

# SGDA WITH SHUFFLING: FASTER CONVERGENCE FOR NONCONVEX-PŁ MINIMAX OPTIMIZATION

Hanseul Cho, Chulhee Yun

Kim Jaechul Graduate School of AI, KAIST

{jhs4015, chulhee.yun}@kaist.ac.kr

## ABSTRACT

Stochastic gradient descent-ascent (SGDA) is one of the main workhorses for solving finite-sum minimax optimization problems. Most practical implementations of SGDA randomly reshuffle components and sequentially use them (*i.e.*, without-replacement sampling); however, there are few theoretical results on this approach for minimax algorithms, especially outside the easier-to-analyze (strongly-)monotone setups. To narrow this gap, we study the convergence bounds of SGDA with random reshuffling (**SGDA-RR**) for smooth nonconvex-nonconcave objectives with Polyak-Łojasiewicz (PŁ) geometry. We analyze both simultaneous and alternating SGDA-RR for *nonconvex-PŁ* and *primal-PŁ-PŁ* objectives, and obtain convergence rates faster than with-replacement SGDA. Our rates extend to mini-batch SGDA-RR, recovering known rates for full-batch gradient descent-ascent (GDA). Lastly, we present a comprehensive lower bound for GDA with an arbitrary step-size ratio, which matches the full-batch upper bound for the *primal-PŁ-PŁ* case.

## 1 INTRODUCTION

A finite-sum minimax optimization problem aims to solve the following:

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}; \mathbf{y}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}; \mathbf{y}), \quad (1)$$

where  $f_i$  denotes the  $i$ -th component function. In plain language, we want to minimize the average of  $n$  component functions for  $\mathbf{x}$ , while maximizing it for  $\mathbf{y}$  given  $\mathbf{x}$ . There are many important areas in modern machine learning that fall within the minimax problem, including generative adversarial networks (GANs) (Goodfellow et al., 2020), adversarial attack and robust optimization (Madry et al., 2018; Sinha et al., 2018), multi-agent reinforcement learning (MARL) (Li et al., 2019), AUC maximization (Ying et al., 2016; Liu et al., 2020; Yuan et al., 2021), and many more. In most cases, the objective  $f$  is usually nonconvex-nonconcave, *i.e.*, neither convex in  $\mathbf{x}$  nor concave in  $\mathbf{y}$ . Since general nonconvex-nonconcave problems are known to be intractable, we would like to tackle the problems with some additional structures, such as smoothness and Polyak-Łojasiewicz (PŁ) condition(s). We elaborate the detailed settings for our analysis, *nonconvex-PŁ* and *primal-PŁ-PŁ* (or, *PŁ(Φ)-PŁ*), in Section 2.

One of the simplest and most popular algorithms to solve the problem (1) would be *stochastic gradient descent-ascent* (**SGDA**). This naturally extends the idea of stochastic gradient descent (SGD) used for minimization problems. Given an initial iterate  $(\mathbf{x}_0; \mathbf{y}_0)$ , at time  $t \in \mathbb{N}$ , SGDA (randomly) chooses an index  $i(t) \in \{1, \dots, n\}$  and accesses the  $i(t)$ -th component to perform a pair of updates

$$\begin{cases} \mathbf{x}_t = \mathbf{x}_{t-1} - \alpha \nabla_1 f_{i(t)}(\mathbf{x}_{t-1}; \mathbf{y}_{t-1}), \\ \mathbf{y}_t = \mathbf{y}_{t-1} + \beta \nabla_2 f_{i(t)}(\mathbf{x}'_t; \mathbf{y}_{t-1}), \end{cases} \quad \text{where } \mathbf{x}'_t = \begin{cases} \mathbf{x}_{t-1}, & (\text{simSGDA}), \text{ or} \\ \mathbf{x}_t, & (\text{altSGDA}). \end{cases}$$

Here,  $\alpha > 0$  and  $\beta > 0$  are the step sizes and  $\nabla_j$  denotes the gradient with respect to  $j$ -th argument for  $f_{i(t)}$  ( $j = 1, 2$ ). As shown in the update equations above, there are two widely used versions of SGDA: *simultaneous SGDA* (**simSGDA**), and *alternating SGDA* (**altSGDA**).

