# **Demystifying SGD with Doubly Stochastic Gradients**

Kyurae Kim<sup>1</sup> Joohwan Ko<sup>2</sup> Yi-An Ma<sup>3</sup> Jacob R. Gardner<sup>1</sup>

## **Abstract**

Optimization objectives in the form of a sum of intractable expectations are rising in importance (e.g., diffusion models, variational autoencoders, and many more), a setting also known as "finite sum with infinite data." For these problems, a popular strategy is to employ SGD with doubly stochastic gradients (doubly SGD): the expectations are estimated using the gradient estimator of each component, while the sum is estimated by subsampling over these estimators. Despite its popularity, little is known about the convergence properties of doubly SGD, except under strong assumptions such as bounded variance. In this work, we establish the convergence of doubly SGD with independent minibatching and random reshuffling under general conditions, which encompasses dependent component gradient estimators. In particular, for dependent estimators, our analysis allows fined-grained analysis of the effect correlations. As a result, under a per-iteration computational budget of  $b \times m$ , where b is the minibatch size and m is the number of Monte Carlo samples, our analysis suggests where one should invest most of the budget in general. Furthermore, we prove that random reshuffling (RR) improves the complexity dependence on the subsampling noise.

## 1. Introduction

Stochastic gradient descent (SGD; Robbins & Monro, 1951; Bottou, 1999; Nemirovski et al., 2009; Shalev-Shwartz et al., 2011) is the *de facto* standard for solving large scale optimization problems of the form of finite sums

Proceedings of the 41<sup>st</sup> International Conference on Machine Learning, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

such as

$$\underset{\boldsymbol{x} \in \mathcal{X} \subseteq \mathbb{R}^d}{\text{minimize}} \left\{ F(\boldsymbol{x}) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\boldsymbol{x}) \right\}. \tag{1}$$

When n is large, SGD quickly converges to low-accuracy solutions by subsampling over components  $f_1, ..., f_n$ . The properties of SGD on the finite sum class have received an immense amount of interest (Bottou et al., 2018) as it includes empirical risk minimization (ERM; Vapnik, 1991).

Unfortunately, for an emerging large set of problems in machine learning, we may not have direct access to the components  $f_1, ..., f_n$ . That is, each  $f_i$  may be defined as an intractable expectation, or an "infinite sum"

$$f_i(\mathbf{x}) = \mathbb{E}_{\mathbf{n} \sim \varphi} f_i(\mathbf{x}; \mathbf{n}), \qquad (2)$$

where we only have access to the noise distribution  $\varphi$ and the integrand  $f_i(x; \eta)$ , and  $\eta$  is a potentially continuous and unbounded source of stochasticity; a setting Zheng & Kwok (2018); Bietti & Mairal (2017) have previously called "finite sum with infinite data." Such problems include the training of diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song & Ermon, 2019), variational autoencoders (Kingma & Welling, 2014; Rezende et al., 2014), solving ERM under differential privacy (Bassily et al., 2014; Song et al., 2013), and also classical problems such as variational inference (Ranganath et al., 2014; Titsias & Lázaro-Gredilla, 2014; Kucukelbir et al., 2017), and variants of empirical risk minimization (Dai et al., 2014; Bietti & Mairal, 2017; Shi et al., 2021; Orvieto et al., 2023; Liu et al., 2021). In contrast to the finite sum setting where SGD has traditionally been applied, our problem takes the form of

$$\underset{\boldsymbol{x} \in \mathcal{X} \subseteq \mathbb{R}^d}{\text{minimize}} \left\{ F\left(\boldsymbol{x}\right) \triangleq \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}_{\boldsymbol{\eta} \sim \varphi} f_i\left(\boldsymbol{x}; \boldsymbol{\eta}\right) \right\}.$$

These optimization problems are typically solved using SGD with *doubly stochastic gradients* (doubly SGD; coined by Dai et al. 2014; Titsias & Lázaro-Gredilla 2014), so-called because, in addition to subsampling over  $f_i$ , stochastic estimates of each component  $f_i$  are used.

Previous studies have relied on strong assumptions to analyze doubly stochastic gradients. For instance, Kulunchakov & Mairal (2020); Bietti & Mairal (2017); Zheng & Kwok (2018) have (i) assumed that the variance of each component estimator is bounded by a constant, which contradicts componentwise strong convexity (Nguyen et al.,

<sup>&</sup>lt;sup>1</sup>Department of Computer and Information Sciences, University of Pennsylvania, Philadelphia, PA, U.S.A. <sup>2</sup>KAIST, Daejeon, South Korea, Republic of <sup>3</sup>Halıcıoğlu Data Science Institute, University of California San Diego, San Diego, CA, U.S.A.. Correspondence to: Kyurae Kim < kyrkim@seas.upenn.edu>.

2018) when  $\mathcal{X} = \mathbb{R}^d$ , (ii) or that the integrand  $\nabla f_i(\mathbf{x}; \boldsymbol{\eta})$ , is L-Lipschitz smooth "uniformly" over  $\eta$ . That is, for any fixed  $\eta$  and i,

$$\|\nabla f_i(\mathbf{x}; \boldsymbol{\eta}) - \nabla f_i(\mathbf{y}; \boldsymbol{\eta})\| \le L\|\mathbf{x} - \mathbf{y}\|_2^2$$

holds for all  $(x, y) \in \mathcal{X}^2$ . Unfortunately, this only holds for additive noise and is otherwise unrealizable when nhas an unbounded support. Therefore, analyses relying on uniform smoothness obscure a lot of interesting behavior. Meanwhile, weaker assumptions such as expected smoothness (ES; Moulines & Bach, 2011; Gower et al., 2021b) have shown to be realizable even for complex gradient estimators (Domke, 2019; Kim et al., 2023). Therefore, a key question is how these ES-type assumptions propagate to doubly stochastic estimators. Among these, we focus on the expected residual (ER; Gower et al., 2019) condition.

Furthermore, in practice, certain applications of doubly SGD share the randomness  $\eta$  across the batch B. (See Section 2.2 for examples.) This introduces dependence between the gradient estimate for each component such that  $\nabla f_i(\mathbf{x}; \boldsymbol{\eta}) \perp \nabla f_i(\mathbf{x}; \boldsymbol{\eta})$  for  $i, j \in B$ . Little is known about the effect of this practice apart from some empirical results (Kingma et al., 2015). For instance, when m Monte Carlo samples of  $\eta$  and a minibatch of size b are used, what is the trade-off between m and b? To answer this question, we provide a theoretical analysis of doubly SGD that encompasses dependent gradient estimators.

## **Technical Contributions**

**Theorem 1**: For doubly stochastic estimators, we establish a general variance bound of the form of

$$\mathcal{O}\left(\frac{\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}^{2}}{mb}+\rho\frac{\left(\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\right)^{2}}{m}+\frac{\tau^{2}}{b}\right),$$

where  $\sigma_i^2$  is the variance of the estimator of  $\nabla f_i$ ,  $\rho \in [0,1]$  is the correlation between the estimators, and  $\tau^2$  is the variance of subsampling.

- Theorems 2 and 3: Using the general variance bound, we show that a doubly stochastic estimator subsampling over correlated estimators satisfying the ER condition and the bounded variance (BV; Definition 2; bounded only on the solution set) condition equally satisfies the ER and BV conditions as well. This is sufficient to guarantee the convergence of doubly SGD on convex, quasar convex, and nonconvex smooth objectives.
- **Theorem 5**: Under similar assumptions, we also prove the convergence of doubly SGD with random reshuffling (doubly SGD-RR), instead of independent subsampling, on a strongly convex objective with strongly convex components.

## **Practical Insights**

- Should I invest in (increase) m or b? When dependent gradient estimators are used, increasing m or b does not have the same impact on the gradient variance as the subsampling strategy also affects the resulting correlation between the estimators. Through Lemma 9, our analysis provides insight into this effect. In particular, we reveal that reducing subsampling variance also reduces Monte Carlo variances. Therefore, for a fixed budget  $m \times b$ , increasing b should always be preferred over increasing m.
- Random Reshuffling Improves Complexity. Our analysis of doubly SGD-RR reveals that, for strongly convex objectives, random reshuffling improves the iteration complexity of doubly SGD from  $\mathcal{O}\left(\frac{1}{\epsilon}\sigma_{\rm mc}^2 + \frac{1}{\epsilon}\sigma_{\rm sub}^2\right)$  to  $\mathcal{O}\left(\frac{1}{\epsilon}\sigma_{\rm mc}^2 + \frac{1}{\sqrt{\epsilon}}\sigma_{\rm sub}\right)$ . Furthermore, for dependent gradient estimators, doubly SGD-RR is "super-efficient": for a batch taking  $\Theta(mb)$  samples to compute, it achieves a n/btighter asymptotic sample complexity compared to full-batch SGD.

## 2. Preliminaries

**Notation** We denote random variables (RVs) in serif (e.g., x, x, X, B), vectors and matrices in bold (e.g., x, X, B) $\boldsymbol{x}$ ,  $\boldsymbol{A}$ ,  $\boldsymbol{A}$ ). For a vector  $\boldsymbol{x}$ , we denote the  $\ell_2$ -norm as  $\|x\|_2 \triangleq \sqrt{\langle x, x \rangle} = \sqrt{x^\top x}$ , where  $\langle x, x \rangle = x^\top x$  is the inner product. Lastly,  $X \perp Y$  denotes independence of X and Y.

Table 1. Nomenclature

Symb.	Description	Ref.
$F(\mathbf{x})$	Objective function	Eq. (1)
$f_i(\mathbf{x})$	<i>i</i> th component of <i>F</i>	Eq. (1)
$\nabla f_B(\mathbf{x})$	Minibatch subsampling estimator of $\nabla F$	Eq. (4)
В	Minibatch of component indices	Eq. (3)
$\pi$	Minibatch subsampling strategy	Eq. (3)
$b_{ m eff}$	Effective sample size of $\pi$	Eq. (5)
$\mathbf{g}_{i}\left(\mathbf{x}\right)$	Unbiased stochastic estimator of $\nabla f_i$	Eq. (7)
$g_i(x; \eta)$	Integrand of estimator $\mathbf{g}_i(\mathbf{x})$	Eq. (7)
$\mathbf{g}_{B}(\mathbf{x})$	Doubly stochastic estimator of $\nabla F$	Eq. (8)
$\mathcal{L}_{ ext{sub}}$	ER constant (Definition 1) of $\pi$	Assu. 5
$\mathcal{L}_i$	ER constant (Definition 1) of $\mathbf{g}_i$	Assu. 6
$ au^2$	BV constant (Definition 2) of $\pi$	Assu. 7
$\sigma_i^2$	BV constant (Definition 2) of $\mathbf{g}_i$	Assu. 7

#### 2.1. Stochastic Gradient Descent on Finite-Sums

Stochastic gradient descent (SGD) is an optimization algorithm that repeats the steps

$$\mathbf{x}_{t+1} = \prod_{\mathcal{X}} (\mathbf{x}_t - \mathbf{y}_t \mathbf{g}(\mathbf{x}_t)).$$

 $\boldsymbol{x}_{t+1} = \Pi_{\mathcal{X}} \left( \boldsymbol{x}_t - \gamma_t \boldsymbol{g} \left( \boldsymbol{x}_t \right) \right),$  where,  $\Pi_{\mathcal{X}}$  is a projection operator onto  $\mathcal{X}, \left( \gamma_t \right)_{i=0}^{T-1}$  is some stepsize schedule,  $\mathbf{g}(\mathbf{x})$  is an unbiased estimate of  $\nabla F(\mathbf{x})$ .

**Finite-Sum Problems.** When the objective can be represented as a "finite sum" it is typical to approximate the gradients of the objective as

$$\nabla F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\mathbf{x}) = \mathbb{E}_{B \sim \pi} \left[ \frac{1}{b} \sum_{i \in B} \nabla f_i(\mathbf{x}) \right], \quad (3)$$

where  $B \sim \pi$  is an index set of cardinality |B| = b, or "minibatch," formed by subsampling over the datapoint indices  $\{1, ..., n\}$ . More formally, we are approximating  $\nabla F$  using the (minibatch) subsampling estimator

$$\nabla f_{B}(\mathbf{x}) \triangleq \frac{1}{b} \sum_{i \in B} \nabla f_{i}(\mathbf{x}), \qquad (4)$$

where the performance of this estimator, or equivalently, of the subsampling strategy  $\pi$ , can be quantified by its variance

$$\operatorname{tr} \mathbb{V}\left[\nabla f_{B}(\mathbf{x})\right] = \frac{1}{b_{\text{eff}}} \underbrace{\frac{1}{n} \sum_{i=1}^{n} \left\|\nabla f_{i}(\mathbf{x}) - \nabla F(\mathbf{x})\right\|_{2}^{2}}_{\text{(unit) subsampling variance}}, (5)$$

where we say  $b_{\rm eff}$  is the "effective sample size" of  $\pi$ . For instance, independent subsampling achieves  $b_{\rm eff} = b$ , and sampling without replacement, also known as "b-nice sampling" (Gower et al., 2019; Richtárik & Takáč, 2016; Csiba & Richtárik, 2018), achieves  $b_{\rm eff} = (n-1)b/n-b$  (Lemma 2).

#### 2.2. Doubly Stochastic Gradients

For problems where the components are defined as intractable expectations as in Eq. (2), we have to rely on an additional Monte Carlo approximation step such as

$$\nabla F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\mathbf{x}) = \mathbb{E}_{B \sim \pi} \left[ \frac{1}{b} \sum_{i \in B} \mathbb{E}_{\boldsymbol{\eta} \sim \varphi} \left[ \nabla f_i(\mathbf{x}; \boldsymbol{\eta}) \right] \right]$$

$$= \mathbb{E}_{\boldsymbol{B} \sim \boldsymbol{\pi}, \ \boldsymbol{\eta}_{j} \sim \varphi} \left[ \frac{1}{mb} \sum_{i \in \boldsymbol{B}} \sum_{j=1}^{m} \nabla f_{i} \left( \boldsymbol{x}; \boldsymbol{\eta}_{j} \right) \right], \tag{6}$$

where  $\eta_j \sim \varphi$  are *m* independently and identically distributed (*i.i.d.*) Monte Carlo samples from  $\varphi$ .

**Doubly Stochastic Gradient** Consider an unbiased estimator of the component gradient  $\nabla f_i$  such that

$$\mathbb{E}\mathbf{g}_{i}(\mathbf{x}) = \mathbb{E}_{\mathbf{\eta} \sim \varphi}\mathbf{g}_{i}(\mathbf{x}; \mathbf{\eta}) = \nabla f_{i}(\mathbf{x}), \tag{7}$$

where  $g_i(x; \eta)$  is the measurable integrand. Using these, we can estimate  $\nabla F$  through the *doubly stochastic* gradient estimator

$$\mathbf{g}_{B}(\mathbf{x}) \triangleq \frac{1}{b} \sum_{i \in B} \mathbf{g}_{i}(\mathbf{x}),$$
 (8)

We separately define the integrand  $g(x; \eta)$  since, in practice, a variety of unbiased estimators of  $\nabla f_i$  can be obtained by appropriately defining the integrand  $g_i$ . For example, one can form the m-sample "naive" Monte Carlo estimator by setting

$$\mathbf{g}_{i}(\mathbf{x}; \boldsymbol{\eta}) = \frac{1}{m} \sum_{j=1}^{m} \nabla f_{i}(\mathbf{x}; \boldsymbol{\eta}_{j}),$$

where  $\boldsymbol{\eta} = [\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_m] \sim \varphi^{\otimes m}$ .

**Dependent Component Gradient Estimators.** Notice that, in Eq. (6), the subcomponents in the batch share the Monte Carlo samples, which may occur in practice. This means  $\mathbf{g}_i$  and  $\mathbf{g}_j$  in the same batch are dependent and, in the worst case, positively correlated, which complicates the analysis. While it is possible to make the estimators independent by sampling m unique Monte Carlo samples for each component (mb Monte Carlo samples in total) as highlighted by Kingma et al. (2015), it is common to use dependent estimators for various practical reasons:

- 1. **ERM with Randomized Smoothing**: In the ERM context, recent works have studied the generalization benefits of randomly perturbing the model weights before computing the gradient (Orvieto et al., 2023; Liu et al., 2021). When subsampling is used, perturbing the weights independently for each datapoint is computationally inefficient. Therefore, the perturbation is shared across the batch, creating dependence.
- Black-Box Variational inference (Titsias & Lázaro-Gredilla, 2014; Kucukelbir et al., 2017): Here, each component can be decomposed as

$$f_i(\mathbf{x}; \boldsymbol{\eta}) = \ell_i(\mathbf{x}; \boldsymbol{\eta}) + r(\mathbf{x}; \boldsymbol{\eta}),$$

where  $\ell_i$  is the log likelihood and r is the log-density of the prior. By sharing  $(\eta_j)_{j=1}^m$ , r only needs to be evaluated m times. To create independent estimators, it needs to be evaluated mb times instead, but r can be expensive to compute.

3. Random feature kernel regression with doubly SGD (Dai et al., 2014): The features are shared across the batch <sup>1</sup>. This reduces the peak memory requirement from  $bmd_{\eta}$ , where  $d_{\eta}$  is the size of the random features, to  $md_{\eta}$ .

One of the analysis goals of this work is to characterize the effect of dependence in the context of SGD.

## 2.3. Technical Assumptions on Gradient Estimators

To establish convergence of SGD, contemporary analyses use the "variance transfer" strategy (Moulines & Bach, 2011; Johnson & Zhang, 2013; Nguyen et al., 2018; Gower et al., 2019; 2021b). That is, by assuming the gradient noise satisfies some condition resembling smoothness, it is possible to bound the gradient noise on some arbitrary point x by the gradient variance on the solution set  $x_* \in \arg\min_{x \in \mathcal{X}} F(x)$ .

**ER Condition.** In this work, we will use the *expected residual* (ER) condition by Gower et al. (2021a):

<sup>&</sup>lt;sup>1</sup>See the implementation at https://github.com/zixu1986/Doubly\_Stochastic\_Gradients

**Definition 1 (Expected Residual; ER).** A gradient estimator  $\mathbf{g}$  of  $F: \mathcal{X} \to \mathbb{R}$  is said to satisfy  $ER(\mathcal{L})$  if

$$\operatorname{tr} \mathbb{V}\left[\boldsymbol{g}\left(\boldsymbol{x}\right) - \boldsymbol{g}\left(\boldsymbol{x}_{*}\right)\right] \leq 2\mathcal{L}\left(F(\boldsymbol{x}) - F(\boldsymbol{x}_{*})\right),$$

for some  $0 < \mathcal{L} < \infty$  and all  $\mathbf{x} \in \mathcal{X}$  and all  $\mathbf{x}_* \in \arg\min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$ .

When f is convex, a weaker form can be used: We will also consider the *convex* variant of the ER condition that uses the Bregman divergence defined as

$$D_{\phi}(\mathbf{y}, \mathbf{x}) \triangleq \phi(\mathbf{y}) - \phi(\mathbf{x}) - \langle \nabla \phi(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle,$$

 $\forall (x, y) \in \mathcal{X}^2$ , where  $\phi : \mathcal{X} \to \mathbb{R}$  is a convex function.

Why the ER condition? A way to think about the ER condition is that it corresponds to the "variance form" equivalent of the expected smoothness (ES) condition by Gower et al. (2021b) defined as

$$\mathbb{E}\|\boldsymbol{g}\left(\boldsymbol{x}\right) - \boldsymbol{g}\left(\boldsymbol{x}_{*}\right)\|_{2}^{2} \leq 2\mathcal{L}\left(F\left(\boldsymbol{x}\right) - F\left(\boldsymbol{x}_{*}\right)\right), \quad (ES)$$

but is slightly weaker, as shown by Gower et al. (2021a). The main advantage of the ER condition is that, due to the properties of the variance, it composes more easily:

**Proposition 1.** Let g satisfy ER  $(\mathcal{L})$ . Then, the m-sample i.i.d. average of g satisfy ER  $(\mathcal{L}/m)$ .

**BV Condition.** From the ER property, the gradient variance on any point  $x \in \mathcal{X}$  can be bounded by the variance on the solution set as long as the following holds:

**Definition 2** (Bounded Gradient Variance). A gradient estimator g of  $F: \mathcal{X} \to \mathbb{R}$  satisfies BV  $(\sigma^2)$  if

$$\mathrm{tr}\mathbb{V}\left[\boldsymbol{g}\left(\boldsymbol{x}_{*}\right)\right]\leq\sigma^{2}$$

for some  $\sigma^2 < \infty$  and all  $\mathbf{x}_* \in \arg\max_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$ .

#### 2.4. Convergence Guarantees for SGD

**Sufficiency of ER and BV.** From the ER and BV conditions, other popular conditions such as ES (Gower et al., 2021b) and ABC (Khaled & Richtárik, 2023) can be established with minimal additional assumptions. As a result, we retrieve the previous convergence results on SGD established for various objective function classes:

- ➤ strongly convex (Gower et al., 2019),
- ➤ quasar convex (+PL) (Gower et al., 2021a),
- ➤ smooth (+PL) (Khaled & Richtárik, 2023).

(Note: quasar convexity is strictly weaker than convexity Guminov et al., 2023; PL: Polyak-Łojasiewicz.) (See also the comprehensive treatment by Garrigos & Gower, 2023.) Therefore, ER and BV are sufficient conditions for SGD to converge on problem classes typically considered in SGD convergence analysis.

In this work, we will specifically focus on smooth and strongly convex objectives:

**Assumption 1.** There exists some  $\mu, L$  satisfying  $0 < \mu \le L < \infty$  suc that the objective function  $F: \mathcal{X} \to \mathbb{R}$  is  $\mu$ -strongly convex and L-smooth as

$$F(\mathbf{y}) - F(\mathbf{x}) \ge \langle \nabla F(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2} ||\mathbf{x} - \mathbf{y}||_2^2$$
$$F(\mathbf{y}) - F(\mathbf{x}) \le \langle \nabla F(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} ||\mathbf{x} - \mathbf{y}||_2^2$$

hold for all  $(x, y) \in \mathcal{X}^2$ .

Also, we will occasionally assume that F is comprised of a finite sum of convex and smooth components:

**Assumption 2.** The objective function  $F: \mathcal{X} \to \mathbb{R}$  is a finite sum as  $F = \frac{1}{n} (f_1 + ... + f_n)$ , where each component is  $L_i$ -smooth and convex such that

$$\left\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\right\|_2^2 \le 2L_i D_{f_i}(\mathbf{x}, \mathbf{y})$$

holds for all  $(x, y) \in \mathcal{X}^2$ .

Note that Assumption 2 alone already implies that F is convex and  $L_{\text{max}}$ -smooth with  $L_{\text{max}} = \max\{L_1, \dots, L_n\}$ .

Why focus on strongly convex functions? We focus on strongly convex objectives as the effect of stochasticity is the most detrimental: in the deterministic setting, one only needs  $\mathcal{O}(\log(1/\epsilon))$  iterations to achieve an  $\epsilon$ -accurate solution. But with SGD, one actually needs  $\mathcal{O}(1/\epsilon)$  iterations due to noise. As such, we can observe a clear contrast between the effect of optimization and noise in this setting.

With that said, for completeness, we provide full proof of convergence on strongly convex-smooth objectives:

**Lemma 1.** Let the objective F satisfy Assumption 1 and the gradient estimator  $\mathbf{g}$  satisfy  $\mathrm{ER}\left(\mathcal{L}\right)$  and  $\mathrm{BV}\left(\sigma^{2}\right)$ . Then, the last iterate of SGD is  $\epsilon$ -close to the global optimum  $\mathbf{x}_{*} = \arg\min_{\mathbf{x} \in \mathcal{X}} F\left(\mathbf{x}\right)$  such that  $\mathbb{E}||\mathbf{x}_{T} - \mathbf{x}_{*}||_{2}^{2} \leq \epsilon$  after a number of iterations at least

$$T \ge 2 \max\left(\frac{\sigma^2}{\mu^2} \frac{1}{\epsilon}, \frac{\mathcal{L} + L}{\mu}\right) \log\left(2||\boldsymbol{x}_0 - \boldsymbol{x}_*||_2^2 \frac{1}{\epsilon}\right)$$

and the fixed stepsize

$$\gamma = \min\left(\frac{\epsilon\mu}{2\sigma^2}, \frac{1}{2(\mathcal{L} + L)}\right).$$

See the *full proof* in page 22.

Note that our complexity guarantee is only  $\mathcal{O}(1/\varepsilon \log (1/\varepsilon))$  due to the use of a fixed stepsize. It is also possible to establish a  $\mathcal{O}(1/\varepsilon)$  guarantee using decreasing stepsize schedules proposed by Gower et al. (2019); Stich (2019). In practice, these schedules are rarely used, and the resulting complexity guarantees are less clear than with fixed stepsizes. Therefore, we will stay on fixed stepsizes.

#### 3. Main Results

## 3.1. Doubly Stochastic Gradients

First, while taming notational complexity, we will prove a general result that holds for combin-

$$\begin{array}{cccc} \textbf{\textit{x}}_i & \leftrightarrow & \textbf{\textit{g}}_i \\ \textbf{\textit{x}}_B & \leftrightarrow & \textbf{\textit{g}}_B \\ \bar{\textbf{\textit{x}}}_i & \leftrightarrow & \nabla f_i \\ \bar{\textbf{\textit{x}}} & \leftrightarrow & \nabla F \end{array}$$

ing unbiased but potentially correlated estimators through subsampling. All of the later results on SGD will fall out as special cases following the correspondence in Table 2.

#### 3.1.1. GENERAL VARIANCE BOUND

**Theoretical Setup.** Consider the problem of estimating the population mean  $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} \bar{x}_i$  with a collection of RVs  $x_1, \dots, x_n$ , each an unbiased estimator of the component  $\bar{x}_i = \mathbb{E}x_i$ . Then, any subsampled ensemble

$$\mathbf{x}_{B} \triangleq \frac{1}{b} \sum_{i \in B} \mathbf{x}_{i} \quad \text{with} \quad B \sim \pi,$$

where  $\pi$  is an unbiased subsampling strategy with an effective sample size of  $b_{\rm eff}$ , is also an unbiased estimator of  $\bar{\mathbf{x}}$ . The goal is to analyze how the variance of the component estimators  ${\rm tr} \mathbb{V} \mathbf{x}_i$  for  $i=1,\ldots,n$  and the variance of  $\pi$  affect the variance of  $\mathbf{x}_B$ . The following condition characterizes the correlation between the component estimators:

**Assumption 3.** The component estimators  $\mathbf{x}_1, \dots, \mathbf{x}_n$  have finite variance  $\operatorname{tr} \mathbb{V} \mathbf{x}_i < \infty$  for all  $i = 1, \dots, n$  and, there exists some  $\rho \in [0, 1]$  for all  $i \neq j$  such that

$$\operatorname{tr} \operatorname{Cov} (\boldsymbol{x}_i, \boldsymbol{x}_j) \leq \rho \sqrt{\operatorname{tr} \mathbb{V} \boldsymbol{x}_i} \sqrt{\operatorname{tr} \mathbb{V} \boldsymbol{x}_j}.$$

**Remark 1.** Assumption 3 always holds with  $\rho = 1$  as a basic consequence of the Cauchy-Schwarz inequality.

**Remark 2.** For a collection of mutually independent estimators  $\mathbf{x}_1, \dots, \mathbf{x}_n$  such that  $\mathbf{x}_i \perp \mathbf{x}_j$  for all  $i \neq j$ , Assumption 3 holds with  $\rho = 0$ .

