient estimator  $\tilde{\nabla}$ .

- 1: **for**  $k = 0, 1, \dots, T - 1$  **do**
- 2:     Compute the stochastic gradient estimate  $\tilde{\nabla}_k$  at the current iterate  $x_k$ .
- 3:     Choose the step size/learning rate  $\eta_k$ .
- 4:     Update  $x_{k+1}$ :

$$x_{k+1} \leftarrow \text{prox}_{\eta_k g}(x_k - \eta_k \tilde{\nabla}_k). \tag{2}$$

- 5: **end for**
- 

Below we summarize several popular stochastic gradient estimators  $\tilde{\nabla}_k$ :

- **SGD** Classic stochastic gradient descent (Robbins and Monro, 1951) uses the following gradient estimator at iteration  $k$ :

$$\left[ \begin{array}{l} \text{sample } j_k \text{ uniformly at random from } \{1, \dots, n\}, \\ \tilde{\nabla}_k^{\text{SGD}} = \nabla f_{j_k}(x_k). \end{array} \right.$$

At each step, SGD uses the gradient of the sampled function  $\nabla f_{j_k}(x_k)$  as a stochastic approximation of the full gradient  $\nabla f(x_k)$ . It is an unbiased estimate as  $\mathbb{E}_k[\nabla f_{j_k}(x_k)] = \nabla f(x_k)$ . It is also *memoryless*: every update of  $x_{k+1}$  depends only upon  $x_k$  and the random variable  $j_k$ . The variance of SGD is does not vanish as the algorithm converges.

- **SAG** To deal with the non-vanishing variance of SGD, in (Roux et al., 2012; Schmidt et al., 2017) the authors introduce the SAG gradient estimator, which uses the gradient history to decrease its variance. With  $\nabla f_i(\varphi_0^i) = 0, i = 1, \dots, n$ , the SAG gradient estimator is computed using the following procedure:

$$\left[ \begin{array}{l} \text{sample } j_k \text{ uniformly at random from } \{1, \dots, n\}, \\ \tilde{\nabla}_k^{\text{SAG}} = \frac{1}{n}(\nabla f_{j_k}(x_k) - \nabla f_{j_k}(\varphi_k^{j_k})) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i), \\ \text{update the gradient history: } \nabla f_i(\varphi_{k+1}^i) = \begin{cases} \nabla f_i(x_k) & \text{if } i = j_k, \\ \nabla f_i(\varphi_k^i) & \text{o.w.} \end{cases} \end{array} \right.$$

Here, for each  $i \in \{1, \dots, n\}$ ,  $\nabla f_i(\varphi_k^i)$  is a stored gradient of  $\nabla f_i$  from a previous iteration. With the help of memory, the variance of the SAG gradient estimator diminishes as the algorithm converges. Estimators that satisfy this property are known as *variance-reduced* estimators. In contrast to the SGD estimator,  $\tilde{\nabla}_k^{\text{SAG}}$  is a *biased* estimate of  $\nabla f(x_k)$ .

- **SAGA** Building on (Roux et al., 2012; Schmidt et al., 2017), Defazio et al. (2014a) propose the unbiased gradient estimator SAGA, which is computed using the procedure below.

$$\left[ \begin{array}{l} \text{Sample } j_k \text{ uniformly at random from } \{1, \dots, n\}, \\ \tilde{\nabla}_k^{\text{SAGA}} = \nabla f_{j_k}(x_k) - \nabla f_{j_k}(\varphi_k^{j_k}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i), \\ \text{update the gradient history : } \nabla f_i(\varphi_{k+1}^i) = \begin{cases} \nabla f_i(x_k) & \text{if } i = j_k, \\ \nabla f_i(\varphi_k^i) & \text{o.w.} \end{cases} \end{array} \right.$$

Compared to  $\tilde{\nabla}^{\text{SAG}}$ , the SAGA estimator gives less weight to stored gradients. With this adjustment,  $\tilde{\nabla}^{\text{SAGA}}$  is unbiased while maintaining the variance reduction property. Similar gradient estimators can be found in Point-SAGA (Defazio, 2016), Finito (Defazio et al., 2014b), MISO (Mairal, 2014), SDCA (Shalev-Shwartz and Zhang, 2013), and those in (Hofmann et al., 2015).

- **SVRG** Another popular variance-reduced estimator is SVRG (Johnson and Zhang, 2013). The SVRG gradient estimator is computed as follows:

$$\left[ \begin{array}{l} \text{For } s = 0, \dots, S \\ \nabla f(\varphi_s) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_s), \\ \text{For } k = 1, \dots, m \\ \left[ \begin{array}{l} \text{Sample } j_k \text{ uniformly at random from } \{1, \dots, n\}, \\ \tilde{\nabla}_k^{\text{SVRG}} = \nabla f_{j_k}(x_k) - \nabla f_{j_k}(\varphi_s) + \nabla f(\varphi_s), \end{array} \right. \end{array} \right.$$

where  $\varphi_s$  is a “snapshot” point updated every  $m$  steps. The algorithms prox-SVRG (Xiao and Zhang, 2014), Natasha (Allen-Zhu, 2017), Katyusha (Allen-Zhu, 2018a),

KatyushaX (Allen-Zhu, 2018b), Natasha2 (Allen-Zhu, 2018c), MiG (Zhou et al., 2018), ASVRG (Shang et al., 2018), and VARAG (Lan et al., 2019) use the SVRG gradient estimator.

- **SARAH** In (Nguyen et al., 2017) the authors propose a recursive modification to SVRG.

$$\left\{ \begin{array}{l} \text{For } s = 0, \dots, S \\ \tilde{\nabla}_{k-1}^{\text{SARAH}} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_s), \\ \text{For } k = 1, \dots, m \\ \left[ \begin{array}{l} \text{Sample } j_k \text{ uniformly at random from } \{1, \dots, n\}, \\ \tilde{\nabla}_k^{\text{SARAH}} = \nabla f_{j_k}(x_k) - \nabla f_{j_k}(x_{k-1}) + \tilde{\nabla}_{k-1}^{\text{SARAH}}, \end{array} \right. \end{array} \right.$$

Like SAG, SARAH is a biased gradient estimator. It is also used in prox-SARAH (Pham et al., 2019), SPIDER (Fang et al., 2018), SPIDERBoost (Wang et al., 2018) and SPIDER-M (Wang et al., 2019).

We refer to algorithms employing (un)biased gradient estimators as (un)biased stochastic algorithms, respectively. The body of work on biased algorithms is stunted compared to the enormous literature on unbiased algorithms, leading to several gaps in the development of biased stochastic gradient methods. We list a few below.

- **Complex convergence proofs.** It is often difficult to analyze biased stochastic algorithms. The convergence proof of the biased algorithm SAG is especially complex, requiring computational verification (Roux et al., 2012; Schmidt et al., 2017).
- **Sub-optimal convergence rates.** In the convex setting with  $g \equiv 0$ , SARAH achieves a complexity bound of  $\mathcal{O}(\frac{\log(1/\epsilon)}{\epsilon})$  (Nguyen et al., 2017) for finding a point  $\bar{x}_k$  such that  $\mathbb{E}[F(\bar{x}_k) - F(x^*)] \leq \epsilon$ . In comparison, SAGA and SVRG achieve a complexity bound of  $\mathcal{O}(1/\epsilon)$  which is the same as deterministic proximal gradient descent.
- **Lack of proximal support.** Bias also makes it difficult to handle non-smooth functions. To the best of our knowledge, there are no theoretical guarantees for biased algorithms to solve (1) with  $g \not\equiv 0$  that take advantage of convexity when it is present.

Despite the above shortcomings, there are notable exceptions that suggest biased algorithms are worth further consideration. Recently, (Pham et al., 2019; Fang et al., 2018; Wang et al., 2018, 2019) proved that algorithms using the SARAH gradient estimator require  $\mathcal{O}(\sqrt{n}/\epsilon^2)$  stochastic gradient evaluations to find an  $\epsilon$ -first-order stationary point. This matches the complexity lower-bound for non-convex, finite-sum optimization for smooth functions  $f_i$  and  $n \leq \mathcal{O}(\epsilon^{-4})$  (Fang et al., 2018). For comparison, the best complexity bound obtained for SAGA and SVRG in this setting is  $\mathcal{O}(n^{2/3}/\epsilon^2)$  (Reddi et al., 2016a; Allen-Zhu and Hazan, 2016), and this performance requires large mini-batches of size  $\mathcal{O}(n^{2/3})$ .

A detailed summary of existing complexity bounds for the variance-reduced gradient estimators mentioned above is provided in Tables 1 and 2 for convex and non-convex objectives, respectively. The complexity bound for SAGA (Defazio et al., 2014a), SVRG (Johnson and Zhang, 2013) and SAG (Schmidt et al., 2017) on convex objectives is  $\mathcal{O}(\frac{n+L}{\epsilon})$ , which can be improved to  $\mathcal{O}((n+\kappa)\log(1/\epsilon))$  when strong convexity is present. For smooth objectives (where  $g \equiv 0$ ), this complexity is  $\mathcal{O}(\frac{\sqrt{nL}}{\epsilon})$  for SAGA and SVRG (Reddi et al., 2016a), while SAG requires  $g \equiv 0$  for any convergence results. Convergence results for SARAH

	Convex	Convex and $g \equiv 0$	Strongly Convex	Prox-Support?
SAGA	$\mathcal{O}(\frac{n+L}{\epsilon})$	$\mathcal{O}(\frac{\sqrt{nL}}{\epsilon})^*$	$\mathcal{O}((n + \kappa) \log(1/\epsilon))$	Yes
SVRG	$\mathcal{O}(\frac{n+L}{\epsilon})$	$\mathcal{O}(\frac{\sqrt{nL}}{\epsilon})^*$	$\mathcal{O}((n + \kappa) \log(1/\epsilon))$	Yes
SAG	<b>None</b>	$\mathcal{O}(\frac{n+L}{\epsilon})$	$\mathcal{O}((n + \kappa) \log(1/\epsilon))$	<b>No</b>
SARAH	<b>None</b>	$\mathcal{O}(\frac{n+L \log(1/\epsilon)}{\epsilon})^*$	$\mathcal{O}((n + \kappa) \log(1/\epsilon))$	<b>No</b>

Table 1: Existing complexity bounds for stochastic gradient estimators for convex objectives. Complexities with a “\*” represent the number of stochastic gradient oracle calls required to find an  $\epsilon$ -approximate stationary point (as in Definition 7). The other complexities are for finding a point satisfying  $\mathbb{E}[F(\bar{x}_k) - F(x^*)] \leq \epsilon$  in the convex case and  $\mathbb{E}[\|x_k - x^*\|^2] \leq \epsilon$  in the strongly convex case. The parameter  $\mu$  is the strong convexity constant, and  $\kappa = L/\mu$  is the condition number.

on convex objectives also require  $g \equiv 0$ , and the proven complexity is worse than similar results for SAGA, SVRG, and SAG by a logarithmic factor (Nguyen et al., 2017).

There are several accelerated algorithms that achieve better convergence rates than those in Table 1. SVRG++ achieves a complexity of  $\mathcal{O}(n \log(1/\epsilon) + L/\epsilon)$  on convex objectives using an epoch-doubling procedure (Allen-Zhu and Yuan, 2018). Katyusha, an accelerated variant of SVRG, has complexities of  $\mathcal{O}(n \log(1/\epsilon) + \sqrt{nL/\epsilon})$  on convex objectives and  $\mathcal{O}(n + \sqrt{n\kappa} \log(1/\epsilon))$  with strong convexity. Combining a variance-reduced algorithm with the Catalyst acceleration scheme produces algorithms with the same convergence rates up to logarithmic factors (Lin et al., 2015). These accelerated algorithms are not directly comparable to the non-accelerated algorithms in this paper, so we leave these rates out of Table 1. The algorithms considered in this work can be accelerated using momentum schemes as well, and this is the subject of a related work (Driggs et al., 2020).

On non-convex objectives, SAGA and SVRG achieve complexities of  $\mathcal{O}(\frac{nL}{\epsilon^2})$ , and this rate can be improved to  $\mathcal{O}(\frac{n^{2/3}L}{\epsilon^2})$  using large mini-batches of size  $\mathcal{O}(n^{2/3})$  (Reddi et al., 2016b). Although we do not consider mini-batching in this work, using large mini-batch sizes could similarly improve the presented complexities for B-SAGA and B-SVRG. SAG has not been previously analyzed in the non-convex setting, so this work presents the first convergence results for SAG in this setting as a special case of our results for B-SAGA. SARAH achieves the oracle complexity lower-bound of  $\mathcal{O}(\frac{\sqrt{nL}}{\epsilon^2})$  (Pham et al., 2019).

Table 3 summarises the convergence rates provided in this work. Our results provide proximal support to biased algorithms such as SARAH and SAG for the first time, prove state-of-the-art convergence rates for all algorithms in the non-convex setting, and improve the best-known convergence rates for SARAH on convex objectives. These strong results for non-smooth, non-convex problems and biased estimators comes at the cost of recovering suboptimal convergence rates for SAGA and SVRG on convex problems.

## 1.2 Contributions

This work provides three main contributions:

	Non-Convex, No Mini-Batching	With Mini-Batching	Prox-Support?
SAGA	$\mathcal{O}(\frac{nL}{\epsilon^2})$	$\mathcal{O}(\frac{n^{2/3}L}{\epsilon^2})$	Yes
SVRG	$\mathcal{O}(\frac{nL}{\epsilon^2})$	$\mathcal{O}(\frac{n^{2/3}L}{\epsilon^2})$	Yes
SAG	<b>None</b>	<b>None</b>	<b>No</b>
SARAH	$\mathcal{O}(\frac{\sqrt{n}L}{\epsilon^2})$	$\mathcal{O}(\frac{\sqrt{n}L}{\epsilon^2})$	Yes

Table 2: Existing complexity bounds for stochastic gradient estimators for non-convex optimization. These complexities represent the number of stochastic gradient oracle calls required to find an  $\epsilon$ -approximate stationary point (as in Definition 7). Using mini-batches of size  $n^{2/3}$  optimizes the complexity of SAGA and SVRG in this setting. While we do not consider mini-batching in this work, this improvement likely extends to B-SAGA and B-SVRG as well.

1. We introduce a framework for the systematic analysis of a large class of stochastic gradient methods and investigate a bias-variance tradeoff arising from our analysis. Our analysis provides proximal support to biased algorithms for the first time in the convex setting.
2. We apply our framework to derive convergence rates for SARAH and biased versions of SAGA and SVRG on convex, strongly convex, and non-convex problems.
3. We design a new recursive gradient estimator, Stochastic Average Recursive GradiEnt (SARGE), that achieves the same convergence rates as SARAH but never computes a full gradient, giving it a strictly smaller per-iteration complexity than SARAH. In particular, we show that SARGE achieves the oracle complexity lower bound for non-convex finite-sum optimization.

To study the effects of bias on the SAGA and SVRG estimators, we introduce Biased SAGA (B-SAGA) and Biased SVRG (B-SVRG). For  $\theta > 0$ , these two gradient estimators read

- B-SAGA:  $\tilde{\nabla}_k^{\text{B-SAGA}} \stackrel{\text{def}}{=} \frac{1}{\theta}(\nabla f_{j_k}(x_k) - \nabla f_{j_k}(\varphi_k^{j_k})) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i)$ .
- B-SVRG:  $\tilde{\nabla}_k^{\text{B-SVRG}} \stackrel{\text{def}}{=} \frac{1}{\theta}(\nabla f_{j_k}(x_k) - \nabla f_{j_k}(\varphi_s)) + \nabla f(\varphi_s)$ .

In both B-SAGA and B-SVRG, the bias parameter  $\theta$  adjusts how much weight is given to stored gradient information. When  $\theta = n$ ,  $\tilde{\nabla}_k^{\text{B-SAGA}}$  recovers the SAG gradient estimator.

Motivated by the desirable properties of SARAH, we propose a new gradient estimator, Stochastic Average Recursive GradiEnt (SARGE), which is defined below

$$\tilde{\nabla}_k^{\text{SARGE}} \stackrel{\text{def}}{=} \nabla f_{j_k}(x_k) - \psi_k^{j_k} + \frac{1}{n} \sum_{i=1}^n \psi_k^i - (1 - \frac{1}{n})(\nabla f_{j_k}(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}),$$

where the variables  $\psi_k^i$  follow the update rule  $\psi_{k+1}^{j_k} = \nabla f_{j_k}(x_k) - (1 - \frac{1}{n})\nabla f_{j_k}(x_{k-1})$ . Similar to SAGA, SARGE uses stored gradient information to avoid having to compute the full gradient, a computational burden that SVRG and SARAH require for variance reduction.

A summary of the complexity results obtained from our analysis for SAG/B-SAGA, B-SVRG, SARAH, and SARGE are provided in Table 3. Note that the result for SAG is included in B-SAGA.

	Convex	Strongly Convex	Non-Convex	Prox-Support?
B-SAGA	$\mathcal{O}(\frac{nL}{\epsilon})$	$\mathcal{O}(n\kappa \log(1/\epsilon))$	$\mathcal{O}(\frac{nL}{\epsilon^2})$	Yes
B-SVRG	$\mathcal{O}(\frac{nL}{\epsilon})$	$\mathcal{O}(n\kappa \log(1/\epsilon))$	$\mathcal{O}(\frac{nL}{\epsilon^2})$	Yes
SARAH	$\mathcal{O}(\frac{\sqrt{n}L + \sqrt{Ln}^{3/4}}{\epsilon})$	$\mathcal{O}(\max\{\sqrt{n}\kappa, n\} \log(1/\epsilon))$	$\mathcal{O}(\frac{\sqrt{n}L}{\epsilon^2})$	Yes
SARGE	$\mathcal{O}(\frac{\sqrt{n}L + \sqrt{Ln}^{3/4}}{\epsilon})$	$\mathcal{O}(\max\{\sqrt{n}\kappa, n\} \log(1/\epsilon))$	$\mathcal{O}(\frac{\sqrt{n}L}{\epsilon^2})$	Yes

Table 3: Complexity bounds obtained from our analytical framework. These complexities represent the number of stochastic gradient oracle calls required to find a point  $\bar{x}_k$  satisfying  $\mathbb{E}[F(\bar{x}_k) - F(x^*)] \leq \epsilon$  for the convex case,  $\mathbb{E}[\|x_k - x^*\|^2] \leq \epsilon$  for the strongly convex case, and an  $\epsilon$ -approximate stationary point in the non-convex case. While we do not recover the best-known rates for (unbiased) SAGA and SVRG in the convex setting, our rates for B-SAGA, B-SVRG, SARAH, and SARGE are the first known for this problem class, and our rates for SARAH and SARGE are better than the best-known rates for SAGA and SVRG in the convex setting.

**Paper organization** Preliminary results and notations are provided in Section 2. A discussion on the tradeoff between bias and variance in stochastic optimization is included in Section 3. Our main convergence results are presented in Section 4. In Section 5, we substantiate our theoretical results using numerical experiments involving several classic regression tasks arising from machine learning. All the proofs of the main results are collected in the appendix.

## 2. Preliminaries and notations

Throughout the paper,  $\mathbb{R}^p$  is a  $p$ -dimensional Euclidean space equipped with scalar inner product  $\langle \cdot, \cdot \rangle$  and associated norm  $\|\cdot\|$ . The sub-differential of a proper closed convex function  $g$  is the set-valued operator defined by  $\partial g(x) \stackrel{\text{def}}{=} \{v \in \mathbb{R}^n | g(x') \geq g(x) + \langle v, x' - x \rangle, \forall x' \in \mathbb{R}^n\}$ , the *proximity*, or *proximal map* of  $g$  is defined as

$$\text{prox}_{\eta g}(y) \stackrel{\text{def}}{=} \arg \min_{x \in \mathbb{R}^n} \eta g(x) + \frac{1}{2} \|x - y\|^2, \quad (3)$$

where  $\eta > 0$  and  $y \in \mathbb{R}^p$ . Below we summarize some useful results in convex analysis.

**Lemma 1 (Nesterov (2004, Thm 2.1.5))** *Suppose  $f$  is convex with an  $L$ -Lipschitz continuous gradient. We have for every  $x, u \in \mathbb{R}^p$ ,*

$$\|\nabla f(x) - \nabla f(u)\|^2 \leq 2L(f(x) - f(u) - \langle \nabla f(u), x - u \rangle).$$

When  $f$  is a finite sum as in (1), Lemma 1 is equivalent to the following result.

**Lemma 2** *Let  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ , where each  $f_i$  is convex with an  $L$ -Lipschitz gradient. Then for every  $x, u \in \mathbb{R}^p$ ,*

$$\frac{1}{2Ln} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f_i(u)\|^2 \leq f(x) - f(u) - \langle \nabla f(u), x - u \rangle.$$

Lemma 2 is obtained by applying Lemma 1 to each  $f_i$  and averaging.

**Lemma 3** *Suppose  $g$  is  $\mu$ -strongly convex with  $\mu \geq 0$ , and let  $z = \text{prox}_{\eta g}(x - \eta d)$  for some  $x, d \in \mathbb{R}^p$  and  $\eta > 0$ . Then, for any  $y \in \mathbb{R}^p$ ,*

$$\eta \langle d, z - y \rangle \leq \frac{1}{2} \|x - y\|^2 - \frac{1+\mu\eta}{2} \|z - y\|^2 - \frac{1}{2} \|z - x\|^2 - \eta g(z) + \eta g(y).$$

**Proof** By the strong convexity of  $g$ , we have  $g(z) - g(y) \leq \langle \xi, z - y \rangle - \frac{\mu}{2} \|z - y\|^2$ ,  $\forall \xi \in \partial g(z)$ . From the definition of the proximal operator, we have that  $\frac{1}{\eta}(x - z) - d \in \partial g(z)$ . Therefore,

$$\begin{aligned} g(z) - g(y) &\leq \langle \xi, z - y \rangle - \frac{\mu}{2} \|z - y\|^2 \\ &= \frac{1}{\eta} \langle x - z - \eta d, z - y \rangle - \frac{\mu}{2} \|z - y\|^2 \\ &= -\langle d, z - y \rangle + \frac{1}{\eta} \langle x - z, z - y \rangle - \frac{\mu}{2} \|z - y\|^2 \\ &= -\langle d, z - y \rangle - \frac{1}{2\eta} \|x - z\|^2 - \frac{1}{2\eta} \|z - y\|^2 + \frac{1}{2\eta} \|x - y\|^2 - \frac{\mu}{2} \|z - y\|^2. \end{aligned}$$

Multiplying by  $\eta$  and rearranging yields the assertion.  $\blacksquare$

The next lemma is an analogue of the descent lemma for gradient descent when the gradient is replaced with an arbitrary vector  $d$ .