In such stochastic gradient methods, there are two main categories of sampling schemes for the component indices  $i(t)$ . One way is to sample  $i(t)$  independently (in time) and uniformly at random

from  $\{1, \dots, n\}$ , which is called *with-replacement sampling*. This scheme is widely adopted in theory papers because it makes analysis of stochastic methods amenable: the noisy gradients  $\nabla f_{i(t)}$  are independent over time  $t$  and are unbiased estimators of the full-batch gradient  $\nabla f$ . In contrast, the vast majority of practical implementations employ *without-replacement sampling*, indicating a huge theory-practice gap. In without-replacement sampling, we sample each index precisely once at each epoch. Perhaps the most popular of such schemes is *random reshuffling (RR)*, which uniformly randomly shuffles the order of indices at the beginning of every epoch. Unfortunately, it is well-known that without-replacement methods are much more difficult to analyze theoretically, largely because the sampled indices in each epoch are no longer independent of each other.

Interestingly, for minimization problems, several recent works overcome this obstacle and show that SGD using without-replacement sampling leads to faster convergence, given that the number of epochs is large enough (Nagaraj et al., 2019; Ahn et al., 2020; Mishchenko et al., 2020; Rajput et al., 2020; Nguyen et al., 2021; Yun et al., 2021; 2022). On the other hand, for minimax problems like (1), the majority of the studies still assume with-replacement sampling and/or rely on independent unbiased gradient oracles (Nouiehed et al., 2019; Guo et al., 2020; Lin et al., 2020; Yan et al., 2020; Yang et al., 2020; Loizou et al., 2021; Beznosikov et al., 2022). There are very few results on minimax algorithms using without-replacement sampling; even most of the existing ones take advantage of (strong-)convexity (in  $\mathbf{x}$ ) and/or (strong-)concavity (in  $\mathbf{y}$ ) (Das et al., 2022; Maheshwari et al., 2022; Yu et al., 2022). Detailed comparative analysis of these works is conducted in Section 4.

Putting all these issues into consideration, our main question is the following.

*Does SGDA using without-replacement component sampling provably converge fast, even on smooth nonconvex-nonconcave objective  $f$  with PL structures?*

## 1.1 SUMMARY OF OUR CONTRIBUTIONS

To answer the question, we analyze the convergence of SGDA with random reshuffling (**SGDA-RR**, Algorithm 1). We analyze both the simultaneous and alternating versions of SGDA-RR and prove convergence theorems for the following two regimes. Here we denote the step size ratio as  $r = \beta/\alpha$ .

- When  $-f(\mathbf{x}; \mathbf{y})$  satisfies  $\mu_2$ -PL condition in  $\mathbf{y}$  (**nonconvex-PL**) and component function  $f_i$ 's are  $L$ -smooth, we prove that SGDA-RR with  $r \gtrsim (L/\mu_2)^2$  converges to  $\varepsilon$ -stationarity in expectation after  $\mathcal{O}(nrL\varepsilon^{-2} + \sqrt{nr}^{1.5}L\varepsilon^{-3})$  gradient evaluations (Theorem 1).
- Further assuming  $\mu_1$ -PL condition on  $\Phi(\cdot) := \max_{\mathbf{y}} f(\cdot; \mathbf{y})$  (**primal-PL-PL**, or **PL( $\Phi$ )-PL**), we prove that SGDA-RR with  $r \gtrsim (L/\mu_2)^2$  converges within  $\varepsilon$ -accuracy in expectation after  $\tilde{\mathcal{O}}\left(\frac{nLr}{\mu_1} \log(\varepsilon^{-1}) + \sqrt{n}L\left(\frac{r}{\mu_1}\right)^{1.5}\varepsilon^{-1}\right)$  gradient evaluations (Theorem 2).

As will be discussed in Section 4, the rates shown above are *faster* than existing results on with-replacement SGDA. In fact, Theorems 1 & 2 are special cases ( $b = 1$ ) of our extended theorems (Theorems 4 & 5 in Appendix A) that analyze *mini-batch* SGDA-RR of batch size  $b$ ; by setting  $b = n$ , we also recover known convergence rates for full-batch gradient descent ascent (GDA). Hence, our analysis covers the entire spectrum between vanilla SGDA-RR ( $b = 1$ ) and GDA ( $b = n$ ).