**Remark 3.** The equality in Assumption 3 holds with  $\rho = 0$  for independent estimators, while it holds with  $\rho = 1$  when they are perfectly positively correlated such that, for all  $i \neq j$ , there exists some constant  $\alpha_{ij} \geq 0$  such that  $\mathbf{x}_i = \alpha_{i,j} \mathbf{x}_j$ 

**Theorem 1.** Let the component estimators  $\mathbf{x}_1, \dots, \mathbf{x}_n$  satisfy Assumption 3. Then, the variance of the doubly stochastic estimator  $\mathbf{x}_B$  is bounded as

$$\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{B} \right] \leq V_{\operatorname{com}} + V_{\operatorname{cor}} + V_{\operatorname{sub}},$$

where

$$\begin{split} V_{\text{com}} &= \left(\frac{\rho}{b_{\text{eff}}} + \frac{1-\rho}{b}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \text{tr} \mathbb{V}\left[\boldsymbol{x}_{i}\right]\right), \\ V_{\text{cor}} &= \rho \left(1 - \frac{1}{b_{\text{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\text{tr} \mathbb{V}\left[\boldsymbol{x}_{i}\right]}\right)^{2}, \text{ and} \\ V_{\text{sub}} &= \frac{1}{b_{\text{off}}} \frac{1}{n} \sum_{i=1}^{n} \left\|\bar{\boldsymbol{x}}_{i} - \bar{\boldsymbol{x}}\right\|_{2}^{2}. \end{split}$$

Equality holds when the equality in Assumption 3 holds.

*Proof.* We start from the law of total (co)variance,

$$\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{B} \right] = \underbrace{\mathbb{E}_{\pi} \left[ \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{B} \mid B \right] \right]}_{\text{Variance of ensemble}} + \underbrace{\operatorname{tr} \mathbb{V}_{\pi} \left[ \mathbb{E} \left[ \mathbf{x}_{B} \mid B \right] \right]}_{\text{Variance of subsampling}}.$$

This splits the variance into the variance of the specific ensemble of *B* and subsampling variance. The main challenge is to relate the variance of the ensemble of *B* with the variance of the individual estimators in the sum

$$\mathbb{E}_{\pi}\left[\operatorname{tr}\mathbb{V}\left[\mathbf{x}_{B}\mid B\right]\right] = \mathbb{E}_{\pi}\left[\operatorname{tr}\mathbb{V}\left[\frac{1}{h}\sum_{i\in B}\mathbf{x}_{i}\right]\right].\tag{9}$$

Since the individual estimators may not be independent, analyzing the variance of the sum can be tricky. However, the following lemma holds generally:

**Lemma 9.** Let  $\mathbf{x}_1, ..., \mathbf{x}_b$  be a collection of vector-variate RVs dependent on some random variable B satisfying Assumption 3. Then, the expected variance of the sum of  $\mathbf{x}_1, ..., \mathbf{x}_b$  conditioned on B is bounded as

$$\mathbb{E}\left[\operatorname{tr}\mathbb{V}\left[\sum_{i=1}^{b}\boldsymbol{x}_{i}\mid\boldsymbol{B}\right]\right]\leq\rho\mathbb{V}\left[\boldsymbol{S}\right]+\rho(\mathbb{E}\boldsymbol{S})^{2}+\left(1-\rho\right)\mathbb{E}\left[\boldsymbol{V}\right],$$

where

$$S = \sum_{i=1}^b \sqrt{\operatorname{tr} \mathbb{V}\left[\boldsymbol{x}_i \mid \boldsymbol{B}\right]} \ \ and \ \ \boldsymbol{V} = \sum_{i=1}^b \operatorname{tr} \mathbb{V}\left[\boldsymbol{x}_i \mid \boldsymbol{B}\right].$$

Equality holds when the equality in Assumption 3 holds.

Here, S is the sum of conditional standard deviations, while V is the sum of conditional variances. Notice that the "variance of the variances" is playing a role: if we reduce the subsampling variance, then the variance of the ensemble,  $V_{\rm com}$ , also decreases.

The rest of the proof, along with the proof of Lemma 9, can be found in Appendix B.3 page 23.  $\Box$ 

In Theorem 1,  $V_{\rm com}$  is the contribution of the variance of the component estimators, while  $V_{\rm cor}$  is the contribution of the correlation between component estimators, and  $V_{\rm sub}$  is the subsampling variance.

## Monte Carlo with Subampling Without Replacement.

Theorem 1 is very general: it encompasses both the correlated and uncorrelated cases and matches the constants of all of the important special cases. We will demonstrate this in the following corollary along with variance reduction by Monte Carlo averaging of m i.i.d. samples. That is, we subsample over  $\mathbf{x}_1^m, \dots, \mathbf{x}_n^m$ , where each estimator is an m-sample Monte Carlo estimator:

$$\mathbf{x}_i^m \triangleq \frac{1}{m} \sum_{j=1}^m \mathbf{x}_i^{(j)},$$

where  $\mathbf{x}_i^{(1)}, \dots, \mathbf{x}_i^{(m)}$  are i.i.d replications with mean  $\bar{\mathbf{x}}_i = \mathbb{E}\mathbf{x}_i^{(j)}$ . Then, the variance of the doubly stochastic estimator  $\mathbf{x}_B$  of the mean  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{x}}_i$  defined as

$$\mathbf{x}_{B}^{m} \triangleq \frac{1}{h} \sum_{i \in B} \mathbf{x}_{i}^{m}$$
 with  $B \sim \pi$ ,

can be bounded as follows:

**Corollary 1.** For each j=1,...,m, let  $\mathbf{x}_1^{(j)},...,\mathbf{x}_n^{(j)}$  satisfy Assumption 3. Then, the variance of the doubly stochastic estimator  $\mathbf{x}_B^m$  of the mean  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{x}}_i$ , where  $\pi$  is b-minibatch sampling without replacement, satisfy the following corollaries:

(i) 
$$\rho = 1$$
 and  $1 < b < n$ :

$$\operatorname{tr} \mathbb{V}\left[\mathbf{x}_{B}^{m}\right] \leq \frac{n-b}{(n-1)mb} \left(\frac{1}{n} \sum_{i=1}^{n} \sigma_{i}^{2}\right) + \frac{n(b-1)}{(n-1)mb} \left(\frac{1}{n} \sum_{i=1}^{n} \sigma_{i}\right)^{2} + \frac{n-b}{(n-1)b} \tau^{2}$$

(ii) 
$$\rho = 1$$
 and  $b = 1$ :  

$$\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_B^m \right] \le \frac{1}{m} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i^2 \right) + \tau^2$$

(iii) 
$$\rho = 1$$
 and  $b = n$ :  

$$\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{B}^{m} \right] \leq \frac{1}{m} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \right)^{2}$$

(iv) 
$$\sigma_i = 0$$
 for all  $i = 1, ..., n$ :

$$\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{B}^{m} \right] \leq \frac{n-b}{(n-1)b} \tau^{2},$$

$$\begin{array}{c} (v) \ \rho = 0: \\ \operatorname{tr} \mathbb{V} \left[ \boldsymbol{x}_{B}^{m} \right] \leq \frac{1}{mb} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_{i}^{2} \right) + \frac{n-b}{(n-1)b} \tau^{2} \end{array}$$

where, for all i = 1, ..., n and any j = 1, ..., m,  $\sigma_i^2 = \text{tr} \mathbb{V} \, \mathbf{x}_i^{(j)} \qquad \text{is invidual variance and}$   $\tau^2 = \frac{1}{n} \sum_{i=1}^n \left\| \bar{\mathbf{x}}_i - \bar{\mathbf{x}} \right\|_2^2 \quad \text{is the subsampling variance.}$ 

Remark 4 (For dependent estimators, increasing b also reduces component variance.). Notice that, for case of  $\rho = 1$ , Corollary 1 (i), the term with  $\frac{1}{n} \sum_{i=1}^{n} \sigma_i^2$  is reduced in a rate of  $\mathcal{O}(1/mb)$ . This means reducing the subsampling noise by increasing b also reduces the noise of estimating each component. Furthermore, the first term dominates the second term as

$$\left(\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\right)^{2} \leq \frac{1}{n}\sum_{i=1}^{n}\sigma_{i}^{2},$$

as stated by Jensen's inequality. Therefore, despite correlations, increasing b will have a more significant effect since it reduces both dominant terms  $\frac{1}{n} \sum_{i=1}^{n} \sigma_i^2$  and  $\tau^2$ .

**Remark 5.** When independent estimators are used, Corollary 1 (v) shows that increasing b reduces the full variance in a O(1/b) rate, but increasing m does not.

**Remark 6.** Corollary 1 achieves all known endpoints in the context of SGD: For b = n (full batch), doubly SGD reduces to SGD with a Monte Carlo estimator, where there is no subsampling noise (no  $\tau^2$ ). When the Monte Carlo noise is 0, then doubly SGD reduces to SGD with a subsampling estimator (no  $\sigma_i$ ), retrieving the result of Gower et al. (2019).

#### 3.1.2. Gradient Variance Conditions for SGD

From Theorem 1, we can establish the ER and BV conditions (Section 2.3) of the doubly stochastic gradient estimators. Following the notation in Section 2.2, we will denote the doubly stochastic gradient estimator as  $\mathbf{g}_B$ , which combines the estimators  $\mathbf{g}_1, \ldots, \mathbf{g}_n$  according to the subsampling strategy  $B \sim \pi$ , which achieves an effective sample size of  $b_{\rm eff}$ . We will also use the corresponding minibatch subsampling estimator  $\nabla f_B$  for the analysis.

**Assumption 4.** For all  $x \in \mathcal{X}$ , the component gradient estimators  $\mathbf{g}_1(x), \dots, \mathbf{g}_n(x)$  satisfy Assumption 3 with some  $\rho \in [0,1]$ .

Again, this assumption is always met with  $\rho = 1$  and holds with  $\rho = 0$  if the estimators are independent.

**Assumption 5.** The subsampling estimator  $\nabla f_B$  satisfies the ER ( $\mathcal{L}_{\text{sub}}$ ) condition in Definition 1.

This is a classical assumption used to analyze SGD on finite sums and is automatically satisfied by Assumption 2. (See Lemma 10 in Appendix B.4.3 for a proof.)

**Assumption 6.** For all i = 1, ..., n and  $x \in \mathcal{X}$  and global minimizers  $x_* \in \arg\min_{x \in \mathcal{X}} F(x_*)$ , the component gradient estimator  $g_i$  satisfies at least one of the following variants of the ER condition:

$$(\mathbf{A}^{\text{CVX}}) \text{ tr} \mathbb{V} \left[ \mathbf{g}_{i}(\mathbf{x}) - \mathbf{g}_{i}(\mathbf{y}) \right] \leq 2\mathcal{L}_{i} \mathbf{D}_{f_{i}}(\mathbf{x}, \mathbf{y}),$$
 where  $f_{i}$  is convex.

(A<sup>ITP</sup>) tr
$$\mathbb{V}\left[\mathbf{g}_{i}\left(\mathbf{x}\right) - \mathbf{g}_{i}\left(\mathbf{y}\right)\right] \leq 2\mathcal{L}_{i}\left(f_{i}\left(\mathbf{x}\right) - f_{i}\left(\mathbf{x}_{*}\right)\right)$$
 where  $f_{i}\left(\mathbf{x}\right) \geq f_{i}\left(\mathbf{x}_{*}\right)$ .

(B) 
$$\operatorname{tr} \mathbb{V} \left[ \mathbf{g}_{i}(\mathbf{x}) - \mathbf{g}_{i}(\mathbf{y}) \right] \leq 2\mathcal{L}_{i} \left( F(\mathbf{x}) - F(\mathbf{x}_{*}) \right).$$

Each of these assumptions holds under different assumptions and problem setups. For instance,  $A^{CVX}$  holds only under componentwise convexity, while  $A^{ITP}$  requires majorization  $f_i(\mathbf{x}) \geq f_i(\mathbf{x}_*)$ , which is essentially assuming "interpolation" (Vaswani et al., 2019; Ma et al., 2018; Gower et al., 2021a) in the ERM context. Among these, (B) is the strongest since it directly relates the individual components  $f_1, \ldots, f_n$  with the full objective F.

We now state our result establishing the ER condition:

**Theorem 2.** Let Assumption 4 to 6 hold. Then, we have:

(i) If 
$$(A^{CVX})$$
 or  $(A^{ITP})$  hold,  $\mathbf{g}_B$  satisfies  $ER(\mathcal{L}_A)$ .

(ii) If (B) holds,  $\mathbf{g}_{B}$  satisfies ER ( $\mathcal{L}_{B}$ ).

where 
$$\mathcal{L}_{\max} = \max \{\mathcal{L}_1, \dots, \mathcal{L}_n\}$$
,

$$\mathcal{L}_{\mathrm{A}} = \left(\frac{\rho}{b_{\mathrm{eff}}} + \frac{1-\rho}{b}\right) \mathcal{L}_{\mathrm{max}} + \rho \left(1 - \frac{1}{b_{\mathrm{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i}\right) + \frac{\mathcal{L}_{\mathrm{sub}}}{b_{\mathrm{eff}}}$$

$$\begin{split} \mathcal{L}_{\mathrm{B}} &= \left(\frac{\rho}{b_{\mathrm{eff}}} + \frac{1-\rho}{b}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i}\right) \\ &+ \rho \left(1 - \frac{1}{b_{\mathrm{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\mathcal{L}_{i}}\right)^{2} + \frac{\mathcal{L}_{\mathrm{sub}}}{b_{\mathrm{eff}}}. \end{split}$$

See the *full proof* in page 25.

**Remark 7.** Assuming the conditions in Assumption 6 hold with the same value of  $\mathcal{L}_i$ , the inequality

$$\left(\frac{1}{n}\sum_{i=1}^{n}\sqrt{\mathcal{L}_{i}}\right)^{2} \leq \frac{1}{n}\sum_{i=1}^{n}\mathcal{L}_{i} \leq \mathcal{L}_{\max},$$

implies  $\mathcal{L}_{B} \leq \mathcal{L}_{A}$ .

Meanwhile, The BV condition follows by assuming equivalent conditions on each component estimator:

**Assumption 7.** Variance is bounded for all  $x_* \in \arg\min_{x \in \mathcal{X}} F(x)$  such that the following hold:

1. 
$$\frac{1}{n} \sum_{i=1}^{n} \left\| \nabla f_i(\mathbf{x}_*) \right\|_2^2 \le \tau^2$$
 for some  $\tau^2 < \infty$  and,

2. 
$$\operatorname{tr} \mathbb{V} [\mathbf{g}_i(\mathbf{x}_*)] \leq \sigma_i^2$$
 for some  $\sigma_i^2 < \infty$ , for all  $i = 1, \dots, n$ .

Based on these, Theorem 1 immediately yields the result:

**Theorem 3.** Let Assumption 4 and 7 hold. Then,  $\mathbf{g}_B$  satisfies BV  $(\sigma^2)$ , where

$$\begin{split} \sigma^2 &= \left(\frac{\rho}{b_{\mathrm{eff}}} + \frac{1-\rho}{b}\right) \left(\frac{1}{n} \sum_{i=1}^n \sigma_i^2\right) \\ &+ \rho \left(1 - \frac{1}{b_{\mathrm{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^n \sigma_i\right)^2 + \frac{\tau^2}{b_{\mathrm{eff}}}. \end{split}$$

Equality in Definition 2 holds if equality in Assumption 4 holds.

See the *full proof* in page 27.

As discussed in Section 2.4, Theorems 2 and 3 are sufficient to guarantee convergence of doubly SGD. For completeness, let us state a specific result for  $\rho = 1$ :

**Corollary 2.** Let the objective F satisfy Assumption 1 and 2, the global optimum  $\mathbf{x}_* = \arg\min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$  be a stationary point of F, the component gradient estimators  $\mathbf{g}_1, \dots, \mathbf{g}_n$  satisfy Assumption 6 (B) and 7, and  $\pi$  be b-minibatch sampling without replacement. Then the last iterate of SGD with  $\mathbf{g}_B$  is  $\epsilon$ -close to  $\mathbf{x}_*$  as  $\mathbb{E}\|\mathbf{x}_T - \mathbf{x}_*\|_2^2 \le \epsilon$  after a number of iterations of at least

$$T \ge 2 \max \left( C_{\text{var}} \frac{1}{\epsilon}, \ C_{\text{bias}} \right) \log \left( 2 || \boldsymbol{x}_0 - \boldsymbol{x}_* ||_2^2 \frac{1}{\epsilon} \right)$$

for some fixed stepsize where

$$C_{\text{var}} = \frac{2}{b} \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_i^2}{\mu^2} \right) + 2 \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_i}{\mu} \right)^2 + \frac{2}{b} \frac{\tau^2}{\mu^2},$$

$$C_{\text{bias}} = \frac{2}{b} \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\mathcal{L}_i}{\mu} \right) + 2 \left( \frac{1}{n} \sum_{i=1}^{n} \sqrt{\frac{\mathcal{L}_i}{\mu}} \right)^2 + \frac{2}{b} \frac{L}{\mu}.$$

See the *full proof* in page 28.

#### 3.2. Random Reshuffling of Stochastic Gradients

We now move to our analysis of SGD with random reshuffling (SGD-RR). In the doubly stochastic setting, this corresponds to reshuffling over stochastic estimators instead of gradients, which we will denote as doubly SGD-RR. In practice, doubly SGD-RR is often observed to converge faster than doubly SGD, even when dependent estimators are used.

#### 3.2.1. ALGORITHM

**Doubly SGD-RR** The algorithm is stated as follows:

- Reshuffle and partition the gradient estimators into minibatches of size b as  $P = \{P_1, ..., P_p\}$ , where p = n/b is the number of partitions or minibatches.
- **2** Perform gradient descent for i = 1, ..., p steps as

$$\boldsymbol{x}_{k}^{i+1} = \Pi_{\mathcal{X}} \left( \boldsymbol{x}_{k}^{i} - \gamma \boldsymbol{g}_{P_{i}} \left( \boldsymbol{x}_{k}^{i} \right) \right)$$

**3**  $k \leftarrow k + 1$  and go back to step **1**.

(We assume n is an integer multiple of b for clarity.) Here,  $i=1,\ldots,p$  denotes the step within the epoch,  $k=1,\ldots,K$  denotes the epoch number.

## 3.2.2. PROOF SKETCH

Why SGD-RR is Faster A key aspect of random reshuffling in the finite sum setting (SGD-RR) is that it uses conditionally biased gradient estimates. Because of this, on strongly convex finite sums, Mishchenko et al. (2020) show that the Lyapunov function for random reshuffling is not the usual  $\|\boldsymbol{x}_k^i - \boldsymbol{x}_*\|_2^2$ , but some *biased* Lyapunov function  $\|\boldsymbol{x}_k^i - \boldsymbol{x}_*^k\|_2^2$ , where the reference point is

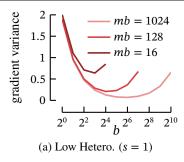
$$\boldsymbol{x}_{*}^{i} \triangleq \Pi_{\mathcal{X}} \left( \boldsymbol{x}_{*} - \gamma \sum_{j=0}^{i-1} \nabla f_{P_{i}}(\boldsymbol{x}_{*}) \right). \tag{10}$$

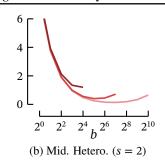
Under this definition, the convergence rate of SGD is not determined by the gradient variance anymore; it is determined by the squared error of the Lyapunov reference point,  $\|\boldsymbol{x}_*^i - \boldsymbol{x}_*\|_2^2$ . There are two key properties of this quantity:

- The peak mean-squared error decreases at a rate of  $\gamma^2$  with respect to the stepsize  $\gamma$ .
- The squared error is 0 at the following two endpoints: beginning of the epoch and at the end of the epoch.

For some stepsize achieving a  $\mathcal{O}(1/T)$  rate on SGD, these two properties combined result in SGD-RR attaining a  $\mathcal{O}(1/T^2)$  rate at exactly the end of each epoch.

**Is doubly SGD-RR as Fast as SGD-RR?** Unfortunately, doubly SGD-RR does not achieve the same rate as SGD-RR. Since stochastic gradients are used in addition to reshuffling, doubly SGD-RR deviates from the path that minimizes the biased Lyapunov function. Still, doubly SGD-RR does have provable benefits.





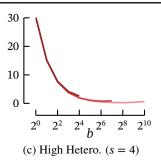


Figure 1. Trade-off between b and m on the gradient variance  $\operatorname{tr} \mathbb{V} g(x_*)$  under varying budgets  $m \times b$ . The problem is a finite sum of d = 10, n = 1024 isotropic quadratics with smoothness constants sampled as  $L_i \sim \operatorname{Inv-Gamma}(1/2, 1/2)$  and stationary points sampled as  $x_i^* \sim \mathcal{N}(\mathbf{0}_d, s^2\mathbf{I}_d)$ , where the gradient has additive noise of  $\eta \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ . Larger s means more heterogeneous data.

#### 3.2.3. COMPLEXITY ANALYSIS

We provide the general complexity guarantee for doubly SGD-RR on strongly convex objectives with  $\mu$ -strongly convex components and fully correlated component estimators ( $\rho = 1$ ):

**Theorem 4.** Let the objective F satisfy Assumption 1 and 2, where each component  $f_i$  is additionally  $\mu$ -strongly convex, and Assumption 6 (A<sup>CVX</sup>), 7 hold. Then, the last iterate  $\mathbf{x}_T$  of doubly SGD-RR is  $\epsilon$ -close to the global optimum  $\mathbf{x}_* = \arg\max_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$  such that  $\mathbb{E}||\mathbf{x}_T - \mathbf{x}_*||_2^2 \le \epsilon$  after a number of iterations of at least

$$T \, \geq \max \left( 4 C_{\mathrm{var}}^{\mathrm{com}} \frac{1}{\epsilon} + C_{\mathrm{var}}^{\mathrm{sub}} \frac{1}{\sqrt{\epsilon}}, \, \, C_{\mathrm{bias}} \right) \log \left( 2 \left\| \boldsymbol{x}_{1}^{0} - \boldsymbol{x}_{*} \right\|_{2}^{2} \frac{1}{\epsilon} \right)$$

for some fixed stepsize, where T = Kp = Kn/b,

$$\begin{split} &C_{\text{bias}} = \left(\mathcal{L}_{\text{max}} + L\right) / \mu \\ &C_{\text{var}}^{\text{com}} = \frac{2}{b} \left(\frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_{i}^{2}}{\mu^{2}}\right) + 2 \left(\frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_{i}}{\mu}\right)^{2}, \\ &C_{\text{var}}^{\text{sub}} = \sqrt{\frac{L_{\text{max}}}{\mu}} \frac{\sqrt{n}}{b} \frac{\tau}{\mu}. \end{split}$$

See the *full proof* in page 34.

**Remark 8.** When  $\sigma_i = 0$  for all i = 1, ..., n, the anytime convergence bound Theorem 5 in the Appendix reduces exactly to Theorem 1 of Mishchenko et al. (2020). Therefore, Theorem 5 is a strict generalization of SGD-RR to the doubly stochastic setting.

Using *m*-sample Monte Carlo improves the constants as follows:

**Corollary 3.** Let the assumptions of Theorem 4 hold. Then, for 1 < b < n and m-sample Monte Carlo, the same guarantees hold with the constant

$$C_{\text{var}}^{\text{com}} = \frac{2}{mb} \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_i^2}{\mu^2} \right) + \frac{2}{m} \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_i}{\mu} \right)^2.$$

**Remark 9.** Compared to doubly SGD, doubly SGD-RR improves the dependence on the subsampling noise  $\tau^2$  from  $\mathcal{O}(1/\varepsilon)$  to  $\mathcal{O}(1/\sqrt{\varepsilon})$ . Therefore, random reshuffling does improve the complexity of doubly SGD. Unfortunately, it also means that it does not achieve a better asymptotic complexity as in the finite sum setting. However, non-asymptotically, if the subsampling noise dominates component estimation noise, doubly SGD-RR will behave closely to an  $\mathcal{O}(1/\sqrt{\varepsilon})$  (or equivalently,  $\mathcal{O}(1/T)$ ) algorithm.

**Remark 10.** As was the case with independent subsampling, increasing *b* also reduces component estimation noise for RR-SGD. However, the impact on the complexity is more subtle. Consider that the iteration complexity is

$$\mathcal{O}\left(\kappa_{\sigma}^{2}\left(\frac{1}{mb} + \frac{1}{m}\right)\frac{1}{\epsilon} + \kappa \kappa_{\tau} \frac{\sqrt{n}}{b} \frac{1}{\sqrt{\epsilon}}\right),\tag{11}$$

where  $\kappa_{\sigma} = \max_{i=1,\dots,n} \sigma_i/\mu$ ,  $\kappa_{\tau} = \tau/\mu$  and  $\kappa = \max_{i=1,\dots,n} L_i/\mu$ . The  $1/\varepsilon$  term decreases the fastest with m. Therefore, it might seem that increasing m is advantageous. However, the  $1/\sqrt{\varepsilon}$  term has a  $\mathcal{O}\left(\sqrt{n}\right)$  dependence on the dataset size, which would be non-negligible for large datasets. As a result, in the large n, large  $\varepsilon$  regime, increasing b over m should be more effective.

**Remark 11.** Eq. (11) also implies that, for dependent estimators, doubly SGD-RR achieves an asymptotic speedup of n/b compared to full-batch SGD with only component estimation noise. Assume that the sample complexity of a single estimate is  $\Theta(mb)$  ( $\Theta(mn)$  for full-batch). Then, the sample complexity of doubly SGD-RR is  $\mathcal{O}(b^1/\varepsilon)$  and  $\mathcal{O}(n^1/\varepsilon)$  for full-batch SGD. However, the n/b seed-up comes from correlations. Therefore, for independent estimators, the asymptotic complexity of the two is equal.