**Lemma 4** *Suppose  $g$  is  $\mu$ -strongly convex for  $\mu \geq 0$ , and let  $z = \text{prox}_{\eta g}(x - \eta d)$ . The following inequality holds for any  $\lambda > 0$ .*

$$0 \leq \eta(F(x) - F(z)) + \frac{\eta}{2L\lambda} \|d - \nabla f(x)\|^2 + \left(\frac{\eta L(\lambda+1)}{2} - \frac{2+\mu\eta}{2}\right) \|z - x\|^2.$$

**Proof** This follows immediately from Lemma 3.

$$\begin{aligned} 0 &= \eta \langle d, x - z \rangle + \eta \langle d, z - x \rangle \\ &\stackrel{\textcircled{1}}{\leq} \eta \langle d, x - z \rangle - \frac{2+\mu\eta}{2} \|z - x\|^2 + \eta(g(x) - g(z)) \\ &= \eta \langle \nabla f(x), x - z \rangle + \eta \langle d - \nabla f(x), x - z \rangle - \frac{2+\mu\eta}{2} \|z - x\|^2 + \eta(g(x) - g(z)) \\ &\stackrel{\textcircled{2}}{\leq} \eta(F(x) - F(z)) + \eta \langle d - \nabla f(x), x - z \rangle + \left(\frac{\eta L}{2} - \frac{2+\mu\eta}{2}\right) \|z - x\|^2 \\ &\stackrel{\textcircled{3}}{\leq} \eta(F(x) - F(z)) + \frac{\eta}{2L\lambda} \|d - \nabla f(x)\|^2 + \left(\frac{\eta L(\lambda+1)}{2} - \frac{2+\mu\eta}{2}\right) \|z - x\|^2. \end{aligned}$$

Inequality  $\textcircled{1}$  is due to Lemma 3 with  $y = x$ ,  $\textcircled{2}$  is due to the Lipschitz continuity of  $\nabla f_i$ , and  $\textcircled{3}$  is Young's.  $\blacksquare$

The previous two lemmas require that  $g$  to be convex. Similar results hold in the non-convex case as well.

**Lemma 5** *Let  $z \in \text{prox}_{\eta g}(x - \eta d)$  for some  $x, d \in \mathbb{R}^p$  and  $\eta > 0$ . Then, for any  $y \in \mathbb{R}^p$ ,*

$$\eta \langle d, z - y \rangle \leq \frac{1}{2} \|x - y\|^2 - \frac{1}{2} \|z - x\|^2 - \eta g(z) + \eta g(y).$$

**Proof** By the definition of  $z$ ,

$$z \in \arg \min_v \left\{ \langle d, v - x \rangle + \frac{1}{2\eta} \|v - x\|^2 + g(v) \right\}.$$

Let  $v = y$ , then

$$g(z) - g(y) \leq \langle d, y - z \rangle + \frac{1}{2\eta} (\|x - y\|^2 - \|x - z\|^2).$$

Multiplying by  $\eta$  completes the proof. ■

**Lemma 6** *Let  $z \in \text{prox}_{\eta g}(x - \eta d)$ . Then*

$$F(z) \leq F(y) + \langle \nabla f(x) - d, z - y \rangle + \left(\frac{L}{2} - \frac{1}{2\eta}\right) \|x - z\|^2 + \left(\frac{L}{2} + \frac{1}{2\eta}\right) \|x - y\|^2.$$

**Proof** By the Lipschitz continuity of  $\nabla f$ , we have the inequalities

$$\begin{aligned} f(x) - f(y) &\leq \langle \nabla f(x), x - y \rangle + \frac{L}{2} \|x - y\|^2, \\ f(z) - f(x) &\leq \langle \nabla f(x), z - x \rangle + \frac{L}{2} \|z - x\|^2. \end{aligned}$$

Furthermore, by Lemma 5,  $g(z) - g(y) \leq \langle d, y - z \rangle + \frac{1}{2\eta} (\|x - y\|^2 - \|x - z\|^2)$ . Adding these inequalities together completes the proof. ■

In the non-convex setting, to measure convergence of the sequence to a first-order stationary point, we use the notion of a generalized gradient (Nesterov, 2004).

**Definition 7 (generalized gradient map)** *The generalized gradient map is defined as*

$$\mathcal{G}_\eta(x_k) \stackrel{\text{def}}{=} \frac{1}{\eta} (x_k - \text{prox}_{\eta_k g}(x_k - \eta \nabla f(x_k))).$$

For any  $\eta > 0$ . A point  $x$  satisfying  $\mathcal{G}_\eta(x) = 0$  is a first-order stationary point of  $f + g$ , and an  $\epsilon$ -first-order stationary point is a point satisfying  $\|\mathcal{G}_\eta(x)\| \leq \epsilon$ .

When  $g \equiv 0$ , we have  $\mathcal{G}_{\eta_k}(x_k) = \nabla f(x_k) \rightarrow 0$  if the sequence  $\{x_k\}$  converges to some  $x^* \in \mathbb{R}^p$  such that  $\nabla f(x^*) = 0$ . For nontrivial  $g$ , suppose  $\inf_k \eta_k > 0$  and  $x_k$  converges to some  $x^*$  such that  $x^* \in \text{prox}_{\eta g}(x^* - \eta \nabla f(x^*))$ , then  $\mathcal{G}_{\eta_k}(x_k) \rightarrow 0$ .

### 3. A bias-variance tradeoff in stochastic gradient methods

In this section, we discuss the effect of bias and variance of a stochastic gradient estimator on the performance of Algorithm 1. It is elementary that the mean-squared error (MSE) of a stochastic estimator can be decomposed into the sum of its variance and squared bias. In our setting, we have

$$\mathbb{E}_k[\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] = \|\mathbb{E}_k[\tilde{\nabla}_k] - \nabla f(x_k)\|^2 + \mathbb{E}_k[\|\tilde{\nabla}_k - \mathbb{E}_k[\tilde{\nabla}_k]\|^2].$$

This decomposition shows that a biased estimator might have a smaller MSE than an unbiased estimator as long as the bias sufficiently diminishes the variance. This is the *bias-variance tradeoff*. As we see below, a bias-variance tradeoff exists in our analysis of stochastic gradient methods, but with a slightly different form.

In what follows, we first discuss the bias-variance tradeoff in the convex settings and then the non-convex setting.

### 3.1 The convex case

Let  $x^*$  be a global minimizer of problem (1). From the update (2), let  $w_{k+1} \in \partial g(x_{k+1})$ . We have the following bound on the suboptimality at  $x_{k+1}$ :

$$\begin{aligned}
 & \mathbb{E}_k[F(x_{k+1}) - F(x^*)] \\
 &= \mathbb{E}_k[f(x_{k+1}) - f(x_k) + f(x_k) - f(x^*) + g(x_{k+1}) - g(x^*)] \\
 &\stackrel{\textcircled{1}}{\leq} \frac{L}{2} \mathbb{E}_k[\|x_{k+1} - x_k\|^2] + \mathbb{E}_k[\langle \nabla f(x_k), x_{k+1} - x_k \rangle] + \langle \nabla f(x_k), x_k - x^* \rangle + \mathbb{E}_k[g(x_{k+1}) - g(x^*)] \\
 &= \frac{L}{2} \mathbb{E}_k[\|x_{k+1} - x_k\|^2] + \mathbb{E}_k[\langle \nabla f(x_k) - \tilde{\nabla}_k, x_{k+1} - x_k \rangle] \\
 &\quad + \mathbb{E}_k[\langle \nabla f(x_k) - \tilde{\nabla}_k, x_k - x^* \rangle] + \mathbb{E}_k[\langle \tilde{\nabla}_k, x_{k+1} - x^* \rangle] + \mathbb{E}_k[g(x_{k+1}) - g(x^*)] \\
 &\stackrel{\textcircled{2}}{\leq} \frac{\epsilon}{2} \mathbb{E}_k[\|\nabla f(x_k) - \tilde{\nabla}_k\|^2] + (\frac{L}{2} + \frac{1}{2\epsilon}) \mathbb{E}_k[\|x_{k+1} - x_k\|^2] \\
 &\quad + \mathbb{E}_k[\langle \nabla f(x_k) - \tilde{\nabla}_k, x_k - x^* \rangle] + \mathbb{E}_k[\langle \tilde{\nabla}_k + w_{k+1}, x_{k+1} - x^* \rangle - \frac{\mu}{2} \|x_{k+1} - x^*\|^2] \\
 &\stackrel{\textcircled{3}}{\leq} \frac{\epsilon}{2} \mathbb{E}_k[\|\nabla f(x_k) - \tilde{\nabla}_k\|^2] + (\frac{L}{2} + \frac{1}{2\epsilon} - \frac{1}{2\eta}) \mathbb{E}_k[\|x_{k+1} - x_k\|^2] \\
 &\quad + \mathbb{E}_k[\langle \nabla f(x_k) - \tilde{\nabla}_k, x_k - x^* \rangle] - \frac{1+\mu\eta}{2\eta} \mathbb{E}_k[\|x_{k+1} - x^*\|^2] + \frac{1}{2\eta} \|x_k - x^*\|^2.
 \end{aligned} \tag{4}$$

Inequality  $\textcircled{1}$  follows from the convexity of  $f$  and Lipschitz continuity of  $\nabla f$ ,  $\textcircled{2}$  follows from the (strong) convexity of  $g$ , and  $\textcircled{3}$  comes from the implicit definition of the proximal operator (3). For the last line of (4), we observe that the inner product term  $\mathbb{E}_k[\langle \nabla f(x_k) - \tilde{\nabla}_k, x_k - x^* \rangle]$  vanishes when  $\tilde{\nabla}_k$  is an unbiased estimate of  $\nabla f(x_k)$ . When the estimator is biased, we must develop a new way to control this term, together with  $\mathbb{E}_k[\|\nabla f(x_k) - \tilde{\nabla}_k\|^2]$ .

Hence, the following terms arise in our convergence analysis from the bias and the variance of the gradient estimator:

$$\begin{aligned}
 & \text{Bias : } \mathbb{E}_k[\langle \nabla f(x_k) - \tilde{\nabla}_k, x_k - x^* \rangle] \quad \text{and} \quad \|\mathbb{E}_k[\tilde{\nabla}_k] - \nabla f(x_k)\|^2, \\
 & \text{Variance : } \mathbb{E}_k[\|\tilde{\nabla}_k - \mathbb{E}_k[\tilde{\nabla}_k]\|^2].
 \end{aligned} \tag{5}$$

An effective gradient estimator minimizes all three of these terms simultaneously. As we later show in our MSE and bias-term bounds, SARAH and SARGE minimize these terms more effectively than biased SAGA and SVRG estimators, leading to better convergence rates. We provide an explicit comparison between SARAH and SVRG in Appendix G.

**Remark 8 (Non-composite case  $g = 0$ )** *When  $g = 0$ , for gradient descent, the descent property of  $f$  yields*

$$f(x_{k+1}) - f(x^*) \leq (\frac{L}{2} - \frac{1}{\eta}) \|x_{k+1} - x_k\|^2 + f(x_k) - f(x^*),$$

where  $\eta \leq 2/L$ . For stochastic gradient descent, we obtain the following relationship:

$$\begin{aligned}
 & \mathbb{E}_k[f(x_{k+1}) - f(x^*)] \\
 &= \mathbb{E}_k[f(x_{k+1}) - f(x_k) + f(x_k) - f(x^*)] \\
 &\leq \mathbb{E}_k[\langle \nabla f(x_k) - \tilde{\nabla}_k, x_{k+1} - x_k \rangle] + (\frac{L}{2} - \frac{1}{\eta}) \mathbb{E}_k[\|x_{k+1} - x_k\|^2] + f(x_k) - f(x^*) \\
 &\leq \frac{\epsilon}{2} \mathbb{E}_k[\|\nabla f(x_k) - \tilde{\nabla}_k\|^2] + (\frac{L}{2} + \frac{1}{2\epsilon} - \frac{1}{\eta}) \mathbb{E}_k[\|x_{k+1} - x_k\|^2] + f(x_k) - f(x^*).
 \end{aligned} \tag{6}$$

Compared to (4), there is no inner product term in (6), which makes the analysis of the non-composite case much simpler. This is one reason why biased algorithms have been successfully studied in non-composite setting, but not in the composite setting.

### 3.2 The non-convex case

The influence of bias is simpler in the non-convex setting and independent of  $g$ , which explains why biased algorithms have recently found success for these problems. To begin, let  $\hat{x}_{k+1} \in \text{prox}_{\eta g/2}(x_k - \frac{\eta}{2}\nabla f(x_k))$ . Applying Lemma 6 with  $z = \hat{x}_{k+1}$ ,  $y = x = x_k$  and  $d = \nabla f(x_k)$ , we have

$$F(\hat{x}_{k+1}) \leq F(x_k) + \left(\frac{L}{2} - \frac{1}{\eta}\right)\|\hat{x}_{k+1} - x_k\|^2.$$

Again, applying Lemma 6 with  $z = x_{k+1}$ ,  $y = \hat{x}_{k+1}$ ,  $x = x_k$ , and  $d = \tilde{\nabla}_k$ , we obtain

$$\begin{aligned} F(x_{k+1}) &\leq F(\hat{x}_{k+1}) + \langle \nabla f(x_k) - \tilde{\nabla}_k, x_{k+1} - \hat{x}_{k+1} \rangle \\ &\quad + \left(\frac{L}{2} - \frac{1}{2\eta}\right)\|x_{k+1} - x_k\|^2 + \left(\frac{L}{2} + \frac{1}{2\eta}\right)\|\hat{x}_{k+1} - x_k\|^2 \end{aligned}$$

Adding these two inequalities together gives

$$\begin{aligned} F(x_{k+1}) &\leq F(x_k) + \left(L - \frac{1}{2\eta}\right)\|\hat{x}_{k+1} - x_k\|^2 + \left(\frac{L}{2} - \frac{1}{2\eta}\right)\|x_{k+1} - x_k\|^2 \\ &\quad + \langle \nabla f(x_k) - \tilde{\nabla}_k, x_{k+1} - \hat{x}_{k+1} \rangle \\ &\stackrel{\textcircled{1}}{\leq} F(x_k) + \left(L - \frac{1}{2\eta}\right)\|\hat{x}_{k+1} - x_k\|^2 + \left(\frac{L}{2} - \frac{1}{2\eta}\right)\|x_{k+1} - x_k\|^2 + 2\eta\|\nabla f(x_k) - \tilde{\nabla}_k\|^2 \\ &\quad + \frac{1}{8\eta}\|\hat{x}_{k+1} - x_{k+1}\|^2 \\ &\stackrel{\textcircled{2}}{\leq} F(x_k) + \left(L - \frac{1}{4\eta}\right)\|\hat{x}_{k+1} - x_k\|^2 + \left(\frac{L}{2} - \frac{1}{4\eta}\right)\|x_{k+1} - x_k\|^2 + 2\eta\|\nabla f(x_k) - \tilde{\nabla}_k\|^2. \end{aligned} \tag{7}$$

Inequality  $\textcircled{1}$  is Young's, and  $\textcircled{2}$  is the standard inequality  $\|a - c\|^2 \leq 2\|a - b\|^2 + 2\|b - c\|^2$ . In the non-convex case, the inner-product bias term does not appear, so the bias-variance tradeoff is the classical one.

### 3.3 General bounds on bias and variance

To ensure convergence of Algorithm 1 using a particular gradient estimator, we must bound the inner-product bias term,  $\mathbb{E}_k[\langle \nabla f(x_k) - \tilde{\nabla}_k, x_k - x^* \rangle]$ , and the MSE,  $\mathbb{E}_k[\|\nabla f(x_k) - \tilde{\nabla}_k\|^2]$ . Below we introduce general bounds on these terms that allow us to establish convergence rates for a variety of gradient estimators. The first of these is a bound on the MSE term.

**Definition 9 (Bounded MSE)** *The stochastic gradient estimator  $\tilde{\nabla}$  is said to satisfy the BMSE( $M_1, M_2, \rho_M, \rho_F, m$ ) property with parameters  $M_1, M_2 \geq 0$ ,  $\rho_M, \rho_F \in (0, 1]$  and  $m \geq 1$  if there exist sequences  $\mathcal{M}_k$  and  $\mathcal{F}_k$  such that*

$$\sum_{k=m_s}^{m(s+1)-1} \mathbb{E}[\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] \leq \mathcal{M}_{m_s},$$

and the following bounds hold:

$$\begin{aligned} \mathcal{M}_{m_s} &\leq (1 - \rho_M)^m \mathcal{M}_{m(s-1)} + \mathcal{F}_{m_s} + \frac{M_1}{n} \sum_{k=m_s}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2], \\ \mathcal{F}_{m_s} &\leq \sum_{\ell=0}^s \frac{M_2(1-\rho_F)^{m(s-\ell)}}{n} \sum_{k=m_s}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2]. \end{aligned}$$

The constant  $m$  is the epoch length of the gradient estimator, hence it is usually set to be  $\mathcal{O}(n)$ . This property is useful in convergence analyses because it bounds the MSE by a geometrically decaying sequence  $\{\mathcal{M}_{mk}\}_{k \in \mathbb{N}}$  and a component that is proportional to the one-iteration progress of gradient descent  $(1/n \sum_{i=1}^n \|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2)$ .

**Remark 10**

- *Most variance-reduced stochastic gradient estimators satisfy the BMSE property, including SAG, SAGA, SVRG, SARAH, and all the estimators considered by Hofmann et al. (2015). SGD does not satisfy this property, as its variance does not decay along the iterations.*
- *Most existing work on the analysis of general stochastic gradient algorithms enforce bounds of this form on either the MSE or the moments of the stochastic estimator, with the crucial difference that existing works require the bounds to be Markovian (that is, dependent on only the previous iteration) (Bottou et al., 2018). In contrast, the BMSE property allows non-Markovian MSE bounds through the sequence  $\mathcal{F}_k$ . This relaxation is crucial for the analysis of our new gradient estimator, SARGE.*

In order to bound the inner-product bias term, we require the gradient estimator to admit a certain structure in its bias. In biased estimators such as SAG, the bias depends on the stored gradient values:

$$\nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k^{\text{SAG}}] = (1 - \frac{1}{n})(\nabla f(x_k) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i)).$$

We call estimators whose bias admits the above structure *memory-biased* gradient estimators. These include SAG, and more generally B-SAGA and B-SVRG.

**Definition 11 (Memory-biased gradient estimator)** *The stochastic gradient estimator  $\tilde{\nabla}$  is memory-biased with parameters  $\theta > 0$ ,  $B_1 \geq 0$ , and  $m \geq 1$  if*

$$\nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k] = (1 - \frac{1}{\theta})(\nabla f(x_k) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i)),$$

for some  $\{\varphi_k^i\}_{i=1}^n \subset \{x_\ell\}_{\ell=0}^{k-1}$ , and for any  $s \in \mathbb{N}_0$ ,

$$\sum_{k=ms}^{m(s+1)-1} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|x_k - \varphi_k^i\|^2] \leq B_1 \sum_{k=ms}^{m(s+1)-1} \mathbb{E}[\|x_k - x_{k-1}\|^2]. \quad (8)$$

B-SAGA is clearly a memory-biased estimator, and so is B-SVRG where  $\varphi_k^i = \varphi_{ms}^i$  for all  $k$  in epoch  $s$ . The parameter  $\theta$  controls the amount of bias in the estimator, and  $B_1$ , in a sense, measures how “stale” the stored gradient information is. For memory-biased gradient estimators, the bias-term can be bounded by terms of the form  $\|x_k - \varphi_k^i\|^2$ .

**Lemma 12** *Suppose  $\tilde{\nabla}$  is memory-biased with parameter  $\theta \geq 1$  and that  $F$  is  $\mu$ -strongly convex with  $\mu \geq 0$ . For any  $\lambda > 0$ , the following inequality holds:*

$$\begin{aligned} \eta \mathbb{E}_k[F(x_{k+1}) - F(x^*)] &\leq \frac{\eta}{2L\lambda} \mathbb{E}_k[\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] - \frac{1+\mu\eta}{2} \mathbb{E}_k[\|x_{k+1} - x^*\|^2] + \frac{1}{2} \|x_k - x^*\|^2 \\ &\quad + (\frac{\eta L(\lambda+1)}{2} - \frac{1}{2}) \mathbb{E}_k[\|x_{k+1} - x_k\|^2] + \frac{\eta L}{2n} (1 - \frac{1}{\theta}) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2. \end{aligned}$$

The proof of Lemma 12 can be found in Appendix A. The bound of Lemma 12 is analogous to the bound in (4), but the inner-product bias term is replaced with  $\frac{\eta L}{2n}(1-\frac{1}{\theta})\sum_{i=1}^n\|x_k-\varphi_k^i\|^2$ . This term is proportional to the progress of gradient descent (by (8)), so this provides the necessary control over the inner-product bias term.

**Remark 13** *Lemma 12 requires that  $\theta \geq 1$ , so the rates we derive for the convex setting hold only for  $\theta \geq 1$ . However, the convergence rate we prove in Section 4.1.2 for the non-convex setting, which allow  $\theta \in (0, 1)$ , also hold for convex problems. In summary, our results guarantee convergence for all  $\theta > 0$ , but they suggest different rates for the parameter settings  $\theta < 1$  and  $\theta \geq 1$ .*

For estimators such as SARAH, the bias depends on the error in the previous gradient estimate, rather than previous stochastic gradients:

$$\nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k^{\text{SARAH}}] = \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARAH}}.$$

We refer to estimators of this type as *recursively biased*.