- Additionally, we provide complexity lower bounds for solving strongly-convex-strongly-concave (SC-SC) minimax problems using full-batch simultaneous GDA with an arbitrarily fixed step size ratio  $r = \beta/\alpha$ . Perhaps surprisingly, we find that the lower bound for SC-SC functions matches the convergence upper bound for a much larger class of primal-PL-PL functions when the step size ratio satisfies  $r \gtrsim L^2/\mu_2^2$  (Theorem 3).

## 2 PROBLEM SETUP

### 2.1 NOTATION

In our problem (1), the domain of every  $f_i$  is  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ , where  $\mathcal{X} = \mathbb{R}^{d_x}$ ,  $\mathcal{Y} = \mathbb{R}^{d_y}$ , and  $\mathcal{Z} = \mathbb{R}^d$ : we concern unconstrained problems for simplicity. We denote the Euclidean norm and the standard inner product by  $\|\cdot\|$  and  $\langle \cdot, \cdot \rangle$ , respectively. We often use an abbreviated notation  $\mathbf{z} = (\mathbf{x}; \mathbf{y}) \in \mathcal{Z}$  for  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{y} \in \mathcal{Y}$ . Even when  $\mathbf{z}$  or  $(\mathbf{x}; \mathbf{y})$  is followed by superscripts and/or subscripts, we use the symbols interchangeably; e.g.,  $\mathbf{z}_i^k = (\mathbf{x}_i^k; \mathbf{y}_i^k)$ . Note that we split the arguments  $\mathbf{x}$  (for minimization) and  $\mathbf{y}$  (for maximization) by a semicolon (;). We use  $\nabla_1$  and  $\nabla_2$  to denote the

**Algorithm 1** simSGDA/altSGDA-RR

---

```

1: Given: The number of components  $n$ ; the number of epochs  $K$ ; step sizes  $\alpha, \beta > 0$ 
2: Initialize:  $(\mathbf{x}_0^1; \mathbf{y}_0^1) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ 
3: for  $k \in [K]$  do
4:   Sample  $\sigma_k \sim \text{Unif}(\mathbb{S}_n)$  ▷ RR: uniformly randomly shuffle the indices every epoch
5:   for  $i \in [n]$  do
6:      $\mathbf{x}_i^k = \mathbf{x}_{i-1}^k - \alpha \nabla_1 f_{\sigma_k(i)}(\mathbf{x}_{i-1}^k; \mathbf{y}_{i-1}^k)$ 
7:     if simSGDA-RR then
8:        $\mathbf{y}_i^k = \mathbf{y}_{i-1}^k + \beta \nabla_2 f_{\sigma_k(i)}(\mathbf{x}_{i-1}^k; \mathbf{y}_{i-1}^k)$  ▷ simultaneous update:  $\mathbf{x}$  &  $\mathbf{y}$ 
9:     else if altSGDA-RR then
10:       $\mathbf{y}_i^k = \mathbf{y}_{i-1}^k + \beta \nabla_2 f_{\sigma_k(i)}(\mathbf{x}_i^k; \mathbf{y}_{i-1}^k)$  ▷ alternating update:  $\mathbf{x} \rightarrow \mathbf{y}$ 
11:    $(\mathbf{x}_0^{k+1}; \mathbf{y}_0^{k+1}) = (\mathbf{x}_n^k; \mathbf{y}_n^k)$ 

```

---

gradients with respect to first and second arguments, respectively. Accordingly, we can write the full gradient as, e.g.,  $\nabla g = [\nabla_1 g^\top; \nabla_2 g^\top]^\top$ . For a positive integer  $N$ , we denote  $[N] := \{1, \dots, N\}$ . Let the set  $\mathbb{S}_N$  be a symmetric group of degree  $N$ . That is, each *permutation*  $\sigma \in \mathbb{S}_N$  is a bijection from  $[N]$  to itself, or equivalently, a re-arrangement of  $[N]$ . Lastly, we use the usual  $\mathcal{O}/\Omega/\Theta$  notation for bounds, where  $\tilde{\mathcal{O}}/\tilde{\Omega}/\tilde{\Theta}$  are used for hiding some logarithmic factors, respectively.