## 4. Simulation

**Setup** We evaluate the insight on the tradeoff between b and m for correlated estimators on a synthetic problem. In particular, we set

$$f_i(\mathbf{x}; \boldsymbol{\eta}) = \frac{L_i}{2} \|\mathbf{x} - \mathbf{x}_i^* + \boldsymbol{\eta}\|_2^2,$$

where the smoothness constants  $L_i \sim \text{Inv-Gamma}(1/2, 1/2)$  and the stationary points  $\boldsymbol{x}_i^* \sim \mathcal{N}(\mathbf{0}_d, s^2\mathbf{I}_d)$  are sampled randomly, where  $\mathbf{0}_d$  is a vector of d zeros and  $\mathbf{I}_d$  is a  $d \times d$  identity matrix. Then, we compute the gradient variance on the global optimum, corresponding to computing the BV (Definition 2) constant. Note that  $s^2$  here corresponds to the "heterogeneity" of the data. We make the estimators dependent by sharing  $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_m$  across the batch.

**Results** The results are shown in Fig. 1. At low heterogeneity, there exists a "sweet spot" between m and b. However, this sweet spot moves towards large values of b, where, at high heterogeneity levels, the largest values of b are more favorable. Especially in the low budget regime where  $mb \ll n$ , the largest b values appear to achieve the lowest variance. This confirms our theoretical results that a large b should be preferred on challenging (large number of datapoints, high heterogeneity) problems.

## 5. Discussions

## 5.1. Applications

In Appendix C, we establish Assumption 6 and 7 on the following applications:

- ERM with Randomized Smoothing: In this problem, we consider ERM, where the model weights are perturbed by noise. This variant of ERM has recently gathered interest as it is believed to improve generalization performance (Orvieto et al., 2023; Liu et al., 2021). In Appendix C.1, we establish Assumption 6 (A<sup>ITP</sup>) under the interpolation assumption.
- **Reparameterization Gradient**: In certain applications, *e.g.*, variational inference, generative modeling, and reinforcement learning (see Mohamed et al., 2020, §5), the optimization problem is over the parameters of some distribution, which is taken expectation over. Among gradient estimators for this problem, the reparameterization gradient is widely used due to lower variance (Xu et al., 2019). For this, in Appendix C.2, we establish Assumption 6 (A<sup>CVX</sup>) and (B) by assuming a convexity and smooth integrand.

#### 5.2. Related Works

Unlike SGD in the finite sum setting, doubly SGD has received little interest. Previously, Bietti & Mairal (2017); Zheng & Kwok (2018); Kulunchakov & Mairal (2020) have studied the convergence of variance-reduced gradients (Gower et al., 2020) specific to the doubly stochastic setting under the uniform Lipchitz integrand assumption ( $g_i(\cdot; \eta)$ ) is L-Lipschitz for all  $\eta$ ). Although this assumption has often been used in the stochastic optimization literature (Nemirovski et al., 2009; Moulines & Bach, 2011; Shalev-Shwartz et al., 2009; Nguyen et al., 2018), it is easily shown to be restrictive: for some L-smooth  $\hat{f}_i(x)$ ,

 $\nabla f_i(\mathbf{x}; \boldsymbol{\eta}) = \nabla \widehat{f}_i(\mathbf{x}) + x_1 \boldsymbol{\eta}$  is not *L*-Lipschitz unless the support of  $\boldsymbol{\eta}$  is compact. In contrast, we established results under weaker conditions. We also provide a discussion on the relationships of different conditions in Appendix A.

Furthermore, we extended doubly SGD to the case where random reshuffling is used in place of sampling independent batches. In the finite-sum setting, the fact that SGD-RR converges faster than independent subsampling (SGD) has been empirically known for a long time (Bottou, 2009). While Gürbüzbalaban et al. (2021) first demonstrated that SGD-RR can be fast for quadratics, a proof under general conditions was demonstrated recently (Haochen & Sra, 2019): In the strongly convex setting, Mishchenko et al. (2020) Ahn et al. (2020); Nguyen et al. (2021) establish a  $\mathcal{O}\left(1/\sqrt{\varepsilon}\right)$  complexity to be  $\varepsilon$ -accurate, which is tight in terms of asymptotic complexity (Safran & Shamir, 2020; Cha et al., 2023; Safran & Shamir, 2021).

Lastly, Dai et al. (2014); Xie et al. (2015); Shi et al. (2021) provided convergence guarantees for doubly SGD for ERM of random feature kernel machines. However, these analyses are based on concentration arguments that doubly SGD does not deviate too much from the optimization path of finite-sum SGD. Unfortunately, concentration arguments require stronger assumptions on the noise, and their analysis is application-specific. In contrast, we provide a general analysis under the general ER assumption.

## 5.3. Conclusions

In this work, we analyzed the convergence of SGD with doubly stochastic and dependent gradient estimators. In particular, we showed that if the gradient estimator of each component satisfies the ER and BV conditions, the doubly stochastic estimator also satisfies both conditions; this implies convergence of doubly SGD.

Practical Recommendations An unusual conclusion of our analysis is that when Monte Carlo is used with minibatch subsampling, it is generally more beneficial to increase the minibatch size b instead of the number of Monte Carlo samples m. That is, for both SGD and SGD-RR, increasing b decreases the variance in a rate close to 1/b when (i) the gradient variance of the component gradient estimators varies greatly such that  $\left(\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}\right)^{2}\ll\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}^{2}$  or when (ii) the estimators are independent as  $\rho=0$ . Surprisingly, such a benefit persists even in the interpolation regime  $\tau^2 = 0$ . On the contrary, when the estimators are both dependent and have similar variance, it is necessary to increase both m and b, where a sweet spot between the two exists. However, such a regime is unlikely to occur in practice; in statistics and machine learning applications, the variance of the gradient estimators tends to vary greatly due to the heterogeneity of data.

## Acknowledgements

The authors would like to thank the anonymous reviewers for their comments and Jason Altschuler (UPenn) for numerous suggestions that strengthened the work.

K. Kim was supported by a gift from AWS AI to Penn Engineering's ASSET Center for Trustworthy AI; Y.-A. Ma was funded by the NSF Grants [NSF-SCALE MoDL-2134209] and [NSF-CCF-2112665 (TILOS)], the U.S. Department Of Energy, Office of Science, as well as the DARPA AIE program; J. R. Gardner was supported by NSF award [IIS-2145644].

## **Impact Statement**

This paper presents a theoretical analysis of stochastic gradient descent under doubly stochastic noise to broaden our understanding of the algorithm. The work itself is theoretical, and we do not expect direct societal consequences, but SGD with doubly stochastic gradients is widely used in various aspects of machine learning and statistics. Therefore, we inherit the societal impact of the downstream applications of SGD.

### References

- Ahn, K., Yun, C., and Sra, S. SGD with shuffling: Optimal rates without component convexity and large epoch requirements. In *Advances in Neural Information Processing Systems*, volume 33, pp. 17526–17535. Curran Associates, Inc., 2020. (page 9)
- Bassily, R., Smith, A., and Thakurta, A. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *Proceedings of the IEEE Annual Symposium on Foundations of Computer Science*, FOCS '14, pp. 464–473, USA, October 2014. IEEE Computer Society. (page 1)
- Bietti, A. and Mairal, J. Stochastic optimization with variance reduction for infinite datasets with finite sum structure. In *Advances in Neural Information Processing Systems*, volume 30, pp. 1623–1633. Curran Associates, Inc., 2017. (pages 1, 9)
- Bottou, L. On-line learning and stochastic approximations. In *On-Line Learning in Neural Networks*, pp. 9–42. Cambridge University Press, 1 edition, January 1999. (page 1)
- Bottou, L. Curiously fast convergence of some stochastic gradient descent algorithms. Unpublished open problem at the International Symposium on Statistical Learning and Data Sciences (SLDS), 2009. (page 9)

- Bottou, L., Curtis, F. E., and Nocedal, J. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311, January 2018. (page 1)
- Cha, J., Lee, J., and Yun, C. Tighter lower bounds for shuffling SGD: Random permutations and beyond. In *Proceedings of the International Conference on Machine Learning*, volume 202 of *PMLR*, pp. 3855–3912. JMLR, July 2023. (page 9)
- Csiba, D. and Richtárik, P. Importance sampling for minibatches. *Journal of Machine Learning Research*, 19(27): 1–21, 2018. (pages 3, 37)
- Dai, B., Xie, B., He, N., Liang, Y., Raj, A., Balcan, M.-F. F., and Song, L. Scalable kernel methods via doubly stochastic gradients. In *Advances in Neural Information Processing Systems*, volume 27, pp. 3041–3049. Curran Associates, Inc., 2014. (pages 1, 3, 9)
- Domke, J. Provable gradient variance guarantees for black-box variational inference. In *Advances in Neural Information Processing Systems*, volume 32, pp. 329–338. Curran Associates, Inc., 2019. (pages 2, 37)
- Domke, J. Provable smoothness guarantees for black-box variational inference. In *Proceedings of the International Conference on Machine Learning*, volume 119 of *PMLR*, pp. 2587–2596. JMLR, July 2020. (page 38)
- Domke, J., Gower, R., and Garrigos, G. Provable convergence guarantees for black-box variational inference. In *Advances in Neural Information Processing Systems*, volume 36, pp. 66289–66327. Curran Associates, Inc., 2023. (page 16)
- Duchi, J. C., Bartlett, P. L., and Wainwright, M. J. Randomized smoothing for stochastic optimization. *SIAM Journal on Optimization*, 22(2):674–701, January 2012. (page 35)
- Garrigos, G. and Gower, R. M. Handbook of convergence theorems for (stochastic) gradient methods. Preprint arXiv:2301.11235, arXiv, February 2023. (pages 4, 19)
- Gorbunov, E., Hanzely, F., and Richtarik, P. A unified theory of SGD: Variance reduction, sampling, quantization and coordinate descent. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 108 of *PMLR*, pp. 680–690. JMLR, June 2020. (page 37)
- Gower, R., Sebbouh, O., and Loizou, N. SGD for structured nonconvex functions: Learning rates, minibatching and interpolation. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 130 of *PMLR*, pp. 1315–1323. JMLR, March 2021a. (pages 3, 4, 6, 15, 16, 21, 35)

- Gower, R. M., Loizou, N., Qian, X., Sailanbayev, A., Shulgin, E., and Richtárik, P. SGD: General analysis and improved rates. In *Proceedings of the International Conference on Machine Learning*, volume 97 of *PMLR*, pp. 5200–5209. JMLR, June 2019. (pages 2, 3, 4, 6, 16, 21, 37)
- Gower, R. M., Schmidt, M., Bach, F., and Richtarik, P. Variance-reduced methods for machine learning. *Proceedings of the IEEE*, 108(11):1968–1983, November 2020. (page 9)
- Gower, R. M., Richtárik, P., and Bach, F. Stochastic quasi-gradient methods: Variance reduction via Jacobian sketching. *Mathematical Programming*, 188(1): 135–192, July 2021b. (pages 2, 3, 4)
- Guminov, S., Gasnikov, A., and Kuruzov, I. Accelerated methods for weakly-quasi-convex optimization problems. *Computational Management Science*, 20(1): 36, December 2023. (pages 4, 16)
- Gürbüzbalaban, M., Ozdaglar, A., and Parrilo, P. A. Why random reshuffling beats stochastic gradient descent. *Mathematical Programming*, 186(1):49–84, March 2021. (page 9)
- Haochen, J. and Sra, S. Random shuffling beats SGD after finite epochs. In *Proceedings of the International Conference on Machine Learning*, volume 97 of *PMLR*, pp. 2624–2633. JMLR, May 2019. (page 9)
- Hinder, O., Sidford, A., and Sohoni, N. Near-optimal methods for minimizing star-convex functions and beyond. In *Proceedings of Conference on Learning Theory*, volume 125 of *PMLR*, pp. 1894–1938. JMLR, July 2020. (page 16)
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pp. 6840–6851. Curran Associates, Inc., 2020. (page 1)
- Johnson, R. and Zhang, T. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances* in *Neural Information Processing Systems*, volume 26, pp. 315–323. Curran Associates, Inc., 2013. (page 3)
- Karimi, H., Nutini, J., and Schmidt, M. Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In *Machine Learning and Knowledge Discovery in Databases*, Lecture Notes in Computer Science, pp. 795–811, Cham, 2016. Springer International Publishing. (pages 15, 38)
- Khaled, A. and Richtárik, P. Better theory for SGD in the nonconvex world. *Transactions of Machine Learning Research*, 2023. (pages 4, 16)

- Kim, K., Oh, J., Wu, K., Ma, Y., and Gardner, J. R. On the convergence of black-box variational inference. In *Advances in Neural Information Processing Systems*, volume 36, pp. 44615–44657, New Orleans, LA, USA, December 2023. Curran Associates Inc. (pages 2, 38)
- Kingma, D. P. and Welling, M. Auto-encoding variational Bayes. In *Proceedings of the International Conference* on *Learning Representations*, Banff, AB, Canada, April 2014. (pages 1, 37)
- Kingma, D. P., Salimans, T., and Welling, M. Variational dropout and the local reparameterization trick. In *Advances in Neural Information Processing Systems*, volume 28, pp. 2575–2583. Curran Associates, Inc., 2015. (pages 2, 3)
- Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., and Blei, D. M. Automatic differentiation variational inference. *Journal of Machine Learning Research*, 18(14): 1–45, 2017. (pages 1, 3, 37)
- Kulunchakov, A. and Mairal, J. Estimate sequences for stochastic composite optimization: Variance reduction, acceleration, and robustness to noise. *Journal of Machine Learning Research*, 21(155):1–52, 2020. (pages 1, 9)
- Liu, T., Li, Y., Wei, S., Zhou, E., and Zhao, T. Noisy gradient descent converges to flat minima for nonconvex matrix factorization. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 130 of *PMLR*, pp. 1891–1899. JMLR, March 2021. (pages 1, 3, 9, 35)
- Ma, S., Bassily, R., and Belkin, M. The power of interpolation: Understanding the effectiveness of SGD in modern over-parametrized learning. In *Proceedings of the International Conference on Machine Learning*, volume 80 of *PMLR*, pp. 3325–3334. JMLR, July 2018. (pages 6, 35)
- Mishchenko, K., Khaled, A., and Richtarik, P. Random reshuffling: Simple analysis with vast improvements. In *Advances in Neural Information Processing Systems*, volume 33, pp. 17309–17320. Curran Associates, Inc., 2020. (pages 7, 8, 9, 30, 33)
- Mohamed, S., Rosca, M., Figurnov, M., and Mnih, A. Monte Carlo gradient estimation in machine learning. *Journal of Machine Learning Research*, 21(132):1–62, 2020. (pages 9, 37)
- Moulines, E. and Bach, F. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems*, volume 24, pp. 451–459. Curran Associates, Inc., 2011. (pages 2, 3, 9, 16)

- Needell, D. and Ward, R. Batched stochastic gradient descent with weighted sampling. In *Approximation Theory XV: San Antonio 2016*, Springer Proceedings in Mathematics & Statistics, pp. 279–306, Cham, 2017. Springer International Publishing. (page 37)
- Needell, D., Srebro, N., and Ward, R. Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm. *Mathematical Programming*, 155(1): 549–573, January 2016. (page 37)
- Nemirovski, A., Juditsky, A., Lan, G., and Shapiro, A. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4): 1574–1609, January 2009. (pages 1, 9)
- Nesterov, Y. and Polyak, B. Cubic regularization of Newton method and its global performance. *Mathematical Programming*, 108(1):177–205, August 2006. (page 16)
- Nesterov, Yu. Smooth minimization of non-smooth functions. *Mathematical Programming*, 103(1):127–152, May 2005. (page 35)
- Nguyen, L., Nguyen, P. H., van Dijk, M., Richtarik, P., Scheinberg, K., and Takac, M. SGD and Hogwild! Convergence without the bounded gradients assumption. In *Proceedings of the International Conference on Machine Learning*, volume 80 of *PMLR*, pp. 3750–3758. JMLR, July 2018. (pages 1, 3, 9, 15)
- Nguyen, L. M., Tran-Dinh, Q., Phan, D. T., Nguyen, P. H., and Van Dijk, M. A unified convergence analysis for shuffling-type gradient methods. *The Journal of Machine Learning Research*, 22(1):207:9397–207:9440, January 2021. (page 9)
- Orvieto, A., Raj, A., Kersting, H., and Bach, F. Explicit regularization in overparametrized models via noise injection. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 206 of *PMLR*, pp. 7265–7287. JMLR, April 2023. (pages 1, 3, 9, 35)
- Polyak, B. T. and d Aleksandr Borisovich, T. Optimal order of accuracy of search algorithms in stochastic optimization. *Problemy Peredachi Informatsii*, 26(2):45–53, 1990. (page 35)
- Ranganath, R., Gerrish, S., and Blei, D. Black box variational inference. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 33 of *PMLR*, pp. 814–822. JMLR, April 2014. (page 1)
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the International*

- Conference on Machine Learning, volume 32 of PMLR, pp. 1278–1286. JMLR, June 2014. (pages 1, 37)
- Richtárik, P. and Takáč, M. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, 156(1-2):433–484, March 2016. (page 3)
- Robbins, H. and Monro, S. A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3): 400–407, September 1951. (page 1)
- Safran, I. and Shamir, O. How good is SGD with random shuffling? In *Proceedings of Conference on Learning Theory*, volume 125 of *PMLR*, pp. 3250–3284. JMLR, July 2020. (page 9)
- Safran, I. and Shamir, O. Random shuffling beats SGD only after many epochs on ill-conditioned problems. In *Advances in Neural Information Processing Systems*, volume 34, pp. 15151–15161. Curran Associates, Inc., 2021. (page 9)
- Schmidt, M. and Roux, N. L. Fast convergence of stochastic gradient descent under a strong growth condition. arXiv Preprint arXiv:1308.6370, arXiv, August 2013. (page 16)
- Shalev-Shwartz, S., Shamir, O., Srebro, N., and Sridharan, K. Stochastic convex optimization. In *Proceedings* of the Conference on Computational Learning Theory, June 2009. (page 9)
- Shalev-Shwartz, S., Singer, Y., Srebro, N., and Cotter, A. Pegasos: Primal estimated sub-gradient solver for SVM. *Mathematical Programming*, 127(1):3–30, March 2011. (page 1)
- Shi, W., Gu, B., Li, X., Deng, C., and Huang, H. Triply stochastic gradient method for large-scale nonlinear similar unlabeled classification. *Machine Learning*, 110(8): 2005–2033, August 2021. (pages 1, 9)
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the International Conference on Machine Learning*, volume 37 of *PMLR*, pp. 2256–2265. JMLR, June 2015. (page 1)
- Song, S., Chaudhuri, K., and Sarwate, A. D. Stochastic gradient descent with differentially private updates. In *Proceedings of the IEEE Global Conference on Signal and Information Processing*, pp. 245–248, Austin, TX, USA, December 2013. IEEE. (page 1)
- Song, Y. and Ermon, S. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*, volume 32, pp. 11918–11930. Curran Associates, Inc., 2019. (page 1)

- Stich, S. U. Unified optimal analysis of the (stochastic) gradient method. Preprint arXiv:1907.04232, arXiv, December 2019. (page 4)
- Titsias, M. and Lázaro-Gredilla, M. Doubly stochastic variational Bayes for non-conjugate inference. In *Proceedings of the International Conference on Machine Learning*, volume 32 of *PMLR*, pp. 1971–1979. JMLR, June 2014. (pages 1, 3, 37)
- Vapnik, V. Principles of risk minimization for learning theory. In *Advances in Neural Information Processing Systems*, volume 4, pp. 831–838. Morgan-Kaufmann, 1991. (page 1)
- Vaswani, S., Bach, F., and Schmidt, M. Fast and faster convergence of SGD for over-parameterized models and an accelerated perceptron. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 89 of *PMLR*, pp. 1195–1204. JMLR, April 2019. (pages 6, 16, 35)
- Wright, S. J. and Recht, B. *Optimization for Data Analysis*. Cambridge University Press, New York, 2021. (page 16)
- Xie, B., Liang, Y., and Song, L. Scale up nonlinear component analysis with doubly stochastic gradients. In *Advances in Neural Information Processing Systems*, volume 28, pp. 2341–2349. Curran Associates, Inc., 2015. (page 9)
- Xu, M., Quiroz, M., Kohn, R., and Sisson, S. A. Variance reduction properties of the reparameterization trick. In Proceedings of the International Conference on Artificial Intelligence and Statistics, volume 89 of PMLR, pp. 2711–2720. JMLR, April 2019. (pages 9, 37)
- Zheng, S. and Kwok, J. T.-Y. Lightweight stochastic optimization for minimizing finite sums with infinite data. In *Proceedings of the International Conference on Machine Learning*, volume 80 of *PMLR*, pp. 5932–5940. JMLR, July 2018. (pages 1, 9)

# Demystifying SGD with Doubly Stochastic Gradients

# \_TABLE OF CONTENTS\_

1	Introduction	1
2	Preliminaries  2.1 Stochastic Gradient Descent on Finite-Sums  2.2 Doubly Stochastic Gradients  2.3 Technical Assumptions on Gradient Estimators  2.4 Convergence Guarantees for SGD	2 2 3 3 4
3	Main Results         3.1 Doubly Stochastic Gradients       3.1.1 General Variance Bound         3.1.2 Gradient Variance Conditions for SGD         3.2 Random Reshuffling of Stochastic Gradients         3.2.1 Algorithm         3.2.2 Proof Sketch         3.2.3 Complexity Analysis	5 5 6 7 7 7 8
4	Simulation	8
5	Discussions5.1 Applications5.2 Related Works5.3 Conclusions	<b>9</b> 9 9
A	Gradient Variance Conditions  A.1 Definitions	15 15 15 16 16
В		18 18 21 23 25 25 27 27 29 29 31 34
C	Applications  C.1 ERM with Randomized Smoothing C.1.1 Description C.1.2 Preliminaries C.1.3 Theoretical Analysis  C.2 Reparameterization Gradient C.2.1 Description C.2.2 Preliminaries C.2.3 Theoretical Analysis	35 35 35 35 36 37 37 37 38

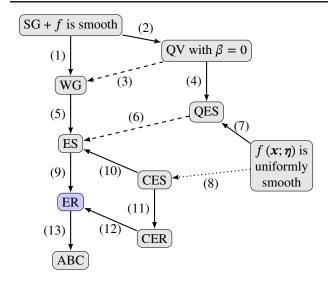


Figure 2. Implications between general gradient variance conditions for some unbiased estimator  $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}; \boldsymbol{\eta})$  of  $\nabla f(\mathbf{x}) = \mathbb{E}\mathbf{g}(\mathbf{x})$ . The dashed arrows  $(-- \rightarrow)$  hold if f is further assumed to be QFG; the dotted arrow  $(\cdots \rightarrow)$  holds if the integrand  $f(\mathbf{x}; \boldsymbol{\eta})$  is uniformly convex such that it is convex with respect to  $\mathbf{x}$  for any fixed  $\boldsymbol{\eta}$ . (1), (5), (9), (13) are established by Gower et al. (2021a, Theorem 3.4); (2) is proven in Proposition 3; (3) is proven in Proposition 7; (4) is proven in Proposition 6; (7) is proven in Proposition 4; (8) is proven in Proposition 5; (6) is proven by Nguyen et al. (2018, Lemma 2) but we restate the proof in Proposition 8; (11) is proven in Proposition 2; (10), (12) hold trivially if  $\mathbf{x}_* \in \arg\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$  are all stationary points.

# **A. Gradient Variance Conditions**

In this section, we will discuss some additional aspects of the ER and ES conditions introduced in Section 2.3. We will also look into alternative gradient variance conditions that have been proposed in the literature and their relationship with the ER condition.

## A.1. Definitions

For this section, we will use the following additional definitions:

**Definition 3** (Quadratic Functional Growth; QFG). We say  $f: \mathcal{X} \to \mathbb{R}$  satisfies  $\mu$ -quadratic functional growth if there exists some  $\mu > 0$  such that

$$\frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_*\|_2^2 \le f(\mathbf{x}) - f(\mathbf{x}_*)$$

holds for all  $\mathbf{x} \in \mathcal{X}$ , where  $\mathbf{x}_* \in \arg\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ .

This condition implies that f grows at least as fast as some quadratic and is weaker than the Polyak-Łojasiewicz. However, for any convex function f that satisfies this condition also means that f is  $\mu$ -strongly convex (Karimi et al.,

2016).

**Definition 4** (Uniform Smoothness). For the unbiased estimator  $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}; \boldsymbol{\eta})$  of  $\nabla F(\mathbf{x}) = \mathbb{E} \nabla f(\mathbf{x}; \boldsymbol{\eta}) = \mathbb{E} \nabla f(\mathbf{x}; \boldsymbol{\eta})$ , we say the integrand  $\nabla f_i(\mathbf{x}; \boldsymbol{\eta})$  satisfies uniform L-smoothness if there exist some  $L < \infty$  such that, for any fixed  $\boldsymbol{\eta}$ ,

$$\left\| \nabla f(\mathbf{x}; \boldsymbol{\eta}) - \nabla f(\mathbf{x}'; \boldsymbol{\eta}) \right\|_{2} \le L \|\mathbf{x} - \mathbf{x}'\|_{2}$$

holds for all  $(x, x') \in \mathcal{X}^2$  simultaneously.

As discussed in Sections 1 and 5.2, this condition is rather strong: it does not hold for multiplicative noise unless the support is bounded.

**Definition 5** (Uniform Convexity). For the unbiased estimator  $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}; \boldsymbol{\eta})$  of  $\nabla F(\mathbf{x}) = \mathbb{E}\nabla f(\mathbf{x}; \boldsymbol{\eta})$ , we say the integrand  $f(\mathbf{x}; \boldsymbol{\eta})$  is uniformly convex if it is convex for any  $\boldsymbol{\eta}$  such that, for any fixed  $\boldsymbol{\eta}$ ,

$$f(\mathbf{x}; \boldsymbol{\eta}) - f(\mathbf{x}'; \boldsymbol{\eta}) \le \langle \nabla f(\mathbf{x}; \boldsymbol{\eta}), \mathbf{x} - \mathbf{x}' \rangle$$

holds for all  $(x, x') \in \mathcal{X}^2$  simultaneously.