**Definition 14 (Recursively biased gradient estimator)** *For any sequence  $\{x_k\}$ , let  $\tilde{\nabla}_k$  be a stochastic gradient estimator generated from the points  $\{x_\ell\}_{\ell=0}^k$ . This estimator is recursively biased with parameters  $\rho_B \in (0, 1]$  and  $\nu \geq 1$  if*

$$\nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k] = \begin{cases} 0 & \text{for } k \in \nu\mathbb{N}_0, \\ (1 - \rho_B)(\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}) & \text{o.w.} \end{cases}$$

The parameter  $\nu$  represents how many steps occur between full gradient evaluations. For SARGE,  $\nu = \infty$  because the full gradient is never computed. Recursively biased estimators admit a bound on the inner-product bias term that involves the estimator's MSE.

**Lemma 15** *Suppose  $\tilde{\nabla}$  is a recursively biased gradient estimator with parameters  $\nu \geq 1$  and  $\rho_B \in (0, 1]$ . Then, for any epoch  $s \in \mathbb{N} \cup \{0\}$  and  $\epsilon > 0$ ,*

$$\begin{aligned} & \sum_{k=\nu s+1}^{\nu(s+1)-1} |\mathbb{E}\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle| \\ & \leq \min\left\{\nu, \frac{1}{\rho_B}\right\} \sum_{k=\nu s+1}^{\nu(s+1)-1} \mathbb{E}\left[\frac{\epsilon}{2}\|\nabla f(x_k) - \tilde{\nabla}_k\|^2 + \frac{1}{2\epsilon}\|x_{k+1} - x_k\|^2\right]. \end{aligned}$$

The proof of Lemma 15 is in Appendix B. Lemma 15 shows that, for recursively biased estimators, the inner-product bias term  $\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle$  is bounded from above by the MSE, implying that introducing bias to decrease the MSE is a reasonable approach to design improved gradient estimators.

**Remark 16** *When  $\nu = \infty$ , which is true for SARGE, there is only a single epoch,  $s = 0$ . In this case, we adopt the convention  $\infty \cdot 0 = 0$ , so that the sums appearing in Lemma 15 are well-defined.*

## 4. Convergence rates

In this section, we analyze the convergence rates for the stochastic gradient methods. We first provide very general convergence rates based on the bounds from the last section. Then, we specify the result to specific gradient estimators including memory-biased B-SAGA/B-SVRG, and recursively biased SARAH and SARGE.

### 4.1 General convergence rates

For Algorithm 1, we consider a constant step size  $\eta_k \equiv \eta > 0$ . Given  $T$  iterations of Algorithm 1, define the average iterate  $\bar{x}_T \stackrel{\text{def}}{=} 1/T \sum_{k=1}^T x_k$ .

#### 4.1.1 CONVEX AND STRONGLY CONVEX CASES

The following theorem establishes convergence rates for memory-biased estimators in the convex regime.

**Theorem 17 (Memory-biased estimators)** *Let  $\tilde{\nabla}$  be a memory-biased gradient estimator parameterized by  $\theta \geq 1$  and  $B_1 \geq 0$ , which satisfies the BMSE( $M_1, M_2, \rho_M, \rho_F, m$ ) property. Let  $\Theta = \frac{M_1 \rho_F + 2M_2}{\rho_M \rho_F}$  and  $\rho = \min\{\rho_M, \rho_F\}$ .*

- In the convex setting, let  $\eta = \frac{1}{L(1+3\sqrt{2\Theta})}$ , then

$$\mathbb{E}[F(\bar{x}_T) - F(x^*)] \leq \frac{1}{T} \left( \frac{L(1+3\sqrt{2\Theta})\|x_0 - x^*\|^2}{2} + \max\left\{ \frac{B_1(1-1/\theta)}{\sqrt{2\Theta}} - 1, 0 \right\} (F(x_0) - F(x^*)) \right).$$

- When  $g$  is additionally  $\mu$ -strongly convex with  $\mu > 0$ , let  $\eta = \min\left\{ \frac{1}{3L(1+3\sqrt{2\Theta})}, \frac{\sqrt{2\Theta}}{B_1\mu(1-1/\theta)}, \frac{\rho}{2\mu} \right\}$ . The iterate  $x_T$  satisfies

$$\mathbb{E}[\|x_T - x^*\|^2] \leq (1 + \mu\eta)^{-T} \left( \frac{2}{\mu} (F(x_0) - F(x^*)) + \|x_0 - x^*\|^2 \right).$$

The proof of Theorem 17 is provided in Appendix A. The next result establishes convergence rates for recursively biased gradient estimators whose proof is in Appendix B.

**Theorem 18 (Recursively biased estimators)** *Let  $\tilde{\nabla}$  be a recursively biased gradient estimator parameterized by  $\rho_B \in (0, 1)$  and  $\nu \geq 1$ , which satisfies the BMSE( $M_1, M_2, \rho_M, \rho_F, m$ ) property. Let  $B_2 \stackrel{\text{def}}{=} \min\{\nu, 1/\rho_B\}$ ,  $\Theta = \frac{M_1 \rho_F + 2M_2}{\rho_M \rho_F}$  and  $\rho = \min\{\rho_M, \rho_F\}$ .*

- In the convex setting, let  $\eta = \frac{1}{2\sqrt{2\Theta}L(1 + \frac{(1-\rho_B)B_2}{1+\delta}) + L}$  with  $\delta = \max\{\sqrt{L\Theta^{1/2}(1-\rho_B)B_2} - 1, 0\}$ . Then

$$\mathbb{E}[F(\bar{x}_T) - F(x^*)] \leq \frac{1}{T} \left( \frac{1}{2\eta} \|x_0 - x^*\|^2 + \delta (F(x_0) - F(x^*)) \right).$$

- When  $g$  is additionally  $\mu$ -strongly convex with  $\mu > 0$ , let  $\eta = \min\left\{ \frac{1}{3L(4\sqrt{2\Theta}+1)}, \frac{1}{\mu(1-\rho_B)B_2}, \frac{\rho}{2\mu} \right\}$ , then

$$\mathbb{E}[\|x_T - x^*\|^2] \leq (1 + \mu\eta)^{-T} \left( \frac{2}{\mu} (F(x_0) - F(x^*)) + \|x_0 - x^*\|^2 \right).$$

**Remark 19**

- Both theorems hold true for smaller  $\eta$ ; the choices in the theorems are the largest ones allowed by our analysis.
- For B-SAGA and B-SVRG,  $\Theta = \mathcal{O}(n^2)$ , while for SARAH and SARGE,  $\Theta = \mathcal{O}(n)$ . This gives these recursive gradient estimators improved convergence rates and suggests that the bias in these estimators is more effective than the bias in SAGA and SVRG.

## 4.1.2 NON-CONVEX CASE

The analysis of biased gradient estimators is simpler for the non-convex setting than the convex ones due to the absence of the inner-product bias term in (7). Below we provide a uniform convergence guarantee for all gradient estimators satisfying the BMSE property, regardless of their bias. This suggests that in the non-convex setting, a large-bias, small-MSE gradient estimator is favorable over an estimator with small bias and large MSE.

**Theorem 20** *Let  $\tilde{\nabla}$  be a gradient estimator that satisfies the BMSE( $M_1, M_2, \rho_M, \rho_F, m$ ) property, let  $\Theta = \frac{M_1 \rho_F + 2M_2}{\rho_M \rho_F}$ , and let  $\alpha$  be a chosen uniformly at random from the set  $\{0, 1, \dots, T-1\}$ . If  $F$  is non-convex, set  $\eta = \frac{\sqrt{16\Theta+1}-1}{16L\Theta}$  in Algorithm 1, and the point  $x_\alpha$  satisfies the following bound on its generalized gradient:*

$$\mathbb{E}[\|\mathcal{G}_{\eta/2}(x_\alpha)\|^2] \leq \frac{16(F(x_0) - F(x^*))}{T\eta(1-4\eta L)}.$$

The proof of this result is provided in Appendix C.

**Remark 21** *The convergence result of Theorem 20 does not depend on the bias except through the MSE of the gradient estimator, which implies that incorporating arbitrary amounts of bias for a smaller MSE improves the convergence rate. This fact is what allows the recursively biased estimators SARAH and SARGE to achieve the oracle complexity lower bound for non-convex optimization when they are used in Algorithm 1.*

**4.2 Convergence rates for specific gradient estimators**

In this section, we specialize the general convergence rates to analyze the performance of B-SAGA, B-SVRG, SARAH, and SARGE.

## 4.2.1 BIASED SAGA AND SVRG

B-SAGA and B-SVRG are examples of memory-biased gradient estimators, as their biases take the form

$$\nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k] = (1 - \frac{1}{\theta})(\nabla f(x_k) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i)),$$

for some previous iterates  $\varphi_k^i$ . To establish convergence rates for B-SAGA and B-SVRG, we only need to show these estimators satisfy the BMSE property with suitable constants.

**Lemma 22** *The B-SAGA gradient estimator is memory-biased with  $B_1 = 2n(2n+1)$ , and it satisfies the BMSE property with parameters  $\rho_M = \frac{1}{2n}$ ,  $m = 1$ ,  $M_2 = 0$ ,  $\rho_F = 1$ , and*

$$M_1 = \begin{cases} \frac{2n+1}{\theta^2} & \theta \in (0, 2], \\ (2n+1)(1 - \frac{1}{\theta})^2 & \theta > 2. \end{cases}$$

The proof of Lemma 22 uses a slight modification of existing variance bounds for the SAGA estimator, appearing in (Defazio et al., 2014a) for example. We include the proof in Appendix D. The B-SVRG gradient estimator satisfies the BMSE property with similar constants.

**Lemma 23** *The B-SVRG gradient estimator is memory-biased with  $B_1 = 3m(m+1)$ , and it satisfies the BMSE property with parameters  $\rho_M = 1$ ,  $M_2 = 0$ ,  $\rho_F = 1$ , and*

$$M_1 = \begin{cases} \frac{3m(m+1)}{\theta^2} & \theta \in (0, 2], \\ 3m(m+1)(1 - \frac{1}{\theta})^2 & \theta > 2. \end{cases}$$

With these constants established, Theorem 17 provides rates of convergence.<sup>1</sup>

**Corollary 24 (Convergence rates for B-SAGA)** *Algorithm 1 achieves the following convergence guarantees using the B-SAGA gradient estimator:*

- *In the convex setting, depending on the choice of  $\theta$ , set the step size to*

$$\eta = \eta_\theta \stackrel{\text{def}}{=} \begin{cases} \frac{1}{L(1 + \frac{6}{\theta}\sqrt{n(2n+1)})} : \theta \in [1, 2], \\ \frac{1}{L(1 + 6(1 - \frac{1}{\theta})\sqrt{n(2n+1)})} : \theta > 2, \end{cases}$$

*and  $\bar{x}_T$  satisfies  $\mathbb{E}[F(\bar{x}_T) - F(x^*)] = \mathcal{O}(\frac{Ln}{T})$ .*

- *If additionally  $g$  is  $\mu$ -strongly convex, set  $\eta = \min\{\eta_\theta, \frac{1}{4\mu}\}$ . Then  $x_T$  satisfies  $\mathbb{E}[\|x_T - x^*\|^2] = \mathcal{O}((1 + \mu\eta)^{-T})$ .*
- *In the non-convex setting, after  $T$  iterations, the generalized gradient at  $x_\alpha$  satisfies*

$$\mathbb{E}[\|\mathcal{G}_{\eta/2}(x_\alpha)\|^2] = \begin{cases} \mathcal{O}(\frac{Ln}{T\theta}) : \eta = \frac{\theta}{2L\sqrt{n(2n+1)}}, \theta \in (0, 2], \\ \mathcal{O}(\frac{Ln}{T(1 - \frac{1}{\theta})}) : \eta = \frac{1}{2L(1 - \frac{1}{\theta})\sqrt{n(2n+1)}}, \theta > 2. \end{cases}$$

**Corollary 25 (Convergence rates for B-SVRG)** *Algorithm 1 achieves the following convergence guarantees using the B-SVRG gradient estimator:*

- *In the convex setting, depending on the choice of  $\theta$ , set the step size to*

$$\eta = \eta_\theta = \begin{cases} \frac{1}{L(1 + \frac{3}{\theta}\sqrt{6m(m+1)})} : \theta \in [1, 2], \\ \frac{1}{L(1 + 3(1 - \frac{1}{\theta})\sqrt{6m(m+1)})} : \theta > 2. \end{cases}$$

*After  $S$  epochs, the point  $\bar{x}_{mS}$  satisfies  $\mathbb{E}[F(\bar{x}_{mS}) - F(x^*)] = \mathcal{O}(L/S)$ .*

- *If additionally  $g$  is  $\mu$ -strongly convex, let  $\eta = \min\{\eta_\theta, \frac{1}{2\mu}\}$ . After  $S$  epochs,  $x_{mS}$  satisfies  $\mathbb{E}[\|x_{mS} - x^*\|^2] = \mathcal{O}((1 + \mu\eta)^{-mS})$ .*
- *In the non-convex setting, after  $S$  epochs, the generalized gradient at  $x_\alpha$  satisfies*

$$\mathbb{E}[\|\mathcal{G}_{\eta/2}(x_\alpha)\|^2] = \begin{cases} \mathcal{O}(\frac{L}{S\theta}) : \eta = \frac{\sqrt{2}\theta}{2L\sqrt{3m(m+1)}}, \theta \in (0, 2], \\ \mathcal{O}(\frac{L}{S(1 - \frac{1}{\theta})}) : \eta = \frac{\sqrt{2}\theta}{2L(1 - \frac{1}{\theta})\sqrt{3m(m+1)}}, \theta > 2. \end{cases}$$

**Remark 26**

- Our MSE bounds and convergence rates are optimized when  $\theta = 2$ . Numerical experiments (including those in Section 5) suggest that setting  $\theta$  in the range  $1 < \theta \ll n$  gives the best performance for convex problems, and B-SAGA prefers larger values of  $\theta$  than B-SVRG.
- The convergence guarantees for  $\theta \in (0, 1)$  still hold for convex objectives, but in this setting, the rates we obtain for  $\theta \geq 1$  are superior, suggesting that for convex problems,  $\theta$  should be larger than or equal to one for best performance. Our numerical experiments in Section 5 support this; setting  $\theta < 1$  can be beneficial for non-convex problems, but we do not observe this for convex problems.
- In the special case  $\theta = 1$ , Corollaries 24 and 25 recover the state-of-the-art rates for SAGA and SVRG in the non-convex regime. For strongly convex problems, these rates are worse than existing convergence rates of  $\mathcal{O}((1 + \min\{\frac{\mu}{L}, \frac{1}{n}\})^{-T})$  proven for SAGA and SVRG (Defazio et al., 2014a; Xiao and Zhang, 2014). This difference is due to the generality of Theorem 17, as some memory-biased estimators, including B-SVRG, exhibit poor performance on strongly convex problems when the bias is large.
- Corollaries 24 and 25 require step sizes that decrease with  $n$ , while existing results for SAG, SAGA, and SVRG allow step sizes that are independent of  $n$ . This is also due to the generality of Theorem 17. In practice, we find B-SAGA converges with step sizes that are independent of  $n$ , but B-SVRG requires step sizes to decrease when the epoch length is larger.

4.2.2 SARAH AND SARGE

The SARAH and SARGE gradient estimators are recursively biased, with

$$\begin{aligned} \nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k^{\text{SARAH}}] &= \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARAH}} \\ \text{and } \nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k^{\text{SARGE}}] &= (1 - \frac{1}{n})(\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}). \end{aligned}$$

As we shall see, these biased estimators admit smaller MSE bounds than unbiased and memory-biased estimators, and this is reflected in their improved convergence rates. The following two lemmas establish the constants appearing in Theorem 18 for these estimators.

**Lemma 27** *The SARAH gradient estimator is recursively biased with parameters  $\rho_B = 0$  and  $\nu = m$ , and it satisfies the BMSE property with  $M_1 = m$ ,  $\rho_M = 1$ ,  $\rho_F = 1$ , and  $M_2 = 0$ .*

**Lemma 28** *The SARGE gradient estimator is recursively biased with parameters  $\rho_B = 1/n$  and  $\nu = \infty$ , and it satisfies the BMSE property with  $M_1 = 12$ ,  $M_2 = 39/n$ ,  $\rho_M = \frac{1}{4n}$ ,  $\rho_F = \frac{1}{2n}$ , and  $m = 1$ .*

Proofs of these results are included in Appendices E and F, respectively. It is enlightening to compare these BMSE constants to those of B-SVRG and B-SAGA.  $M_1$  is a factor of  $n$  smaller for the SARAH and SARGE estimators than for the B-SVRG and B-SAGA estimators (as long as  $m = \mathcal{O}(n)$  in SARAH and B-SVRG). This translates to an  $\mathcal{O}(\sqrt{n})$  improvement in the convergence rates for SARAH and SARGE when  $L$  is  $\mathcal{O}(\sqrt{n})$  or larger.

**Corollary 29 (Convergence rates for SARAH)** *When using the SARAH gradient estimator in Algorithm 1,*

- If  $F$  is convex, set  $\eta = \frac{1}{L(2\sqrt{2m+1})+\sqrt{Lm}^{3/4}}$ . After  $T$  iterations,  $\bar{x}_T$  satisfies  $\mathbb{E}[F(\bar{x}_T) - F(x^*)] = \mathcal{O}(\frac{L\sqrt{m}+\sqrt{Lm}^{3/4}}{T})$ .
- If  $F$  is  $\mu$ -strongly convex, set  $\eta = \min\{\frac{1}{3L(4\sqrt{2m+1})}, \frac{1}{\mu m}\}$ , then  $\mathbb{E}[\|x_T - x^*\|^2] = \mathcal{O}((1 + \mu\eta)^{-T})$ .
- If  $F$  is non-convex, set  $\eta = \frac{1}{L\sqrt{2m}}$ , then  $\mathbb{E}[\|\mathcal{G}_{\eta/2}(x_\alpha)\|^2] \leq \mathcal{O}(L\sqrt{m}/T)$ .

**Corollary 30 (Convergence rates for SARGE)** *When using the SARGE gradient estimator in Algorithm 1,*

- If  $F$  is convex, set  $\eta = \frac{1}{L(74\sqrt{n}+15\sqrt{Ln}^{3/4})+1}$ , then  $\mathbb{E}[F(\bar{x}_T) - F(x^*)] = \mathcal{O}(\frac{L\sqrt{n}+\sqrt{Ln}^{3/4}}{T})$ .
- If  $F$  is  $\mu$ -strongly convex, set  $\eta = \min\{\frac{1}{3L(16\sqrt{3(n+13)+1})}, \frac{1}{4\mu n}\}$ , then  $\mathbb{E}[\|x_T - x^*\|^2] = \mathcal{O}((1 + \mu\eta)^{-T})$ .
- If  $F$  is non-convex, set  $\eta = \frac{1}{4L\sqrt{3(n+13)}}$ , then  $\mathbb{E}[\|\mathcal{G}_{\eta/2}(x_\alpha)\|^2] \leq \mathcal{O}(L\sqrt{n}/T)$ .

**Remark 31** *Our theoretical results suggest the step sizes for SARAH and SARGE should decrease with  $n$ , and such step sizes lead to optimal convergence guarantees for non-convex problems. However this is not true in practice, as we find that using larger step sizes that are independent of  $n$  leads to better performance. We provide examples in Section 5.*

These convergence rates for convex objectives represent a significant improvement over the performance of SAGA, SVRG, and full-gradient methods. Each of these algorithms require  $\mathcal{O}(\frac{nL}{\epsilon})$  stochastic gradient evaluations to find a point satisfying  $F(x_T) - F(x^*) \leq \epsilon$ , while SARAH and SARGE require only  $\mathcal{O}(\frac{\sqrt{nL}}{\epsilon})$ . These rates do not require the epoch-doubling procedure of (Allen-Zhu and Yuan, 2018), although epoch-doubling can potentially be used to improve the performance of SARAH just as it improves the performance of SVRG on non-strongly convex objectives.

This square-root dependence on  $n$  is present in the convergence rates for strongly convex and non-convex objectives as well, which is a significant improvement over the dependence on  $n$  in the convergence rates of B-SAGA and B-SVRG. This better dependence on  $n$  is most significant in the non-convex regime, where these convergence rates imply that the SARAH and SARGE estimators require only  $\mathcal{O}(\frac{L\sqrt{n}+\sqrt{Ln}^{3/4}}{\epsilon^2})$  stochastic gradient evaluations to find an  $\epsilon$ -approximate stationary point, which is the oracle-complexity lower bound (Fang et al., 2018). Similar results already exist for algorithms using the SARAH estimator (Fang et al., 2018; Wang et al., 2019, 2018; Pham et al., 2019). Our results for SARGE show that achieving this complexity is possible without ever computing the full gradient.

## 5. Numerical Experiments

In this section, we present numerical experiments testing B-SAGA, B-SVRG, SARAH, and SARGE for minimizing convex, strongly convex, and non-convex objectives. We include one set of experiments comparing different values of  $\theta$  in B-SAGA and B-SVRG with a

fixed step size and one set comparing SARAH and SARGE to B-SAGA and B-SVRG with the best values of  $\theta$ .<sup>2</sup>

### 5.1 Convex and strongly convex objectives

Let  $(h_i, l_i) \in \mathbb{R}^p \times \{\pm 1\}$ ,  $i = 1, \dots, n$  be the training set, where  $h_i \in \mathbb{R}^p$  is the feature vector of each data sample, and  $l_i$  is the binary label. Let  $\beta > 0$  be a tuning parameter. The ridge regression problem takes the form

$$\min_{x \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (h_i^\top x - l_i)^2 + \frac{\beta}{2} \|x\|_2^2.$$

LASSO is similar, but with the regularizer  $\|x\|_1$  replacing  $\|x\|_2^2$ . These problems are of the form (1), where we set  $f_i = (h_i^\top x - l_i)^2$  and  $g$  equal to the regularizer. In ridge regression,  $g$  is strongly convex, and in LASSO,  $g$  is only convex.