## 2.2 ALGORITHMS: SIMSGDA-RR &amp; ALTSFDA-RR

As we explained in Section 1, we consider simSGDA and altSGDA combined with RR, a without-replacement sampling scheme. We call them **simSGDA-RR** and **altSGDA-RR**, respectively. We present a detailed description of the methods in Algorithm 1. For completeness, we also provide an extended version that uses mini-batches of size  $\geq 1$  (Algorithm 2) in Appendix A. For comparison, we call the SGDA algorithms using *with-replacement* sampling by just **simSGDA** and **altSGDA**.

The quantities  $\alpha, \beta > 0$  are step sizes associated with  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. We use two separate symbols  $\alpha$  and  $\beta$  to allow the two step sizes to be different. Such algorithms are sometimes called *two-time-scale* algorithms, in a broader sense, and they are adopted in nonconvex minimax optimization problems (Heusel et al., 2017; Lin et al., 2020; Yang et al., 2020). In fact, a recent result (Li et al., 2022) shows that having  $\alpha \neq \beta$  is sometimes *necessary* for convergence.

## 2.3 ASSUMPTIONS AND DEFINITIONS

To define the function classes that we are interested in solving, we introduce a few assumptions.

**Assumption 1** (Component smoothness). *Every  $i$ -th component  $f_i : \mathcal{Z} \rightarrow \mathbb{R}$  is  $L$ -smooth, i.e.,  $f_i$  is differentiable and  $\nabla f_i$  is  $L$ -Lipschitz continuous:  $\|\nabla f_i(\mathbf{z}) - \nabla f_i(\bar{\mathbf{z}})\| \leq L \|\mathbf{z} - \bar{\mathbf{z}}\|$ . As a result,  $f_i(\bar{\mathbf{z}}) - f_i(\mathbf{z}) \leq \langle \nabla f_i(\mathbf{z}), \bar{\mathbf{z}} - \mathbf{z} \rangle + \frac{L}{2} \|\bar{\mathbf{z}} - \mathbf{z}\|^2$  ( $\forall \mathbf{z}, \bar{\mathbf{z}}$ ) and the average  $f$  of  $f_i$ 's is also  $L$ -smooth.<sup>1</sup>*

**Assumption 2** (Component gradient variance). *There exist constants  $A, B \geq 0$  such that, for any  $\mathbf{z} = (\mathbf{x}; \mathbf{y}) \in \mathcal{Z}$  and  $j \in \{1, 2\}$ , we have  $\frac{1}{n} \sum_{i=1}^n \|\nabla_j f_i(\mathbf{z}) - \nabla_j f(\mathbf{z})\|^2 \leq A \|\nabla_j f(\mathbf{z})\|^2 + B$ .*

**Assumption 3.** *For a function  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ , the **primal function**  $\Phi : \mathcal{X} \rightarrow \mathbb{R}$  is well-defined as  $\Phi(\mathbf{x}) := \max_{\mathbf{y}' \in \mathcal{Y}} f(\mathbf{x}; \mathbf{y}')$ . For each  $\mathbf{x} \in \mathcal{X}$ , the set  $\mathcal{Y}_x^* := \arg \max_{\mathbf{y}' \in \mathcal{Y}} f(\mathbf{x}; \mathbf{y}')$  is non-empty and closed. Moreover, we assume  $\Phi(\mathbf{x})$  is bounded below by  $\Phi^* = \inf_{\mathbf{x}' \in \mathcal{X}} \Phi(\mathbf{x}') > -\infty$ .*

Note that Assumption 2 controls the discrepancy between the objective function  $f$  and its components  $f_i$ 's; it is similar to Assumption 2 of Nguyen et al. (2021), adapted to minimax problems. Letting  $A = 0$  recovers a common assumption of the uniformly bounded variance of component gradients; thus, our assumption is a relaxation. Also, note that  $A = B = 0$  when  $n = 1$ .

We now add an additional structure to our objective function, which is called Polyak-Łojasiewicz (PŁ) condition. A function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is said to be  $\mu$ -PŁ if it has a minimum value  $g^*$  and satisfies

$$\|\nabla g(\mathbf{t})\|^2 \geq 2\mu(g(\mathbf{t}) - g^*). \quad (\forall \mathbf{t} \in \mathbb{R}^d) \quad (\mu\text{-PŁ})$$

<sup>1</sup>As we noted, Assumption 1 directly implies the *average smoothness* which is a common requirement in the analysis with unbiased gradient oracles. Nevertheless, we claim that Assumption 1 is not more crucial than without-replacement sampling to obtain faster convergence rates: see Appendix F for details and proofs.