#### A.2. Additional Gradient Variance Conditions

For some estimator  $\mathbf{g}$  of  $\nabla f$ , the following conditions have been considered in the literature:

• Strong growth condition (SG):

$$\mathbb{E}\|\boldsymbol{g}\left(\boldsymbol{x}\right)\|_{2}^{2} \leq \rho \|\nabla f\left(\boldsymbol{x}\right)\|_{2}^{2}$$

• Weak growth condition (WG):

$$\mathbb{E}\|\boldsymbol{g}\left(\boldsymbol{x}\right)\|_{2}^{2} \leq \rho\left(f\left(\boldsymbol{x}\right) - f\left(\boldsymbol{x}_{*}\right)\right)$$

• Quadratic variance (QV):

$$\mathbb{E}\|\boldsymbol{g}\left(\boldsymbol{x}\right)\|_{2}^{2} \leq \alpha \left\|\boldsymbol{x} - \boldsymbol{x}_{*}\right\|_{2}^{2} + \beta$$

• Convex expected smoothness (CES):

$$\mathbb{E}\|\boldsymbol{g}\left(\boldsymbol{x}\right) - \boldsymbol{g}\left(\boldsymbol{y}\right)\|_{2}^{2} \le 2\mathcal{L}D_{f}\left(\boldsymbol{x}, \boldsymbol{y}\right)$$

• Convex expected residual (CER):

$$\operatorname{tr} \mathbb{V} \left[ \mathbf{g} \left( \mathbf{x} \right) - \mathbf{g} \left( \mathbf{y} \right) \right] \leq 2 \mathcal{L} D_{f} \left( \mathbf{x}, \mathbf{y} \right)$$

• Quadratic expected smoothness (QES):

$$\mathbb{E}||\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})||_{2}^{2} \le \mathcal{L}^{2}||\mathbf{x} - \mathbf{y}||_{2}^{2}$$

• ABC:

$$\mathbb{E}\|\mathbf{g}(\mathbf{x})\|_{2}^{2} \le A(f(\mathbf{x}) - f(\mathbf{x}_{*})) + B\|\nabla f(\mathbf{x})\|_{2}^{2} + C$$

Here,  $x_* \in \operatorname{arg\,min}_{x \in \mathcal{X}} f(x)$  is any stationary point of f and the stated conditions should hold for all  $(x, y) \in \mathcal{X}^2$ .

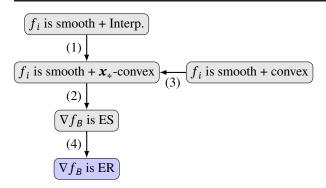


Figure 3. Implications of assumptions on the components  $f_1, ..., f_n$  to the minibatch subsampling gradient estimator  $\nabla f_B$  of  $F = \frac{1}{n}(f_1 + ... + f_n)$ . (1), (4) are established by Gower et al. (2021a, Theorem 3.4), while (3) trivially follows from the fact that  $x_*$ -convexity is strictly weaker than (global) convexity, and (2) was established by Gower et al. (2019, Proposition 3.10).

SG was used by Schmidt & Roux (2013) to establish the linear convergence of SGD for strongly convex objectives, and O(1/T) convergence for convex objectives; **WG** was proposed by Vaswani et al. (2019) to establish similar guarantees to SG under a verifiable condition; QV was used to establish the non-asymptotic convergence on strongly convex functions by Wright & Recht (2021), while convergence on general convex functions was established by Domke et al. (2023), including stochastic proximal gradient descent; QES was used by (Moulines & Bach, 2011) to establish one of the earliest general non-asymptotic convergence results for SGD on strongly convex objectives; ABC was used by Khaled & Richtárik (2023) to establish convergence of SGD for non-convex smooth functions. (See also Khaled & Richtárik (2023) for a comprehensive overview of these conditions.) The relationship of these conditions with the ER condition are summarized in Fig. 2.

As demonstrated in Fig. 2 and discussed by Khaled & Richtárik (2023), the ABC condition is the weakest of all. However, the convergence guarantees for problems that exclusively satisfy the ABC condition are weaker than others. (For instance, the number of iterations *T* has to be fixed *a priori*.) On the other hand, the ER condition retrieves most of the strongest known guarantees for SGD; some of which were listed in Section 2.4.

# A.3. Establishing the ER Condition

For subsampling estimators, it is possible to establish some of the gradient variance conditions through general assumptions on the components. See some examples in Fig. 3. Here, we use the following definitions:

**Definition 6.** For the finite sum objective  $F = \frac{1}{n}(f_1 + ... + f_n)$ , we say interpolation holds if, for all i = 1, ..., n,

$$f_i(\mathbf{x}_*) \leq f_i(\mathbf{x}),$$

holds for all  $\mathbf{x} \in \mathcal{X}$ , where  $\mathbf{x}_* \in \arg\min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$ .

**Definition 7.** For the finite sum objective  $F = \frac{1}{n}(f_1 + ... + f_n)$ , we say the components are  $x_*$ -convex if, for all i = 1, ..., n,

$$f_i(\mathbf{x}_*) - f_i(\mathbf{x}) \leq \langle \nabla f_i(\mathbf{x}_*), \mathbf{x}_* - \mathbf{x} \rangle$$

holds for all  $x \in \mathcal{X}$ , where  $x_* \in \arg\min_{x \in \mathcal{X}} F(x)$ .

This assumption is a weaker version of convexity; convexity needs to hold with respect to  $x_*$  only. It is closely related to star (Nesterov & Polyak, 2006) and quasar convexity (Hinder et al., 2020; Guminov et al., 2023).

## A.4. Proofs of Implications in Fig. 2

We prove new implication results between some of the gradient variance conditions discussed in Appendix A.2. In particular, the relationship between the QES and QV against other conditions has not been considered before.

**Proposition 2.** Let g be an unbiased estimator of  $\nabla f$ . Then,

$$(\mathbf{g} \text{ is CES}) \longrightarrow (\mathbf{g} \text{ is CER})$$

*Proof.* The result immediately follows from the fact that

$$\operatorname{tr} \mathbb{V}\left[\boldsymbol{g}\left(\boldsymbol{x}\right) - \boldsymbol{g}\left(\boldsymbol{x}'\right)\right] \leq \mathbb{E}\left\|\boldsymbol{g}\left(\boldsymbol{x}\right) - \boldsymbol{g}\left(\boldsymbol{x}'\right)\right\|_{2}^{2}$$

holds for all  $x, x' \in \mathcal{X}$ .

**Proposition 3.** Let g be an unbiased estimator of  $\nabla f$ . Then,

$$\begin{array}{c}
\textbf{g is SG} \\
+ \\
f is L\text{-smooth}
\end{array}$$

*Proof.* Notice that, by definition,  $\nabla f(\mathbf{x}_*) = 0$ . Then,

$$\mathbb{E}\|\boldsymbol{g}(\boldsymbol{x})\|_{2}^{2} \leq \rho \|\nabla f(\boldsymbol{x})\|_{2}^{2}$$
$$= \rho \|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{x}_{*})\|_{2}^{2},$$

applying L-smoothness of f,

$$\leq L^2 \rho \|\boldsymbol{x} - \boldsymbol{x}_*\|_2^2$$

**Proposition 4.** Let  $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}; \boldsymbol{\eta})$  be an unbiased estimator of  $\nabla f(\mathbf{x}) = \mathbb{E} \nabla f(\mathbf{x}; \boldsymbol{\eta})$ . Then,

Integrand is uniformly L-smooth 
$$\Longrightarrow QES$$

*Proof.* The result immediately follows from the fact that the integrand  $f(x; \eta)$  is L-smooth with respect to x uniformly over  $\eta$  as

$$\mathbb{E} \left\| \boldsymbol{g} \left( \boldsymbol{x} \right) - \boldsymbol{g} \left( \boldsymbol{x}' \right) \right\|_{2}^{2} = \mathbb{E} \left\| \nabla f \left( \boldsymbol{x}; \boldsymbol{\eta} \right) - \nabla f \left( \boldsymbol{x}'; \boldsymbol{\eta} \right) \right\|_{2}^{2}$$

$$\leq L^{2} \left\| \boldsymbol{x} - \boldsymbol{x}' \right\|_{2}^{2}.$$

**Proposition 5.** Let  $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}; \boldsymbol{\eta})$  be an unbiased estimator of  $\nabla f(\mathbf{x}) = \mathbb{E} \nabla f(\mathbf{x}; \boldsymbol{\eta})$ . Then,

*Proof.* Since the integrand  $f(x; \eta)$  is both uniformly smooth and convex with respect to x for a any fixed  $\eta$ , we have

$$\begin{aligned} \left\| \nabla f\left( \boldsymbol{x}; \boldsymbol{\eta} \right) - \nabla f\left( \boldsymbol{x}'; \boldsymbol{\eta} \right) \right\|_{2} \\ &\leq 2L \left( f\left( \boldsymbol{x}; \boldsymbol{\eta} \right) - f\left( \boldsymbol{x}'; \boldsymbol{\eta} \right) - \left\langle \nabla f\left( \boldsymbol{x}'; \boldsymbol{\eta} \right), \boldsymbol{x} - \boldsymbol{x}' \right\rangle \right). \end{aligned}$$

Then,

$$\mathbb{E}\|\boldsymbol{g}(\boldsymbol{x}) - \boldsymbol{g}(\boldsymbol{x}_*)\|_2^2$$

$$= \mathbb{E}\|\nabla f(\boldsymbol{x}; \boldsymbol{\eta}) - \nabla f(\boldsymbol{x}_*; \boldsymbol{\eta})\|_2^2$$

$$\leq 2L \mathbb{E}\left(f(\boldsymbol{x}; \boldsymbol{\eta}) - f(\boldsymbol{x}_*; \boldsymbol{\eta}) - \left\langle \nabla f(\boldsymbol{x}_*; \boldsymbol{\eta}), \boldsymbol{x} - \boldsymbol{x}' \right\rangle \right)$$

$$= 2L \left(f(\boldsymbol{x}) - f(\boldsymbol{x}_*) - \left\langle \nabla f(\boldsymbol{x}_*), \boldsymbol{x} - \boldsymbol{x}' \right\rangle \right)$$

$$= 2L \left(f(\boldsymbol{x}) - f(\boldsymbol{x}_*)\right)$$

holds for all  $x \in \mathcal{X}$ .

**Proposition 6.** Let g be an unbiased estimator of  $\nabla F$ . Then,

$$g \text{ is } QV \text{ with } \beta = 0 \Longrightarrow QES$$

*Proof.* From the classic inequality  $(a + b)^2 \le 2a^2 + 2b^2$ , we have

$$\mathbb{E}\|\boldsymbol{g}(x) - \boldsymbol{g}(x_*)\|_2^2 \le 2 \mathbb{E}\|\boldsymbol{g}(x)\|_2^2 + 2 \mathbb{E}\|\boldsymbol{g}(x_*)\|_2^2.$$

Now, since QV holds with  $\beta = 0$ , we have  $\mathbb{E}||\boldsymbol{g}(\boldsymbol{x}_*)||_2^2 = 0$ . Therefore,

$$\mathbb{E}\|\boldsymbol{g}(\boldsymbol{x}) - \boldsymbol{g}(\boldsymbol{x}_*)\|_2^2 \le 2 \mathbb{E}\|\boldsymbol{g}(\boldsymbol{x})\|_2^2 \le 2 \alpha \|\boldsymbol{x} - \boldsymbol{x}_*\|_2^2,$$

where we have applied QV at the last inequality.

**Proposition 7.** Let g be an unbiased estimator of  $\nabla f$ . Then,

$$\begin{cases}
\mathbf{g} \text{ is } QV \text{ with } \beta = 0 \\
+ \\
f \text{ is } \mu\text{-}QFG
\end{cases}$$

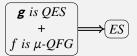
*Proof.* The result immediately follows from QV as

$$\mathbb{E}||\mathbf{g}(\mathbf{x})||_{2}^{2} \leq \alpha ||\mathbf{x} - \mathbf{x}_{*}||_{2}^{2},$$

applying  $\mu$ -quadratic functional growth,

$$\leq \frac{2\alpha}{\mu} \left( f\left( \boldsymbol{x} \right) - f\left( \boldsymbol{x}_* \right) \right).$$

**Proposition 8.** Let g be an unbiased estimator of  $\nabla f$ . Then.



Proof. From QV, we have

$$\mathbb{E}\|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}_*)\|_2^2 \le \mathcal{L}^2 \|\mathbf{x} - \mathbf{x}_*\|_2^2$$

and  $\mu$ -quadratic functional growth yields

$$\leq \frac{2\mathcal{L}^2}{\mu} \left( f\left( \boldsymbol{x} \right) - f\left( \boldsymbol{x}_* \right) \right).$$

The strategy applying QFG when proving Propositions 7 and 8 establishes the stronger variant of the ER condition: Assumption 6 (B). However, the price for this is that one has to pay for an excess  $\kappa = \mathcal{L}/\mu$  factor, and this strategy works only works for quadratically growing objectives.

## **B. Proofs**

## B.1. Auxiliary Lemmas (Lemmas 2 to 6)

**Lemma 2.** Consider a finite population of n vector-variate random variables  $x_1, ..., x_n$ . Then, the variance of the average of b samples chosen without replacement is

$$\operatorname{tr} \mathbb{V} \left[ \frac{1}{b} \sum_{i=1}^{b} \mathbf{x}_{B_i} \right] = \frac{n-b}{b(n-1)} \sigma^2,$$

where  $B = \{B_1, ..., B_b\}$  is the collection of random indices of the samples and  $\sigma^2$  is the variance of independently choosing a single sample.

*Proof.* From the variance of the sum of random variables, we have

$$\operatorname{tr} \mathbb{V} \left[ \sum_{i=1}^{b} \mathbf{x}_{\mathcal{B}_{i}} \right] = \sum_{i=1}^{b} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{\mathcal{B}_{i}} \right] + \sum_{i=1}^{b} \sum_{i \neq j}^{b} \operatorname{Cov} \left( \mathbf{x}_{\mathcal{B}_{i}}, \mathbf{x}_{\mathcal{B}_{j}} \right),$$

and noticing that the covariance is independent of the index in the batch.

$$= b \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{B_i} \right] + b(b-1)C, \tag{12}$$

where  $C = \text{Cov}(\mathbf{x}_{B_i}, \mathbf{x}_{B_j})$ . Using the fact that the variance is 0 for b = n, we can solve for C such that

$$C=-\frac{1}{n-1}\operatorname{tr}\mathbb{V}\left[\mathbf{x}_{B_i}\right],$$

which is negative, and a negative correlation is always great. Plugging this expression to Eq. (12), we have

$$\operatorname{tr} \mathbb{V}\left[\sum_{i=1}^{b} \mathbf{x}_{\mathcal{B}_{i}}\right] = b \operatorname{tr} \mathbb{V}\left[\mathbf{x}_{\mathcal{B}_{i}}\right] - b(b-1) \frac{1}{n-1} \operatorname{tr} \mathbb{V}\left[\mathbf{x}_{\mathcal{B}_{i}}\right],$$

$$= b \left(1 - \frac{b-1}{n-1}\right) \operatorname{tr} \mathbb{V}\left[\mathbf{x}_{\mathcal{B}_{i}}\right]$$

$$= b \left(\frac{n-b}{n-1}\right) \operatorname{tr} \mathbb{V}\left[\mathbf{x}_{\mathcal{B}_{i}}\right].$$

Dividing both sides by  $b^2$  yields the result.

**Lemma 3.** Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be vector-variate random variables. Then, the variance of the sum is upper-bounded as

$$\operatorname{tr} \mathbb{V}\left[\sum_{i=1}^{n} \mathbf{x}_{i}\right] \leq \left(\sum_{i=1}^{n} \sqrt{\operatorname{tr} \mathbb{V}\left[\mathbf{x}_{i}\right]}\right)^{2}$$
 (13)

$$\leq n \sum_{i=1}^{n} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right]. \tag{14}$$

The equality in Eq. (13) holds if and only if  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are constant multiples such that there exists some  $\alpha_{ij} \geq 0$  such that

$$\mathbf{x}_i = \alpha_{ij} \mathbf{x}_j$$

for all i, j.

Proof. The variance of a sum is

$$\mathrm{tr} \mathbb{V} \left[ \sum_{i=1}^{n} \boldsymbol{x}_{i} \right] = \sum_{i=1}^{n} \sum_{j=1}^{n} \mathrm{tr} \, \mathrm{Cov} \left( \boldsymbol{x}_{i}, \boldsymbol{x}_{j} \right).$$

From the Cauchy-Schwarz inequality for expectations,

$$\operatorname{tr}\operatorname{Cov}\left(\mathbf{x}_{i},\mathbf{x}_{j}\right) = \mathbb{E}(\mathbf{x}_{i} - \mathbb{E}\mathbf{x}_{i})^{\top}(\mathbf{x}_{j} - \mathbb{E}\mathbf{x}_{j})$$

$$\leq \mathbb{E}\|\mathbf{x}_{i} - \mathbb{E}\mathbf{x}_{i}\|_{2}\mathbb{E}\|\mathbf{x}_{j} - \mathbb{E}\mathbf{x}_{j}\|_{2}$$

$$= \sqrt{\operatorname{tr}\mathbb{V}\left[\mathbf{x}_{i}\right]}\sqrt{\operatorname{tr}\mathbb{V}\left[\mathbf{x}_{j}\right]}.$$

This implies

$$\operatorname{tr} \mathbb{V}\left[\sum_{i=1}^{n} \mathbf{x}_{i}\right] = \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{tr} \operatorname{Cov}\left(\mathbf{x}_{i}, \mathbf{x}_{j}\right)$$

$$\leq \sum_{i=1}^{n} \sum_{j=1}^{n} \sqrt{\operatorname{tr} \mathbb{V}\left[\mathbf{x}_{i}\right]} \sqrt{\operatorname{tr} \mathbb{V}\left[\mathbf{x}_{j}\right]}$$

$$= \left(\sum_{i=1}^{n} \sqrt{\operatorname{tr} \mathbb{V}\left[\mathbf{x}_{i}\right]}\right)^{2}.$$
(15)

The equality statement comes from the property of the Cauchy-Schwarz inequality. Lastly, Eq. (14) follows from additionally applying Jensen's inequality as

$$\left(\sum_{i=1}^{n} \sqrt{\operatorname{tr}\mathbb{V}[\mathbf{x}_{i}]}\right)^{2} = n^{2} \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\operatorname{tr}\mathbb{V}[\mathbf{x}_{i}]}\right)^{2}$$

$$\leq n^{2} \frac{1}{n} \sum_{i=1}^{n} \left(\sqrt{\operatorname{tr}\mathbb{V}[\mathbf{x}_{i}]}\right)^{2}$$

$$= n \sum_{i=1}^{n} \operatorname{tr}\mathbb{V}[\mathbf{x}_{i}].$$

An equivalent proof strategy is to expand the quadratic in Eq. (15) and apply the arithmetic mean-geometric mean inequality to the cross terms.

**Lemma 4** (Lemma A.2; Garrigos & Gower, 2023). For a recurrence relation given as

$$r_T \le \left(1 - \gamma \mu\right)^T r_0 + B \gamma,$$

for some constant  $0 < \gamma < 1/C$ ,

$$r_T \leq \epsilon$$

can be guaranteed by setting

$$\gamma = \min\left(\frac{\epsilon}{2B}, \frac{1}{C}\right) \text{ and}$$

$$T \ge \frac{1}{\mu} \max\left(2B\frac{1}{\epsilon}, C\right) \log\left(2\frac{r_0}{\epsilon}\right),$$

where  $\mu$ , B > 0 and  $0 < C < \mu$  are some finite constants.

*Proof.* First, notice that the recurrence

$$r_T \leq \underbrace{(1 - \gamma \mu)^T r_0}_{\text{bias}} + \underbrace{B \gamma}_{\text{variance}},$$

is a sum of monotonically increasing (variance) and decreasing (bias) terms with respect to  $\gamma$ . Therefore, the bound is minimized when both terms are equal. This implies that  $r_t \le \varepsilon$  can be achieved by solving for

$$(1 - \gamma \mu)^T r_0 \le \frac{\epsilon}{2}$$
 and  $B\gamma \le \frac{\epsilon}{2}$ 

First, for the variance term,

$$B\gamma \le \frac{\epsilon}{2} \quad \Leftrightarrow \quad \gamma \le \frac{\epsilon}{2B}.$$

For the bias term, as long as  $\gamma < \frac{1}{\mu}$ ,

$$(1 - \gamma \mu)^{T} r_{0} \leq \frac{\epsilon}{2}$$

$$\Leftrightarrow T \log (1 - \gamma \mu) \leq \log \frac{\epsilon}{2r_{0}}$$

$$\Leftrightarrow T \leq \frac{\log \frac{\epsilon}{2r_{0}}}{\log (1 - \gamma \mu)}$$

$$\Leftrightarrow T \geq \frac{\log \frac{2r_{0}}{\epsilon}}{\log (1 / (1 - \gamma \mu))}$$

Furthermore, using the bound  $\log 1/x \ge 1 - x$  for 0 < x < 1, we can achieve the guarantee with

$$T \ge \frac{1}{\gamma \mu} \log \left( \frac{2r_0}{\epsilon} \right).$$

Therefore,  $1/\gamma$  determines the iteration complexity. Plugging in the minimum over the constraints on  $\gamma$  yields the iteration complexity.

**Lemma 5.** For a recurrence relation given as

$$r_T \le (1 - \gamma \mu)^T r_0 + A \gamma^2 + B \gamma,$$

for some constant  $0 < \gamma < 1/C$ ,

$$r_T \leq \epsilon$$

can be guaranteed by setting

$$\begin{split} \gamma &= \min\left(\frac{-B + \sqrt{B^2 + 2A\varepsilon}}{2A}, \frac{1}{C}\right) \ and \\ T &\geq \frac{1}{\mu} \max\left(2B \frac{1}{\varepsilon} + \sqrt{2A} \frac{1}{\sqrt{\varepsilon}}, C\right) \log\left(2\frac{r_0}{\varepsilon}\right), \end{split}$$

where  $\mu$ , A, B > 0 and  $0 < C < \mu$  are some finite constants.

*Proof.* This theorem is a generalization of Lemma A.2 by Garrigos & Gower (2023). First, notice that the recurrence

$$r_T \le \underbrace{(1 - \gamma \mu)^T r_0}_{\text{bias}} + \underbrace{A \gamma^2 + B \gamma}_{\text{variance}},$$

is a sum of monotonically increasing (variance) and decreasing (bias) terms with respect to  $\gamma$ . Therefore, the bound is minimized when both terms are equal. This implies that  $r_t \le \varepsilon$  can be achieved by solving for

$$(1 - \gamma \mu)^T r_0 \le \frac{\epsilon}{2}$$
 and  $A\gamma^2 + B\gamma \le \frac{\epsilon}{2}$ 

First, for the variance term,

$$A\gamma^2 + B\gamma \le \frac{\epsilon}{2}$$
 
$$\Leftrightarrow \qquad A\gamma^2 + B\gamma - \frac{\epsilon}{2} \le 0$$

The solution to this equation is given by the positive solution of the quadratic equation as

$$0 < \gamma \le \frac{-B + \sqrt{B^2 + 2A\epsilon}}{2A}.$$

For the bias term, as long as  $\gamma < \frac{1}{\mu}$ , the solution is identical to Lemma 4. Therefore,

$$T \ge \frac{1}{\gamma \mu} \log \left( \frac{2r_0}{\epsilon} \right) \tag{16}$$

can guarantee the bias term to be smaller than  $\epsilon/2$ , while  $1/\gamma$  determines the iteration complexity. Plugging in the minimum over the constraints on  $\gamma$ ,

$$\gamma = \min\left(\frac{-B + \sqrt{B^2 + 2A\varepsilon}}{2A}, \frac{1}{C}\right) \tag{17}$$

yields the iteration complexity.

Now, since the quadratic formula is not very interpretable, let us simplify the expression for  $1/\gamma$  using the bound

$$\frac{a}{2\sqrt{b^2 + a}} \le -b + \sqrt{b^2 + a},$$

which holds for any a, b > 0 and is tight for  $\epsilon \to 0$ . With our constants, this reads

$$\frac{A\epsilon}{\sqrt{B^2 + 2A\epsilon}} \le -B + \sqrt{B^2 + 2A\epsilon},$$

and therefore

$$\frac{2A}{-B + \sqrt{B^2 + 2A\epsilon}} \le \frac{2\sqrt{B^2 + 2A\epsilon}}{\epsilon}$$
$$\le \frac{2B + \sqrt{2A\epsilon}}{\epsilon}$$
$$= 2B\frac{1}{\epsilon} + \sqrt{2A}\frac{1}{\sqrt{\epsilon}}.$$

Therefore, for the stepsize choice of Eq. (17),

$$\frac{1}{\gamma} \le \min\left(2B\frac{1}{\epsilon} + \sqrt{2A}\frac{1}{\sqrt{\epsilon}}, \frac{1}{C}\right).$$

Plugging this into Eq. (16) yields the statement.

**Lemma 6.** Let  $F: \mathcal{X} \to \mathbb{R}$  be a finite sum of convex functions as  $F = \frac{1}{n}(f_1 + ... + f_n)$ , where  $f_i: \mathcal{X} \to \mathbb{R}$ . Then,

$$\frac{1}{n}\sum_{i=1}^{n}D_{f_{i}}\left(\boldsymbol{x},\boldsymbol{x}'\right)=D_{F}\left(\boldsymbol{x},\boldsymbol{x}'\right),$$

for any  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ .

*Proof.* The result immediately follows from the definition of Bregman divergences as

$$\frac{1}{n} \sum_{i=1}^{n} D_{f_{i}}(\mathbf{x}, \mathbf{x}')$$

$$= \frac{1}{n} \sum_{i=1}^{n} (f_{i}(\mathbf{x}) - f_{i}(\mathbf{x}') - \langle \nabla f_{i}(\mathbf{x}'), \mathbf{x} - \mathbf{x}' \rangle)$$

$$= \left(\frac{1}{n} \sum_{i=1}^{n} f_{i}(\mathbf{x})\right) - \left(\frac{1}{n} \sum_{i=1}^{n} f_{i}(\mathbf{x}')\right)$$

$$- \left\langle\frac{1}{n} \sum_{i=1}^{n} \nabla f_{i}(\mathbf{x}'), \mathbf{x} - \mathbf{x}'\right\rangle$$

$$= F(\mathbf{x}) - F(\mathbf{x}') - \langle \nabla F(\mathbf{x}'), \mathbf{x} - \mathbf{x}' \rangle$$

$$= D_{F}(\mathbf{x}, \mathbf{x}').$$

## B.2. Convergence of SGD (Lemmas 1, 7 and 8)

**Lemma 7.** Let  $F: \mathcal{X} \to \mathbb{R}$  be L-smooth function. Then, the expected squared norm of a gradient estimator  $\mathbf{g}$  satisfying both  $\mathrm{ER}(\mathcal{L})$  and  $\mathrm{BV}(\sigma^2)$  is bounded as

$$\mathbb{E}\|\mathbf{g}(\mathbf{x})\|_{2}^{2} \le 4(\mathcal{L} + L)(F(\mathbf{x}) - F(\mathbf{x}_{*})) + 2\sigma^{2},$$

for any  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{x}_* \in \arg\max_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$ .

*Proof.* The proof is a minor modification of Lemma 2.4 by Gower et al. (2019) and Lemma 3.2 by Gower et al. (2021a).