We consider four binary classification data sets: `australian`, `mushrooms`, `phishing`, and `ijcnn1` from LIBSVM<sup>3</sup>. We rescale the value of the data to  $[-1, 1]$ , set  $\beta = 1/n$ , and set the step size to  $\eta = \frac{1}{5L}$ . To compare performance, we use the objective function value  $F(x_k) - F(x^*)$  is considered.

**Comparison of B-SAGA** We first compare the performance of B-SAGA under different choices of  $\theta$  for solving ridge regression and LASSO problems. Four choices of  $\theta$  are considered:  $\theta \in \{1, 10, 100, n\}$ , the results are provided below in Figure 1, from which we observe that B-SAGA consistently performs better with moderate amounts of bias (that is,  $\theta \in (1, n)$ ). For the considered datasets, overall  $\theta = 10$  provides the best performance.

**Comparison of B-SVRG** We also consider four choices of  $\theta$  for B-SVRG, which are  $\theta \in \{0.5, 0.8, 1, 1.5\}$ . The results are shown below in Figure 2. We observe that B-SVRG is more sensitive to the choice of  $\theta$ ; only small amounts of bias (that is,  $\theta \in [0.8, 1.5]$ ) can occasionally improve performance.

**Comparison of different gradient estimators** Finally, we provide comparison of SAGA, B-SAGA with  $\theta = 10$ , SVRG, SARAH and SARGE, the results are provided below in Figure 3 from which we observe that

- SARAH performs similarly to SVRG, but is occasionally slower in early epochs.
- SARGE consistently outperforms all other methods except for B-SAGA with  $\theta = 10$ .

The above observations indicate that, depending on the data, biased schemes can benefit from their biased gradient estimates. The free parameter  $\theta$  reduces the MSE of the B-SAGA and B-SVRG gradient estimators leading to better performance, and the bias in SARAH and SARGE has a similar effect.

### 5.2 Non-convex objectives

To test the effect of bias in the non-convex setting, we consider the non-negative principal component analysis (NN-PCA) problems, which can be formulated as (Reddi et al., 2016b):

$$\min_{x \in \mathbb{R}^p} \left\{ F(x) - \frac{1}{n} \sum_{i=1}^n (h_i^\top x)^2 + \iota_C(x) \right\},$$

2. See <https://github.com/derekdriggs/StochOpt> for MATLAB scripts reproducing these experiments.

3. <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

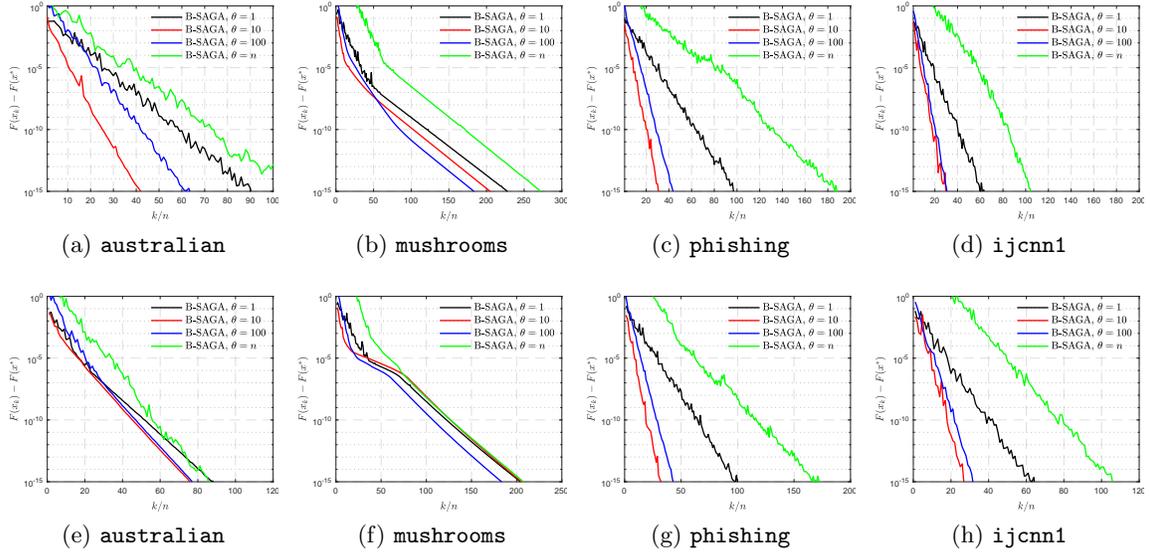


Figure 1: First row, performance comparison fitting a LASSO model for different choices of  $\theta$  in B-SAGA. Second row, performance comparison fitting a ridge regression model for different choices of  $\theta$  in B-SAGA. The step size for each case is set to  $\eta = \frac{1}{5L}$ .

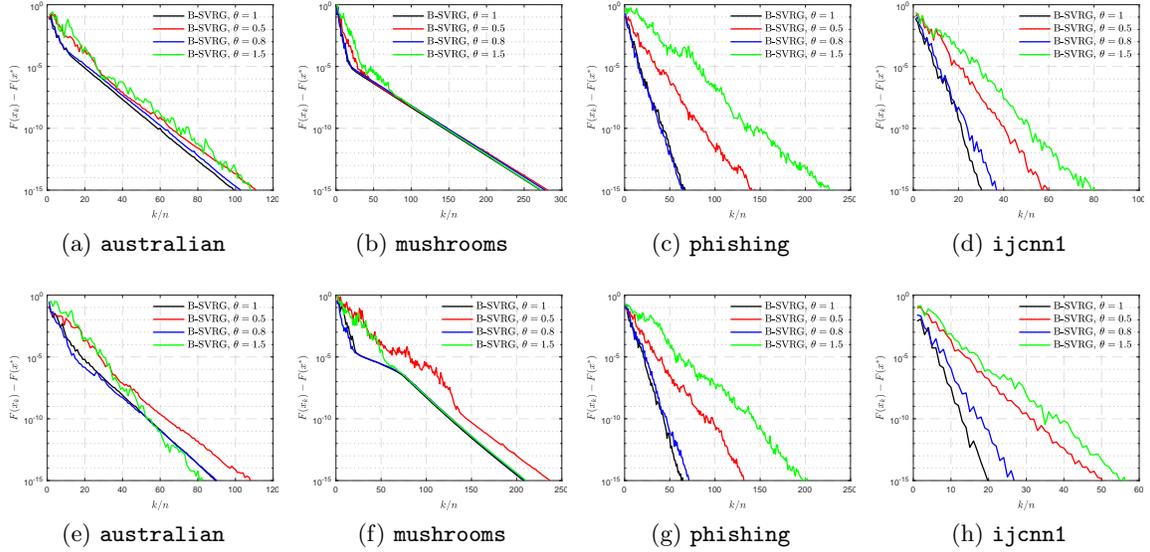


Figure 2: First row, performance comparison fitting a ridge regression model for different choices of  $\theta$  in B-SVRG. Second row, performance comparison fitting a LASSO model for different choices of  $\theta$  in B-SVRG. The step size for each case is set to  $\eta = \frac{1}{5L}$ .

where  $C \stackrel{\text{def}}{=} \{x \in \mathbb{R}^p : \|x\| \leq 1, x \geq 0\}$  is a convex set and  $\iota_C(x) = \begin{cases} 0 & : x \in C \\ +\infty & : x \notin C \end{cases}$  is the indicator function of  $C$ . Letting  $g = \iota_C$ , the operator  $\text{prox}_{\eta g}$  is the projection onto  $C$ , which can be computed efficiently.

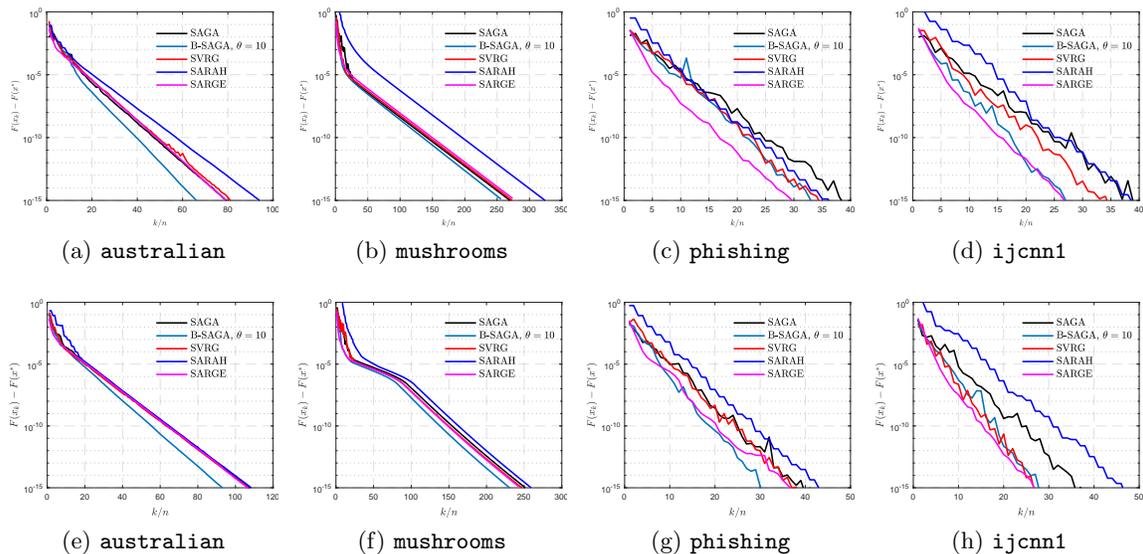


Figure 3: First row, performance comparison for solving ridge regression among different algorithms. Second row, performance comparison for solving LASSO regression among different algorithms. For both examples, step sizes are tuned automatically to minimize the number of iterations required to reach a suboptimality of  $10^{-15}$ .

As the problem is non-convex, we cannot measure convergence with respect to the global optimum  $x^*$ , so we use many iterations of proximal gradient descent with a small step size ( $\eta = \frac{1}{10Ln}$ ) to find a reference point  $x^*$ . Every test is initialized using a random vector with normally distributed i.i.d. entries, and the same starting point is used for testing each value of  $\theta$ . We found that small step sizes generally lead to stationary points with smaller objective values, so we set  $\eta = \frac{1}{5n}$  for all our experiments. We report  $F(x_k) - F(x^*)$  averaged over every  $n$  iterations. These experiments show that the performance of B-SAGA and B-SVRG varies significantly with  $\theta$ , with smaller values leading to better performance. SARAH and SARGE perform similarly to SAGA and SVRG in these experiments, see Figure 4.

For the comparison of all algorithms, B-SAGA and B-SVRG provides the best performance with B-SVRG being slightly faster.

## 6. Conclusion

The complicated convergence proofs of biased stochastic gradient methods have restricted researchers to studying unbiased estimators almost exclusively. Our simple framework for proving convergence rates for biased algorithms overcomes this limitation. Our analysis allows for the study of biased algorithms with proximal support for minimizing convex, strongly convex, and non-convex objectives for the first time.

We also show that biased gradient estimators can offer improvements over unbiased estimators in theory and in practice. Most notably, we find that biased recursive gradient estimators, such as SARAH and SARGE, admit smaller bounds on their MSEs and faster convergence rates than SAGA and SVRG.

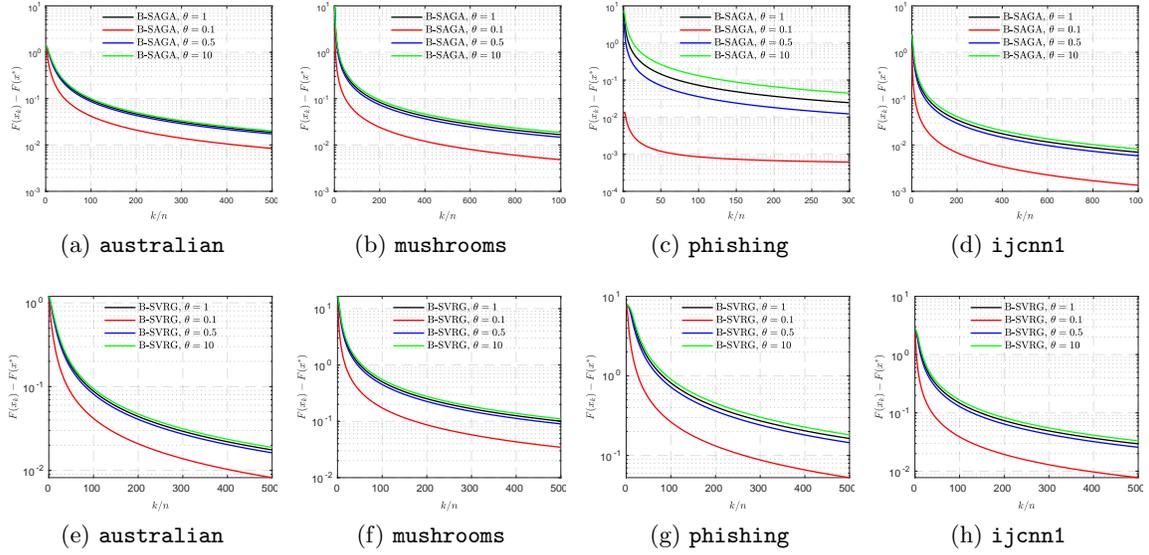


Figure 4: First row, performance comparison for solving NN-PCA with different choices of  $\theta$  in B-SAGA. Second row, performance comparison for solving NN-PCA with different choices of  $\theta$  in B-SVRG. The step size for each case is set to  $\eta = \frac{1}{5Ln}$ . The point  $x^*$  is found by solving the problem using proximal gradient descent.

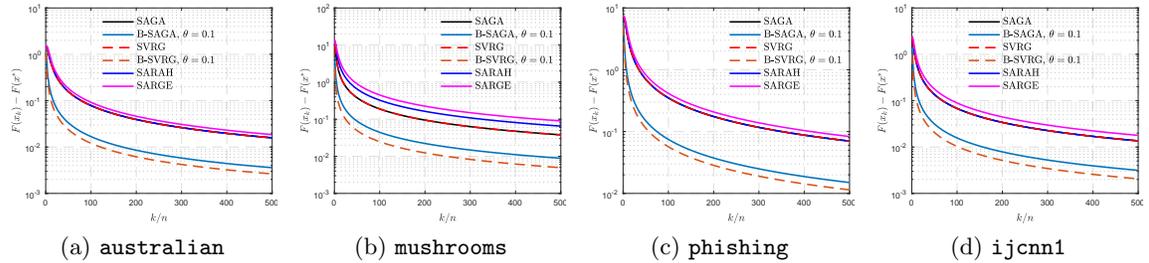


Figure 5: Performance comparison for solving NN-PCA among different algorithms. All step sizes are set to  $\frac{1}{5Ln}$ . Objective values are averaged over each epoch ( $n$  steps).

## Acknowledgements

CBS acknowledges support from the Leverhulme Trust project on Breaking the Non-Convexity Barrier and on Unveiling the Invisible, the Philip Leverhulme Prize, the EPSRC grant No. EP/M00483X/1, the EPSRC Centre No. EP/N014588/1, the European Union Horizon 2020 research and innovation programs under the Marie Skłodowska-Curie grant agreement No. 691070, the Cantab Capital Institute for the Mathematics of Information, and the Alan Turing Institute. JL acknowledges the support from Leverhulme Trust and Newton Trust.

## Appendix

The organization of the appendix is as follows: we prove Theorems 17 and 18 in Appendices A and B, respectively, and we prove Theorem 20 in Appendix C. We provide convergence rates for B-SAGA and B-SVRG as special cases of Theorem 17 in Appendix D, and we provide convergence rates for SARAH and SARGE as special cases of Theorem 18 in Appendices E and F, respectively.

### Appendix A. Proof of Theorem 17

To prove Theorem 17, we begin by showing that the BMSE property (Definition 9) implies the MSE of the gradient estimator over  $T$  iterations is proportional to  $\sum_{k=0}^{T-1} \mathbb{E} \|x_{k+1} - x_k\|^2$ .

**Lemma 32 (MSE bound)** *Suppose that the stochastic gradient estimator  $\tilde{\nabla}$  satisfies the  $\text{BMSE}(M_1, M_2, \rho_M, \rho_F, m)$  property, let  $\rho = \min\{\rho_M, \rho_F\}$ , and let  $\sigma_s$  be any sequence satisfying  $\sigma_s(1 - \rho)^{ms} \leq \sigma_{s-1}(1 - \frac{\rho}{2})^{ms}$ . For convenience, define  $\Theta = \frac{M_1\rho_F + 2M_2}{\rho_M\rho_F}$ . The MSE of the gradient estimator is bounded as*

$$\sum_{s=0}^S \sigma_s \sum_{k=ms}^{m(s+1)-1} \mathbb{E} [\|\nabla f(x_k) - \tilde{\nabla}_k\|^2] \leq 2\Theta L^2 \sum_{s=0}^S \sigma_s \sum_{k=ms}^{m(s+1)-1} \mathbb{E} [\|x_{k+1} - x_k\|^2].$$

**Proof** First, we derive a bound on the sequence  $\mathcal{F}_{ms}$  arising in the BMSE property. Summing this sequence from  $s = 0$  to  $s = S$ ,

$$\begin{aligned} \sum_{s=0}^S \sigma_s \mathcal{F}_{ms} &\leq \sum_{s=0}^S \sum_{\ell=0}^s \frac{M_2 \sigma_s (1 - \rho_F)^{m(s-\ell)}}{n} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2] \\ &\stackrel{\textcircled{1}}{\leq} \sum_{s=0}^S \sum_{\ell=0}^s \frac{M_2 \sigma_\ell (1 - \frac{\rho_F}{2})^{m(s-\ell)}}{n} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2] \\ &\leq \sum_{s=0}^S \left( \sum_{\ell=0}^\infty (1 - \frac{\rho_F}{2})^\ell \right) \frac{M_2 \sigma_s}{n} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2] \\ &= \sum_{s=0}^S \frac{2M_2 \sigma_s}{n\rho_F} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2]. \end{aligned} \quad (9)$$

Inequality  $\textcircled{1}$  uses the fact that  $\sigma_s(1 - \rho_F)^{ms} \leq \sigma_{s-1}(1 - \frac{\rho_F}{2})^{ms}$ . With this bound on  $\mathcal{F}_{ms}$ , we proceed to bound  $\mathcal{M}_{ms}$  similarly.

$$\begin{aligned} &\sum_{s=0}^S \sigma_s \mathcal{M}_{ms} \\ &\stackrel{\textcircled{1}}{\leq} \sum_{s=0}^S \sigma_s \left( \mathcal{F}_{ms} + \frac{M_1}{n} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2] \right) \\ &\quad + (1 - \rho_M)^m \sum_{s=1}^S \sigma_s \mathcal{M}_{m(s-1)} \\ &\stackrel{\textcircled{2}}{\leq} \sum_{s=0}^S \sigma_s \left( \frac{M_1\rho_F + 2M_2}{n\rho_F} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2] \right) \\ &\quad + (1 - \frac{\rho_M}{2})^m \sum_{s=1}^S \sigma_{s-1} \mathcal{M}_{m(s-1)} \\ &= \sum_{s=0}^S \sigma_s \left( \frac{M_1\rho_F + 2M_2}{n\rho_F} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2] \right) \\ &\quad + (1 - \frac{\rho_M}{2})^m \sum_{s=1}^S \sigma_{s-1} \left( \frac{M_1\rho_F + 2M_2}{n\rho_F} \sum_{k=m(s-1)}^{m(s)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2] \right) + \dots \\ &\leq \left( \sum_{\ell=0}^\infty (1 - \frac{\rho_M}{2})^{m\ell} \right) \sum_{s=0}^S \sigma_s \left( \frac{M_1\rho_F + 2M_2}{n\rho_F} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2] \right) \\ &\stackrel{\textcircled{3}}{\leq} \sum_{s=0}^S \frac{2\sigma_s \Theta}{n} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2] \\ &\stackrel{\textcircled{4}}{\leq} 2\Theta L^2 \sum_{s=0}^S \sigma_s \sum_{k=ms}^{m(s+1)-1} \mathbb{E} [\|x_{k+1} - x_k\|^2]. \end{aligned}$$

Inequality ① uses the fact that  $\mathcal{M}_m \leq (1 - \rho_M)^m \mathcal{M}_{m(s-1)}$ . Inequality ② uses  $\sigma_s(1 - \rho_M)^{ms} \leq \sigma_{s-1}(1 - \frac{\rho_M}{2})^{ms}$ , ③ uses the same estimate we applied in (9), and ④ uses the Lipschitz continuity of  $\nabla f_i$ .  $\blacksquare$

**Proof of Lemma 12** By assumption,  $1 - \frac{1}{\theta} \geq 0$ , so we can apply convexity to obtain

$$\begin{aligned} & \frac{\eta}{\theta}(f(x_k) - f(x^*)) + \frac{\eta}{n}(1 - \frac{1}{\theta})(\sum_{i=1}^n f_i(\varphi_k^i) - f_i(x^*)) \\ & \leq \frac{\eta}{\theta}\langle \nabla f(x_k), x_k - x^* \rangle + \frac{\eta}{n}(1 - \frac{1}{\theta}) \sum_{i=1}^n \langle \nabla f_i(\varphi_k^i), \varphi_k^i - x^* \rangle \\ & = \frac{\eta}{\theta}\langle \nabla f(x_k), x_k - x^* \rangle + \frac{\eta}{n}(1 - \frac{1}{\theta}) \sum_{i=1}^n \langle \nabla f_i(\varphi_k^i), x_k - x^* \rangle + \frac{\eta}{n}(1 - \frac{1}{\theta}) \sum_{i=1}^n \langle \nabla f_i(\varphi_k^i), \varphi_k^i - x_k \rangle. \end{aligned}$$