Readers could find several studies and applications that the condition involves, in the papers by Karimi et al. (2016); Nouiehed et al. (2019); Yang et al. (2020); Liu et al. (2020), and more. Note that every  $\mu$ -strongly convex<sup>2</sup> function satisfies  $\mu$ -PŁ condition, whereas a PŁ function does not need to be convex. Hence,  $\mu$ -PŁ is a strict generalization of  $\mu$ -strong convexity. In addition, every stationary point of a PŁ function is a global optimum, which is a benign property for optimization.

We are interested in the case where our objective function  $f(\mathbf{x}; \mathbf{y})$  has such a structure in terms of  $\mathbf{y}$  (Assumption 4). Sometimes, we further assume the primal function  $\Phi$  is also PŁ (Assumption 5). We emphasize that we do not necessarily assume the PŁ conditions for the individual  $f_i$ 's.

**Assumption 4 (y-side PŁ).** For each (fixed)  $\mathbf{x} \in \mathcal{X}$ ,  $-f(\mathbf{x}; \cdot)$  is  $\mu_2$ -PŁ, i.e., for every  $(\mathbf{x}; \mathbf{y}) \in \mathcal{Z}$ ,  $\|\nabla_{\mathbf{y}} f(\mathbf{x}; \mathbf{y})\|^2 \geq 2\mu_2(\Phi(\mathbf{x}) - f(\mathbf{x}; \mathbf{y}))$ , where  $\Phi$  is the primal function associated with  $f$ .

**Assumption 5 (Primal PŁ, or PŁ( $\Phi$ )).** The primal function  $\Phi(\cdot) = \max_{\mathbf{y}'} f(\mathbf{x}; \mathbf{y}')$  of  $f$  is  $\mu_1$ -PŁ, i.e., for every  $\mathbf{x} \in \mathcal{X}$ ,  $\|\nabla \Phi(\mathbf{x})\|^2 \geq 2\mu_1(\Phi(\mathbf{x}) - \Phi^*)$ , where  $\Phi^* = \min_{\mathbf{x}} \Phi(\mathbf{x})$  is well-defined.

We say the function  $f$  is **nonconvex-PŁ** when it satisfies Assumption 4. Since we do not assume any convexity/concavity, it is generally hard to reach global optima. Due to the  $\mathbf{y}$ -side PŁ condition, we can guarantee that the primal function  $\Phi$  is differentiable and even  $L_{\Phi}$ -smooth with  $L_{\Phi} \leq L + L^2/\mu_2$  (Proposition 9 in Appendix B). Since the problem (1) can be reformulated as the minimization problem of  $\Phi$  (when we can always find  $\mathbf{y}$  well that maximizes  $f(\mathbf{x}; \mathbf{y})$  given  $\mathbf{x}$ ), we could aim to find an approximate first-order stationary point of  $\Phi$ , by making the norm of the gradient of  $\Phi$  small.

On top of that, if  $f$  satisfies both Assumptions 4 and 5, the function is said to be **primal-PŁ-PŁ**, or **PŁ( $\Phi$ )-PŁ** for short.<sup>3</sup> In this case, we directly aim not only to decrease the primal function  $\Phi$  associated with the objective function  $f$  but also to increase the function value  $f(\mathbf{x}; \mathbf{y})$  in terms of  $\mathbf{y}$ . To evaluate how close we are to our goal, we define a *potential function*  $V_{\lambda}$  later in Section 3. When we attain  $V_{\lambda}(\mathbf{x}^*, \mathbf{y}^*) = 0$ , it implies that we arrive at a global minimax point:  $f(\mathbf{x}^*, \mathbf{y}^*) = \Phi(\mathbf{x}^*) = \Phi^*$ . The function  $V_{\lambda}$  enables us to develop a unified analysis for nonconvex-PŁ and PŁ( $\Phi$ )-PŁ objective functions; we discuss this in greater detail in Section 3.