By applying the bound  $(a + b)^2 \le 2a^2 + 2b^2$ , we can "transfer" the variance on x to the variance of  $x_*$ . That is

$$\mathbb{E}\|\mathbf{g}(\mathbf{x})\|_{2}^{2} = \mathbb{E}\|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}_{*}) + \mathbf{g}(\mathbf{x}_{*})\|_{2}^{2}$$

$$\leq 2 \underbrace{\mathbb{E}\|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}_{*})\|_{2}^{2}}_{V_{1}} + 2 \underbrace{\mathbb{E}\|\mathbf{g}(\mathbf{x}_{*})\|_{2}^{2}}_{V_{2}}$$

The key is to bound  $V_1$ . It is typical to do this using expected-smoothness-type assumptions such as the ER assumption. That is,

$$V_{1} = \mathbb{E}\|\boldsymbol{g}(\boldsymbol{x}) - \boldsymbol{g}(\boldsymbol{x}_{*})\|_{2}^{2}$$
  
= tr\mathbb{V}[\boldsymbol{g}(\boldsymbol{x}) - \boldsymbol{g}(\boldsymbol{x}\_{\*})] + (\nabla F(\boldsymbol{x}) - \nabla F(\boldsymbol{x}),

from the L-smoothness of F,

$$\leq \operatorname{tr} \mathbb{V}\left[\boldsymbol{g}\left(\boldsymbol{x}\right) - \boldsymbol{g}\left(\boldsymbol{x}_{*}\right)\right] + 2L\left(F\left(\boldsymbol{x}\right) - F\left(\boldsymbol{x}_{*}\right)\right),$$

and the ER condition,

$$\leq 2\mathcal{L}\left(F\left(\boldsymbol{x}\right) - F\left(\boldsymbol{x}_{*}\right)\right) + 2L\left(F\left(\boldsymbol{x}\right) - F\left(\boldsymbol{x}_{*}\right)\right)$$
$$= 2\left(L + \mathcal{L}\right)\left(F\left(\boldsymbol{x}\right) - F\left(\boldsymbol{x}_{*}\right)\right).$$

Finally,  $V_2$  immediately follows from the BV condition as

$$V_2 = \mathbb{E}||\boldsymbol{g}(\boldsymbol{x}_*)||_2^2 \leq \sigma^2.$$

**Lemma 8.** Let the objective function F satisfy Assumption 1 and the gradient estimator  $\mathbf{g}$  be unbiased and satisfy both  $ER(\mathcal{L})$  and  $BV(\sigma^2)$ . Then, the last iterate of SGD guarantees

$$\mathbb{E}\left[\left\|\boldsymbol{x}_{T}-\boldsymbol{x}^{*}\right\|_{2}^{2}\right] \leq \left(1-\mu\gamma\right)^{T}\left\|\boldsymbol{x}_{0}-\boldsymbol{x}_{*}\right\|_{2}^{2}+\frac{2\sigma^{2}}{\mu}\gamma$$

where  $\mathbf{x}_* = \arg\min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$  is the global optimum.

*Proof.* Firstly, we have

$$\|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2$$

$$= \|\Pi_{\mathcal{X}} (\mathbf{x}_t - \gamma \mathbf{g}(\mathbf{x}_t)) - \Pi(\mathbf{x}_*)\|_2^2,$$

and since the projection onto a convex set under a Euclidean metric is non-expansive,

$$\leq \|\mathbf{x}_{t} - \gamma \mathbf{g}(\mathbf{x}_{t}) - \mathbf{x}_{*}\|_{2}^{2}$$

$$= \|\mathbf{x}_{t} - \mathbf{x}_{*}\|_{2}^{2} - 2\gamma \langle \mathbf{g}(\mathbf{x}_{t}), \mathbf{x}_{t} - \mathbf{x}_{*} \rangle + \gamma^{2} \|\mathbf{g}(\mathbf{x}_{t})\|_{2}^{2}.$$

Denoting the  $\sigma$ -algebra formed by the randomness and the iterates up to the tth iteration as  $\mathcal{F}_t$  such that  $(\mathcal{F}_t)_{t\geq 1}$  forms a filtration, the conditional expectation is

$$\mathbb{E}\left[\left\|\boldsymbol{x}_{t+1} - \boldsymbol{x}_{*}\right\|_{2}^{2} \mid \mathcal{F}_{t}\right]$$

$$= \left\|\boldsymbol{x}_{t} - \boldsymbol{x}_{*}\right\|_{2}^{2} - 2\gamma \left\langle \mathbb{E}\left[\boldsymbol{g}\left(\boldsymbol{x}_{t}\right) \mid \mathcal{F}_{t}\right], \boldsymbol{x}_{t} - \boldsymbol{x}_{*}\right\rangle$$

$$+ \gamma^{2} \mathbb{E}\left[\left\|\boldsymbol{g}\left(\boldsymbol{x}_{t}\right)\right\|_{2}^{2} \mid \mathcal{F}_{t}\right].$$

$$= \left\|\boldsymbol{x}_{t} - \boldsymbol{x}_{*}\right\|_{2}^{2} - 2\gamma \left\langle \nabla F\left(\boldsymbol{x}_{t}\right), \boldsymbol{x}_{t} - \boldsymbol{x}_{*}\right\rangle$$

$$+ \gamma^{2} \mathbb{E}\left[\left\|\boldsymbol{g}\left(\boldsymbol{x}_{t}\right)\right\|_{2}^{2} \mid \mathcal{F}_{t}\right],$$

applying the  $\mu$ -strong convexity of F,

$$\leq \|\boldsymbol{x}_{t} - \boldsymbol{x}_{*}\|_{2}^{2} - 2\gamma \left(F(\boldsymbol{x}_{t}) - F(\boldsymbol{x}_{*}) + \frac{\mu}{2} \|\boldsymbol{x}_{t} - \boldsymbol{x}_{*}\|_{2}^{2}\right)$$

$$+ \gamma^{2} \mathbb{E} \left[\|\boldsymbol{g}(\boldsymbol{x}_{t})\|_{2}^{2} \mid \mathcal{F}_{t}\right]$$

$$= (1 - \gamma\mu) \|\boldsymbol{x}_{t} - \boldsymbol{x}_{*}\|_{2}^{2} - 2\gamma \left(F(\boldsymbol{x}_{t}) - F(\boldsymbol{x}_{*})\right)$$

$$+ \gamma^{2} \mathbb{E} \left[\|\boldsymbol{g}(\boldsymbol{x}_{t})\|_{2}^{2} \mid \mathcal{F}_{t}\right]$$

From Lemma 7, we have

$$\mathbb{E}\left[\left\|\boldsymbol{g}\left(\boldsymbol{x}_{t}\right)\right\|_{2}^{2}\mid\mathcal{F}_{t}\right]\leq\left(4\left(\mathcal{L}+L\right)\left(F\left(\boldsymbol{x}_{t}\right)-F\left(\boldsymbol{x}_{*}\right)\right)+2\sigma^{2}\right).$$

Therefore,

$$\mathbb{E}\left[\|\mathbf{x}_{t+1} - \mathbf{x}_*\|_2^2 \mid \mathcal{F}_t\right]$$

$$\leq (1 - \gamma \mu) \|\mathbf{x}_t - \mathbf{x}_*\|_2^2 - 2\gamma (F(\mathbf{x}_t) - F(\mathbf{x}_*))$$

$$+ \gamma^2 \left(4(\mathcal{L} + L)(F(\mathbf{x}_t) - F(\mathbf{x}_*)) + 2\sigma^2\right)$$

$$= (1 - \gamma \mu) \|\mathbf{x}_t - \mathbf{x}_*\|_2^2$$

$$- 2\gamma (1 - 2\gamma (\mathcal{L} + L))(F(\mathbf{x}_t) - F(\mathbf{x}_*)) + 2\gamma^2 \sigma^2,$$

and with a small-enough stepsize satisfying  $\gamma < \frac{1}{2(\mathcal{L}+L)}$ , we can guarantee a partial contraction as

$$\leq (1 - \gamma \mu) \|\mathbf{x}_t - \mathbf{x}_*\|_2^2 + 2\gamma^2 \sigma^2.$$

Note that the coefficient  $1 - \gamma \mu$  is guaranteed to be strictly smaller than 1 since  $\mu \le L$ , which means that we indeed have a partial contraction.

Now, taking full expectation, we have

$$\mathbb{E}||\boldsymbol{x}_{t+1} - \boldsymbol{x}_*||_2^2 \le (1 - \gamma \mu) \mathbb{E}||\boldsymbol{x}_t - \boldsymbol{x}_*||_2^2 + 2\gamma^2 \sigma^2.$$

Unrolling the recursion from 0 to T-1, we have

$$\begin{split} \mathbb{E} \|\boldsymbol{x}_{T} - \boldsymbol{x}_{*}\|_{2}^{2} &\leq (1 - \gamma \mu)^{T} \mathbb{E} \|\boldsymbol{x}_{0} - \boldsymbol{x}_{*}\|_{2}^{2} \\ &+ 2\gamma^{2} \sigma^{2} \sum_{t=0}^{T-1} (1 - \gamma \mu)^{t}. \\ &\leq (1 - \gamma \mu)^{T} \mathbb{E} \|\boldsymbol{x}_{0} - \boldsymbol{x}_{*}\|_{2}^{2} + \frac{2\sigma^{2}}{\mu} \gamma. \end{split}$$

where the last inequality follows from the asymptotic bound on geometric sums.  $\Box$ 

**Lemma 1.** Let the objective F satisfy Assumption 1 and the gradient estimator  $\mathbf{g}$  satisfy  $\mathrm{ER}\left(\mathcal{L}\right)$  and  $\mathrm{BV}\left(\sigma^{2}\right)$ . Then, the last iterate of SGD is  $\epsilon$ -close to the global optimum  $\mathbf{x}_{*} = \arg\min_{\mathbf{x} \in \mathcal{X}} F\left(\mathbf{x}\right)$  such that  $\mathbb{E}\|\mathbf{x}_{T} - \mathbf{x}_{*}\|_{2}^{2} \leq \epsilon$  after a number of iterations at least

$$T \geq 2 \max \left( \frac{\sigma^2}{\mu^2} \frac{1}{\epsilon}, \frac{\mathcal{L} + L}{\mu} \right) \log \left( 2 || \boldsymbol{x}_0 - \boldsymbol{x}_* ||_2^2 \frac{1}{\epsilon} \right)$$

and the fixed stepsize

$$\gamma = \min\left(\frac{\epsilon\mu}{2\sigma^2}, \frac{1}{2(\mathcal{L} + L)}\right).$$

*Proof.* We can apply Lemma 4 to the result of Lemma 8 with the constants

$$r_0 = \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2$$
,  $B = \frac{2\sigma^2}{\mu}$ , and  $C = 2(\mathcal{L} + L)$ .

Then, we can guarantee an  $\epsilon$ -accurate solution with the stepsize

$$\gamma = \min\left(\frac{\epsilon\mu}{2\sigma^2}, \frac{1}{2(\mathcal{L} + L)}\right)$$

and a number of iterations of at least

$$T \ge \frac{1}{\mu} \max \left( \frac{2\sigma^2}{\mu}, 2\left(\mathcal{L} + L\right) \right) \log \left( 2\|\boldsymbol{x}_0 - \boldsymbol{x}_*\|_2^2 \frac{1}{\epsilon} \right)$$
$$= 2 \max \left( \frac{\sigma^2}{\mu^2}, \frac{\mathcal{L} + L}{\mu} \right) \log \left( 2\|\boldsymbol{x}_0 - \boldsymbol{x}_*\|_2^2 \frac{1}{\epsilon} \right).$$

#### **B.3.** General Variance Bound (Theorem 1)

**Lemma 9.** Let  $\mathbf{x}_1, ..., \mathbf{x}_b$  be a collection of vector-variate RVs dependent on some random variable B satisfying Assumption 3. Then, the expected variance of the sum of  $\mathbf{x}_1, ..., \mathbf{x}_b$  conditioned on B is bounded as

$$\mathbb{E}\left[\mathrm{tr}\mathbb{V}\left[\sum_{i=1}^{b}\boldsymbol{x}_{i}\mid\boldsymbol{B}\right]\right]\leq\rho\mathbb{V}\left[\boldsymbol{S}\right]+\rho(\mathbb{E}\boldsymbol{S})^{2}+\left(1-\rho\right)\mathbb{E}\left[\boldsymbol{V}\right],$$

where

$$S = \sum_{i=1}^{b} \sqrt{\operatorname{tr} \mathbb{V}\left[\mathbf{x}_{i} \mid B\right]} \ \ and \ \ V = \sum_{i=1}^{b} \operatorname{tr} \mathbb{V}\left[\mathbf{x}_{i} \mid B\right].$$

Equality holds when the equality in Assumption 3 holds.

*Proof.* From the formula for the variance of sums,

$$\operatorname{tr} \mathbb{V} \left[ \sum_{i=1}^{b} \mathbf{x}_{i} \mid B \right]$$

$$= \sum_{i=1}^{b} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \mid B \right] + \sum_{i=1}^{b} \sum_{j \neq i} \operatorname{tr} \operatorname{Cov} \left( \mathbf{x}_{i}, \mathbf{x}_{j} \mid B \right).$$

$$\leq \sum_{i=1}^{b} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \mid B \right] + \sum_{i=1}^{b} \sum_{j \neq i} \rho \sqrt{\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \mid B \right]} \sqrt{\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{j} \mid B \right]}$$

$$= (1 - \rho) \sum_{i=1}^{b} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \mid B \right]$$

$$+ \rho \sum_{i=1}^{b} \sum_{j=1}^{b} \sqrt{\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \mid B \right]} \sqrt{\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{j} \mid B \right]}$$

$$= (1 - \rho) \sum_{i=1}^{b} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \mid B \right] + \rho \left( \sum_{i=1}^{b} \sqrt{\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \mid B \right]} \right)^{2}$$

$$= (1 - \rho) V + \rho S^{2}.$$

Then, it follows that

$$\mathbb{E}\left[\operatorname{tr}\mathbb{V}\left[\sum_{i=1}^{b} \mathbf{x}_{i} \middle| B\right]\right] \leq \rho \mathbb{E}\left[S^{2}\right] + (1-\rho)\mathbb{E}\left[V\right]$$
$$= \rho \mathbb{V}\left[S\right] + \rho(\mathbb{E}S)^{2} + (1-\rho)\mathbb{E}\left[V\right],$$

from the basic property of the variance:

$$\mathbb{V}\left[S\right] = \mathbb{E}\left[S^2\right] - \left(\mathbb{E}S\right)^2.$$

Since Assumption 3 is the only inequality we use, the equality in the statement holds whenever the equality in Assumption 3 holds.

**Theorem 1.** Let the component estimators  $\mathbf{x}_1, ..., \mathbf{x}_n$  satisfy Assumption 3. Then, the variance of the doubly stochastic estimator  $\mathbf{x}_B$  is bounded as

$$\operatorname{tr} \mathbb{V}[\mathbf{x}_B] \leq V_{\operatorname{com}} + V_{\operatorname{cor}} + V_{\operatorname{sub}},$$

wher

$$\begin{split} V_{\text{com}} &= \left(\frac{\rho}{b_{\text{eff}}} + \frac{1-\rho}{b}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \text{tr} \mathbb{V}\left[\boldsymbol{x}_{i}\right]\right), \\ V_{\text{cor}} &= \rho \left(1 - \frac{1}{b_{\text{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\text{tr} \mathbb{V}\left[\boldsymbol{x}_{i}\right]}\right)^{2}, \text{ and} \\ V_{\text{sub}} &= \frac{1}{b_{\text{eff}}} \frac{1}{n} \sum_{i=1}^{n} \left\|\bar{\boldsymbol{x}}_{i} - \bar{\boldsymbol{x}}\right\|_{2}^{2}. \end{split}$$

Equality holds when the equality in Assumption 3 holds.

*Proof.* Starting from the law of total covariance, we have

$$\mathbb{V}\left[\mathbf{x}_{B}\right] = \underbrace{\mathbb{E}_{B \sim \pi}\left[\operatorname{tr}\mathbb{V}\left[\mathbf{x}_{B} \mid B\right]\right]}_{\text{Ensemble Variance}} + \underbrace{\operatorname{tr}\mathbb{V}_{B \sim \pi}\left[\mathbb{E}\left[\mathbf{x}_{B} \mid B\right]\right]}_{\text{Subsampling Variance}}.$$
 (18)

**Ensemble Variance** Bounding the variance of each ensemble is key. From Lemma 9, we have

$$\mathbb{E}\left[\operatorname{tr}\mathbb{V}\left[\mathbf{x}_{B}\mid B\right]\right] = \mathbb{E}\left[\operatorname{tr}\mathbb{V}\left[\frac{1}{b}\sum_{i\in B}\mathbf{x}_{i}\left|B\right|\right]\right]$$

$$= \mathbb{E}\left[\operatorname{tr}\mathbb{V}\left[\sum_{i\in B}\left(\frac{1}{b}\mathbf{x}_{i}\right)\left|B\right|\right]\right]$$

$$<\rho\mathbb{V}S + \rho(\mathbb{E}S)^{2} + (1-\rho)\mathbb{E}V, \quad (19)$$

where

$$S \triangleq \sum_{i \in B} \sqrt{\operatorname{tr} \mathbb{V} \left[ \frac{1}{b} \boldsymbol{x}_{i} \right]} = \frac{1}{b} \sum_{i \in B} \sqrt{\operatorname{tr} \mathbb{V} \left[ \boldsymbol{x}_{i} \right]},$$
$$V \triangleq \sum_{i \in B} \operatorname{tr} \mathbb{V} \left[ \frac{1}{b} \boldsymbol{x}_{i} \right] = \frac{1}{b^{2}} \sum_{i \in B} \operatorname{tr} \mathbb{V} \left[ \boldsymbol{x}_{i} \right].$$

In our context, S is the batch average of the standard deviations, and V is the batch average of the variance (scaled with a factor of 1/b).

Notice that *S* is an *b*-sample average of the standard deviations. Therefore, if  $\pi$  is an unbiased subsampling strategy, we retrieve the population average standard deviation as

$$\mathbb{E}_{B \sim \pi} \left[ S \right] = \mathbb{E}_{B \sim \pi} \left[ \frac{1}{b} \sum_{i \in B} \sqrt{\text{tr} \mathbb{V} \left[ \mathbf{x}_i \right]} \right] = \frac{1}{n} \sum_{i=1}^n \sqrt{\text{tr} \mathbb{V} \left[ \mathbf{x}_i \right]}. \tag{20}$$

Under a similar reasoning, the variance of the standard deviations follows as

$$\begin{aligned} \mathbb{V}_{B \sim \pi} \left[ S \right] \\ &= \mathbb{V}_{B \sim \pi} \left[ \frac{1}{b} \sum_{i \in B} \sqrt{\text{tr} \mathbb{V} \left[ \mathbf{x}_i \right]} \right] \end{aligned}$$

$$= \frac{1}{b_{\text{eff}}} \mathbb{V}_{i \sim \text{Uniform}\{1,\dots,n\}} \left[ \sqrt{\text{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right]} \right]$$

$$= \frac{1}{b_{\text{eff}}} \left( \frac{1}{n} \sum_{i=1}^{n} \text{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right] - \left( \frac{1}{n} \sum_{i=1}^{n} \sqrt{\text{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right]} \right)^{2} \right), \quad (21)$$

where the last identity is the well-known formula for the variance:  $\mathbb{V}X = \mathbb{E}X^2 - (\mathbb{E}X)^2$ . Likewise, the average variance follows as

$$\mathbb{E}_{B \sim \pi} V = \frac{1}{b^2} \mathbb{E}_{B \sim \pi} \left[ \sum_{i \in B} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_i \right] \right]$$

$$= \frac{1}{b} \mathbb{E}_{B \sim \pi} \left[ \frac{1}{b} \sum_{i \in B} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_i \right] \right]$$

$$= \frac{1}{b} \left( \frac{1}{n} \sum_{i=1}^{n} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_i \right] \right)$$
(22)

Plugging Eqs. (20) to (22) into Eq. (19), we have

$$\mathbb{E}_{B \sim \pi} \left[ \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{B} \mid B \right] \right]$$

$$\leq \rho \mathbb{V} S + \rho (\mathbb{E} S)^{2} + (1 - \rho) \mathbb{E} V$$

$$= \frac{\rho}{b_{\text{eff}}} \left( \frac{1}{n} \sum_{i=1}^{n} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right] - \left( \frac{1}{n} \sum_{i=1}^{n} \sqrt{\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right]} \right)^{2} \right)$$

$$+ \rho \left( \frac{1}{n} \sum_{i=1}^{n} \sqrt{\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right]} \right)^{2}$$

$$+ \frac{1 - \rho}{b} \left( \frac{1}{n} \sum_{i=1}^{n} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right] \right)$$

$$= \left( \frac{\rho}{b_{\text{eff}}} + \frac{1 - \rho}{b} \right) \left( \frac{1}{n} \sum_{i=1}^{n} \operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right] \right)$$

$$+ \rho \left( 1 - \frac{1}{b_{\text{eff}}} \right) \left( \frac{1}{n} \sum_{i=1}^{n} \sqrt{\operatorname{tr} \mathbb{V} \left[ \mathbf{x}_{i} \right]} \right)^{2} .$$

$$(23)$$

**Subsampling Variance** The subsampling noise is straightforward. For this, we will denote the minibatch subsampling estimator of the component means as

$$\bar{\mathbf{x}}_B \triangleq \frac{1}{b} \sum_{i \in B} \bar{\mathbf{x}}_i.$$

Since each component estimator  $x_i$  is unbiased, the expectation conditional on the minibatch B is

$$\mathbb{E}\left[\bar{\boldsymbol{x}}_{B}\mid B\right] = \frac{1}{b}\sum_{i\in B}\bar{\boldsymbol{x}}_{i}.$$

Therefore,

$$\operatorname{tr} \mathbb{V}_{B \sim \pi} \left[ \mathbb{E} \left[ \mathbf{x}_B \mid B \right] \right] = \operatorname{tr} \mathbb{V}_{B \sim \pi} \left[ \bar{\mathbf{x}}_B \right]$$

$$= \frac{1}{b_{\text{eff}}} \left( \frac{1}{n} \sum_{i=1}^{n} ||\bar{x}_i - \bar{x}||_2^2 \right). \tag{24}$$

Combining Eqs. (23) and (24) into Eq. (18) yields the result. Notice that the only inequality we used is Eq. (19), Lemma 9, in which equality holds if the equality in Assumption 3 holds.

## **B.4. Doubly Stochastic Gradients**

# B.4.1. EXPECTED RESIDUAL CONDITION (THEOREM 2)

**Theorem 2.** Let Assumption 4 to 6 hold. Then, we have:

- (i) If  $(A^{CVX})$  or  $(A^{ITP})$  hold,  $\mathbf{g}_B$  satisfies  $ER(\mathcal{L}_A)$ .
- (ii) If (B) holds,  $\mathbf{g}_{B}$  satisfies ER ( $\mathcal{L}_{B}$ ).

where 
$$\mathcal{L}_{\max} = \max \{\mathcal{L}_1, \dots, \mathcal{L}_n\}$$

$$\begin{split} \mathcal{L}_{\mathrm{A}} &= \left(\frac{\rho}{b_{\mathrm{eff}}} + \frac{1-\rho}{b}\right) \mathcal{L}_{\mathrm{max}} + \rho \left(1 - \frac{1}{b_{\mathrm{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i}\right) + \frac{\mathcal{L}_{\mathrm{sub}}}{b_{\mathrm{eff}}} \\ \mathcal{L}_{\mathrm{B}} &= \left(\frac{\rho}{b_{\mathrm{eff}}} + \frac{1-\rho}{b}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i}\right) \\ &+ \rho \left(1 - \frac{1}{b_{\mathrm{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\mathcal{L}_{i}}\right)^{2} + \frac{\mathcal{L}_{\mathrm{sub}}}{b_{\mathrm{eff}}}. \end{split}$$

*Proof.* From Theorem 1, we have

$$\operatorname{tr} \mathbb{V}\left[\mathbf{g}_{\mathcal{B}}\left(\mathbf{x}\right) - \mathbf{g}_{\mathcal{B}}\left(\mathbf{x}_{*}\right)\right]$$

$$\leq \left(\frac{\rho}{b_{\mathrm{eff}}} + \frac{1 - \rho}{b}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \operatorname{tr} \mathbb{V}\left[\mathbf{g}_{i}\left(\mathbf{x}\right) - \mathbf{g}_{i}\left(\mathbf{x}_{*}\right)\right]\right)$$

$$+ \rho \left(1 - \frac{1}{b_{\mathrm{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\operatorname{tr} \mathbb{V}\left[\mathbf{g}_{i}\left(\mathbf{x}\right) - \mathbf{g}_{i}\left(\mathbf{x}_{*}\right)\right]}\right)^{2}$$

+  $\frac{1}{b_{\text{eff}}} \text{tr} \mathbb{V} \left[ \nabla f_B(\mathbf{x}) - \nabla F(\mathbf{x}) \right],$ where Assumption 5 yields

$$\leq \left(\frac{\rho}{b_{\text{eff}}} + \frac{1-\rho}{b}\right) \underbrace{\left(\frac{1}{n} \sum_{i=1}^{n} \text{tr} \mathbb{V}\left[\mathbf{g}_{i}\left(\mathbf{x}\right) - \mathbf{g}_{i}\left(\mathbf{x}_{*}\right)\right]\right)}_{\triangleq \mathbf{T}_{corr}}$$

$$+\rho\left(1-\frac{1}{b_{\text{eff}}}\right)\left(\underbrace{\frac{1}{n}\sum_{i=1}^{n}\sqrt{\text{tr}\mathbb{V}\left[\boldsymbol{g}_{i}\left(\boldsymbol{x}\right)-\boldsymbol{g}_{i}\left(\boldsymbol{x}_{*}\right)\right]}}_{\triangleq T...}\right)^{2}$$

$$+ \frac{2\mathcal{L}_{\text{sub}}}{b_{\text{eff}}} (F(\mathbf{x}) - F(\mathbf{x}_*))$$

$$= \left(\frac{\rho}{b_{\text{eff}}} + \frac{1 - \rho}{b}\right) T_{\text{var}} + \rho \left(1 - \frac{1}{b_{\text{eff}}}\right) T_{\text{cov}}$$

$$+ \frac{2\mathcal{L}_{\text{sub}}}{b_{\text{off}}} (F(\mathbf{x}) - F(\mathbf{x}_*)). \tag{25}$$