Because  $\tilde{\nabla}_k$  is memory-biased, hence  $\frac{1}{\theta}\nabla f(x_k) + \frac{1}{n}(1 - \frac{1}{\theta})\sum_{i=1}^n \nabla f_i(\varphi_k^i) = \mathbb{E}_k[\tilde{\nabla}_k]$ . Therefore,

$$\begin{aligned} & \frac{\eta}{\theta}\langle \nabla f(x_k), x_k - x^* \rangle + \frac{\eta}{n}(1 - \frac{1}{\theta}) \sum_{i=1}^n \langle \nabla f_i(\varphi_k^i), x_k - x^* \rangle \\ & = \mathbb{E}_k[\eta\langle \tilde{\nabla}_k, x_k - x^* \rangle] \\ & = \mathbb{E}_k[\eta\langle \tilde{\nabla}_k, x_k - x_{k+1} \rangle + \eta\langle \tilde{\nabla}_k, x_{k+1} - x^* \rangle] \\ & \leq \mathbb{E}_k[\eta\langle \tilde{\nabla}_k, x_k - x_{k+1} \rangle - \frac{1}{2}\|x_{k+1} - x_k\|^2 + \frac{1}{2}\|x_k - x^*\|^2 - \frac{1+\mu\eta}{2}\|x_{k+1} - x^*\|^2 - \eta g(x_{k+1}) + \eta g(x^*)]. \end{aligned}$$

The inequality is due to Lemma 3 with  $z = x_{k+1}$ ,  $x = x_k$ ,  $d = \tilde{\nabla}_k$ , and  $y = x^*$ . Combining these two inequalities, we have shown

$$\begin{aligned} & \frac{\eta}{\theta}(f(x_k) - f(x^*)) + \frac{\eta}{n}(1 - \frac{1}{\theta}) \sum_{i=1}^n (f_i(\varphi_k^i) - f_i(x^*)) \\ & \leq \mathbb{E}_k[\eta\langle \tilde{\nabla}_k, x_k - x_{k+1} \rangle - \frac{1}{2}\|x_{k+1} - x_k\|^2 - \eta g(x_{k+1}) + \eta g(x^*)] \\ & \quad + \frac{1}{2}\|x_k - x^*\|^2 - \frac{1+\mu\eta}{2}\|x_{k+1} - x^*\|^2 + \frac{\eta}{n}(1 - \frac{1}{\theta}) \sum_{i=1}^n \langle \nabla f_i(\varphi_k^i), \varphi_k^i - x_k \rangle. \end{aligned} \tag{10}$$

We bound the first three terms on the right further.

$$\begin{aligned} & \eta\langle \tilde{\nabla}_k, x_k - x_{k+1} \rangle - \frac{1}{2}\|x_{k+1} - x_k\|^2 - \eta g(x_{k+1}) \\ & = \eta(\langle \nabla f(x_k), x_k - x_{k+1} \rangle - g(x_{k+1})) + \eta\langle \tilde{\nabla}_k - \nabla f(x_k), x_k - x_{k+1} \rangle - \frac{1}{2}\|x_{k+1} - x_k\|^2 \\ & \stackrel{\text{①}}{\leq} \eta(f(x_k) - F(x_{k+1})) + \eta\langle \tilde{\nabla}_k - \nabla f(x_k), x_k - x_{k+1} \rangle + (\frac{\eta L}{2} - \frac{1}{2})\|x_{k+1} - x_k\|^2 \\ & \stackrel{\text{②}}{\leq} \eta(f(x_k) - F(x_{k+1})) + \frac{\eta}{2L\lambda}\|\tilde{\nabla}_k - \nabla f(x_k)\|^2 + (\frac{\eta L(\lambda+1)}{2} - \frac{1}{2})\|x_{k+1} - x_k\|^2. \end{aligned}$$

Inequality ① is due to the Lipschitz continuity of  $\nabla f$ , and inequality ② is Young's. Combining this bound with (10) and rearranging terms, we have shown that

$$\begin{aligned} 0 & \leq -\eta\mathbb{E}_k[F(x_{k+1}) - F(x^*)] + \frac{\eta}{2L\lambda}\mathbb{E}_k[\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] \\ & \quad - \frac{1+\mu\eta}{2}\mathbb{E}_k[\|x_{k+1} - x^*\|^2] + \frac{1}{2}\|x_k - x^*\|^2 + (\frac{\eta L(\lambda+1)}{2} - \frac{1}{2})\mathbb{E}_k[\|x_{k+1} - x_k\|^2] \\ & \quad + \eta(1 - \frac{1}{\theta})(f(x_k) - \frac{1}{n}\sum_{i=1}^n f_i(\varphi_k^i) + \frac{1}{n}\sum_{i=1}^n \langle \nabla f_i(\varphi_k^i), \varphi_k^i - x_k \rangle). \end{aligned}$$

We use Lemma 1 to bound the final term, yielding the desired inequality.  $\blacksquare$

**Proof of Theorem 17 (Convex Case)** We begin with the inequality of Lemma 12 with  $\mu = 0$ . Multiplying the inequality of Lemma 4 with  $z = x_{k+1}$ ,  $x = x_k$ , and  $d = \tilde{\nabla}_k$  by a

non-negative constant  $\delta$  and adding it to the inequality of Lemma 12, we obtain

$$\begin{aligned} & \eta \mathbb{E}_k [F(x_{k+1}) - F(x^*) + \delta(F(x_{k+1}) - F(x_k))] \\ & \leq \frac{\eta(1+\delta)}{2L\lambda} \mathbb{E}_k [\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] - \frac{1}{2} \mathbb{E}_k [\|x_{k+1} - x^*\|^2] + \frac{1}{2} \|x_k - x^*\|^2 \\ & \quad + \left( \frac{\eta L(1+\delta)(\lambda+1)}{2} - \frac{1+2\delta}{2} \right) \mathbb{E}_k [\|x_{k+1} - x_k\|^2] + \frac{\eta L}{2n} \left(1 - \frac{1}{\theta}\right) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2. \end{aligned} \quad (11)$$

Applying the full expectation operator and summing from  $k = 0$  to  $k = T - 1$ , we have

$$\begin{aligned} & \sum_{k=0}^{T-1} \eta \mathbb{E} [F(x_{k+1}) - F(x^*)] + \eta \delta (\mathbb{E} [F(x_T)] - F(x_0)) \\ & \leq -\frac{1}{2} \mathbb{E} [\|x_T - x^*\|^2] + \frac{1}{2} \|x_0 - x^*\|^2 + \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{\eta(1+\delta)}{2L\lambda} \|\tilde{\nabla}_k - \nabla f(x_k)\|^2 \right] \\ & \quad + \left( \frac{\eta L(1+\delta)(\lambda+1)}{2} - \frac{1+2\delta}{2} \right) \mathbb{E} [\|x_{k+1} - x_k\|^2] + \frac{\eta L}{2n} \left(1 - \frac{1}{\theta}\right) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2. \end{aligned}$$

We use Lemma 32 with  $\sigma_s = 1$  to bound the MSE, and we use the fact that the gradient estimator is memory-biased to bound the term  $1/n \sum_{i=1}^n \|x_k - \varphi_k^i\|^2$ . This leaves

$$\begin{aligned} & \eta \sum_{k=0}^{T-1} \mathbb{E} [F(x_{k+1}) - F(x^*)] \\ & \leq -\frac{1}{2} \mathbb{E} [\|x_T - x^*\|^2] + \frac{1}{2} \|x_0 - x^*\|^2 + \eta \delta (F(x_0) - \mathbb{E} [F(x_T)]) \\ & \quad + \left( \frac{\eta L(1+\delta)(\lambda+1)}{2} + \frac{\Theta \eta L(1+\delta)}{\lambda} + \frac{B_1 \eta L}{2} \left(1 - \frac{1}{\theta}\right) - \frac{1+2\delta}{2} \right) \sum_{k=0}^{T-1} \mathbb{E} [\|x_{k+1} - x_k\|^2]. \end{aligned} \quad (12)$$

Setting  $\lambda = \sqrt{2\Theta}$  minimizes the coefficient of the term on the final line. With

$$\eta \leq \frac{1}{L(1+2\sqrt{2\Theta} + \frac{B_1(1-1/\theta)}{1+\delta})},$$

the final term in (12) is non-positive, so we can drop it from the inequality along with the term  $-1/2 \mathbb{E} \|x_T - x^*\|^2$ . Using the fact that  $-F(x_T) \leq -F(x^*)$ , this leaves

$$\sum_{k=0}^{T-1} \mathbb{E} [F(x_{k+1}) - F(x^*)] \leq \frac{1}{2\eta} \|x_0 - x^*\|^2 + \delta (F(x_0) - F(x^*)).$$

We use the convexity of  $F$  to rewrite this inequality as a bound on the suboptimality of the average iterate

$$\mathbb{E} [F(\bar{x}_T) - F(x^*)] \leq \frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E} [F(x_{k+1}) - F(x^*)] \leq \frac{1}{2\eta T} \|x_0 - x^*\|^2 + \frac{\delta}{T} (F(x_0) - F(x^*)).$$

Setting  $\delta = \max\{B_1(1 - 1/\theta)/\sqrt{2\Theta} - 1, 0\}$  approximately minimizes the right side of this inequality, completing the proof.  $\blacksquare$

**Proof of Theorem 17 (Strongly Convex Case)** As in the proof of the convex case, we begin with the inequality of Lemma 12, multiply the inequality of Lemma 4 with  $z = x_{k+1}$ ,  $x = x_k$ , and  $d = \tilde{\nabla}_k$  by a non-negative constant  $\delta$ , and add the two inequalities.

$$\begin{aligned} & \eta \mathbb{E}_k [F(x_{k+1}) - F(x^*) + \delta(F(x_{k+1}) - F(x_k))] \\ & \leq -\frac{1+\mu\eta}{2} \mathbb{E}_k [\|x_{k+1} - x^*\|^2] + \frac{1}{2} \|x_k - x^*\|^2 + \mathbb{E}_k \left[ \frac{\eta(1+\delta)}{2L\lambda} \|\tilde{\nabla}_k - \nabla f(x_k)\|^2 \right] \\ & \quad + \left( \frac{\eta L(1+\delta)(\lambda+1)}{2} - \frac{1+\delta(2+\mu\eta)}{2} \right) \mathbb{E}_k [\|x_{k+1} - x_k\|^2] + \frac{\eta L}{2n} \left(1 - \frac{1}{\theta}\right) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2. \end{aligned}$$

Applying the full expectation operator, multiplying by  $(1 + \mu\eta)^k$ , and summing over the epoch  $k = ms$  to  $k = m(s+1) - 1$  for some  $s \in \mathbb{N}_0$ , we have

$$\begin{aligned} & \eta \sum_{k=ms}^{m(s+1)-1} (1 + \mu\eta)^k \mathbb{E}[F(x_{k+1}) - F(x^*) + \delta(F(x_{k+1}) - F(x_k))] \\ & \leq -\frac{(1+\mu\eta)^{m(s+1)}}{2} \mathbb{E}\|x_{m(s+1)} - x^*\|^2 + \frac{(1+\mu\eta)^{ms}}{2} \mathbb{E}\|x_{ms} - x^*\|^2 \\ & \quad + \sum_{k=ms}^{m(s+1)-1} (1 + \mu\eta)^k \mathbb{E}\left[\frac{\eta(1+\delta)}{2L\lambda} \|\tilde{\nabla}_k - \nabla f(x_k)\|^2 + \left(\frac{\eta L(1+\delta)(\lambda+1)}{2} - \frac{1+\delta(2+\mu\eta)}{2}\right) \|x_{k+1} - x_k\|^2\right] \\ & \quad + \frac{\eta L}{2n} \left(1 - \frac{1}{\theta}\right) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2. \end{aligned}$$

Using the fact that  $\eta \leq \frac{1}{\mu m}$ ,

$$(1 + \mu\eta)^k \leq (1 + \mu\eta)^{m(s+1)} \leq (1 + \mu\eta)^{ms} \lim_{m \rightarrow \infty} \left(1 + \frac{1}{m}\right)^m = e(1 + \mu\eta)^{ms} \leq 3(1 + \mu\eta)^{ms}, \quad (13)$$

where  $e$  is Euler's number. Therefore,

$$\begin{aligned} & \eta \sum_{k=ms}^{m(s+1)-1} (1 + \mu\eta)^k \mathbb{E}[F(x_{k+1}) - F(x^*) + \delta(F(x_{k+1}) - F(x_k))] \\ & \leq -\frac{(1+\mu\eta)^{m(s+1)}}{2} \mathbb{E}\|x_{m(s+1)} - x^*\|^2 + \frac{(1+\mu\eta)^{ms}}{2} \mathbb{E}\|x_{ms} - x^*\|^2 \\ & \quad + (1 + \mu\eta)^{ms} \sum_{k=ms}^{m(s+1)-1} \mathbb{E}\left[\left(\frac{3\eta L(1+\delta)(\lambda+1)}{2} - \frac{1+\delta(2+\mu\eta)}{2}\right) \|x_{k+1} - x_k\|^2\right] \\ & \quad + \frac{3\eta(1+\delta)}{2L\lambda} \mathbb{E}\|\tilde{\nabla}_k - \nabla f(x_k)\|^2 + \frac{3\eta L}{2n} \left(1 - \frac{1}{\theta}\right) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2. \end{aligned} \quad (14)$$

Summing the inequality from epoch  $s = 0$  to  $s = S - 1$ ,

$$\begin{aligned} & \eta \sum_{k=0}^{mS-1} (1 + \mu\eta)^k \mathbb{E}[F(x_{k+1}) - F(x^*) + \delta(F(x_{k+1}) - F(x_k))] \\ & \leq \sum_{s=0}^{S-1} (1 + \mu\eta)^{ms} \sum_{k=ms}^{m(s+1)-1} \mathbb{E}\left[\frac{3\eta(1+\delta)}{2L\lambda} \|\tilde{\nabla}_k - \nabla f(x_k)\|^2 + \left(\frac{3\eta L(1+\delta)(\lambda+1)}{2} - \frac{1+\delta(2+\mu\eta)}{2}\right) \|x_{k+1} - x_k\|^2\right] \\ & \quad + \frac{3\eta L(1+\delta)}{2n} \left(1 - \frac{1}{\theta}\right) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2 - \frac{(1+\mu\eta)^{mS}}{2} \mathbb{E}\|x_{mS} - x^*\|^2 + \frac{1}{2} \|x_0 - x^*\|^2. \end{aligned}$$

We use Lemma 32 with  $\sigma_s = (1 + \mu\eta)^{ms}$  to bound the MSE. Recall  $\rho = \min\{\rho_M, \rho_F\}$  and  $\eta \leq \frac{\rho}{2\mu}$ . This choice for  $\sigma_s$  satisfies the conditions of Lemma 32 because  $(1 + \mu\eta)^{ms} (1 - \rho)^{ms} \leq (1 + \mu\eta)^{m(s-1)} (1 - \rho/2)^{ms}$ . We use the fact that the gradient estimator is memory-biased to bound the term  $1/n \sum_{i=1}^n \|x_k - \varphi_k^i\|^2$ . This leaves

$$\begin{aligned} & \eta \sum_{k=0}^{mS-1} (1 + \mu\eta)^k \mathbb{E}[F(x_{k+1}) - F(x^*) + \delta(F(x_{k+1}) - F(x_k))] \\ & \leq -\frac{(1+\mu\eta)^{mS}}{2} \mathbb{E}\|x_{mS} - x^*\|^2 + \frac{1}{2} \|x_0 - x^*\|^2 + C \sum_{s=0}^{S-1} (1 + \mu\eta)^{ms} \sum_{k=ms}^{m(s+1)-1} \mathbb{E}\|x_{k+1} - x_k\|^2, \end{aligned} \quad (15)$$

where  $C = \frac{3\eta L(1+\delta)(\lambda+1)}{2} + \frac{3\Theta\eta L(1+\delta)}{\lambda} + \frac{3B_1\eta L}{2} \left(1 - \frac{1}{\theta}\right) - \frac{1+\delta(2+\mu\eta)}{2}$ . We must choose  $\eta, \lambda$ , and  $\delta$  so that  $C \leq 0$ . Setting  $\lambda = \sqrt{2\Theta}$  minimizes  $C$  over  $\lambda$ . Using the approximation  $\delta(2 + \mu\eta) \geq \delta$ , we see that  $C$  is non-positive if

$$\eta \leq \frac{1}{3L(1+2\sqrt{2\Theta} + \frac{B_1(1-1/\theta)}{1+\delta})}.$$

Setting  $\delta = \max\{B_1(1 - 1/\theta)/\sqrt{2\Theta} - 1, 0\}$ , we are guaranteed that

$$\frac{1}{3L(1+3\sqrt{2\Theta})} \leq \frac{1}{3L(1+2\sqrt{2\Theta} + \frac{B_1(1-1/\theta)}{1+\delta})},$$

so the step size in the theorem statement ensures  $C \leq 0$ , and the final term in (15) is non-positive. Dropping this non-positive term from the inequality, we have

$$\begin{aligned} & \eta(1 + \delta) \sum_{k=0}^{mS-1} (1 + \mu\eta)^k \mathbb{E}[F(x_{k+1}) - F(x^*)] + \delta\eta \sum_{k=0}^{mS-1} (1 + \mu\eta)^k \mathbb{E}[F(x_k) - F(x^*)] \\ & \leq -\frac{(1+\mu\eta)^{mS}}{2} \mathbb{E}[\|x_{mS} - x^*\|^2] + \frac{1}{2} \|x_0 - x^*\|^2. \end{aligned} \quad (16)$$

We would like to show that  $1 + \delta \geq (1 + \mu\eta)\delta$  so that the terms on the first line telescope.

We use the fact that  $\eta \leq \frac{\sqrt{2\Theta}}{B_1\mu(1-1/\theta)}$  to say

$$\frac{1}{\mu\eta} \geq \frac{B_1(1-1/\theta)}{\sqrt{2\Theta}} \geq \delta$$

Hence  $\frac{1+\delta}{\delta} \geq 1 + \mu\eta$  and inequality (16) simplifies to

$$(1 + \mu\eta)^{mS} \mathbb{E}[\eta\delta(F(x_{mS}) - F(x^*)) + \frac{1}{2}\|x_{mS} - x^*\|^2] \leq \eta\delta(F(x_0) - F(x^*)) + \frac{1}{2}\|x_0 - x^*\|^2,$$

which implies the result.  $\blacksquare$

## Appendix B. Proof of Theorem 18

The following two lemmas establish an analogue of Lemma 12 for recursively biased estimators.

**Lemma 33** *Suppose  $\tilde{\nabla}$  is recursively biased with parameters  $\rho_B$  and  $\nu$ . Suppose  $g$  is  $\mu$ -strongly convex with  $\mu \geq 0$ , and let  $\lambda > 0$  be a constant whose value we determine later. The following inequality holds:*

$$\begin{aligned} 0 & \leq -\eta \mathbb{E}_k[F(x_{k+1}) - F(x^*)] + \frac{\eta}{2L\lambda} \mathbb{E}_k[\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] \\ & \quad - \frac{1+\mu\eta}{2} \mathbb{E}_k[\|x_{k+1} - x^*\|^2] + \frac{1}{2} \|x_k - x^*\|^2 \\ & \quad + \left(\frac{\eta L(\lambda+1)}{2} - \frac{1}{2}\right) \mathbb{E}_k[\|x_{k+1} - x_k\|^2] + \eta(1 - \rho_B) \langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle. \end{aligned}$$

**Proof** Applying the convexity of  $f$  yields

$$\begin{aligned} & \eta(f(x_k) - f(x^*)) \\ & \leq \eta \langle \nabla f(x_k), x_k - x^* \rangle \\ & = \eta \langle \nabla f(x_k) - (1 - \rho_B)(\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}), x_k - x^* \rangle + \eta(1 - \rho_B) \langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle. \end{aligned}$$

Because the estimator is recursively biased,

$$\mathbb{E}_k[\tilde{\nabla}_k] = \nabla f(x_k) - (1 - \rho_B)(\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}).$$

Therefore,

$$\begin{aligned} & \eta \langle \nabla f(x_k) - (1 - \rho_B)(\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}), x_k - x^* \rangle \\ & = \mathbb{E}_k[\eta \langle \tilde{\nabla}_k, x_k - x^* \rangle] \\ & = \mathbb{E}_k[\eta \langle \tilde{\nabla}_k, x_k - x_{k+1} \rangle + \eta \langle \tilde{\nabla}_k, x_{k+1} - x^* \rangle] \\ & \leq \mathbb{E}_k[\eta \langle \tilde{\nabla}_k, x_k - x_{k+1} \rangle - \frac{1}{2} \|x_{k+1} - x_k\|^2 + \frac{1}{2} \|x_k - x^*\|^2 - \frac{1}{2} \|x_{k+1} - x^*\|^2 + \eta g(x_{k+1}) - \eta g(x^*)], \end{aligned}$$

The inequality is due to Lemma 3. The rest of the proof follows the proof of Lemma 12. ■

### Proof of Lemma 15

We begin with a technical lemma.