### 3 MAIN RESULTS

Based on the assumptions stated in the previous section, we present the convergence results for both smooth nonconvex-PŁ objectives and smooth PŁ( $\Phi$ )-PŁ objectives. Before stating the main theorems, we first introduce the most important tool for our analyses: the *potential function*.

#### 3.1 POTENTIAL FUNCTION $V_{\lambda}$

For our convergence analyses, we utilize a function  $V_{\lambda} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  defined as

$$V_{\lambda}(\mathbf{x}; \mathbf{y}) := \lambda(\Phi(\mathbf{x}) - \Phi^*) + (\Phi(\mathbf{x}) - f(\mathbf{x}; \mathbf{y})), \quad (2)$$

where  $\lambda > 0$  is a constant. We borrow inspiration from Yang et al. (2020) and Das et al. (2022) to come up with this function, although the placement of  $\lambda$  of ours is different. In fact, the convergence to a neighborhood of a global minimax point (if it exists) implies the reduction of this function. For each  $\mathbf{x}$ , a non-negative term  $\Phi(\mathbf{x}) - f(\mathbf{x}; \mathbf{y})$  gets smaller as  $\mathbf{y}$  makes  $f(\mathbf{x}; \mathbf{y})$  larger. The term becomes zero when  $\mathbf{y} = \mathbf{y}^*(\mathbf{x})$  for some  $\mathbf{y}^*(\mathbf{x}) \in \mathcal{Y}_{\mathbf{x}}^*$ , since  $\Phi(\mathbf{x}) = f(\mathbf{x}; \mathbf{y}^*(\mathbf{x}))$ . Also, another non-negative term  $\Phi(\mathbf{x}) - \Phi^*$  gets smaller as  $\mathbf{x}$  makes  $\Phi(\mathbf{x})$  smaller. Thus, as  $(\mathbf{x}; \mathbf{y})$  approaches to a minimax optimal point,  $V_{\lambda}(\mathbf{x}; \mathbf{y})$  decreases to near zero. In general,  $V_{\lambda}$  is not guaranteed to attain exact zero, especially when the objective function  $f(\mathbf{x}; \mathbf{y})$  is nonconvex in  $\mathbf{x}$  (e.g.,  $f$  is nonconvex-PŁ). Nevertheless, the potential function is still useful for deriving our convergence results.

#### 3.2 MAIN THEOREMS: UPPER BOUNDS OF CONVERGENCE RATES

Now, we present our main results. We provide a detailed comparison of our theorems against existing results in Section 4. We present the full proof in Appendices C and D. We remark that both Theorems 1 and 2 are special cases (for mini-batch size  $b = 1$ ) of their *mini-batch* extensions: Theorems 4 and 5 in Appendix A.

<sup>2</sup>We say a function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex for some  $\mu > 0$  if it holds  $g(\mathbf{x}') \geq g(\mathbf{x}) + \langle \nabla g(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + (\mu/2) \|\mathbf{x}' - \mathbf{x}\|^2$  ( $\forall \mathbf{x}, \mathbf{x}'$ ); we say  $g$  is  $\mu$ -strongly concave if  $-g$  is  $\mu$ -strongly convex.

<sup>3</sup>The PŁ( $\Phi$ )-PŁ condition is much weaker than **two-sided PŁ** condition assuming “ $\mathbf{x}$ -side” PŁ condition: see Proposition 10. As pointed out by Guo et al. (2020), there exist a PŁ( $\Phi$ )-PŁ function  $g(\mathbf{x}; \mathbf{y})$  that is not  $\mathbf{x}$ -side  $\mu$ -PŁ for any  $\mu > 0$  but even strongly concave in  $\mathbf{x}$ .

**Theorem 1 (Nonconvex-PŁ).** Suppose that  $f$  satisfies Assumptions 1, 2, 3, and 4. Let  $\kappa_2 = L/\mu_2$ , where  $\mu_2$  is PŁ constant of  $-f(\mathbf{x}; \cdot)$  at all  $\mathbf{x}$ . Let  $\lambda = 4$ . Choose the step sizes  $\alpha$  and  $\beta$  such that

$$\beta = \min \left\{ \frac{1}{6L\sqrt{n(n+A)}}, \mathcal{O} \left( \left( \frac{V_\lambda(\mathbf{z}_0^1)}{Bn^2K} \right)^{\frac{1}{3}} \right) \right\} \quad \text{and} \quad \alpha = \frac{\beta}{r},$$

for some  $r \geq 14\kappa_2^2$ . Then, both simSGDA-RR and altSGDA-RR (Algorithm 1) satisfy