**Proof of (i) with (A**<sup>CVX</sup>) Since Assumption 6 (A<sup>CVX</sup>) requires  $f_1, ..., f_n$  to be convex, F is also convex. Therefore, we can use the identity in Lemma 6 and

$$D_F(\mathbf{x}, \mathbf{x}_*) = F(\mathbf{x}) - F(\mathbf{x}_*).$$

With that said, under (ACVX), we have

$$T_{\text{var}} \leq \frac{1}{n} \sum_{i=1}^{n} \operatorname{tr} \mathbb{V} \left[ \boldsymbol{g}_{i} \left( \boldsymbol{x} \right) - \boldsymbol{g}_{i} \left( \boldsymbol{x}_{*} \right) \right]$$

$$=\frac{1}{n}\sum_{i=1}^{n}2\mathcal{L}_{i}D_{f_{i}}\left(\boldsymbol{x},\boldsymbol{x}_{*}\right),$$

applying  $\mathcal{L}_{\max} \geq \mathcal{L}_i$  for all i = 1, ..., n,

$$\leq 2\mathcal{L}_{\max} \frac{1}{n} \sum_{i=1}^{n} D_{f_i}(\boldsymbol{x}, \boldsymbol{x}_*)$$

and Lemma 6,

$$=2\mathcal{L}_{\max}D_F(\boldsymbol{x},\boldsymbol{x}_*). \tag{26}$$

For  $T_{cov}$ , since

$$T_{\text{cov}} = \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\text{tr} \mathbb{V} \left[ \mathbf{g}_{i} \left( \mathbf{x} \right) - \mathbf{g}_{i} \left( \mathbf{x}_{*} \right) \right]} \right)^{2}$$

is monotonic w.r.t. the variance, we can apply (A<sup>CVX</sup>) as

$$\leq \frac{2}{n^2} \left( \sum_{i=1}^n \sqrt{\mathcal{L}_i D_{f_i}(\boldsymbol{x}, \boldsymbol{x}_*)} \right)^2.$$

Now, applying the Cauchy-Schwarz inequality yields

$$\leq \frac{2}{n^2} \left( \sum_{i=1}^n \mathcal{L}_i \right) \left( \sum_{i=1}^n \mathrm{D}_{f_i} \left( \boldsymbol{x}, \boldsymbol{x}_* \right) \right)$$
$$= 2 \left( \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i \right) \left( \frac{1}{n} \sum_{i=1}^n \mathrm{D}_{f_i} \left( \boldsymbol{x}, \boldsymbol{x}_* \right) \right)$$

and by Lemma 6.

$$= 2\left(\frac{1}{n}\sum_{i=1}^{n}\mathcal{L}_{i}\right)D_{F}\left(\boldsymbol{x},\boldsymbol{x}_{*}\right). \tag{27}$$

Plugging Eqs. (26) and (27) into Eq. (25), we have

$$\begin{split} \operatorname{tr} \mathbb{V} \left[ \boldsymbol{g} \left( \boldsymbol{x} \right) - \boldsymbol{g} \left( \boldsymbol{x}_{*} \right) \right] \\ & \leq \left( \frac{\rho}{b_{\text{eff}}} + \frac{1 - \rho}{b} \right) T_{\text{var}} + \rho \left( 1 - \frac{1}{b_{\text{eff}}} \right) T_{\text{cov}} \\ & + \frac{2\mathcal{L}_{\text{sub}}}{b_{\text{eff}}} \left( F \left( \boldsymbol{x} \right) - F \left( \boldsymbol{x}_{*} \right) \right) \\ & \leq \left( \frac{\rho}{b_{\text{eff}}} + \frac{1 - \rho}{b} \right) 2\mathcal{L}_{\text{max}} D_{F} \left( \boldsymbol{x}, \boldsymbol{x}_{*} \right) \\ & + \rho \left( 1 - \frac{1}{b_{\text{eff}}} \right) 2 \left( \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i} \right) D_{F} \left( \boldsymbol{x}, \boldsymbol{x}_{*} \right) \\ & + \frac{1}{b_{\text{eff}}} 2\mathcal{L}_{\text{sub}} D_{F} \left( \boldsymbol{x}, \boldsymbol{x}_{*} \right). \\ & = 2 \left( \left( \frac{\rho}{b_{\text{eff}}} + \frac{1 - \rho}{b} \right) \mathcal{L}_{\text{max}} + \rho \left( 1 - \frac{1}{b_{\text{eff}}} \right) \left( \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i} \right) \\ & + \frac{1}{b_{\text{eff}}} \mathcal{L}_{\text{sub}} \right) D_{F} \left( \boldsymbol{x}, \boldsymbol{x}_{*} \right) \\ & = 2 \left( \left( \frac{\rho}{b_{\text{eff}}} + \frac{1 - \rho}{b} \right) \mathcal{L}_{\text{max}} + \rho \left( 1 - \frac{1}{b_{\text{eff}}} \right) \left( \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i} \right) \right) \end{split}$$

$$+\frac{1}{b_{\mathrm{eff}}}\mathcal{L}_{\mathrm{sub}}\left(F\left(\boldsymbol{x}\right)-F\left(\boldsymbol{x}_{*}\right)\right).$$

**Proof of (i) with (A**<sup>ITP</sup>) From Assumption 6 (A<sup>ITP</sup>), we have

$$T_{\text{var}} = \frac{1}{n} \sum_{i=1}^{n} \text{tr} \mathbb{V} \left[ \mathbf{g}_{i} \left( \mathbf{x} \right) - \mathbf{g}_{i} \left( \mathbf{x}_{*} \right) \right]$$

$$\leq \frac{1}{n} \sum_{i=1}^{n} 2\mathcal{L}_{i} \left( f_{i} \left( \mathbf{x} \right) - f_{i} \left( \mathbf{x}_{*} \right) \right),$$
applying  $\mathcal{L}_{\text{max}} \geq \mathcal{L}_{i}$  for all  $i = 1, ..., n$ ,

$$\underset{\text{definition}}{\text{definition}} \mathcal{L}_{\text{max}} \geq \mathcal{L}_{i} \text{ for all } t = 1, \dots, n, \\
\leq 2\mathcal{L}_{\text{max}} \frac{1}{n} \sum_{i=1}^{n} (f_{i}(\mathbf{x}) - f_{i}(\mathbf{x}_{*})) \\
= 2\mathcal{L}_{\text{max}} (F(\mathbf{x}) - F(\mathbf{x}_{*})). \tag{28}$$

Similarly,

$$T_{\text{cov}} = \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\text{tr} \mathbb{V} \left[ \mathbf{g}_{i} \left( \mathbf{x} \right) - \mathbf{g}_{i} \left( \mathbf{x}_{*} \right) \right]} \right)^{2},$$

applying (AITP)

$$\leq \frac{2}{n^2} \left( \sum_{i=1}^n \sqrt{\mathcal{L}_i \left( f_i \left( \boldsymbol{x} \right) - f_i \left( \boldsymbol{x}_* \right) \right)} \right)^2,$$

and applying the Cauchy-Schwarz inequality

$$\leq \frac{2}{n^2} \left( \sum_{i=1}^n \mathcal{L}_i \right) \left( \sum_{i=1}^n f_i(\mathbf{x}) - f_i(\mathbf{x}_*) \right) 
= 2 \left( \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i \right) \left( \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) - f_i(\mathbf{x}_*) \right) 
= 2 \left( \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i \right) (F(\mathbf{x}) - F(\mathbf{x}_*)).$$
(29)

Plugging Eqs. (28) and (29) into Eq. (25), we have

$$\begin{split} \operatorname{tr} \mathbb{V} \left[ \boldsymbol{g} \left( \boldsymbol{x} \right) - \boldsymbol{g} \left( \boldsymbol{x}_{*} \right) \right] \\ & \leq \left( \frac{\rho}{b_{\mathrm{eff}}} + \frac{1 - \rho}{b} \right) T_{\mathrm{var}} + \rho \left( 1 - \frac{1}{b_{\mathrm{eff}}} \right) T_{\mathrm{cov}} \\ & + \frac{2 \mathcal{L}_{\mathrm{sub}}}{b_{\mathrm{eff}}} \left( F \left( \boldsymbol{x} \right) - F \left( \boldsymbol{x}_{*} \right) \right) \\ & \leq \left( \frac{\rho}{b_{\mathrm{eff}}} + \frac{1 - \rho}{b} \right) 2 \mathcal{L}_{\mathrm{max}} \left( F \left( \boldsymbol{x} \right) F \left( \boldsymbol{x}_{*} \right) \right) \\ & + \rho \left( 1 - \frac{1}{b_{\mathrm{eff}}} \right) 2 \left( \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i} \right) \left( F \left( \boldsymbol{x} \right) - F \left( \boldsymbol{x}_{*} \right) \right) \\ & + \frac{1}{b_{\mathrm{eff}}} 2 \mathcal{L}_{\mathrm{sub}} \left( F \left( \boldsymbol{x} \right) - F \left( \boldsymbol{x}_{*} \right) \right) . \\ & = 2 \left( \left( \frac{\rho}{b_{\mathrm{eff}}} + \frac{1 - \rho}{b} \right) \mathcal{L}_{\mathrm{max}} + \rho \left( 1 - \frac{1}{b_{\mathrm{eff}}} \right) \left( \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i} \right) \right) \end{split}$$

$$+\frac{1}{b_{\mathrm{eff}}}\mathcal{L}_{\mathrm{sub}}\left(F\left(\mathbf{x}\right)-F\left(\mathbf{x}_{*}\right)\right).$$

**Proof of (ii)** From Assumption 6 (B), we have

$$T_{\text{var}} = \frac{1}{n} \sum_{i=1}^{n} \text{tr} \mathbb{V} \left[ \mathbf{g}_{i} \left( \mathbf{x} \right) - \mathbf{g}_{i} \left( \mathbf{x}_{*} \right) \right]$$

$$\leq \frac{1}{n} \sum_{i=1}^{n} 2\mathcal{L}_{i} \left( F\left( \mathbf{x} \right) - F\left( \mathbf{x}_{*} \right) \right)$$

$$= 2 \left( \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_{i} \right) \left( F\left( \mathbf{x} \right) - F\left( \mathbf{x}_{*} \right) \right). \tag{30}$$

And,

$$T_{\text{cov}} = \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\text{tr} \mathbb{V} \left[ \mathbf{g}_{i} \left( \mathbf{x} \right) - \mathbf{g}_{i} \left( \mathbf{x}_{*} \right) \right]} \right)^{2}$$

$$\leq \frac{2}{n^{2}} \left( \sum_{i=1}^{n} \sqrt{\mathcal{L}_{i} \left( F\left( \mathbf{x} \right) - F\left( \mathbf{x}_{*} \right) \right)} \right)^{2}$$

$$= 2 \left( \frac{1}{n} \sum_{i=1}^{n} \sqrt{\mathcal{L}_{i}} \right)^{2} \left( F\left( \mathbf{x} \right) - F\left( \mathbf{x}_{*} \right) \right). \tag{31}$$

Plugging Eqs. (30) and (31) into Eq. (25), we have

$$\operatorname{tr}\mathbb{V}\left[\boldsymbol{g}\left(\boldsymbol{x}\right) - \boldsymbol{g}\left(\boldsymbol{x}_{*}\right)\right]$$

$$\leq \left(\frac{\rho}{b_{\mathrm{eff}}} + \frac{1-\rho}{b}\right) 2\left(\frac{1}{n}\sum_{i=1}^{n}\mathcal{L}_{i}\right) (F\left(\boldsymbol{x}\right) - F\left(\boldsymbol{x}_{*}\right))$$

$$+ \rho\left(1 - \frac{1}{b_{\mathrm{eff}}}\right) 2\left(\frac{1}{n}\sum_{i=1}^{n}\sqrt{\mathcal{L}_{i}}\right)^{2} (F\left(\boldsymbol{x}\right) - F\left(\boldsymbol{x}_{*}\right))$$

$$+ \frac{1}{b_{\mathrm{eff}}} 2\mathcal{L}_{\mathrm{sub}} (F\left(\boldsymbol{x}\right) - F\left(\boldsymbol{x}_{*}\right)),$$

$$= 2\left(\left(\frac{\rho}{b_{\mathrm{eff}}} + \frac{1-\rho}{b}\right) \left(\frac{1}{n}\sum_{i=1}^{n}\mathcal{L}_{i}\right)$$

$$+ \rho\left(1 - \frac{1}{b_{\mathrm{eff}}}\right) \left(\frac{1}{n}\sum_{i=1}^{n}\sqrt{\mathcal{L}_{i}}\right)^{2}$$

$$+ \frac{1}{b_{\mathrm{eff}}}\mathcal{L}_{\mathrm{sub}}\right) (F\left(\boldsymbol{x}\right) - F\left(\boldsymbol{x}_{*}\right)).$$

# B.4.2. BOUNDED VARIANCE CONDITION (THEOREM 3)

**Theorem 3.** Let Assumption 4 and 7 hold. Then,  $\mathbf{g}_B$  satisfies BV  $(\sigma^2)$ , where

$$\sigma^{2} = \left(\frac{\rho}{b_{\text{eff}}} + \frac{1-\rho}{b}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sigma_{i}^{2}\right) + \rho \left(1 - \frac{1}{b_{\text{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sigma_{i}\right)^{2} + \frac{\tau^{2}}{b_{\text{eff}}}.$$

Equality in Definition 2 holds if equality in Assumption 4 holds.

*Proof.* For any element of the solution set  $x_*$  = arg min<sub> $x \in Y$ </sub> F(x), by Theorem 1, we have

$$\begin{split} \operatorname{tr} \mathbb{V}\left[\mathbf{\textit{g}}_{B}\left(\mathbf{\textit{x}}_{*}\right)\right] & \leq \left(\frac{\rho}{b_{\mathrm{eff}}} + \frac{1-\rho}{b}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \operatorname{tr} \mathbb{V}\left[\mathbf{\textit{g}}_{i}\left(\mathbf{\textit{x}}_{*}\right)\right]\right) \\ & + \rho \left(1 - \frac{1}{b_{\mathrm{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\operatorname{tr} \mathbb{V}\left[\mathbf{\textit{g}}_{i}\left(\mathbf{\textit{x}}_{*}\right)\right]}\right)^{2} \\ & + \frac{1}{b_{\mathrm{eff}}} \operatorname{tr} \mathbb{V}\left[\nabla f_{B}\left(\mathbf{\textit{x}}_{*}\right)\right]. \end{split}$$

Applying Assumption 7, we have

$$\leq \left(\frac{\rho}{b_{\text{eff}}} + \frac{1-\rho}{b}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sigma_{i}^{2}\right)$$

$$+ \left(1 - \frac{1}{b_{\text{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sqrt{\sigma_{i}^{2}}\right)^{\frac{1}{2}}$$

$$+ \frac{1}{b_{\text{eff}}} \tau^{2}$$

$$= \left(1 - \frac{1}{b_{\text{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sigma_{i}^{2}\right)$$

$$+ \rho \left(1 - \frac{1}{b_{\text{eff}}}\right) \left(\frac{1}{n} \sum_{i=1}^{n} \sigma_{i}\right)^{2}$$

$$+ \frac{1}{b_{\text{eff}}} \tau^{2}.$$

#### B.4.3. COMPLEXITY ANALYSIS (COROLLARY 2)

**Lemma 10.** Let the objective function F satisfy Assumption 2,  $\pi$  be sampling b samples without replacement, and all elements of the solution set  $\mathbf{x}_* \in \arg\min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$  be stationary points of F. Then, the subsampling estimator  $\nabla f_B$  satisfies the ER condition as

$$\operatorname{tr} \mathbb{V}_{B \sim \pi} \left[ \nabla f_{B}(\boldsymbol{x}) - \nabla f_{B}(\boldsymbol{x}_{*}) \right]$$

$$\leq 2 \frac{n - b}{b(n - 1)} L_{\max} \left( F(\boldsymbol{x}) - F(\boldsymbol{x}_{*}) \right),$$

where  $L_{\max} = \max\{L_1, \dots, L_n\}$ .

*Proof.* Consider that, for any random vector **x**,

$$\operatorname{tr} \mathbb{V} \left[ \boldsymbol{x}^2 \right] \leq \mathbb{E} \| \boldsymbol{x} \|_2^2$$

holds. Also, sampling without replacement achieves  $b_{\rm eff} = \frac{(n-1)b}{n-b}$ . Therefore, we have

$$\operatorname{tr} \mathbb{V}_{B \sim \pi} \left[ \nabla f_{B}(\mathbf{x}) - \nabla f_{B}(\mathbf{x}_{*}) \right]$$

$$= \frac{n - b}{b (n - 1)} \operatorname{tr} \mathbb{V} \left[ \nabla f_{i}(\mathbf{x}) - \nabla f_{i}(\mathbf{x}_{*}) \right]$$

$$\leq \frac{n - b}{b (n - 1)} \left( \frac{1}{n} \sum_{i=1}^{n} \left\| \nabla f_{i}(\mathbf{x}) - \nabla f_{i}(\mathbf{x}_{*}) \right\|_{2}^{2} \right),$$

and from Assumption 2,

$$=\frac{n-b}{b(n-1)}\left(\frac{1}{n}\sum_{i=1}^{n}2L_{i}D_{f_{i}}(\boldsymbol{x},\boldsymbol{x}_{*})\right).$$

Using the bound  $L_{\text{max}} \ge L_i$  for all i = 1, ..., n,

$$\leq 2L_{\max} \frac{n-b}{b(n-1)} \left( \frac{1}{n} \sum_{i=1}^{n} D_{f_i}(\boldsymbol{x}, \boldsymbol{x}_*) \right),$$

applying Lemma 6.

$$=2L_{\max}\frac{n-b}{b(n-1)}D_F(\boldsymbol{x},\boldsymbol{x}_*),$$

and since  $x_*$  is a stationary point of F

$$=2\frac{n-b}{b(n-1)}L_{\max}\left(F\left(\boldsymbol{x}\right)-F\left(\boldsymbol{x}_{*}\right)\right).$$

**Corollary 2.** Let the objective F satisfy Assumption 1 and 2, the global optimum  $\mathbf{x}_* = \arg\min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$  be a stationary point of F, the component gradient estimators  $\mathbf{g}_1, \dots, \mathbf{g}_n$  satisfy Assumption 6 (B) and 7, and  $\pi$  be b-minibatch sampling without replacement. Then the last iterate of SGD with  $\mathbf{g}_B$  is  $\epsilon$ -close to  $\mathbf{x}_*$  as  $\mathbb{E} ||\mathbf{x}_T - \mathbf{x}_*||_2^2 \le \epsilon$  after a number of iterations of at least

$$T \ge 2 \max \left( C_{\text{var}} \frac{1}{\epsilon}, C_{\text{bias}} \right) \log \left( 2 || \boldsymbol{x}_0 - \boldsymbol{x}_* ||_2^2 \frac{1}{\epsilon} \right)$$

for some fixed stepsize where

$$C_{\text{var}} = \frac{2}{b} \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_i^2}{\mu^2} \right) + 2 \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_i}{\mu} \right)^2 + \frac{2}{b} \frac{\tau^2}{\mu^2},$$

$$C_{\text{bias}} = \frac{2}{b} \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\mathcal{L}_i}{\mu} \right) + 2 \left( \frac{1}{n} \sum_{i=1}^{n} \sqrt{\frac{\mathcal{L}_i}{\mu}} \right)^2 + \frac{2}{b} \frac{L}{\mu}.$$

*Proof.* From Assumption 2 and the assumption that  $x_*$  is a stationary point, Lemma 10 establishes that  $\nabla f_B$  satisfies the ER  $(\mathcal{L}_{\text{sub}})$  holds with

$$\mathcal{L}_{\text{sub}} = \frac{n-b}{(n-1)b} L_{\text{max}}.$$

Therefore, Assumption 5 holds. Furthermore, since the component gradient estimators satisfy Assumption 6 (B) and Assumption 3 always hold with  $\rho = 1$ , we can apply Theorem 2 which estblishes that  $\mathbf{g}_B$  satisfies ER ( $\mathcal{L}$ ) with

$$\mathcal{L} = \frac{n-b}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}_i \right) + \frac{n(b-1)}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sqrt{\mathcal{L}_i} \right)^2 + \frac{n-b}{(n-1)b} L_{\text{max}}.$$

Furthermore, under Assumption 7, Theorem 3 shows that BV ( $\sigma^2$ ) holds with

$$\sigma^{2} = \frac{n-b}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_{i}^{2} \right) + \frac{n(b-1)}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \right)^{2} + \frac{n-b}{(n-1)b} \tau^{2}.$$

Since both ER ( $\mathcal{L}$ ) and BV ( $\sigma^2$ ) hold and F satisfies Assumption 1, we can now invoke Lemma 1, which guarantees that we can obtain an  $\epsilon$ -accurate solution after

$$T \ge 2 \max\left(\underbrace{\frac{\sigma^2}{\mu^2}}_{C_{\text{var}}} \underbrace{\frac{1}{\epsilon}}, \underbrace{\frac{\mathcal{L} + L}{\mu}}_{C_{\text{bias}}}\right) \log\left(2\|\boldsymbol{x}_0 - \boldsymbol{x}_*\|_2^2 \frac{1}{\epsilon}\right)$$

iterations and fixed stepsize of

$$\gamma = \min\left(\frac{\epsilon\mu}{2\sigma^2}, \frac{1}{2(\mathcal{L} + L)}\right).$$

The constants in the lower bound on the number of required iterations can be made more precise as

$$C_{\text{var}} = \frac{n - b}{(n - 1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_i^2}{\mu^2} \right)$$

$$+ \frac{n(b - 1)}{(n - 1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_i}{\mu} \right)^2 + \frac{n - b}{(n - 1)b} \frac{\tau^2}{\mu^2}$$

$$C_{\text{bias}} = \frac{n - b}{(n - 1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \frac{\mathcal{L}_i}{\mu} \right)$$

$$+ \frac{n(b - 1)}{(n - 1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sqrt{\frac{\mathcal{L}_i}{\mu}} \right)^2 + \frac{n - b}{(n - 1)b} \frac{L}{\mu}.$$

Using the fact that  $(n-b)/n \le (n-1)/n \le 2$  for all  $n \ge 2$  yields the simplified constants in the statement.  $\Box$ 

## **B.5. Random Reshuffling of Stochastic Gradients**

# B.5.1. GRADIENT VARIANCE CONDITIONS (LEMMA 11, LEMMA 12)

**Lemma 11.** Let the objective function satisfy Assumption 2, B be any b-minibatch of indices such that  $B \subseteq \{1, ..., n\}$  and the component gradient estimators  $g_1, ..., g_n$  satisfy Assumption 6 (A<sup>CVX</sup>). Then,  $g_B$  is convex-smooth in expectation such that

$$\mathbb{E}_{\varphi} \left\| \boldsymbol{g}_{B}\left(\boldsymbol{x}\right) - \boldsymbol{g}_{B}\left(\boldsymbol{x}_{*}\right) \right\|_{2}^{2} \leq 2 \left(\mathcal{L}_{\max} + L_{\max}\right) D_{f_{B}}\left(\boldsymbol{x}, \boldsymbol{x}_{*}\right),$$

for any  $x \in \mathcal{X}$ , where

$$\begin{aligned} \boldsymbol{x}_* &= \operatorname*{arg\,min}_{\boldsymbol{x} \in \mathcal{X}} F\left(\boldsymbol{x}\right), \\ \mathcal{L}_{\max} &= \max\left\{\mathcal{L}_1, \dots, \mathcal{L}_n\right\}, \\ L_{\max} &= \max\left\{L_1, \dots, L_n\right\}. \end{aligned}$$

*Proof.* Notice that, for this Lemma, we do not assume that the minibatch B is a random variable. Therefore, the only randomness is the stochasticity of the component gradient estimators  $g_1, \dots, g_n$ .

Now, from the property of the variance, we can decompose the expected squared norm as

$$\mathbb{E}\|\boldsymbol{g}_{B}\left(\boldsymbol{x}\right)-\boldsymbol{g}_{B}\left(\boldsymbol{x}_{*}\right)\|_{2}^{2}$$

$$=\underbrace{\operatorname{tr}\mathbb{V}_{\varphi}\left[\boldsymbol{g}_{B}\left(\boldsymbol{x}\right)-\boldsymbol{g}_{B}\left(\boldsymbol{x}_{*}\right)\right]}_{V_{\operatorname{com}}}+\underbrace{\left\|\nabla f_{B}\left(\boldsymbol{x}\right)-\nabla f_{B}\left(\boldsymbol{x}_{*}\right)\right\|_{2}^{2}}_{V_{\operatorname{sub}}}.$$

First, the contribution of the variances of the component gradient estimators follows as

$$\begin{split} V_{\text{com}} &= \text{tr} \mathbb{V}_{\varphi} \left[ \boldsymbol{g} \left( \boldsymbol{x} \right) - \boldsymbol{g} \left( \boldsymbol{x}_{*} \right) \right] \\ &= \text{tr} \mathbb{V}_{\varphi} \left[ \frac{1}{b} \sum_{i \in B} \boldsymbol{g}_{i} \left( \boldsymbol{x} \right) - \boldsymbol{g}_{i} \left( \boldsymbol{x}_{*} \right) \right], \end{split}$$

applying Eq. (14) of Lemma 3,

$$\leq \frac{1}{b} \sum_{i \in B} \operatorname{tr} \mathbb{V}_{\varphi} \left[ \mathbf{g}_{i} \left( \mathbf{x} \right) - \mathbf{g}_{i} \left( \mathbf{x}_{*} \right) \right], \tag{32}$$

and then Assumption 6 (ACVX),

$$\leq \frac{1}{b} \sum_{i \in B} 2\mathcal{L}_i \, \mathcal{D}_{f_i}(\boldsymbol{x}, \boldsymbol{x}_*).$$

Now, since  $\mathcal{L}_{\max} \geq \mathcal{L}_i$  for all i = 1, ..., n,

$$\leq 2\mathcal{L}_{\max} \frac{1}{b} \sum_{i \in B} D_{f_i}(\boldsymbol{x}, \boldsymbol{x}_*)$$

$$=2\mathcal{L}_{\max}\mathrm{D}_{f_{B}}\left(\boldsymbol{x},\boldsymbol{x}_{*}\right).$$

On the other hand, the squared error of subsampling (it is not the variance since we do not take expectation over the batches) follows as

$$V_{\text{sub}} = \left\| \nabla f_B(\mathbf{x}) - \nabla f_B(\mathbf{x}_*) \right\|_2^2$$
$$= \left\| \frac{1}{b} \sum_{i \in B} \nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}_*) \right\|_2^2,$$

by Jensen's inequality,

$$\leq \frac{1}{b} \sum_{i \in R} \left\| \nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}_*) \right\|_2^2,$$

from Assumption 2

$$\leq \frac{1}{b}\sum_{i\in B}2L_{i}\mathrm{D}_{f_{i}}\left(\boldsymbol{x},\boldsymbol{x}_{*}\right)$$

and since  $L_{\max} \ge L_i$  for all i = 1, ..., n,

$$\leq 2L_{\max} \frac{1}{b} \sum_{i \in B} D_{f_i}(\boldsymbol{x}, \boldsymbol{x}_*)$$
$$= 2L_{\max} D_{f_B}(\boldsymbol{x}, \boldsymbol{x}_*).$$

Combining the bound on  $V_{\rm com}$  and  $V_{\rm sub}$  immediately yields the result.