**Lemma 34** *Let  $\{E_\ell\}_{\ell=0}^\infty$  be a sequence of non-negative real numbers. For any  $\rho \in (0, 1]$  and  $c \geq 1$ , the following inequality holds:*

$$\sum_{k=cs+2}^{c(s+1)-1} \sum_{\ell=cs+1}^{k-1} (1-\rho)^{k-\ell-1} E_\ell \leq \min\{c, 1/\rho\} \sum_{k=cs+1}^{c(s+1)-1} E_k.$$

**Proof** This inequality can be seen by expanding the sums.

$$\begin{aligned} & \sum_{k=cs+2}^{c(s+1)-1} \sum_{\ell=cs+1}^{k-1} (1-\rho)^{k-\ell-1} E_\ell \\ &= [E_{cs+1}] + [(1-\rho)E_{cs+1} + E_{cs+2}] + [(1-\rho)^2 E_{cs+1} + (1-\rho)E_{cs+2} + E_{cs+3}] + \cdots \\ & \quad + [(1-\rho)^{c-2} E_{cs+1} + \cdots + E_{c(s+1)-2}] \\ &= (1 + (1-\rho) + \cdots + (1-\rho)^{c-2}) E_{cs+1} + (1 + (1-\rho) + \cdots + (1-\rho)^{c-3}) E_{cs+2} + \cdots \end{aligned}$$

The coefficient of each term of the sequence  $E_\ell$  is less than  $c-2$ . Therefore, replacing each coefficient with  $c$ , we obtain the bound

$$\sum_{k=cs+2}^{c(s+1)-1} \sum_{\ell=cs+1}^{k-1} (1-\rho)^{k-\ell-1} E_\ell \leq c \sum_{k=cs+1}^{c(s+1)-1} E_k.$$

This proves the first bound. For the second bound, notice that each coefficient is less than  $\sum_{k=0}^\infty (1-\rho)^k = 1/\rho$ , proving the second bound. ■

First, we use the fact that  $\tilde{\nabla}_k$  is recursively biased.

$$\begin{aligned} & \mathbb{E}\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle \\ & \stackrel{\textcircled{1}}{=} \mathbb{E}[\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x_{k-1} \rangle + \langle \nabla f(x_{k-1}) - \mathbb{E}_{k-1} \tilde{\nabla}_{k-1}, x_{k-1} - x^* \rangle] \\ & \stackrel{\textcircled{2}}{\leq} \mathbb{E}\left[\frac{\epsilon}{2} \|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}\|^2 + \frac{1}{2\epsilon} \|x_k - x_{k-1}\|^2 + \langle \nabla f(x_{k-1}) - \mathbb{E}_{k-1} \tilde{\nabla}_{k-1}, x_{k-1} - x^* \rangle\right] \\ & \stackrel{\textcircled{3}}{=} \mathbb{E}\left[\frac{\epsilon}{2} \|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}\|^2 + \frac{1}{2\epsilon} \|x_k - x_{k-1}\|^2 + (1-\rho_B) \langle \nabla f(x_{k-2}) - \tilde{\nabla}_{k-2}, x_{k-1} - x^* \rangle\right]. \end{aligned}$$

We can pass the conditional expectation  $\mathbb{E}_{k-1}$  into the second inner-product in  $\textcircled{1}$  because  $x_{k-1}$  is independent of  $j_{k-1}$ . Inequality  $\textcircled{2}$  is Young's, and  $\textcircled{3}$  uses the definition of a recursively biased gradient estimator.

This is a recursive inequality, and expanding the recursion gives

$$\begin{aligned} \mathbb{E}\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle & \leq \sum_{\ell=\nu_s+1}^{k-1} (1-\rho_B)^{k-\ell-1} \mathbb{E}\left[\frac{\epsilon}{2} \|\nabla f(x_\ell) - \tilde{\nabla}_\ell\|^2 + \frac{1}{2\epsilon} \|x_{\ell+1} - x_\ell\|^2\right] \\ & \quad + (1-\rho_B) \langle \nabla f(x_{\nu_s}) - \tilde{\nabla}_{\nu_s}, x_{\nu_s+1} - x^* \rangle \\ & \stackrel{\textcircled{1}}{=} \sum_{\ell=\nu_s+1}^{k-1} (1-\rho_B)^{k-\ell-1} \mathbb{E}\left[\frac{\epsilon}{2} \|\nabla f(x_\ell) - \tilde{\nabla}_\ell\|^2 + \frac{1}{2\epsilon} \|x_{\ell+1} - x_\ell\|^2\right]. \end{aligned}$$

Equality ① is due to the fact that  $\tilde{\nabla}_{\nu s} = \nabla f(x_{\nu s})$ . Taking the absolute value and summing this from  $k = \nu s + 1$  to  $k = \nu(s + 1) - 1$ ,

$$\begin{aligned} & \sum_{k=\nu s+1}^{\nu(s+1)-1} |\mathbb{E}\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle| \\ & \leq \sum_{k=\nu s+2}^{\nu(s+1)-1} \sum_{\ell=\nu s+1}^{k-1} (1 - \rho_B)^{k-\ell-1} \mathbb{E} \left[ \frac{\epsilon}{2} \|\nabla f(x_\ell) - \tilde{\nabla}_\ell\|^2 + \frac{1}{2\epsilon} \|x_{\ell+1} - x_\ell\|^2 \right] \\ & \stackrel{\textcircled{2}}{\leq} \min \left\{ \nu, \frac{1}{\rho_B} \right\} \sum_{k=\nu s+1}^{\nu(s+1)-1} \mathbb{E} \left[ \frac{\epsilon}{2} \|\nabla f(x_k) - \tilde{\nabla}_k\|^2 + \frac{1}{2\epsilon} \|x_{k+1} - x_k\|^2 \right]. \end{aligned}$$

Inequality ② follows from the technical Lemma 34. Summing this inequality from  $s = 0$  to  $s = S$  completes the proof.  $\blacksquare$

**Proof of Theorem 18 (Convex Case)** To begin, we sum the inequality of Lemma 33 and the inequality of Lemma 4 scaled by  $\delta > 0$  with  $z = x_{k+1}$ ,  $x = x_k$ , and  $d = \tilde{\nabla}_k$ .

$$\begin{aligned} & \eta \mathbb{E}_k [F(x_{k+1}) - F(x^*) + \delta(F(x_{k+1}) - F(x_k))] \\ & \leq -\frac{1}{2} \mathbb{E}_k [\|x_{k+1} - x^*\|^2] + \frac{1}{2} \|x_k - x^*\|^2 + \mathbb{E}_k \left[ \frac{\eta(1+\delta)}{2L\lambda} \|\tilde{\nabla}_k - \nabla f(x_k)\|^2 \right] \\ & \quad + (1 + \delta) \left( \frac{\eta L(\lambda+1)}{2} - \frac{1}{2} \right) \|x_{k+1} - x_k\|^2 + \eta(1 - \rho_B) \langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle. \end{aligned} \quad (17)$$

Applying the full expectation operator, setting  $\mu = 0$ , and summing from  $k = 0$  to  $k = T - 1$  where  $T = mS$  for some  $S \in \mathbb{N}$ , we have

$$\begin{aligned} & \eta \sum_{k=0}^{T-1} \mathbb{E} [F(x_{k+1}) - F(x^*)] + \eta \delta \mathbb{E} [F(x_T) - F(x_0)] \\ & \leq -\frac{1}{2} \mathbb{E} [\|x_T - x^*\|^2] + \frac{1}{2} \|x_0 - x^*\|^2 + \sum_{k=0}^{T-1} \mathbb{E} \left[ \frac{\eta(1+\delta)}{2L\lambda} \|\tilde{\nabla}_k - \nabla f(x_k)\|^2 \right] \\ & \quad + (1 + \delta) \left( \frac{\eta L(\lambda+1)}{2} - \frac{1}{2} \right) \|x_{k+1} - x_k\|^2 + \eta(1 - \rho_B) \langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle. \end{aligned}$$

We use Lemma 15 to bound the inner-product bias term.

$$\begin{aligned} & \eta \sum_{k=0}^{T-1} \mathbb{E} [F(x_{k+1}) - F(x^*)] + \eta \delta \mathbb{E} [F(x_T) - F(x_0)] \\ & \leq -\frac{1}{2} \mathbb{E} [\|x_T - x^*\|^2] + \frac{1}{2} \|x_0 - x^*\|^2 + \sum_{k=0}^{T-1} \mathbb{E} \left[ \left( \frac{\eta(1+\delta)}{2L\lambda} + \frac{B_2\eta(1-\rho_B)\epsilon}{2} \right) \|\tilde{\nabla}_k - \nabla f(x_k)\|^2 \right] \\ & \quad + (1 + \delta) \left( \frac{\eta L(\lambda+1)}{2} + \frac{B_2\eta(1-\rho_B)}{2\epsilon(1+\delta)} - \frac{1}{2} \right) \|x_{k+1} - x_k\|^2. \end{aligned}$$

To bound the MSE, we use Lemma 32 with  $\sigma_s = 1$ . This leaves

$$\begin{aligned} & \eta \sum_{k=0}^{T-1} \mathbb{E} [F(x_{k+1}) - F(x^*)] + \eta \delta \mathbb{E} [F(x_T) - F(x_0)] \\ & \leq -\frac{1}{2} \mathbb{E} [\|x_T - x^*\|^2] + \frac{1}{2} \|x_0 - x^*\|^2 + w \sum_{k=0}^{T-1} \mathbb{E} [\|x_{k+1} - x_k\|^2], \end{aligned} \quad (18)$$

where  $w = \frac{\eta L(\lambda+1)(1+\delta)}{2} + \frac{B_2\eta(1-\rho_B)}{2\epsilon} + \frac{\Theta\eta L(1+\delta)}{\lambda} + B_2\eta L^2(1 - \rho_B)\epsilon\Theta - \frac{1+\delta}{2}$ . To minimize the coefficient of the final term, we set  $\lambda = \sqrt{2\Theta}$  and  $\epsilon = (2L^2\Theta)^{-1/2}$ . This coefficient is then equal to

$$\sqrt{2\Theta}\eta L(1 + \delta) + \frac{\eta L(1+\delta)}{2} + \sqrt{2}(1 - \rho_B)\eta L B_2 \sqrt{\Theta} - \frac{1+\delta}{2},$$

which is non-positive when  $\eta \leq \frac{1}{2\sqrt{2\Theta}L(1+\frac{(1-\rho_B)B_2}{1+\delta})+L}$ . This ensures that the final term in (18) is non-positive, so we can drop it from the inequality along with the term  $-1/2\mathbb{E}\|x_T - x^*\|^2$ . This leaves

$$\sum_{k=0}^{T-1} \mathbb{E}[F(x_{k+1}) - F(x^*)] \leq \frac{1}{2\eta} \|x_0 - x^*\|^2 + \delta \mathbb{E}[F(x_0) - F(x_T)].$$

By the convexity of  $F$  and the fact that  $-F(x_T) \leq -F(x^*)$

$$\mathbb{E}[F(\bar{x}_T) - F(x^*)] \leq \frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E}[F(x_{k+1}) - F(x^*)] \leq \frac{1}{2\eta T} \|x_0 - x^*\|^2 + \frac{\delta}{T} (F(x_0) - F(x^*)).$$

We now choose  $\delta$  so that this upper-bound is approximately minimized. Assuming  $\|x_0 - x^*\|^2$  and  $F(x_0) - F(x^*)$  are approximately equal to  $K > 0$ , the right side becomes

$$\frac{K}{T} \left( \frac{2\sqrt{2\Theta}L(1-\rho_B)B_2}{1+\delta} + \delta + 2\sqrt{2\Theta}L + L \right).$$

Choosing  $\delta = \max\{\sqrt{L\sqrt{\Theta}(1-\rho_B)B_2} - 1, 0\}$  approximately minimizes the right side of this inequality, completing the proof.  $\blacksquare$

**Proof of Theorem 18 (Strongly Convex Case)** We begin with inequality (17), but without setting  $\mu = 0$ .

$$\begin{aligned} & \eta(1+\delta)\mathbb{E}_k[F(x_{k+1}) - F(x^*)] + \frac{1+\mu\eta}{2}\mathbb{E}_k\|x_{k+1} - x^*\|^2 \\ & \leq \eta\delta(F(x_k) - F(x^*)) + \frac{1}{2}\|x_k - x^*\|^2 + \mathbb{E}_k\left[\frac{\eta(1+\delta)}{2L\lambda}\|\tilde{\nabla}_k - \nabla f(x_k)\|^2\right] \\ & \quad + (1+\delta)\left(\frac{\eta L(\lambda+1)}{2} - \frac{1}{2}\right)\|x_{k+1} - x_k\|^2 + \eta(1-\rho_B)\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle. \end{aligned}$$

Applying the full expectation operator, multiplying by  $(1+\mu\eta)^k$ , and summing over the epoch  $k = ms$  to  $k = m(s+1) - 1$  for some  $s \in \mathbb{N}_0$ , we have

$$\begin{aligned} & \eta(1+\delta) \sum_{k=ms}^{m(s+1)-1} (1+\mu\eta)^k \mathbb{E}[F(x_{k+1}) - F(x^*)] + \frac{(1+\mu\eta)^{m(s+1)}}{2} \mathbb{E}[\|x_{m(s+1)} - x^*\|^2] \\ & \leq \eta\delta \sum_{k=ms}^{m(s+1)-1} (1+\mu\eta)^k \mathbb{E}[F(x_k) - F(x^*)] + \frac{(1+\mu\eta)^{ms}}{2} \mathbb{E}[\|x_{ms} - x^*\|^2] \\ & \quad + \sum_{k=ms}^{m(s+1)-1} (1+\mu\eta)^k \mathbb{E}\left[\frac{\eta(1+\delta)}{2L\lambda}\|\tilde{\nabla}_k - \nabla f(x_k)\|^2\right] + (1+\delta)\left(\frac{\eta L(\lambda+1)}{2} - \frac{1}{2}\right)\|x_{k+1} - x_k\|^2 \\ & \quad + \eta(1-\rho_B)\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle. \end{aligned}$$

Choosing  $\delta$  so that  $1+\delta \geq (1+\mu\eta)\delta$ , the terms involving  $F(x_{k+1}) - F(x^*)$  telescope, giving the inequality

$$\begin{aligned} & (1+\mu\eta)^{m(s+1)} \mathbb{E}[\delta\eta(F(x_{m(s+1)}) - F(x^*)) + \frac{1}{2}\|x_{m(s+1)} - x^*\|^2] \\ & \leq \delta\eta(1+\mu\eta)^{ms} \mathbb{E}[F(x_{ms}) - F(x^*)] + \frac{(1+\mu\eta)^{ms}}{2} \mathbb{E}[\|x_{ms} - x^*\|^2] \\ & \quad + \sum_{k=ms}^{m(s+1)-1} (1+\mu\eta)^k \mathbb{E}\left[\frac{\eta(1+\delta)}{2L\lambda}\|\tilde{\nabla}_k - \nabla f(x_k)\|^2\right] + (1+\delta)\left(\frac{\eta L(\lambda+1)}{2} - \frac{1}{2}\right)\|x_{k+1} - x_k\|^2 \\ & \quad + \eta(1-\rho_B)\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle. \end{aligned}$$

We would like to bound the inner-product bias term using Lemma 15, and we can do this after some manipulation. Because  $\eta \leq \frac{1}{\mu m}$ , we have  $(1 + \mu\eta)^k \leq (1 + \frac{1}{m})^{k-ms} (1 + \mu\eta)^{ms} \leq 3(1 + \mu\eta)^{ms}$ . Using the same estimate as in equations (13) and (14), we can say

$$\begin{aligned} & \sum_{k=ms}^{m(s+1)-1} (1 + \mu\eta)^k \mathbb{E}[\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle] \\ & \leq 3(1 + \mu\eta)^{ms} \sum_{k=ms}^{m(s+1)-1} |\mathbb{E}[\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle]|, \end{aligned}$$

which produces the inequality

$$\begin{aligned} & (1 + \mu\eta)^{m(s+1)} \mathbb{E}[\delta\eta(F(x_{m(s+1)}) - F(x^*)) + \frac{1}{2}\|x_{m(s+1)} - x^*\|^2] \\ & \leq \delta\eta(1 + \mu\eta)^{ms} \mathbb{E}[F(x_{ms}) - F(x^*)] + \frac{(1 + \mu\eta)^{ms}}{2} \mathbb{E}\|x_{ms} - x^*\|^2 \\ & \quad + (1 + \mu\eta)^{ms} \left( \sum_{k=ms}^{m(s+1)-1} \mathbb{E} \left[ \frac{3\eta(1+\delta)}{2L\lambda} \|\tilde{\nabla}_k - \nabla f(x_k)\|^2 + (1 + \delta) \left( \frac{3\eta L(\lambda+1)}{2} - \frac{1}{2} \right) \|x_{k+1} - x_k\|^2 \right] \right. \\ & \quad \left. + 3\eta(1 - \rho_B) |\mathbb{E}\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle| \right). \end{aligned}$$

Summing this inequality from  $s = 0$  to  $s = S - 1$ ,

$$\begin{aligned} & (1 + \mu\eta)^{mS} \mathbb{E}[\delta\eta(F(x_{mS}) - F(x^*)) + \frac{1}{2}\|x_{mS} - x^*\|^2] \\ & \leq \delta\eta(F(x_0) - F(x^*)) + \frac{1}{2}\|x_0 - x^*\|^2 \\ & \quad + \sum_{s=0}^{S-1} (1 + \mu\eta)^{ms} \left( \sum_{k=ms}^{m(s+1)-1} \mathbb{E} \left[ (1 + \delta) \left( \frac{3\eta L(\lambda+1)}{2} - \frac{1}{2} \right) \|x_{k+1} - x_k\|^2 \right] \right. \\ & \quad \left. + \frac{3\eta(1+\delta)}{2L\lambda} \|\tilde{\nabla}_k - \nabla f(x_k)\|^2 + 3\eta(1 - \rho_B) |\mathbb{E}\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle| \right). \end{aligned}$$

We use Lemma 32 with  $\sigma_s = (1 + \mu\eta)^{ms}$  to bound the MSE and Lemma 15 to bound the inner-product bias term.

$$\begin{aligned} & (1 + \mu\eta)^{mS} \mathbb{E}[\delta\eta(F(x_{mS}) - F(x^*)) + \frac{1}{2}\|x_{mS} - x^*\|^2] \\ & \leq \delta\eta(F(x_0) - F(x^*)) + \frac{1}{2}\|x_0 - x^*\|^2 + w \sum_{s=0}^{S-1} \sum_{k=ms}^{m(s+1)-1} (1 + \mu\eta)^{ms} \mathbb{E}\|x_{k+1} - x_k\|^2, \end{aligned} \tag{19}$$

where  $w = \frac{3\eta L(\lambda+1)(1+\delta)}{2} + \frac{3B_2\eta(1-\rho_B)}{2\epsilon} + \frac{3\Theta\eta L(1+\delta)}{\lambda} + 3B_2\eta L^2(1 - \rho_B)\epsilon\Theta - \frac{1+\delta}{2}$ . To minimize the coefficient of the final term, we set  $\lambda = \sqrt{2\Theta}$  and  $\epsilon = (2L^2\Theta)^{-1/2}$ . This coefficient is then equal to

$$3\sqrt{2\Theta}\eta L(1 + \delta) + \frac{3\eta L(1+\delta)}{2} + 3\sqrt{2}(1 - \rho_B)\eta LB_2\sqrt{\Theta} - \frac{1+\delta}{2}.$$

With

$$\eta \leq \frac{1}{6\sqrt{2\Theta}L(1 + \frac{(1-\rho_B)B_2}{1+\delta}) + L}$$

this term is non-positive. Setting  $\delta = \max\{(1 - \rho_B)B_2 - 1, 0\}$ , we are assured that

$$\eta \leq \frac{1}{3L(1+4\sqrt{2\Theta})} \leq \frac{1}{6\sqrt{2\Theta}L(1 + \frac{(1-\rho_B)B_2}{1+\delta}) + L},$$

so the final term in (19) is non-positive, and we can drop it from the inequality. The resulting inequality is

$$(1 + \mu\eta)^T \mathbb{E}[\delta\eta(F(x_T) - F(x^*)) + \frac{1}{2}\|x_T - x^*\|^2] \leq \delta\eta(F(x_0) - F(x^*)) + \frac{1}{2}\|x_0 - x^*\|^2.$$

All that remains is to show that our choice for  $\delta$  satisfies  $(1 + \delta) \geq (1 + \mu\eta)\delta$ . Using the fact that  $\eta \leq \frac{1}{(1-\rho_B)B_2\mu}$  we can say

$$\frac{1}{\mu\eta} \geq (1 - \rho_B)B_2 \geq \delta$$

which ensures that  $(1 + \delta) \geq (1 + \mu\eta)\delta$  and concludes the proof.  $\blacksquare$

## Appendix C. Proof of Theorem 20

Theorem 20 follows immediately from inequality (7) and the MSE bound of Lemma 32.

### Proof of Theorem 20

Summing inequality (7) from  $k = 0$  to  $k = T - 1$  and applying the full expectation operator, we obtain

$$\begin{aligned} 0 \leq & -\mathbb{E}[F(x_T)] + F(x_0) + (L - \frac{1}{4\eta}) \sum_{k=0}^{T-1} \mathbb{E}[\|\hat{x}_{k+1} - x_k\|^2] \\ & + (\frac{L}{2} - \frac{1}{4\eta}) \sum_{k=0}^{T-1} \mathbb{E}[\|x_{k+1} - x_k\|^2] + 2\eta \sum_{k=0}^{T-1} \mathbb{E}[\|\nabla f(x_k) - \tilde{\nabla}_k\|^2]. \end{aligned}$$

We bound the MSE using Lemma 32 with  $\sigma_s = 1$ .

$$\begin{aligned} 0 \leq & -\mathbb{E}[F(x_T)] + F(x_0) + (L - \frac{1}{4\eta}) \sum_{k=0}^{T-1} \mathbb{E}\|\hat{x}_{k+1} - x_k\|^2 \\ & + (\frac{L}{2} + 4\Theta\eta L^2 - \frac{1}{4\eta}) \sum_{k=0}^{T-1} \mathbb{E}\|x_{k+1} - x_k\|^2. \end{aligned}$$

With  $\eta \leq \frac{\sqrt{16\Theta+1}-1}{16L\Theta}$ , the final term is non-positive, so we can drop it from the inequality. Using the fact that  $-F(x_T) \leq -F(x^*)$ , our inequality simplifies to

$$-(L - \frac{1}{4\eta}) \sum_{k=0}^{T-1} \mathbb{E}[\|\hat{x}_{k+1} - x_k\|^2] \leq F(x_0) - F(x^*).$$

Writing the left side in terms of the generalized gradient, we have the bound

$$\sum_{k=0}^{T-1} \mathbb{E}[\|\mathcal{G}_{\eta/2}(x_k)\|^2] \leq \frac{16(F(x_0) - F(x^*))}{\eta(1 - 4\eta L)}.$$

With  $x_\alpha$  chosen uniformly at random from the set  $\{x_k\}_{k=0}^{T-1}$ , this is equivalent to

$$\mathbb{E}[\|\mathcal{G}_{\eta/2}(x_\alpha)\|^2] \leq \frac{16(F(x_0) - F(x^*))}{\eta(1 - 4\eta L)T}.$$

This completes the proof.  $\blacksquare$

## Appendix D. Proofs of convergence rates for B-SAGA and B-SVRG

The following lemma establishes an MSE bound on the B-SAGA and B-SVRG gradient estimators. For the unbiased case  $\theta = 1$ , this result was essentially first proved in (Defazio et al., 2014a), but the authors ultimately use a looser variance bound.