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla \Phi(\mathbf{x}_0^k)\|^2 \leq \mathcal{O} \left( \frac{rLV_\lambda(\mathbf{z}_0^1)}{K} \sqrt{1 + \frac{A}{n}} + r \left( \frac{L^2BV_\lambda(\mathbf{z}_0^1)^2}{nK^2} \right)^{1/3} \right).$$

**Upper bound on gradient complexity.** To achieve  $\varepsilon$ -stationarity of the primal function, i.e.,  $\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla \Phi(\mathbf{x}_0^k)\|^2 \leq \varepsilon^2$ , a sufficient number of gradient evaluations (denoted by  $T_\varepsilon = nK$ ) is

$$T_\varepsilon = \mathcal{O} \left( \frac{rLV_\lambda(\mathbf{z}_0^1)}{\varepsilon^2} \max \left\{ \sqrt{n^2 + nA}, \frac{\sqrt{rnB}}{\varepsilon} \right\} \right).$$

**Theorem 2 (PŁ(Φ)-PŁ).** Suppose that  $f$  satisfies Assumptions 1, 2, 3, 4, and 5. Let  $\kappa_1 = L/\mu_1$  and  $\kappa_2 = L/\mu_2$ , where  $\mu_1$  and  $\mu_2$  are PŁ constants of  $\Phi(\cdot)$  and  $-f(\mathbf{x}; \cdot)$  (at all  $\mathbf{x}$ ), respectively. Let  $\lambda = 4$ . Choose appropriate step sizes  $\alpha$  and  $\beta$  such that

$$\beta = \min \left\{ \frac{1}{6L\sqrt{n(n+A)}}, \tilde{\mathcal{O}} \left( \frac{\kappa_2^2}{\mu_1 n K} \right) \right\} \quad \text{and} \quad \alpha = \frac{\beta}{r},$$

for some  $r \geq 14\kappa_2^2$ . Then, both simSGDA-RR and altSGDA-RR (Algorithm 1) satisfy

$$\mathbb{E}[V_\lambda(\mathbf{z}_0^{K+1})] \leq \mathcal{O} \left( V_\lambda(\mathbf{z}_0^1) \cdot \exp \left( -\frac{K}{12\kappa_1 r \sqrt{1 + \frac{A}{n}}} \right) \right) + \tilde{\mathcal{O}} \left( \frac{\kappa_1^2 r^3 B}{\mu_1 n K^2} \right).$$

**Upper bound on gradient complexity.** To achieve  $\varepsilon^2$ -accuracy on expectation of  $V_\lambda(\mathbf{z}_n^K)$ , i.e.,  $\mathbb{E}[V_\lambda(\mathbf{z}_n^K)] \leq \varepsilon^2$ , a sufficient number of gradient evaluations (denoted by  $T'_\varepsilon = nK$ ) is

$$T'_\varepsilon = \max \left\{ \mathcal{O} \left( \kappa_1 r \sqrt{n^2 + nA} \cdot \log \left( \frac{V_\lambda(\mathbf{z}_0^1)}{\varepsilon} \right) \right), \tilde{\mathcal{O}} \left( \frac{\kappa_1 r^{3/2} \sqrt{nB}}{\varepsilon \mu_1} \right) \right\}.$$

**Remark on step size ratio.** In both theorems, we use the step sizes of ratio  $r = \beta/\alpha \gtrsim \kappa_2^2$ . It is common to use such a step size scheme  $r = \Theta(\kappa_2^2)$  to analyze two-time-scale (S)GDA for nonconvex minimax problems (Jin et al., 2020; Lin et al., 2020; Yang et al., 2020).

**Remark on the parameter  $\lambda$ .** In our convergence analyses, we arbitrarily choose  $\lambda = 4$  which makes the numerical calculations easier. The value of  $\lambda > 0$  does not matter for the equivalence between the equation  $V_\lambda(\mathbf{x}^*; \mathbf{y}^*) = 0$  and global minimax condition (Proposition 11 in Appendix B). Also, the choice of  $\lambda$  in both theorems can be arbitrary as long as  $\lambda > 1$ ; our logic does not fall apart if other appropriate step sizes for that  $\lambda$  are chosen. That is to say, we can show that the sequence  $V_\lambda(\mathbf{z}_0^k)$  almost monotonically decreases, ignoring some small variance terms.