**Lemma 12.** For any b-minibatch reshuffling strategy, the squared error of the reference point of the Lyapunov function (Eq. (10)) under reshuffling is bounded as

$$\mathbb{E}||\boldsymbol{x}_{*}^{i} - \boldsymbol{x}_{*}||_{2}^{2} \le \frac{\gamma^{2}n}{4b^{2}}\tau^{2}$$

for all i = 1, ..., p, where  $\mathbf{x}_* \in \arg\min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$ .

*Proof.* The proof is a generalization of Mishchenko et al. (2020, Proposition 1), where we sample b-minibatches instead of single datapoints. Recall that P denotes the (possibly random) partitioning of the n datapoints into b-minibatches  $P_1, \ldots, P_p$ . From the definition of the squared error of the Lyapunov function in Eq. (10), we have

$$\mathbb{E}\left[\left\|\boldsymbol{x}_{*}^{i}-\boldsymbol{x}_{*}\right\|_{2}^{2}\right]$$

$$=\mathbb{E}\left[\left\|\Pi_{\mathcal{X}}\left(\boldsymbol{x}_{*}-\sum_{k=0}^{i-1}\gamma\nabla f_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right)-\Pi_{\mathcal{X}}\left(\boldsymbol{x}_{*}\right)\right\|_{2}^{2}\right],$$

and since the projection onto a convex set under a Euclidean metric is non-expansive,

$$\leq \mathbb{E}\left[\left\|\boldsymbol{x}_{*} - \sum_{k=0}^{i-1} \gamma \nabla f_{P_{i}}\left(\boldsymbol{x}_{*}\right) - \boldsymbol{x}_{*}\right\|_{2}^{2}\right]$$

$$= \mathbb{E}\left[\left\|\sum_{k=0}^{i-1} \gamma \nabla f_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right\|_{2}^{2}\right],$$

introducing a factor of i in and out of the squared norm,

$$= \frac{i^2}{2} \mathbb{E} \left[ \left\| \frac{1}{i} \sum_{k=0}^{i-1} \gamma \nabla f_{P_i} (\boldsymbol{x}_*) \right\|_2^2 \right]$$
$$= \frac{\gamma^2 i^2}{2} \mathbb{E} \left[ \left\| \frac{1}{i} \sum_{k=0}^{i-1} \nabla f_{P_i} (\boldsymbol{x}_*) \right\|_2^2 \right].$$

Now notice that  $\frac{1}{i} \sum_{j=0}^{i-1} \nabla f_{P_i}(\boldsymbol{x}_*)$  is a sample average of ib samples drawn without replacement. Therefore, it is an unbiased estimate of  $\nabla F(\boldsymbol{x}_*)$ . This implies

$$\mathbb{E}\left[\left\|\boldsymbol{x}_{*}^{i}-\boldsymbol{x}_{*}\right\|_{2}^{2}\right] = \frac{\gamma^{2}i^{2}}{2}\,\mathbb{E}\left[\left\|\frac{1}{i}\sum_{k=0}^{i-1}\nabla f_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right\|_{2}^{2}\right]$$
$$=\frac{\gamma^{2}i^{2}}{2}\,\mathrm{tr}\mathbb{V}\left[\frac{1}{i}\sum_{k=0}^{i-1}\nabla f_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right],$$

and from Lemma 2 with a sample size of ib,

$$= \frac{\gamma^{2} i^{2}}{2} \frac{n - ib}{(n - 1) ib} \frac{1}{n} \sum_{i=1}^{n} \|\nabla f_{i}(\mathbf{x}_{*})\|_{2}^{2}$$

$$=\frac{\gamma^2 i \left(\frac{n}{b}-i\right)}{2(n-1)} \tau^2.$$

Notice that this is a quadratic with respect to i, where the maximum is obtained by i = n/2b. Then,

$$\mathbb{E}\left[\left\|\boldsymbol{x}_{*}^{i}-\boldsymbol{x}_{*}\right\|_{2}^{2}\right] \leq \frac{\gamma^{2}\left(\frac{n}{2b}\right)^{2}}{2(n-1)}\tau^{2}$$
$$=\frac{\gamma^{2}n^{2}}{8b^{2}(n-1)}\tau^{2},$$

and using the bound  $n/(n-1) \le 2$  for all  $n \ge 2$ ,  $\le \frac{\gamma^2 n}{4h^2} \tau^2.$ 

#### B.5.2. Convergence Analysis (Theorem 5)

**Theorem 5.** Let the objective F satisfy Assumption 1 and 2, where, each component  $f_i$  is additionally  $\mu$ -strongly convex and Assumption 6 (A<sup>CVX</sup>), 7 hold. Then, the last iterate  $\mathbf{x}_T$  of doubly SGD-RR with a stepsize satisfying  $\gamma < 1/(\mathcal{L}_{max} + L_{max})$  guarantees

$$\mathbb{E}\|\boldsymbol{x}_{K+1}^{0}-\boldsymbol{x}_{*}\|_{2}^{2} \leq r^{Kp}\|\boldsymbol{x}_{1}^{0}-\boldsymbol{x}_{*}\|_{2}^{2} + C_{\text{var}}^{\text{sub}} \gamma^{2} + C_{\text{var}}^{\text{com}} \gamma$$

where p = n/b is the number of epochs,  $\mathbf{x}_* = \arg\min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$ ,  $r = 1 - \gamma \mu$  is the contraction coefficient,

$$C_{\text{var}}^{\text{com}} = \frac{4}{\mu b} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_i^2 \right) + \frac{4}{\mu} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_i \right)^2, \text{ and}$$

$$C_{\text{var}}^{\text{sub}} = \frac{1}{4} \frac{L_{\text{max}}}{\mu} \frac{n}{b^2} \left( \frac{1}{n} \sum_{i=1}^{n} \left\| \nabla f_i \left( \boldsymbol{x}_* \right) \right\|_2^2 \right).$$

*Proof.* The key element of the analysis of random reshuffling is that the Lyapunov function that achieves a fast convergence is  $\|\boldsymbol{x}_k^{i+1} - \boldsymbol{x}_*^{i+1}\|_2^2$  not  $\|\boldsymbol{x}_k^{i+1} - \boldsymbol{x}_*\|_2^2$ . This stems from the well-known fact that random reshuffling results in a conditionally biased gradient estimator.

Recall that P denotes the partitioning of the n datapoints into b-minibatches  $P_1, \dots, P_p$ . As usual, we first expand the Lyapunov function as

$$\begin{aligned} & \|\boldsymbol{x}_{k}^{i+1} - \boldsymbol{x}_{*}^{i+1}\|_{2}^{2} \\ &= \|\Pi_{\mathcal{X}}(\boldsymbol{x}_{k}^{i} - \gamma \, \boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i})) - \Pi_{\mathcal{X}}(\boldsymbol{x}_{*}^{i} - \gamma \, \nabla f_{P_{i}}(\boldsymbol{x}_{*}))\|_{2}^{2} \end{aligned}$$

and since the projection onto a convex set under a Euclidean metric is non-expansive,

$$\leq \|(\boldsymbol{x}_{k}^{i} - \gamma \, \boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i})) - (\boldsymbol{x}_{*}^{i} - \gamma \, \nabla f_{P_{i}}(\boldsymbol{x}_{*}))\|_{2}^{2}$$

$$= \|\boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}\|_{2}^{2} - 2\gamma \left\langle \boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i}, \boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i}) - \nabla f_{P_{i}}(\boldsymbol{x}_{*}) \right\rangle$$

$$+ \gamma^{2} \|\boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i}) - \nabla f_{P_{i}}(\boldsymbol{x}_{*})\|_{2}^{2}.$$

Taking expectation over the Monte Carlo noise conditional on the partitioning P,

$$\begin{split} & \mathbb{E}_{\varphi} \| \boldsymbol{x}_{k}^{i+1} - \boldsymbol{x}_{*}^{i+1} \|_{2}^{2} \\ & = \left\| \boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i} \right\|_{2}^{2} - 2\gamma \left\langle \boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i} , \mathbb{E}_{\varphi} [\boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i})] - \nabla f_{P_{i}}(\boldsymbol{x}_{*}) \right\rangle \\ & + \gamma^{2} \mathbb{E}_{\varphi} \| \boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i}) - \nabla f_{P_{i}}(\boldsymbol{x}_{*}) \| \\ & = \left\| \boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i} \right\|_{2}^{2} - 2\gamma \left\langle \boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i} , \nabla f_{P_{i}}(\boldsymbol{x}_{k}^{i}) - \nabla f_{P_{i}}(\boldsymbol{x}_{*}) \right\rangle \\ & + \gamma^{2} \mathbb{E}_{\varphi} \| \boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i}) - \nabla f_{P_{i}}(\boldsymbol{x}_{*}) \|_{2}^{2}. \end{split}$$

From the three-point identity, we can more precisely characterize the effect of the conditional bias such that

$$\begin{split} &\left\langle \boldsymbol{x}_{k}^{i}-\boldsymbol{x}_{*}^{i}\;,\nabla f_{P_{i}}(\boldsymbol{x}_{k}^{i})-\nabla f_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right\rangle \\ &=\mathrm{D}_{f_{P_{i}}}(\boldsymbol{x}_{*}^{i},\boldsymbol{x}_{k}^{i})+\mathrm{D}_{f_{P_{i}}}(\boldsymbol{x}_{k}^{i},\boldsymbol{x}_{*})-\mathrm{D}_{f_{P_{i}}}(\boldsymbol{x}_{*}^{i},\boldsymbol{x}_{*}). \end{split}$$

For the gradient noise,

$$\begin{split} &\mathbb{E}_{\varphi} \| \boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i}) - \nabla f_{P_{i}}\left(\boldsymbol{x}_{*}\right) \|_{2}^{2} \\ &= \mathbb{E}_{\varphi} \| \boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i}) - \boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right) + \boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right) - \nabla f_{P_{i}}\left(\boldsymbol{x}_{*}\right) \|_{2}^{2} \\ &\leq 2\mathbb{E}_{\varphi} \| \boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i}) - \boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right) \|_{2}^{2} + 2\mathbb{E}_{\varphi} \| \boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right) - \nabla f_{P_{i}}\left(\boldsymbol{x}_{*}\right) \|_{2}^{2} \\ &= 2\mathbb{E}_{\varphi} \| \boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{k}^{i}) - \boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right) \|_{2}^{2} + 2\operatorname{tr} \mathbb{V}_{\varphi} \left[\boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right], \\ \text{and from Lemma 11,} \\ &\leq 4\left(\mathcal{L}_{\max} + L_{\max}\right) D_{f_{P_{i}}}(\boldsymbol{x}_{k}^{i}, \boldsymbol{x}_{*}) + 2\operatorname{tr} \mathbb{V}_{\varphi} \left[\boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right] \end{split}$$

 $= (\mathsf{Sinax} \setminus \mathsf{Sinax}) \setminus f_{P_i}(\mathsf{C}_k, \mathsf{C}_k) + \mathsf{Cor}_{\varphi}(\mathsf{C}_k, \mathsf{C}_k)$ 

Notice the variance term  $\operatorname{tr} \mathbb{V}_{\varphi} \left[ \mathbf{g}_{P_i} \left( \mathbf{x}_* \right) \right]$ . This quantifies the amount of deviation from the trajectory of singly stochastic random reshuffling. As such, it quantifies how slower we will be compared to its fast rate.

Now, we will denote the  $\sigma$ -algebra formed by the randomness and the iterates up to the *i*th step of the *k*th epoch as  $\mathcal{F}_k^i$  such that  $(\mathcal{F}_k^i)_{k\geq 1, i\geq 1}$  is a filtration. Then,

$$\begin{split} &\mathbb{E}_{\boldsymbol{\eta}_{k}^{i} \sim \varphi} \left[ \left\| \boldsymbol{x}_{k}^{i+1} - \boldsymbol{x}_{*}^{i+1} \right\|_{2}^{2} \middle| \mathcal{F}_{k}^{i} \right] \\ & \leq \left\| \boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i} \right\|_{2}^{2} \\ & - 2 \gamma \left( \mathbf{D}_{f_{P_{i}}}(\boldsymbol{x}_{*}^{i}, \boldsymbol{x}_{k}^{i}) + \mathbf{D}_{f_{P_{i}}}(\boldsymbol{x}_{k}^{i}, \boldsymbol{x}_{*}) - \mathbf{D}_{f_{P_{i}}}(\boldsymbol{x}_{*}^{i}, \boldsymbol{x}_{*}) \right) \\ & + 4 \gamma^{2} \left( \mathcal{L}_{\max} + L_{\max} \right) \mathbf{D}_{f_{P_{i}}}(\boldsymbol{x}_{k}^{i}, \boldsymbol{x}_{*}) \\ & + 2 \gamma^{2} \mathrm{tr} \mathbb{V}_{\varphi} \left[ \boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right) \right]. \end{split}$$

Now, the  $\mu$ -strong convexity of the component functions imply  $D_{f_{P_i}}\left(\boldsymbol{x}_*^i, \boldsymbol{x}_k^i\right) \leq \frac{\mu}{2} \|\boldsymbol{x}_k^i - \boldsymbol{x}_*^i\|_2^2$ . Therefore,

$$\leq \|\boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i}\|_{2}^{2}$$

$$-2\gamma \left(\frac{\mu}{2} \|\boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i}\|_{2}^{2} + D_{f_{P_{i}}}(\boldsymbol{x}_{k}^{i}, \boldsymbol{x}_{*}) - D_{f_{P_{i}}}(\boldsymbol{x}_{*}^{i}, \boldsymbol{x}_{*})\right)$$

$$+4\gamma^{2} \left(\mathcal{L}_{\max} + L_{\max}\right) D_{f_{P_{i}}}(\boldsymbol{x}_{k}^{i}, \boldsymbol{x}_{*})$$

$$+2\gamma^{2} \operatorname{tr} \mathbb{V}_{\varphi} \left[\boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right],$$

and reorganizing the terms,

$$= (1 - \gamma \mu) \|\boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i}\|_{2}^{2}$$

$$- 2\gamma (1 - 2\gamma (\mathcal{L}_{\max} + L_{\max})) D_{f_{P_{i}}}(\boldsymbol{x}_{k}^{i}, \boldsymbol{x}_{*})$$

$$+ 2\gamma D_{f_{P_{i}}}(\boldsymbol{x}_{*}^{i}, \boldsymbol{x}_{*})$$

$$+ \gamma^{2} 2 \operatorname{tr} \mathbb{V}_{\varphi} \left[\boldsymbol{g}_{P_{i}}(\boldsymbol{x}_{*})\right].$$

Taking full expectation,

$$\mathbb{E}\|\boldsymbol{x}_{k}^{i+1} - \boldsymbol{x}_{*}^{i+1}\|_{2}^{2}$$

$$\leq (1 - \gamma\mu) \mathbb{E}\|\boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i}\|_{2}^{2}$$

$$- 2\gamma (1 - 2\gamma (\mathcal{L}_{\max} + L_{\max})) \mathbb{E}\left[D_{f_{P_{i}}}\left(\boldsymbol{x}_{k}^{n}, \boldsymbol{x}_{*}\right)\right]$$

$$+ 2\gamma \mathbb{E}\left[D_{f_{P_{i}}}\left(\boldsymbol{x}_{*}^{i}, \boldsymbol{x}_{*}\right)\right]$$

$$+ 2\gamma^{2}\mathbb{E}\left[\text{tr}\mathbb{V}_{\varphi}\left[\boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right]\right],$$
and as long as  $\gamma < 1/(2(\mathcal{L}_{\max} + L_{\max}))$ 

$$\leq (1 - \gamma\mu) \mathbb{E}\|\boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i}\|_{2}^{2} + 2\gamma \mathbb{E}\left[D_{f_{P_{i}}}\left(\boldsymbol{x}_{*}^{i}, \boldsymbol{x}_{*}\right)\right]$$

$$+ 2\gamma^{2}\mathbb{E}\left[\text{tr}\mathbb{V}_{\varphi}\left[\boldsymbol{g}_{P_{i}}\left(\boldsymbol{x}_{*}\right)\right]\right].$$

$$(33)$$

**Bounding**  $T_{\text{err}}$  From the definition of the Bregman divergence and L-smoothness, for all j = 1, ..., n, notice that we have

$$D_{f_{j}}(\mathbf{y}, \mathbf{x}) = f_{j}(\mathbf{y}) - f_{j}(\mathbf{x}) - \langle \nabla f_{j}(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$$

$$\leq \frac{L}{2} ||\mathbf{y} - \mathbf{x}||_{2}^{2}.$$
(34)

for all  $(x, x') \in \mathcal{X}^2$ . Given this, the Lyapunov error term

$$\mathbb{E}\left[D_{f_{P_i}}\left(\boldsymbol{x}_k^i, \boldsymbol{x}_*\right)\right] = \mathbb{E}\left[\frac{1}{b}\sum_{j \in P_i}D_{f_j}\left(\boldsymbol{x}_k^i, \boldsymbol{x}_*\right)\right]$$

can be bounded using L-smoothness by Eq. (34),

$$\leq \mathbb{E}\left[\frac{1}{b}\sum_{j\in\mathcal{P}_i}\frac{L_j}{2}\|\boldsymbol{x}_k^i-\boldsymbol{x}_*\|_2^2\right]$$

and  $L_{\max} \ge L_i$  for all i = 1, ..., n,

$$\leq \frac{L_{\max}}{2} \mathbb{E} \left[ \frac{1}{b} \sum_{j \in P_i} \left\| \boldsymbol{x}_k^i - \boldsymbol{x}_* \right\|_2^2 \right]$$
$$= \frac{L_{\max}}{2} \mathbb{E} \left\| \boldsymbol{x}_k^i - \boldsymbol{x}_* \right\|_2^2. \tag{35}$$

The squared error  $\left\| oldsymbol{x}_k^i - oldsymbol{x}_* 
ight\|_2^2$  is bounded in Lemma 12 as

$$\mathbb{E}\|\boldsymbol{x}_{k}^{i}-\boldsymbol{x}_{*}\|_{2}^{2} \leq \epsilon_{\mathrm{sfl}}^{2} \triangleq \frac{\gamma^{2}n}{4h^{2}} \tau^{2} < \infty. \tag{36}$$

**Bounding**  $T_{\text{var}}$  Now, let's take a look at the variance term. First, notice that, by the Law of Total Expectation,

$$\mathbb{E}\left[\mathrm{tr}\mathbb{V}_{\varphi}\left[\mathbf{\textit{g}}_{\textit{P}_{i}}\left(\mathbf{\textit{x}}_{*}\right)\right]\right] = \mathbb{E}\left[\mathbb{E}\left[\mathrm{tr}\mathbb{V}_{\varphi}\left[\mathbf{\textit{g}}_{\textit{P}_{i}}\left(\mathbf{\textit{x}}_{*}\right)\right]\mid\textit{P}\right]\right].$$

Here,

$$\mathbb{E}\left[\mathrm{tr}\mathbb{V}_{\varphi}\left[\mathbf{g}_{P_{i}}\left(\mathbf{x}_{*}\right)\right]\mid P\right]$$

is the variance from selecting b samples without replacement. We can thus apply Lemma 9 with  $b_{\rm eff} = \frac{(n-1)b}{n-b}$  such that

$$\mathbb{E}\left[\operatorname{tr}\mathbb{V}_{\varphi}\left[\mathbf{g}_{P_{i}}\left(\mathbf{x}_{*}\right)\right]\mid P\right]$$

$$\leq \frac{n-b}{(n-1)\,b}\left(\frac{1}{n}\sum_{j=1}^{n}\sigma_{j}^{2}\right) + \frac{n\,(b-1)}{(n-1)\,b}\left(\frac{1}{n}\sum_{j=1}^{n}\sigma_{j}\right)^{2},$$

which we will denote as

$$=\sigma^2\tag{37}$$

for clarity. Also, notice that  $\sigma^2$  no longer depends on the partitioning.

**Per-step Recurrence Equation** Applying Eqs. (35) and (37) to Eq. (33), we now have the recurrence equation

$$\begin{split} \mathbb{E} \|\boldsymbol{x}_{k}^{i+1} - \boldsymbol{x}_{*}^{i+1}\|_{2}^{2} &\leq (1 - \gamma \mu) \, \mathbb{E} \|\boldsymbol{x}_{k}^{i} - \boldsymbol{x}_{*}^{i}\|_{2}^{2} \\ &+ L_{\max} \varepsilon_{\text{sfl}}^{2} \, \gamma + 2\sigma^{2} \, \gamma^{2}. \end{split}$$

Now that we have a contraction of the Lyapunov function  $\mathbb{E}\|\boldsymbol{x}_k^{i+1}-\boldsymbol{x}_*^{i+1}\|_2^2$ , it remains to convert this that the Lyapunov function bounds our objective  $\mathbb{E}\|\boldsymbol{x}_k^{i+1}-\boldsymbol{x}_*\|_2^2$ . This can be achieved by noticing that, at the end of each epoch, we have  $\boldsymbol{x}_{k+1}-\boldsymbol{x}_*=\boldsymbol{x}_k^p-\boldsymbol{x}_*^p$ , and equivalently, we have  $\boldsymbol{x}_k-\boldsymbol{x}_*=\boldsymbol{x}_k^0-\boldsymbol{x}_*^0$  at the beginning of the epoch. The fact that the relationship with the original objective is only guaranteed at the endpoints (beginning and end of the epoch) is related to the fact that the bias of random reshuffling starts increasing at the beginning of the epoch and starts decreasing near the end.

**Per-Epoch Recurrence Equation** Nevertheless, this implies that by simply unrolling the recursion as in the analysis of regular SGD, we obtain a per-epoch contraction of

$$\begin{split} \mathbb{E} \| \boldsymbol{x}_{k+1}^{0} - \boldsymbol{x}_{*} \|_{2}^{2} &\leq (1 - \gamma \mu)^{p} \mathbb{E} \| \boldsymbol{x}_{k}^{0} - \boldsymbol{x}_{*} \|_{2}^{2} \\ &+ \left( L_{\max} \epsilon_{\text{sfl}}^{2} \gamma + 2 \sigma^{2} \gamma^{2} \right) \left( \sum_{i=0}^{p-1} \left( 1 - \mu \gamma \right)^{i} \right). \end{split}$$

And after K epochs,

$$\begin{split} & \mathbb{E} \left\| \boldsymbol{x}_{K+1}^{0} - \boldsymbol{x}_{*} \right\|_{2}^{2} \leq \left(1 - \gamma \mu\right)^{pK} \mathbb{E} \left\| \boldsymbol{x}_{0}^{0} - \boldsymbol{x}_{*} \right\|_{2}^{2} \\ & + \left(L_{\max} \epsilon_{\mathrm{sfl}}^{2} \gamma + 2\sigma^{2} \gamma^{2}\right) \left(\sum_{i=0}^{p-1} \left(1 - \mu \gamma\right)^{i}\right) \left(\sum_{j=0}^{pK-1} \left(1 - \mu \gamma\right)^{pj}\right). \end{split}$$

Note that T = pK.