**Lemma 35** *The MSEs of the B-SAGA and B-SVRG gradient estimators satisfy*

$$\mathbb{E}_k[\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] \leq \frac{1}{n\theta^2} \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f_i(\varphi_k^i)\|^2 + (1 - \frac{2}{\theta}) \|\nabla f(x_k) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i)\|^2. \quad (20)$$

**Proof** Let  $\tilde{\nabla}_k \equiv \tilde{\nabla}_k^{\text{B-SAGA}}$  or  $\tilde{\nabla}_k^{\text{B-SVRG}}$ . The proof amounts to computing the expectation of the estimator and applying the Lipschitz continuity of  $\nabla f_i$ .

$$\begin{aligned} \mathbb{E}_k[\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] &= \mathbb{E}_k[\|\frac{1}{\theta}(\nabla f_{j_k}(x_k) - \nabla f_{j_k}(\varphi_k^{j_k})) - \nabla f(x_k) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i)\|^2] \\ &= \frac{1}{\theta^2} \mathbb{E}_k[\|\nabla f_{j_k}(x_k) - \nabla f_{j_k}(\varphi_k^{j_k})\|^2] + \|\nabla f(x_k) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i)\|^2 \\ &\quad - \frac{2}{\theta} \mathbb{E}_k[\langle \nabla f_{j_k}(x_k) - \nabla f_{j_k}(\varphi_k^{j_k}), \nabla f(x_k) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i) \rangle] \\ &= \frac{1}{n\theta^2} \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f_i(\varphi_k^i)\|^2 + (1 - \frac{2}{\theta}) \|\nabla f(x_k) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i)\|^2, \end{aligned}$$

which is the desired result.  $\blacksquare$

The following two lemmas establish the constants in the BMSE property for the B-SAGA and B-SVRG estimators.

**Proof of Lemma 22** We begin with the inequality of Lemma 35 and consider two cases.

**Case 1.** Suppose  $\theta \in [1, 2]$ . In this case the second term in (20) is non-positive, so we drop it from the inequality. For the remaining term, we use the following bound.

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(\varphi_k^i)\|^2] \\ &\stackrel{\textcircled{1}}{\leq} \frac{1+2n}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] + \frac{1}{n} (1 + \frac{1}{2n}) \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-1}) - \nabla f_i(\varphi_k^i)\|^2] \\ &\stackrel{\textcircled{2}}{\leq} \frac{1+2n}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] + \frac{1}{n} (1 + \frac{1}{2n}) (1 - \frac{1}{n}) \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-1}) - \nabla f_i(\varphi_{k-1}^i)\|^2] \\ &\stackrel{\textcircled{3}}{\leq} \frac{1+2n}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] + \frac{1}{n} (1 - \frac{1}{2n}) \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-1}) - \nabla f_i(\varphi_{k-1}^i)\|^2]. \end{aligned} \quad (21)$$

Inequality  $\textcircled{1}$  is the standard inequality  $\|a - c\|^2 \leq (1 + \delta)\|a - b\|^2 + (1 + \delta^{-1})\|b - c\|^2$  (where we let  $\delta = \frac{1}{2n}$ ). Inequality  $\textcircled{2}$  follows from the definition of  $\varphi_k^i$  and computing the expectation over  $j_{k-1}$ , and  $\textcircled{3}$  uses the fact that  $(1 + \frac{1}{2n})(1 - \frac{1}{n}) \leq (1 - \frac{1}{2n})$ . Altogether, this gives

$$\begin{aligned} &\mathbb{E}[\|\tilde{\nabla}_k^{\text{SAGA}} - \nabla f(x_k)\|^2] \\ &\leq \frac{1}{n\theta^2} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(\varphi_k^i)\|^2] \\ &\leq \frac{2n+1}{n\theta^2} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] + \frac{1}{n\theta^2} (1 - \frac{1}{2n}) \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-1}) - \nabla f_i(\varphi_{k-1}^i)\|^2]. \end{aligned}$$

With  $\mathcal{M}_k = \frac{1}{n\theta^2} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(\varphi_k^i)\|^2]$ , it is clear that the SAGA estimator satisfies the BMSE property with  $M_1 = \frac{2n+1}{\theta^2}$ ,  $\rho_M = \frac{1}{2n}$ ,  $M_2 = 0$ ,  $\rho_F = 1$ , and  $m = 1$ .

**Case 2.** Suppose  $\theta > 2$ , so that the second term in (20) is non-negative. Jensen's inequality gives

$$\mathbb{E}_k[\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] \leq \frac{1}{n}(1 - \frac{1}{\theta})^2 \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f_i(\varphi_k^i)\|^2.$$

Following the argument of Case 1, it is easy to see that the B-SAGA gradient estimator satisfies the BMSE property with  $\mathcal{M}_k = \frac{1}{n}(1 - \frac{1}{\theta})^2 \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f_i(\varphi_k^i)\|^2$ ,  $M_1 = (2n + 1)(1 - \frac{1}{\theta})^2$ ,  $\rho_M = \frac{1}{2n}$ ,  $M_2 = 0$ ,  $\rho_F = 1$ , and  $m = 1$ .

To prove that the B-SAGA is memory-biased, we need computing its expectation.

$$\begin{aligned} \nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k^{\text{B-SAGA}}] &= \nabla f(x_k) - \frac{1}{\theta} \mathbb{E}_k[\nabla f_{j_k}(x_k) - f_i(\varphi_k^{j_k})] - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i) \\ &= (1 - \frac{1}{\theta})(\nabla f(x_k) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i)). \end{aligned}$$

To compute a value for  $B_1$ , we follow (21) to obtain

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|x_k - \varphi_k^i\|^2] &\leq (2n + 1)\|x_k - x_{k-1}\|^2 + \frac{1}{n}(1 - \frac{1}{2n}) \sum_{i=1}^n \mathbb{E}[\|x_{k-1} - \varphi_{k-1}^i\|^2] \\ &\leq (2n + 1) \sum_{\ell=1}^k (1 - \frac{1}{2n})^{k-\ell} \|x_\ell - x_{\ell-1}\|^2. \end{aligned}$$

Summing this inequality from  $k = 0$  to  $k = T - 1$ , we obtain

$$\begin{aligned} \frac{1}{n} \sum_{k=0}^{T-1} \sum_{i=1}^n \mathbb{E}[\|x_k - \varphi_k^i\|^2] &\leq (2n + 1) \sum_{k=0}^{T-1} \sum_{\ell=1}^k (1 - \frac{1}{2n})^{k-\ell} \|x_\ell - x_{\ell-1}\|^2 \\ &\leq (2n + 1) (\sum_{\ell=0}^{\infty} (1 - \frac{1}{2n})^\ell) \sum_{k=0}^{T-1} \|x_{k+1} - x_k\|^2 \\ &= 2n(2n + 1) \sum_{k=0}^{T-1} \|x_{k+1} - x_k\|^2, \end{aligned}$$

which completes the proof.  $\blacksquare$

**Proof of Lemma 23** Suppose  $k \in \{ms, ms + 1, \dots, m(s + 1) - 1\}$  for some  $s \in \mathbb{N}_0$ . As in the proof of Lemma 22, we begin with the inequality of Lemma 35 and consider two cases.

**Case 1.** Suppose  $\theta \in [1, 2]$ , so that we may drop the second term in (20). We can bound the remaining term as follows.

$$\begin{aligned} &\frac{1}{n\theta^2} \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f_i(\varphi_s)\|^2 \\ &\stackrel{\textcircled{1}}{\leq} \frac{1+m}{n\theta^2} \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2 + \frac{1+1/m}{n\theta^2} \sum_{i=1}^n \|\nabla f_i(x_{k-1}) - \nabla f_i(\varphi_s)\|^2 \\ &\stackrel{\textcircled{2}}{\leq} \frac{1+m}{n\theta^2} \sum_{\ell=ms}^k (1 + \frac{1}{m})^{k-\ell} \sum_{i=1}^n \|\nabla f_i(x_{\ell+1}) - \nabla f_i(x_\ell)\|^2. \end{aligned}$$

Inequality  $\textcircled{1}$  uses the inequality  $\|u - w\|^2 \leq (1 + 1/m)\|u - v\|^2 + (1 + m)\|v - w\|^2$ , and  $\textcircled{2}$  follows from the fact that  $x_{ms} = \varphi_s$ . Summing this inequality from  $k = ms$  to  $k = m(s + 1) - 1$  gives

$$\begin{aligned} \frac{1}{n\theta^2} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f_i(\varphi_s)\|^2 &\leq \frac{m+1}{n\theta^2} (1 + \frac{1}{m})^m \sum_{k=ms}^{m(s+1)-1} \sum_{\ell=ms}^k \sum_{i=1}^n \|\nabla f_i(x_{\ell+1}) - \nabla f_i(x_\ell)\|^2 \\ &\leq \frac{m(m+1)}{n\theta^2} (1 + \frac{1}{m})^m \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2 \\ &\leq \frac{3m(m+1)}{n\theta^2} \sum_{k=ms}^{m(s+1)-1} \|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2. \end{aligned}$$

The final inequality uses the fact that  $(1 + \frac{1}{m})^m < \lim_{m \rightarrow \infty} (1 + \frac{1}{m})^m = e < 3$ . From this inequality, it is clear that the B-SVRG gradient estimator satisfies the BMSE property with  $M_1 = \frac{3m(m+1)}{\theta^2}$ ,  $\rho_M = 1$ ,  $M_2 = 0$ , and  $\rho_F = 1$ .

**Case 2.** If  $\theta > 2$ , then applying Jensen's inequality to (20) produces

$$\mathbb{E}_k[\|\tilde{\nabla}_k^{\text{B-SVRG}} - \nabla f(x_k)\|^2] \leq \frac{1}{n}(1 - \frac{1}{\theta})^2 \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f_i(\varphi_s)\|^2.$$

A similar argument to the one in Case 1 shows that  $M_1 = 3m(m+1)(1 - \frac{1}{\theta})^2$ ,  $\rho_M = 1$ ,  $M_2 = 0$ , and  $\rho_F = 1$ .

All that is left is to prove the stated value for  $B_1$ . Following the proof in Case 1,

$$\begin{aligned} \sum_{k=ms}^{m(s+1)-1} \|x_k - \varphi_s\|^2 &\leq \sum_{k=ms}^{m(s+1)-1} \sum_{\ell=ms}^k (1+m)(1 + \frac{1}{m})^m \|x_{\ell+1} - x_\ell\|^2 \\ &\leq 3m(m+1) \sum_{k=ms}^{m(s+1)-1} \|x_{k+1} - x_k\|^2. \end{aligned}$$

Summing over the epochs  $s = 0$  to  $s = S$  shows  $B_1 = 3m(m+1)$ . ■

Combining Lemmas 22 and 23 with Theorems 17 and 20 proves convergence rates for B-SAGA and B-SVRG.

## Appendix E. Proof of convergence rates for SARAH

Lemma 27 establishes the BMSE constants for the SARAH estimator. The convergence rates of Corollary 29 then follow immediately from Theorem 18.

**Proof of Lemma 27** Let  $k \in \{ms+1, ms+2, \dots, m(s+1)-1\}$ . The claim follows immediately from the well-known bound on the MSE of the SARAH gradient estimator

$$\|\tilde{\nabla}_k^{\text{SARAH}} - \nabla f(x_k)\|^2 \leq \frac{1}{n} \sum_{\ell=ms}^k \sum_{i=1}^n \|\nabla f_i(x_{\ell+1}) - \nabla f_i(x_\ell)\|^2.$$

We refer to Fang et al. (2018) for a proof of this inequality. Summing over an epoch and applying the estimate

$$\frac{1}{n} \sum_{k=ms}^{m(s+1)-1} \sum_{\ell=ms}^k \sum_{i=1}^n \|\nabla f_i(x_{\ell+1}) - \nabla f_i(x_\ell)\|^2 \leq \frac{m}{n} \sum_{k=ms}^{m(s+1)-1} \sum_{i=1}^n \|\nabla f_i(x_{k+1}) - \nabla f_i(x_k)\|^2$$

completes the proof. ■

## Appendix F. Proof of convergence rates for SARGE

For our analysis, we write the SARGE gradient estimator in terms of the SAGA estimator. Define the estimator

$$\tilde{\nabla}_k^{\xi\text{-SARGE}} \stackrel{\text{def}}{=} \nabla f_{j_k}(x_{k-1}) - \nabla f_{j_k}(\xi_k^{j_k}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\xi_k^i),$$

where the variables  $\{\xi_k^i\}_{i=1}^n$  follow the update rules  $\xi_{k+1}^{j_k} = x_{k-1}$  and  $\xi_{k+1}^i = \xi_k^i$  for all  $i \neq j_k$ . The SARGE estimator is equal to

$$\widetilde{\nabla}_k^{\text{SARGE}} = \widetilde{\nabla}_k^{\text{SAGA}} - (1 - \frac{1}{n})(\widetilde{\nabla}_k^{\xi\text{-SAGA}} - \widetilde{\nabla}_{k-1}^{\text{SARGE}}).$$

Before we prove Lemma 28, we require a bound on the MSE of the  $\xi$ -SAGA gradient estimator that follows immediately from Lemma 35.

**Lemma 36** *The MSE of the  $\xi$ -SAGA gradient estimator satisfies the following bound:*

$$\mathbb{E}[\|\widetilde{\nabla}_k^{\xi\text{-SAGA}} - \nabla f(x_{k-1})\|^2] \leq 3 \sum_{\ell=1}^{k-1} (1 - \frac{1}{2n})^{k-\ell-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2].$$

**Proof** Following the proof of Lemma 35,

$$\begin{aligned} \mathbb{E}_k[\|\widetilde{\nabla}_k^{\xi\text{-SAGA}} - \nabla f(x_{k-1})\|^2] &= \mathbb{E}_k[\|\nabla f_{j_k}(x_{k-1}) - \nabla f_{j_k}(\xi_k^{j_k}) - \nabla f(x_{k-1}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\xi_k^i)\|^2] \\ &\stackrel{\textcircled{1}}{=} \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_{k-1}) - \nabla f_i(\xi_k^i)\|^2 - \|\nabla f(x_{k-1}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\xi_k^i)\|^2 \\ &\leq \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_{k-1}) - \nabla f_i(\xi_k^i)\|^2. \end{aligned}$$

Equality  $\textcircled{1}$  is the standard variance decomposition. To continue, we follow the proof of Lemma 35.

$$\begin{aligned} &\mathbb{E}[\|\widetilde{\nabla}_k^{\xi\text{-SAGA}} - \nabla f(x_{k-1})\|^2] \\ &\leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-1}) - \nabla f_i(\xi_k^i)\|^2] \\ &\leq \frac{1+2n}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-1}) - \nabla f_i(x_{k-2})\|^2] + \frac{1}{n}(1 + \frac{1}{2n}) \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-2}) - \nabla f_i(\xi_k^i)\|^2] \\ &\stackrel{\textcircled{2}}{=} \frac{(1+2n)}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-1}) - \nabla f_i(x_{k-2})\|^2] \\ &\quad + \frac{1}{n}(1 + \frac{1}{2n})(1 - \frac{1}{n}) \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-2}) - \nabla f_i(\xi_{k-1}^i)\|^2] \\ &\stackrel{\textcircled{3}}{\leq} 3 \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-1}) - \nabla f_i(x_{k-2})\|^2] + \frac{1}{n}(1 - \frac{1}{2n}) \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_{k-2}) - \nabla f_i(\xi_{k-1}^i)\|^2] \\ &\leq 3 \sum_{\ell=1}^{k-1} (1 - \frac{1}{2n})^{k-\ell-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2]. \end{aligned}$$

Equality  $\textcircled{2}$  follows from computing expectations, and  $\textcircled{3}$  uses the estimate  $(1 - \frac{1}{n})(1 + \frac{1}{2n}) \leq (1 - \frac{1}{2n})$ .  $\blacksquare$

Due to the recursive nature of the SARGE gradient estimator, its MSE depends on the difference between the current estimate and the estimate from the previous iteration. The next lemma provides a bound on  $\mathbb{E}\|\widetilde{\nabla}_k^{\text{SARGE}} - \widetilde{\nabla}_{k-1}^{\text{SARGE}}\|^2$ .

**Lemma 37** *The SARGE gradient estimator satisfies the following bound:*

$$\begin{aligned} \mathbb{E}[\|\widetilde{\nabla}_k^{\text{SARGE}} - \widetilde{\nabla}_{k-1}^{\text{SARGE}}\|^2] &\leq \frac{12}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] + \frac{3}{2n^2} \mathbb{E}\|\nabla f(x_{k-1}) - \widetilde{\nabla}_{k-1}^{\text{SARGE}}\|^2 \\ &\quad + \frac{39}{n^2} \sum_{\ell=1}^k (1 - \frac{1}{2n})^{k-\ell} \sum_{i=1}^n \mathbb{E}\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2. \end{aligned}$$

**Proof** To begin, we use the standard inequality  $\|a - c\|^2 \leq (1 + \delta)\|a - b\|^2 + (1 + \delta^{-1})\|b - c\|^2$  for any  $\delta > 0$  twice. For simplicity, we set  $\delta = \sqrt{3/2} - 1$  and use the fact that  $1 + \frac{1}{\sqrt{3/2-1}} \leq 6$

for both applications of this inequality.

$$\begin{aligned}
 & \mathbb{E}[\|\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2] \\
 &= \mathbb{E}[\|\tilde{\nabla}_k^{\text{SAGA}} - (1 - \frac{1}{n})(\tilde{\nabla}_k^{\xi\text{-SAGA}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2] \\
 &\leq 6\mathbb{E}[\|\tilde{\nabla}_k^{\text{SAGA}} - \tilde{\nabla}_k^{\xi\text{-SAGA}}\|^2] + \frac{\sqrt{3}}{\sqrt{2n^2}}\mathbb{E}[\|\tilde{\nabla}_k^{\xi\text{-SAGA}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2] \\
 &\leq 6\mathbb{E}[\|\tilde{\nabla}_k^{\text{SAGA}} - \tilde{\nabla}_k^{\xi\text{-SAGA}}\|^2] + \frac{6\sqrt{3}}{\sqrt{2n^2}}\mathbb{E}[\|\tilde{\nabla}_k^{\xi\text{-SAGA}} - \nabla f(x_{k-1})\|^2] + \frac{3}{2n^2}\mathbb{E}[\|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2].
 \end{aligned} \tag{22}$$

We use  $\frac{6\sqrt{3}}{\sqrt{2n^2}} \leq \frac{9}{n^2}$  to simplify the coefficient of the second term. We now bound the first two of these three terms separately. Consider the first term.

$$\begin{aligned}
 & 6\mathbb{E}[\|\tilde{\nabla}_k^{\text{SAGA}} - \tilde{\nabla}_k^{\xi\text{-SAGA}}\|^2] \\
 &= 6\mathbb{E}[\|\nabla f_{j_k}(x_k) - \nabla f_{j_k}(\varphi_k^{j_k}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i) - \nabla f_{j_k}(x_{k-1}) - \nabla f_{j_k}(\xi_k^{j_k}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\xi_k^i)\|^2] \\
 &\leq 12\mathbb{E}[\|\nabla f_{j_k}(x_k) - \nabla f_{j_k}(x_{k-1})\|^2] \\
 &\quad + 12\mathbb{E}[\|\nabla f_{j_k}(\varphi_k^{j_k}) - \nabla f_{j_k}(\xi_k^{j_k}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\xi_k^i)\|^2] \\
 &\stackrel{\textcircled{1}}{=} 12\mathbb{E}[\|\nabla f_{j_k}(x_k) - \nabla f_{j_k}(x_{k-1})\|^2] \\
 &\quad + 12\mathbb{E}[\|\nabla f_{j_k}(\varphi_k^{j_k}) - \nabla f_{j_k}(\xi_k^{j_k})\|^2] - 12\mathbb{E}[\|\frac{1}{n} \sum_{i=1}^n \nabla f_i(\varphi_k^i) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\xi_k^i)\|^2] \\
 &\leq \frac{12}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] + 12\mathbb{E}[\|\nabla f_{j_k}(\varphi_k^{j_k}) - \nabla f_{j_k}(\xi_k^{j_k})\|^2] \\
 &\leq \frac{12}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] + 12\mathbb{E}[\|\nabla f_{j_k}(\varphi_k^{j_k}) - \nabla f_{j_k}(\xi_k^{j_k})\|^2].
 \end{aligned}$$

Equality  $\textcircled{1}$  is the standard variance decomposition, which states that for any random variable  $X$ ,  $\mathbb{E}[\|X - \mathbb{E}X\|^2] = \mathbb{E}[\|X\|^2] - \|\mathbb{E}X\|^2$ . The second term can be reduced further by computing the expectation. The probability that  $\nabla f_{j_k}(\varphi_k^{j_k}) = \nabla f_{j_{k-1}}(x_{k-1})$  is equal to the probability that  $j_k = j_{k-1}$ , which is  $1/n$ . The probability that  $\nabla f_{j_k}(\varphi_k^{j_k}) = \nabla f_{j_{k-2}}(x_{k-2})$  is equal to the probability that  $j_k \neq j_{k-1}$  and  $j_k = j_{k-2}$ , which is  $\frac{1}{n}(1 - \frac{1}{n})$ . Continuing in this way,

$$\begin{aligned}
 & \mathbb{E}[\|\nabla f_{j_k}(\varphi_k^{j_k}) - \nabla f_{j_k}(\xi_k^{j_k})\|^2] \\
 &= \frac{1}{n}\mathbb{E}[\|\nabla f_{j_{k-1}}(x_{k-1}) - \nabla f_{j_{k-1}}(x_{k-2})\|^2] + \frac{1}{n}(1 - \frac{1}{n})\mathbb{E}[\|\nabla f_{j_{k-2}}(x_{k-2}) - \nabla f_{j_{k-3}}(x_{k-2})\|^2] + \dots \\
 &= \frac{1}{n} \sum_{\ell=1}^{k-1} (1 - \frac{1}{n})^{k-\ell-1} \mathbb{E}[\|\nabla f_{j_\ell}(x_\ell) - \nabla f_{j_\ell}(x_{\ell-1})\|^2].
 \end{aligned}$$

This implies that

$$\begin{aligned}
 12\mathbb{E}[\|\nabla f_{j_k}(\varphi_k^{j_k}) - \nabla f_{j_k}(\xi_k^{j_k})\|^2] &\leq \frac{12}{n^2} \sum_{\ell=1}^{k-1} (1 - \frac{1}{n})^{k-\ell-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2] \\
 &\leq \frac{12}{n^2} \sum_{\ell=1}^{k-1} (1 - \frac{1}{2n})^{k-\ell-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2].
 \end{aligned}$$

We include the second inequality to simplify later arguments. This completes our bound for the first term of (22).