As done by Mishchenko et al. (2020), the product of sums can be bounded as

$$\begin{split} &\left(\sum_{i=0}^{p-1} (1 - \mu \gamma)^{i}\right) \left(\sum_{j=0}^{T-1} (1 - \mu \gamma)^{pj}\right) \\ &= \sum_{i=0}^{p-1} \sum_{j=0}^{T-1} (1 - \mu \gamma)^{i} (1 - \mu \gamma)^{pj} \\ &= \sum_{i=0}^{p-1} \sum_{j=0}^{T-1} (1 - \mu \gamma)^{i} (1 - \mu \gamma)^{pj} \\ &= \sum_{i=0}^{T-1} (1 - \mu \gamma)^{i} \\ &\leq \sum_{i=0}^{\infty} (1 - \mu \gamma)^{i} \\ &\leq \frac{1}{\gamma \mu}. \end{split}$$

Then,

$$\begin{split} & \mathbb{E} \| \boldsymbol{x}_{K+1}^{0} - \boldsymbol{x}_{*} \|_{2}^{2} \\ & \leq (1 - \gamma \mu)^{pK} \mathbb{E} \| \boldsymbol{x}_{0}^{0} - \boldsymbol{x}_{*} \|_{2}^{2} + \frac{1}{\gamma \mu} \left( L_{\max} \varepsilon_{\text{sfl}}^{2} \gamma + 2\sigma^{2} \gamma^{2} \right) \\ & = (1 - \gamma \mu)^{pK} \mathbb{E} \| \boldsymbol{x}_{0}^{0} - \boldsymbol{x}_{*} \|_{2}^{2} + \frac{\varepsilon_{\text{sfl}}^{2}}{\mu} + \frac{2\sigma^{2}}{\mu} \gamma. \end{split}$$

Plugging in the value of  $\varepsilon_{sfl}^2$  from Eq. (36), we have

$$\begin{split} \mathbb{E} \| \boldsymbol{x}_{K+1}^{0} - \boldsymbol{x}_{*} \|_{2}^{2} & \leq (1 - \gamma \mu)^{pK} \mathbb{E} \| \boldsymbol{x}_{0}^{0} - \boldsymbol{x}_{*} \|_{2}^{2} \\ & + \frac{L_{\max} n \sigma_{\text{sub}}^{2}}{4b^{2} \mu} \gamma^{2} + \frac{2\sigma^{2}}{\mu} \gamma. \end{split}$$

This implies

$$\mathbb{E}\|\boldsymbol{x}_{K+1}^{0}-\boldsymbol{x}_{*}\|_{2}^{2} \leq r^{Kn/b}\|\boldsymbol{x}_{1}^{0}-\boldsymbol{x}_{*}\|_{2}^{2} + C_{\text{var}}^{\text{sub}}\gamma^{2} + C_{\text{var}}^{\text{com}}\gamma,$$

where  $r = 1 - \gamma \mu$ ,

$$C_{\text{var}}^{\text{sub}} = \frac{1}{4} \frac{L_{\text{max}}}{\mu} \frac{n}{b^2} \left( \frac{1}{n} \sum_{i=1}^{n} \|\nabla f_i(\mathbf{x}_*)\|_2^2 \right), \text{ and}$$

$$C_{\text{var}}^{\text{com}} = \frac{2}{\mu} \frac{n-b}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_i^2 \right) + \frac{2}{\mu} \frac{n(b-1)}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_i \right)^2.$$

Applying the fact that  $(n - b)/n \le (n - 1)/n \le 2$  for all  $n \ge 2$  yields the simplified constants in the statement.  $\square$ 

## B.5.3. COMPLEXITY ANALYSIS (THEOREM 4)

**Theorem 4.** Let the objective F satisfy Assumption 1 and 2, where each component  $f_i$  is additionally  $\mu$ -strongly convex, and Assumption 6 (A<sup>CVX</sup>), 7 hold. Then, the last iterate  $\mathbf{x}_T$  of doubly SGD-RR is  $\epsilon$ -close to the global optimum  $\mathbf{x}_* = \arg\max_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x})$  such that  $\mathbb{E}||\mathbf{x}_T - \mathbf{x}_*||_2^2 \le \epsilon$  after a number of iterations of at least

$$T \, \geq \max \left( 4 C_{\mathrm{var}}^{\mathrm{com}} \frac{1}{\epsilon} + C_{\mathrm{var}}^{\mathrm{sub}} \frac{1}{\sqrt{\epsilon}}, \; C_{\mathrm{bias}} \right) \log \left( 2 \left\| \boldsymbol{x}_{1}^{0} - \boldsymbol{x}_{*} \right\|_{2}^{2} \frac{1}{\epsilon} \right)$$

for some fixed stepsize, where T = Kp = Kn/b,

$$\begin{split} &C_{\text{bias}} = \left(\mathcal{L}_{\text{max}} + L\right) / \mu \\ &C_{\text{var}}^{\text{com}} = \frac{2}{b} \left(\frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_{i}^{2}}{\mu^{2}}\right) + 2 \left(\frac{1}{n} \sum_{i=1}^{n} \frac{\sigma_{i}}{\mu}\right)^{2}, \\ &C_{\text{var}}^{\text{sub}} = \sqrt{\frac{L_{\text{max}}}{\mu}} \frac{\sqrt{n}}{b} \frac{\tau}{\mu}. \end{split}$$

*Proof.* From the result of Theorem 5, we can invoke Lemma 5 with

$$\begin{split} A &= \frac{L_{\max} n}{4b^2 \mu} \tau^2, \\ B &= \frac{2}{\mu} \left( \frac{n-b}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i^2 \right) + \frac{n(b-1)}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i \right)^2 \right), \\ C &= \mathcal{L}_{\max} + L_{\max}. \end{split}$$

Then, an  $\epsilon$  accurate solution in expectation can be obtained after

$$T \geq \max\left(\underbrace{\frac{2B}{\mu}}_{\triangleq C_1} \frac{1}{\epsilon} + \underbrace{\frac{\sqrt{2A}}{\mu}}_{\triangleq C_2} \frac{1}{\sqrt{\epsilon}}, \frac{\mathcal{L}_{\max} + L_{\max}}{\mu}\right) \log\left(2r_0^2 \frac{1}{\epsilon}\right)$$

iterations with a stepsize of

$$\gamma = \min\left(\frac{-B + \sqrt{B^2 + 2A\epsilon}}{2A}, \frac{1}{C}\right).$$

To make the iteration complexity more precise, the terms  $C_1$ ,  $C_2$  can be organized as

$$C_{1} = \frac{2B}{\mu} = \frac{2}{\mu} \left( \frac{2}{\mu} \left\{ \frac{n-b}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_{i}^{2} \right) + \frac{n(b-1)}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \right)^{2} \right\} \right)$$

$$= \frac{4}{\mu^{2}} \left( \frac{n-b}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_{i}^{2} \right) + \frac{n(b-1)}{(n-1)b} \left( \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \right)^{2} \right)$$

$$C_2 = \frac{\sqrt{2A}}{\mu}$$

$$= \sqrt{2\frac{L_{\text{max}}n}{4b^2\mu}\tau^2} \frac{1}{\mu^2}$$

$$= \frac{\sqrt{L_{\text{max}}}\tau\sqrt{n}}{\sqrt{2b\mu^{3/2}}}$$

$$\leq \frac{\sqrt{L_{\text{max}}}\sqrt{n}}{\mu^{3/2}} \frac{\sqrt{n}}{b}\tau.$$

Applying the fact that  $(n-b)/n \le (n-1)/n \le 2$  for all  $n \ge 2$  yields the simplified constants in the statement.  $\square$ 

# C. Applications

## C.1. ERM with Randomized Smoothing

#### C.1.1. DESCRIPTION

Randomized smoothing was originally considered by Polyak & d Aleksandr Borisovich (1990); Nesterov (2005); Duchi et al. (2012) in the nonsmooth convex optimization context, where the function is "smoothed" through random perturbation. This scheme has recently renewed interest in the non-convex ERM context as it has been found to improve generalization performance (Orvieto et al., 2023; Liu et al., 2021). Here, we will focus on the computational aspect of this scheme. In particular, we will see if we can obtain similar computational guarantees already established in the finite-sum ERM setting, such as those by Gower et al. (2021a, Lemma 5.2).

Consider the canonical ERM problem, where we are given a dataset  $\mathcal{D} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^n \in (\mathcal{X} \times \mathcal{Y})^n \text{ and solve}$ 

$$\underset{\boldsymbol{w} \in \mathcal{W}}{\text{minimize}} L(\boldsymbol{w}) = \frac{1}{n} \sum_{i=1}^{n} \ell(f_{\boldsymbol{w}}(\boldsymbol{x}_i), y_i) + h(\boldsymbol{w}),$$

where  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$  are the feature and label of the ith instance,  $f_w : \mathcal{X} \to \mathcal{Y}$  is the model,  $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_{\geq 0}$  is a non-negative loss function, and  $h : \mathcal{W} \to \mathbb{R}$  is a regularizer.

For randomized smoothing, we instead minimize

$$L(\boldsymbol{w}) = \frac{1}{n} \sum_{i=1}^{n} R_i(\boldsymbol{w}),$$

where the instance risk is defined as

$$r_i(\mathbf{w}) = \mathbb{E}_{\boldsymbol{\epsilon} \sim \boldsymbol{\sigma}} \ell\left(f_{\mathbf{w} + \boldsymbol{\epsilon}}(\mathbf{x}_i), y_i\right)$$

for some noise distribution  $\boldsymbol{\varepsilon} \sim \varphi$ . The goal is to obtain a solution  $\boldsymbol{w}^* = \arg\min_{\boldsymbol{w} \in \mathcal{W}} L(\boldsymbol{w})$  that is robust to such perturbation.

The integrand of the gradient estimator of the instance risk is defined as

$$g_{i}(\boldsymbol{w};\boldsymbol{\eta}) = \nabla_{\boldsymbol{w}} \ell \left( f_{\boldsymbol{w}+\boldsymbol{\epsilon}}(\boldsymbol{x}_{i}), y_{i} \right)$$

$$= \frac{\partial f_{\boldsymbol{w}+\boldsymbol{\epsilon}}(\boldsymbol{x}_{i})}{\partial \boldsymbol{w}} \ell' \left( f_{\boldsymbol{w}+\boldsymbol{\epsilon}}(\boldsymbol{x}_{i}), y_{i} \right),$$

where it is an unbiased estimate of the instance risk such that

$$\mathbb{E} \mathbf{g}_i(\mathbf{w}) = \nabla R_i(\mathbf{w}).$$

The key challenge in analyzing the convergence of SGD in the ERM setting is dealing with the Jacobian  $\frac{\partial f_{w+\epsilon}(x_i)}{\partial w}$ . Even for simple toy models, analyzing the Jacobian without relying on strong assumptions is hard. In this work, we will assume that it is bounded by an instance-dependent constant.

#### C.1.2. PRELIMINARIES

We use the following assumptions:

## **Assumption 8.**

- (a) Let the mapping  $\hat{y} \mapsto \ell(\hat{y}, y)$  is convex and L-smooth for any  $y_i \ \forall i = 1, ..., n$ .
- (b) The Jacobian of the model with respect to its parameters for all i = 1, ..., n is bounded almost surely as

$$\left\| \frac{f_{\boldsymbol{w}+\boldsymbol{\epsilon}}(\boldsymbol{x}_i)}{\partial \boldsymbol{w}} \right\|_2 \le G_i$$

for all  $w \in \mathcal{W}$ .

(c) Interpolation holds on the solution set such that, for all  $\mathbf{w}_* \in \arg\min_{\mathbf{w} \in \mathcal{W}} L(\mathbf{w})$ , the loss minimized as

$$\ell\left(f_{\boldsymbol{w}_{*}+\boldsymbol{\epsilon}}\left(\boldsymbol{x}_{i}\right),y_{i}\right)=\ell'\left(f_{\boldsymbol{w}_{*}+\boldsymbol{\epsilon}}\left(\boldsymbol{x}_{i}\right),y_{i}\right)=0$$

for all  $(\mathbf{x}_i, y_i) \in \mathcal{D}$ .

(a) holds for the squared loss, (c) basically assumes that the model is overparameterized and there exists a set of optimal weights that are robust with respect to perturbation. The has recently gained popularity as it qualitatively explains some of the empirical phenomenons of non-convex SGD (Vaswani et al., 2019; Gower et al., 2021a; Ma et al., 2018). (b) is a strong assumption but is commonly used to establish convergence guarantees of ERM (Gower et al., 2021a).

**Remark 12.** Under Assumption 8 (c), Assumption 7 holds with arbitrarily small  $\sigma_i^2$ ,  $\tau^2$ .

#### C.1.3. THEORETICAL ANALYSIS

**Proposition 9.** Let Assumption 8 hold. Then, Assumption 6 ( $A^{ITP}$ ) holds.

Proof.

$$\mathbb{E}\|\mathbf{g}_{i}(\mathbf{w}) - \mathbf{g}_{i}(\mathbf{w}_{*})\|_{2}^{2}$$

$$= \mathbb{E}\left\|\frac{\partial f_{\mathbf{w}+\boldsymbol{\epsilon}}(\mathbf{x}_{i})}{\partial \mathbf{w}} \ell'(f_{\mathbf{w}+\boldsymbol{\epsilon}}(\mathbf{x}_{i}), y_{i}) - \frac{\partial f_{\mathbf{w}+\boldsymbol{\epsilon}}(\mathbf{x}_{i})}{\partial \mathbf{w}} \ell'(f_{\mathbf{w}_{*}+\boldsymbol{\epsilon}}(\mathbf{x}_{i}), y_{i})\right\|_{2}^{2},$$

from the interpolation assumption (Assumption 8 (c)),

$$= \mathbb{E} \left\| \frac{\partial f_{\boldsymbol{w}+\boldsymbol{\epsilon}}(\boldsymbol{x}_{i})}{\partial \boldsymbol{w}} \ell' (f_{\boldsymbol{w}+\boldsymbol{\epsilon}}(\boldsymbol{x}_{i}), y_{i}) \right\|_{2}^{2}$$

$$\leq \mathbb{E} \left\| \frac{\partial f_{\boldsymbol{w}+\boldsymbol{\epsilon}}(\boldsymbol{x}_{i})}{\partial \boldsymbol{w}} \right\|_{2}^{2} \left| \ell' (f_{\boldsymbol{w}+\boldsymbol{\epsilon}}(\boldsymbol{x}_{i}), y_{i}) \right|^{2},$$

applying Assumption 8 (b),

$$\leq G_i^2 \mathbb{E} \left| \ell' \left( f_{\boldsymbol{w} + \boldsymbol{\epsilon}} \left( \boldsymbol{x}_i \right), y_i \right) \right|^2.$$

and then the interpolation assumption (Assumption 8 (c)),

$$=G_{i}^{2}\mathbb{E}\left|\ell'\left(f_{\boldsymbol{w}+\boldsymbol{\epsilon}}\left(\boldsymbol{x}_{i}\right),y_{i}\right)-\ell'\left(f_{\boldsymbol{w}_{*}+\boldsymbol{\epsilon}}\left(\boldsymbol{x}_{i}\right),y_{i}\right)\right|.$$

From Assumption 8 (a),

$$\leq 2LG_{i}^{2}\mathbb{E}\left(\ell\left(f_{\boldsymbol{w}+\boldsymbol{\epsilon}}\left(\boldsymbol{x}_{i}\right),y_{i}\right)-\ell\left(f_{\boldsymbol{w}_{*}+\boldsymbol{\epsilon}}\left(\boldsymbol{x}_{i}\right),y_{i}\right)\right.\\\left.\left.-\left\langle \ell'\left(f_{\boldsymbol{w}_{*}+\boldsymbol{\epsilon}}\left(\boldsymbol{x}_{i}\right),y_{i}\right),f_{\boldsymbol{w}+\boldsymbol{\epsilon}}\left(\boldsymbol{x}_{i}\right)-f_{\boldsymbol{w}_{*}+\boldsymbol{\epsilon}}\left(\boldsymbol{x}_{i}\right)\right\rangle\right)$$

and interpolation (Assumption 8 (c)),

$$= 2LG_i^2 \left( \mathbb{E}\ell \left( f_{\boldsymbol{w}+\boldsymbol{\epsilon}} \left( \boldsymbol{x}_i \right), y_i \right) - \mathbb{E}\ell \left( f_{\boldsymbol{w}_*+\boldsymbol{\epsilon}} \left( \boldsymbol{x}_i \right), y_i \right) \right)$$
  
=  $2LG_i^2 \left( R_i \left( \boldsymbol{w} \right) - R_i \left( \boldsymbol{w}_* \right) \right).$ 

**Proposition 10.** Let Assumption 8 hold. Then, Assumption 5 holds.

Proof.

$$\frac{1}{n} \sum_{i=1}^{n} \left\| \nabla R_i(\boldsymbol{w}) - \nabla R_i(\boldsymbol{w}_*) \right\|_2^2 
= \frac{1}{n} \sum_{i=1}^{n} \left\| \mathbb{E} \boldsymbol{g}_i(\boldsymbol{w}) - \mathbb{E} \boldsymbol{g}_i(\boldsymbol{w}_*) \right\|_2^2,$$

and from Jensen's inequality,

$$\leq \frac{1}{n} \sum_{i=1}^{n} \mathbb{E} \| \boldsymbol{g}_{i} \left( \boldsymbol{w} \right) - \boldsymbol{g}_{i} \left( \boldsymbol{w}_{*} \right) \|_{2}^{2}.$$

We can now reuse Proposition 9 as

$$\leq \frac{2}{n} \sum_{i=1}^{n} 2LG_i^2 \left( R_i \left( \boldsymbol{w} \right) - R_i \left( \boldsymbol{w}_* \right) \right)$$

and taking  $G_{\text{max}} > G_i$  for all i = 1, ..., n as

$$\leq 2LG_{\max}^2 \frac{1}{n} \sum_{i=1}^n (R_i(\boldsymbol{w}) - R_i(\boldsymbol{w}_*))$$
  
=  $2LG_{\max}^2 (L(\boldsymbol{w}) - L(\boldsymbol{w}_*)).$ 

#### C.2. Reparameterization Gradient

#### C.2.1. DESCRIPTION

The reparameterization gradient estimator (Kingma & Welling, 2014; Rezende et al., 2014; Titsias & Lázaro-Gredilla, 2014) is a gradient estimator for problems of the form of

$$f_i(\boldsymbol{w}) = \mathbb{E}_{\boldsymbol{z} \sim q_{\boldsymbol{w}}} \ell_i(\boldsymbol{z}),$$

where  $\ell$ :  $\mathbb{R}^{d_z \to \mathbb{R}}$  is some integrand, such that the derivative is taken with respect to the parameters of the distribution  $q_w$  we are integrating over. It was independently proposed by Kingma & Welling (2014); Rezende et al. (2014) in the context of variational expectation maximization of deep latent variable models (a setup commonly known as variational autoencoders) and by Titsias & Lázaro-Gredilla (2014) for variational inference of Bayesian models.

Consider the case where the generative process of  $q_w$  can be represented as

$$\boldsymbol{z} \sim q_{\boldsymbol{w}} \quad \Leftrightarrow \quad \boldsymbol{z} \stackrel{d}{=} \mathcal{T}_{\boldsymbol{w}}(\boldsymbol{u}); \quad \boldsymbol{u} \sim \varphi,$$

where  $\stackrel{d}{=}$  is equivalence in distribution,  $\varphi$  is some *base distribution* independent of  $\boldsymbol{w}$ , and  $\mathcal{T}_{\boldsymbol{w}}$  is a *reparameterization function* measurable with respect to  $\varphi$  and differentiable with respect to all  $\boldsymbol{w} \in \mathcal{W}$ . Then, the reparameterization gradient is given by the integrand

$$g_i(\mathbf{w}; \mathbf{u}) = \nabla_{\mathbf{w}} \ell_i(\mathcal{T}_{\mathbf{w}}(\mathbf{u})),$$

which is unbiased, and often results in lower variance (Kucukelbir et al., 2017; Xu et al., 2019) compared to alternatives such as the score gradient estimator. (See Mohamed et al. (2020) for an overview of such estimators.)

The reparameterization gradient is primarily used to solve problems in the form of

$$\underset{\boldsymbol{w} \in \mathcal{W}}{\text{minimize}} \quad F(\boldsymbol{w}) = \sum_{i=1}^{n} f_i(\boldsymbol{w}) + h(\boldsymbol{w}) \\
= \mathbb{E}_{\boldsymbol{z} \sim q_{\boldsymbol{w}}} \ell_i(\boldsymbol{z}) + h(\boldsymbol{w}),$$

where h is some convex regularization term.

Previously, Domke (2019, Theorem 6) established a bound on the gradient variance of the reparameterization gradient (Kingma & Welling, 2014; Rezende et al., 2014; Titsias & Lázaro-Gredilla, 2014) under the doubly stochastic setting. This bound also incorporates more advanced subsampling strategies such as importance sampling Gower et al. (2019); Gorbunov et al. (2020); Csiba & Richtárik (2018); Needell et al. (2016); Needell & Ward (2017). However, he did not extend the analysis to a complexity analysis of SGD and left out the effect of correlation between components.

#### C.2.2. Preliminaries

The properties of the reparameterization gradient for when  $q_{\boldsymbol{w}}$  is in the location-scale family were studied by Domke (2019).

**Assumption 9.** We assume the variational family

$$Q \triangleq \{q_{\boldsymbol{w}} \mid \boldsymbol{w} \in \mathcal{W}\}$$

satisfies the following:

- (a)  $\mathcal{Q}$  is part of the location-scale family such that  $\mathcal{F}_{w}(u) = Cu + m$ .
- (b) The scale matrix is positive definite such that C > 0.
- (c)  $\boldsymbol{u}=(u_1,\ldots,u_{d_z})$  constitute of *i.i.d.* components, where each component is standardized, symmetric, and finite kurtosis such that  $\mathbb{E}u_i=0$ ,  $\mathbb{E}u_i^2=1$ ,  $\mathbb{E}u_i^3=0$ , and  $\mathbb{E}u_i^4=k_{\varphi}$ , where  $k_{\varphi}$  is the kurtosis.

Under these conditions, Domke (2019) proves the following:

**Lemma 13** (Domke, 2019; Theorem 3). Let Assumption 9 hold and  $\ell_i$  be  $L_i$ -smooth. Then, the squared norm of the reparameterization gradient is bounded:

$$\mathbb{E}\|\mathbf{g}_{i}(\mathbf{w})\|_{2}^{2} \leq (d+1)\|\mathbf{m} - \bar{\mathbf{z}}_{i}\|_{2}^{2} + (d+k_{\varphi})\|\mathbf{C}\|_{F}^{2}$$

for all  $\mathbf{w} = (\mathbf{m}, \mathbf{C}) \in \mathcal{W}$  and all stationary points of  $\ell_i$  denoted with  $\bar{\mathbf{z}}_i$ .

Similarly, Kim et al. (2023) establish the QES condition as part of Lemma 3 (Kim et al., 2023). We refine this into statement we need:

**Lemma 14.** Let Assumption 9 hold and  $\ell_i$  be  $L_i$ -smooth. Then, the squared norm of the reparameterization gradient is bounded:

$$\mathbb{E}\left\|\boldsymbol{g}_{i}\left(\boldsymbol{w}\right)-\boldsymbol{g}_{i}\left(\boldsymbol{w}'\right)\right\|_{2}^{2} \leq L_{i}^{2}\left(d+k_{\varphi}\right)\left\|\boldsymbol{w}-\bar{\boldsymbol{w}}_{i}\right\|_{2}^{2}$$

for all  $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$ .

Proof.

$$\begin{split} & \mathbb{E} \left\| \boldsymbol{g}_{i}\left(\boldsymbol{w}\right) - \boldsymbol{g}_{i}\left(\boldsymbol{w}'\right) \right\| \\ & = \mathbb{E} \left\| \nabla_{\boldsymbol{w}} \ell_{i}\left(\mathcal{T}_{\boldsymbol{w}}\left(\boldsymbol{u}\right)\right) - \nabla_{\boldsymbol{w}} \ell_{i}\left(\mathcal{T}_{\boldsymbol{w}'}\left(\boldsymbol{u}\right)\right) \right\|_{2}^{2} \\ & = \mathbb{E} \left\| \frac{\partial \mathcal{T}_{\boldsymbol{w}}\left(\boldsymbol{u}\right)}{\partial \boldsymbol{w}} \nabla \ell_{i}\left(\mathcal{T}_{\boldsymbol{w}}\left(\boldsymbol{u}\right)\right) - \frac{\partial \mathcal{T}_{\boldsymbol{w}'}\left(\boldsymbol{u}\right)}{\partial \boldsymbol{w}'} \nabla \ell_{i}\left(\mathcal{T}_{\boldsymbol{w}'}\left(\boldsymbol{z}\right)\right) \right\|_{2}^{2} \\ & = \mathbb{E} \left(\nabla \ell_{i}\left(\mathcal{T}_{\boldsymbol{w}}\left(\boldsymbol{u}\right)\right) - \nabla \ell_{i}\left(\mathcal{T}_{\boldsymbol{w}'}\left(\boldsymbol{u}\right)\right)\right)^{\mathsf{T}} \left(\frac{\partial \mathcal{T}_{\boldsymbol{w}}\left(\boldsymbol{u}\right)}{\partial \boldsymbol{w}}\right)^{\mathsf{T}} \frac{\partial \mathcal{T}_{\boldsymbol{w}'}\left(\boldsymbol{u}\right)}{\partial \boldsymbol{w}'} \\ & \times \left(\nabla \ell_{i}\left(\mathcal{T}_{\boldsymbol{w}}\left(\boldsymbol{u}\right)\right) - \nabla \ell_{i}\left(\mathcal{T}_{\boldsymbol{w}'}\left(\boldsymbol{u}\right)\right)\right). \end{split}$$

As shown by Kim et al. (2023, Lemma 6), the squared Jacobian is an identity matrix scaled with a scalar-valued function independent of  $\boldsymbol{w}$ ,  $J_{\mathcal{T}}(\boldsymbol{u}) = \|\boldsymbol{u}\|_{2}^{2} + 1$ , such that

$$= \mathbb{E}J_{\mathcal{T}}(\boldsymbol{u}) \|\nabla \ell_{i}(\mathcal{T}_{\boldsymbol{w}}(\boldsymbol{u})) - \nabla \ell_{i}(\mathcal{T}_{\boldsymbol{w}'}(\boldsymbol{u}))\|_{2}^{2},$$
 applying the  $L_{i}$ -smoothness of  $\ell_{i}$ , 
$$= L_{i}^{2} \mathbb{E}J_{\mathcal{T}}(\boldsymbol{u}) \|\mathcal{T}_{\boldsymbol{w}}(\boldsymbol{u}) - \mathcal{T}_{\boldsymbol{w}'}(\boldsymbol{u})\|_{2}^{2},$$
 and Kim et al. (2023, Corollary 2) show that, 
$$\leq L_{i}^{2} (d + k_{\varphi}) \|\boldsymbol{w} - \boldsymbol{w}'\|_{2}^{2}.$$

Lastly, the properties of  $\ell_i$  are known to transfer to the expectation  $f_i$  as follows:

**Lemma 15.** Let Assumption 9 hold. Then we have the following:

- (i) Let  $\ell_i$  be  $L_i$  smooth. Then,  $f_i$  is also  $L_i$ -smooth
- (ii) Let  $\ell_i$  be convex. Then,  $f_i$  and F are also convex.
- (iii) Let  $\ell_i$  be  $\mu$ -strongly convex. Then,  $f_i$  and F are also  $\mu$ -strongly convex.

*Proof.* (i) is proven by Domke (2020, Theorem 1), while a more general result is provided by Kim et al. (2023, Theorem 1); (ii) and (iii) are proven by Domke (2020, Theorem 9) and follow from the fact that h is convex.

#### C.2.3. THEORETICAL ANALYSIS

We now conclude that the reparameterization gradient fits the framework of this work:

**Proposition 11.** Let Assumption 9 hold and  $\ell_i$  be convex and  $L_i$ -smooth. Then, Assumption 5 holds.

*Proof.* The result follows from combining Lemma 15 and Lemma 10. □

From Lemma 13, we satisfy 7.

**Proposition 12.** Let Assumption 9 hold,  $\ell_i$  be  $L_i$ -smooth, the solutions  $\mathbf{w}_* \in \arg\min_{\mathbf{w} \in \mathcal{W}} F(\mathbf{w})$  and the stationary points of  $\ell_i$ ,  $\bar{\mathbf{z}}$ , be bounded such that  $\|\mathbf{w}_*\|_2 < \infty$  and  $\|\bar{\mathbf{z}}\|_2 < \infty$ . Then, Assumption 7 holds.

*Proof.* Lemma 13 implies that, as long as the  $\mathbf{w}_*$  and  $\bar{\mathbf{z}}$  are bounded, we satisfy the component gradient estimator part of Assumption 7, where the constant is given as

$$\sigma_i^2 = L_i^2 (d+1) \|\boldsymbol{m}_* - \bar{\boldsymbol{z}}_i\|_2^2 + L_i^2 (d+k_\varphi) \|\boldsymbol{C}_*\|_F^2,$$
 where  $\boldsymbol{w}_* = (\boldsymbol{m}_*, \boldsymbol{C}_*)$ .

From Lemma 14, we can conclude that the reparameterization gradient satisfies Assumption 6:

**Proposition 13.** Let Assumption 9 hold and  $\ell_i$  be  $L_i$ -smooth and  $\mu$ -strongly convex. Then, Assumption 6 (A<sup>CVX</sup>) and Assumption 6 (B) hold.

*Proof.* Notice the following:

- 1. Assumption 4 always holds for  $\rho = 1$ .
- 2. From the stated conditions, Lemma 15 establishes that both  $f_i$  and F are  $\mu$ -strongly convex.
- 3. μ-strong convexity of f and F implies that both are μ-QFG (Karimi et al., 2016, Appendix A).
- The reparameterization gradient satisfies the QES condition by Lemma 14.

Item 1, 2 and 3 combined imply the ES condition by Proposition 8, which immediately implies the ER condition with the same constant. Therefore, we satisfy both Assumption 6 (A<sup>CVX</sup>), Assumption 6 (B) where the ER constant  $\mathcal{L}_i$  is given as

$$\mathcal{L}_i = \frac{L_i^2}{\mu} \left( d + k_{\varphi} \right).$$