For the second term of (22), we recall Lemma 36.

$$\mathbb{E}[\|\tilde{\nabla}_k^{\xi\text{-SAGA}} - \nabla f(x_{k-1})\|^2] \leq 3 \sum_{\ell=1}^{k-1} (1 - \frac{1}{2n})^{k-\ell-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2].$$

Combining all of these bounds, we obtain

$$\begin{aligned}
 \mathbb{E}[\|\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2] &\leq \frac{12}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] + \frac{3}{2n^2} \mathbb{E}[\|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2] \\
 &\quad + \frac{39}{n^2} \sum_{\ell=1}^{k-1} (1 - \frac{1}{2n})^{k-\ell-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2],
 \end{aligned}$$

which completes the proof. ■

Lemma 37 allows us to take advantage of the recursive structure of our gradient estimate. With this lemma established, we can prove a bound on the MSE.

**Lemma 38** *The SARGE gradient estimator satisfies the following recursive bound:*

$$\begin{aligned} & \mathbb{E}[\|\tilde{\nabla}_k^{\text{SARGE}} - \nabla f(x_k)\|^2] \\ & \leq (1 - \frac{1}{n} + \frac{3}{2n^2})\mathbb{E}[\|\tilde{\nabla}_{k-1}^{\text{SARGE}} - \nabla f(x_{k-1})\|^2] + \frac{12}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] \\ & \quad + \frac{39}{n^2} \sum_{\ell=1}^{k-1} (1 - \frac{1}{2n})^{k-\ell-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2]. \end{aligned}$$

**Proof** The beginning of our proof is similar to the proof of the variance bound for the SARAH gradient estimator in (Nguyen et al., 2017, Lem. 2).

$$\begin{aligned} & \mathbb{E}_k[\|\tilde{\nabla}_k^{\text{SARGE}} - \nabla f(x_k)\|^2] \\ & = \mathbb{E}_k[\|\tilde{\nabla}_{k-1}^{\text{SARGE}} - \nabla f(x_{k-1}) + \nabla f(x_{k-1}) - \nabla f(x_k) + \tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2] \\ & = \|\tilde{\nabla}_{k-1}^{\text{SARGE}} - \nabla f(x_{k-1})\|^2 + \|\nabla f(x_{k-1}) - \nabla f(x_k)\|^2 + \mathbb{E}_k[\|\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2] \\ & \quad + 2\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}, \nabla f(x_k) - \nabla f(x_{k-1}) \rangle \\ & \quad - 2\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}, \mathbb{E}_k[\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}] \rangle \\ & \quad - 2\langle \nabla f(x_k) - \nabla f(x_{k-1}), \mathbb{E}_k[\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}] \rangle. \end{aligned}$$

We consider each inner product separately. The first inner product is equal to

$$\begin{aligned} & 2\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}, \nabla f(x_k) - \nabla f(x_{k-1}) \rangle \\ & = -\|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2 - \|\nabla f(x_k) - \nabla f(x_{k-1})\|^2 + \|\nabla f(x_k) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2. \end{aligned}$$

For the next two inner products, we use the fact that

$$\begin{aligned} \mathbb{E}_k[\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}] & = \mathbb{E}_k[\tilde{\nabla}_k^{\text{SAGA}} - (1 - \frac{1}{n})\tilde{\nabla}_k^{\xi\text{-SAGA}} + (1 - \frac{1}{n})\tilde{\nabla}_{k-1}^{\text{SARGE}}] - \tilde{\nabla}_{k-1}^{\text{SARGE}} \\ & = \nabla f(x_k) - (1 - \frac{1}{n})\nabla f(x_{k-1}) - \frac{1}{n}\tilde{\nabla}_{k-1}^{\text{SARGE}} \\ & = \nabla f(x_k) - \nabla f(x_{k-1}) + \frac{1}{n}(\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}). \end{aligned}$$

With this equality established, we see that the second inner product is equal to

$$\begin{aligned} & -2\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}, \mathbb{E}_k[\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}] \rangle \\ & = -2\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}, \nabla f(x_k) - \nabla f(x_{k-1}) \rangle - \frac{2}{n}\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}, \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}} \rangle \\ & = \|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2 + \|\nabla f(x_k) - \nabla f(x_{k-1})\|^2 \\ & \quad - \|\nabla f(x_k) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2 - \frac{2}{n}\|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2 \\ & = (1 - \frac{2}{n})\|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2 + \|\nabla f(x_k) - \nabla f(x_{k-1})\|^2 - \|\nabla f(x_k) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2. \end{aligned}$$

The third inner product can be bounded using a similar procedure.

$$\begin{aligned} & -2\langle \nabla f(x_k) - \nabla f(x_{k-1}), \mathbb{E}_k[\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}] \rangle \\ & = -2\langle \nabla f(x_k) - \nabla f(x_{k-1}), \nabla f(x_k) - \nabla f(x_{k-1}) \rangle - \frac{2}{n}\langle \nabla f(x_k) - \nabla f(x_{k-1}), \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}} \rangle \\ & \leq -2\|\nabla f(x_k) - \nabla f(x_{k-1})\|^2 + \frac{1}{n}\|\nabla f(x_k) - \nabla f(x_{k-1})\|^2 + \frac{1}{n}\|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2 \\ & = -(2 - \frac{1}{n})\|\nabla f(x_k) - \nabla f(x_{k-1})\|^2 + \frac{1}{n}\|\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2, \end{aligned}$$

where the inequality is Young's. Altogether and after applying the full expectation operator, we have

$$\begin{aligned} \mathbb{E}[\|\tilde{\nabla}_k^{\text{SARGE}} - \nabla f(x_k)\|^2] &\leq (1 - \frac{1}{n})\mathbb{E}[\|\tilde{\nabla}_{k-1}^{\text{SARGE}} - \nabla f(x_{k-1})\|^2] \\ &\quad - (1 - \frac{1}{n})\mathbb{E}[\|\nabla f(x_k) - \nabla f(x_{k-1})\|^2] + \mathbb{E}[\|\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2] \\ &\leq (1 - \frac{1}{n})\mathbb{E}[\|\tilde{\nabla}_{k-1}^{\text{SARGE}} - \nabla f(x_{k-1})\|^2] + \mathbb{E}[\|\tilde{\nabla}_k^{\text{SARGE}} - \tilde{\nabla}_{k-1}^{\text{SARGE}}\|^2]. \end{aligned}$$

Finally, we bound the last term on the right using Lemma 37.

$$\begin{aligned} &\mathbb{E}\|\tilde{\nabla}_k^{\text{SARGE}} - \nabla f(x_k)\|^2 \\ &\leq (1 - \frac{1}{n} + \frac{3}{2n^2})\mathbb{E}[\|\tilde{\nabla}_{k-1}^{\text{SARGE}} - \nabla f(x_{k-1})\|^2] + \frac{12}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_k) - \nabla f_i(x_{k-1})\|^2] \\ &\quad + \frac{39}{n^2} \sum_{\ell=1}^{k-1} (1 - \frac{1}{2n})^{k-\ell-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2] \end{aligned}$$

This completes the proof.  $\blacksquare$

**Proof of Lemma 28** It is easy to see that  $\rho_B = 1/n$  by computing the expectation of the SARGE gradient estimator.

$$\begin{aligned} \nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k^{\text{SARGE}}] &= \nabla f(x_k) - \mathbb{E}_k[\tilde{\nabla}_k^{\text{SAGA}} - (1 - \frac{1}{n})(\tilde{\nabla}_k^{\xi\text{-SAGA}} - \tilde{\nabla}_{k-1}^{\text{SARGE}})] \\ &= (1 - \frac{1}{n})(\nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}^{\text{SARGE}}). \end{aligned}$$

The result of Lemma 38 makes it clear that  $M_1 = 12$ . To determine  $\rho_M$ , we must first choose a suitable sequence  $\mathcal{M}_k$ . Let  $\mathcal{M}_k = \mathbb{E}[\|\tilde{\nabla}_k^{\text{SARGE}} - \nabla f(x_k)\|^2]$ . If  $n = 1$ , then  $\mathcal{M}_k = 0$  for all  $k$ , so it holds trivially that  $\mathcal{M}_k \leq (1 - \rho_M)\mathcal{M}_{k-1}$ . If  $n \geq 2$ , then  $1 - \frac{1}{n} + \frac{3}{2n^2} \leq 1 - \frac{1}{4n}$ , so Lemma 38 ensures that with  $\rho_M = \frac{1}{4n}$ ,  $\mathcal{M}_k \leq (1 - \rho_M)\mathcal{M}_{k-1}$ .

Finally, we must compute  $M_2$  and  $\rho_F$  with respect to some sequence  $\mathcal{F}_k$ . Lemma 38 motivates the choice

$$\mathcal{F}_k = \sum_{\ell=1}^{k-1} (1 - \frac{1}{2n})^{k-\ell-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x_\ell) - \nabla f_i(x_{\ell-1})\|^2],$$

and the choices  $M_2 = \frac{39}{n}$  and  $\rho_F = \frac{1}{2n}$  are clear.  $\blacksquare$

## Appendix G. Incorporating Bias to Lower the MSE: An Example

From (5), we argue that biased stochastic gradient estimators can admit better convergence guarantees than unbiased estimators if the bias reduces the total effect of the estimator's MSE and inner-product bias term. In this section, we compare the effects of these terms for the SARAH and SVRG gradient estimators, revealing why SARAH admits better convergence rates.

Beginning with the SARAH estimator, equation (5) reduces to (17), which we reproduce:

$$\begin{aligned} &\eta \mathbb{E}_k[F(x_{k+1}) - F(x^*) + \delta(F(x_{k+1}) - F(x_k))] \\ &\leq -\frac{1}{2}\mathbb{E}_k[\|x_{k+1} - x^*\|^2] + \frac{1}{2}\|x_k - x^*\|^2 + \mathbb{E}_k\left[\frac{\eta(1+\delta)}{2L\lambda}\|\tilde{\nabla}_k - \nabla f(x_k)\|^2\right] \\ &\quad + (1 + \delta)\left(\frac{\eta L(\lambda+1)}{2} - \frac{1}{2}\right)\|x_{k+1} - x_k\|^2 + \eta(1 - \rho_B)\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle. \end{aligned}$$

Following the proof of Theorem 18, we set  $\lambda = \sqrt{2\Theta}$ , which optimizes the step size. The effect of the MSE and bias term over an epoch is then

$$\begin{aligned}
 & \sum_{k=ms+1}^{m(s+1)-1} \frac{\eta(1+\delta)}{2L\sqrt{2\Theta}} \mathbb{E}[\|\nabla f(x_k) - \tilde{\nabla}_k\|^2 + \eta(1-\rho_B)\langle \nabla f(x_{k-1}) - \tilde{\nabla}_{k-1}, x_k - x^* \rangle] \\
 & \stackrel{\textcircled{1}}{\leq} \sum_{k=ms+1}^{m(s+1)-1} \mathbb{E}[(\frac{\eta(1+\delta)}{L\sqrt{2\Theta}} + \frac{m\epsilon\eta(1-\rho_B)}{2})\|\nabla f(x_k) - \tilde{\nabla}_k\|^2 + \frac{m\eta(1-\rho_B)}{2\epsilon}\|x_{k+1} - x_k\|^2] \\
 & \stackrel{\textcircled{2}}{\leq} \sum_{k=ms+1}^{m(s+1)-1} \mathbb{E}[(\frac{\eta(1+\delta)}{L\sqrt{2\Theta}} + \frac{m\epsilon\eta(1-\rho_B)}{2})\frac{m}{n} \sum_{i=1}^n \|\nabla f_i(x_{k+1}) - f_i(x_k)\|^2 + \frac{m\eta(1-\rho_B)}{2\epsilon}\|x_{k+1} - x_k\|^2] \\
 & \stackrel{\textcircled{3}}{\leq} \sum_{k=ms+1}^{m(s+1)-1} \mathbb{E}[L^2 m(\frac{\eta(1+\delta)}{L\sqrt{2\Theta}} + \frac{m\epsilon\eta(1-\rho_B)}{2})\|x_{k+1} - x_k\|^2 + \frac{m\eta(1-\rho_B)}{2\epsilon}\|x_{k+1} - x_k\|^2]
 \end{aligned}$$

Inequality ① is an application of Lemma 15, ② is Lemma 27, and ③ uses the Lipschitz continuity of  $\nabla f_i$ . Setting  $\epsilon = (L^2 n)^{-1/2}$  to minimize this bound, setting  $\delta = m - 1$ , and choosing  $m = \mathcal{O}(n)$  gives a coefficient of  $\mathcal{O}(n^{3/2}L\eta)$ .

In contrast, B-SVRG admits a larger bound on these terms. For memory-biased estimators, equation (5) reduces to (c.f. (11))

$$\begin{aligned}
 & \eta \mathbb{E}_k[F(x_{k+1}) - F(x^*) + \delta(F(x_{k+1}) - F(x_k))] \\
 & \leq \frac{\eta(1+\delta)}{2L\lambda} \mathbb{E}_k[\|\tilde{\nabla}_k - \nabla f(x_k)\|^2] - \frac{1}{2} \mathbb{E}_k[\|x_{k+1} - x^*\|^2] + \frac{1}{2} \|x_k - x^*\|^2 \\
 & \quad + (\frac{\eta L(1+\delta)(\lambda+1)}{2} - \frac{1+2\delta}{2}) \mathbb{E}_k[\|x_{k+1} - x_k\|^2] + \frac{\eta L}{2n} (1 - \frac{1}{\theta}) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2.
 \end{aligned}$$

The bias term leads to terms of the form  $\|x_k - \varphi_k^i\|^2$ . Therefore, using  $\lambda = \sqrt{2\Theta}$ , our goal is to minimize

$$\begin{aligned}
 & \sum_{k=ms}^{m(s+1)-1} \mathbb{E}[\frac{\eta(1+\delta)}{2L\sqrt{2\Theta}} \|\nabla f(x_k) - \tilde{\nabla}_k\|^2 + \frac{\eta L}{2n} (1 - \frac{1}{\theta}) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2] \\
 & \stackrel{\textcircled{1}}{\leq} \frac{\eta(1+\delta)}{2L\sqrt{2\Theta}} \mathcal{M}_{ms} + \sum_{k=ms}^{m(s+1)-1} \mathbb{E}[\frac{\eta L}{2n} (1 - \frac{1}{\theta}) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2] \\
 & \stackrel{\textcircled{2}}{\leq} \sum_{k=ms}^{m(s+1)-1} \mathbb{E}[\frac{M_1 \eta L(1+\delta)}{2\sqrt{2\Theta}} \|x_{k+1} - x_k\|^2 + \frac{\eta L}{2n} (1 - \frac{1}{\theta}) \sum_{i=1}^n \|x_k - \varphi_k^i\|^2] \\
 & \stackrel{\textcircled{3}}{\leq} \sum_{k=ms}^{m(s+1)-1} \mathbb{E}[\frac{M_1 \eta L(1+\delta)}{2\sqrt{2\Theta}} \|x_{k+1} - x_k\|^2 + \frac{\eta L B_1}{2} (1 - \frac{1}{\theta}) \|x_{k+1} - x_k\|^2].
 \end{aligned}$$

Inequalities ① and ② use the fact that the SVRG estimator satisfies the BMSE property with  $\rho_M = 1$ . Inequality ③ uses the definition of a memory-biased estimator. With  $\delta = \max\{B_1(1 - 1/\theta)/\sqrt{2\Theta} - 1, 0\}$ ,  $B_1 = 3m(m+1)$ ,  $\Theta = M_1$ , and

$$M_1 = \begin{cases} \frac{3m(m+1)}{\theta^2} & \theta \in (0, 2], \\ 3m(m+1)(1 - \frac{1}{\theta})^2 & \theta > 2, \end{cases}$$

this gives a coefficient of  $\mathcal{O}(Lm^2\eta)$ .

This difference of  $\mathcal{O}(n^{1/2})$  between the two bounds is significant. It manifests as an improvement of  $\sqrt{n}$  in the convergence rate of SARAH over the rate of B-SVRG.

## References

Zeyuan Allen-Zhu. Natasha: Faster non-convex stochastic optimization via strongly non-convex parameter. In *Proceedings of the 34<sup>th</sup> International Conference on Machine Learning*, 2017.

- Zeyuan Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. *Journal of Machine Learning Research*, pages 1–51, 2018a.
- Zeyuan Allen-Zhu. Katyusha X: Practical momentum method for stochastic sum-of-nonconvex optimization. In *Proceedings of the 35<sup>th</sup> International Conference on Machine Learning*, 2018b.
- Zeyuan Allen-Zhu. Natasha 2: Faster non-convex optimization than SGD. In *Proceedings of the 32<sup>nd</sup> Annual Conference on Neural Information Processing Systems*, 2018c.
- Zeyuan Allen-Zhu and Elad Hazan. Variance reduction for faster non-convex optimization. In *Proceedings of the 33<sup>rd</sup> International Conference on Machine Learning*, 2016.
- Zeyuan Allen-Zhu and Yang Yuan. Improved SVRG for non-strongly-convex or sum-of-non-convex objectives. In *Proceedings of the 35<sup>th</sup> International Conference on Machine Learning*, 2018.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- Christopher M. Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- Léon Bottou, Frank E. Curtis, , and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60:223–311, 2018.
- Kristian Bredies and Dirk Lorenz. *Mathematical Image Processing*. Springer, 2018.
- Antonin Chambolle, Matthias J. Ehrhardt, Peter Richtárik, and Carola-Bibiane Schönlieb. Stochastic primal-dual hybrid gradient algorithm with arbitrary sampling and imaging applications. *SIAM J. Optim.*, 28(4):2783–2808, 2018.
- Aaron Defazio. A simple practical accelerated method for finite sums. In *Proceedings of the 30<sup>th</sup> Annual Conference on Neural Information Processing Systems*, 2016.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2014a.
- Aaron Defazio, Tiberio Caetano, and Justin Domke. Finito: A faster, permutable incremental gradient method for big data problems. In *Proceedings of the 31st International Conference on Machine Learning*, 2014b.
- Derek Driggs, Matthias Ehrhardt, and Carola-Bibiane Schönlieb. Accelerating variance-reduced gradient methods. *Mathematical Programming*, 2020.
- Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex optimization via stochastic path integrated differential estimator. In *Proceedings of the 32<sup>nd</sup> Annual Conference on Neural Information Processing Systems*, 2018.
- Dan Garber and Elad Hazan. Fast and simple PCA via convex optimization. *arXiv:1509.05647v4*, 2015.

- Thomas Hofmann, Aurelien Lucchi, Simon Lacoste-Julien, and Brian McWilliams. Variance reduced stochastic gradient descent with neighbors. In *Advances in Neural Information Processing Systems*, volume 28, pages 2296–2304, 2015.
- Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems*, pages 315–323, 2013.
- Guanghui Lan, Zhize Li, and Yi Zhou. A unified variance-reduced accelerated gradient method for convex optimization. *arXiv:1905.12412*, 2019.
- Jingwei Liang, Jalal Fadili, and Gabriel Peyré. Activity identification and local linear convergence of Forward–Backward-type methods. *SIAM Journal on Optimization*, 27(1): 408–437, 2017.
- Hongzhou Lin, Julien Mairal, and Zaid Harchaoui. A universal catalyst for first-order optimization. In *Advances In Neural Information Processing Systems*, 2015.
- Pierre-Louis Lions and Bertrand Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM Journal on Numerical Analysis*, 16(6):964–979, 1979.
- Julien Mairal. Incremental majorization-minimization optimization with application to large-scale machine learning. *Technical report*, 2014.
- Yurii Nesterov. *Introductory lectures on convex programming*. Springer, 2004.
- Lam M. Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 2613–2621, 2017.
- Nhan H. Pham, Lam M. Nguyen, Dzung T. Phan, and Quoc Tran-Dinh. ProxSARAH: An efficient algorithmic framework for stochastic composite nonconvex optimization. *arXiv:1902.05679*, 2019.
- Sashank J. Reddi, Ahmed Hefny, Suvrit Sra, Barnabás Póczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. In *Proc. 33rd International Conference on Machine Learning*, 2016a.
- Sashank J. Reddi, Suvrit Sra, Barnabás Póczos, and Alex Smola. Fast stochastic methods for nonsmooth nonconvex optimization. In *Proc. 30th Annual Conference on Neural Information Processing Systems*, 2016b.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22(3):400–407, 1951.
- Nicolas Le Roux, Mark Schmidt, and Francis Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In *Advances in neural information processing systems*, pages 2663–2671, 2012.
- Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162:83–112, 2017.

- Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized loss minimization. *Journal of Machine Learning Research*, 14:567–599, 2013.
- Fanhua Shang, Licheng Jiao, Kaiwen Zhou, James Cheng, Yan Ren, and Yufei Jin. ASVRG: Accelerated proximal SVRG. In *Asian Conference on Machine Learning*, volume 95, pages 1–32, 2018.
- Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh. SpiderBoost: A class of faster variance-reduced algorithms for nonconvex optimization. *arXiv:1810.10690*, 2018.
- Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh. Spiderboost and momentum: Faster variance reduction algorithms. In *Proceedings of the 33<sup>rd</sup> Annual Conference on Neural Information Processing Systems*, 2019.
- Lin Xiao and Tong Zhang. A proximal stochastic gradient method with progressive variance reduction. *Technical report, Microsoft Research*, 2014.
- Kaiwen Zhou, Fanhua Shang, and James Cheng. A simple stochastic variance reduced algorithm with fast convergence rates. In *Proceedings of the 35<sup>th</sup> International Conference on Machine Learning*, pages 5975–5984, 2018.