## 4 COMPARISON WITH RELATED WORKS

### 4.1 COMPARISON WITH STOCHASTIC WITH-REPLACEMENT SETTING

First of all, we confirm that SGDA with random reshuffling (RR) has *faster* convergence rates (i.e., fewer gradient computations) than SGDA based on with-replacement sampling. In particular, we compare our results with the analyses on the purely stochastic minimax settings which assume that every stochastic gradient oracle is *independently sampled and unbiased*: this assumption is naturally satisfied by with-replacement sampling for the finite-sum settings we consider. To make the comparisons fair and easy, we simply let  $r = \beta/\alpha = \Theta(\kappa_2^2)$ ,  $A = 0$ , and  $B = \tau^2$ .

Lin et al. (2020, Theorem 4.5) present a convergence rate for with-replacement simSGDA with  $r = \Theta(\kappa_2^2)$  run on nonconvex  $\mu_2$ -strongly-concave problems with a convex *bounded* constraint set  $\mathcal{Y}$  for dual variable  $\mathbf{y}$ . Their gradient complexity to achieve  $\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla \Phi(\mathbf{x}_t)\|^2 \leq \varepsilon^2$  (where  $T$  is the

number of iterations) is written as  $T_\varepsilon = \mathcal{O}\left(\frac{\kappa_2^2 L \Delta_\Phi + \kappa_2 L^2 D^2}{\varepsilon^2} \max\left\{1, \frac{\kappa_2 \tau^2}{\varepsilon^2}\right\}\right)$ , where  $\kappa_2 = L/\mu_2$ ,  $\Delta_\Phi = \Phi(\mathbf{x}_0) - \Phi^*$ ,  $D = \text{diam } \mathcal{Y}$ , and  $\tau^2$  is the variance of the (unbiased) stochastic gradient oracles. Their complexity can be simplified as  $\mathcal{O}(\kappa_2^3 \tau^2 \varepsilon^{-4})$ , treating other factors as constants. In contrast, our Theorem 1 has a better gradient complexity in terms of  $\varepsilon$  and  $\tau$ , thanks to shuffling:

$$\mathcal{O}\left(\frac{\kappa_2^2 L V_\lambda(\mathbf{z}_0^1)}{\varepsilon^2} \max\left\{n, \frac{\kappa_2 \tau \sqrt{n}}{\varepsilon}\right\}\right), \quad (\text{Ours, from Theorem 1})$$

or simply  $\mathcal{O}(\kappa_2^3 \tau \sqrt{n} \varepsilon^{-3})$ . Thus, our gradient complexity for both simSGDA-RR and altSGDA-RR is better than that of with-replacement simSGDA when  $\varepsilon$  is small as  $\varepsilon \leq \mathcal{O}(\tau/\sqrt{n})$ . Our rate has three more strengths: (i) we do not require strong concavity in  $\mathbf{y}$ , which is a strictly stronger assumption than requiring  $\mathbf{y}$ -side PL condition; (ii) we do not require the constraint set  $\mathcal{Y}$  to be bounded; (iii) our result can easily extend to the case of *any* mini-batch sizes, whereas Lin et al. (2020) need a particular choice of mini-batch size  $M = \mathcal{O}(\kappa_2 \tau^2/\varepsilon)$  to ensure convergence.

For nonconvex-PL objectives, Yang et al. (2022, Theorem 3.1) provide a convergence rate for with-replacement altSGDA with  $r = \Theta(\kappa_2^2)$ . Their rate can be translated to a gradient complexity for achieving  $\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla \Phi(\mathbf{x}_t)\|^2 \leq \varepsilon^2$ , written as  $\mathcal{O}\left(\frac{\kappa_2^2 L V_\lambda(\mathbf{z}_0)}{\varepsilon^2} \left(1 + \frac{\kappa_2^2 V_\lambda(\mathbf{z}_0)^2 \tau^2}{\Delta_\Phi \varepsilon^2}\right)\right)$  or simply  $\mathcal{O}(\kappa_2^4 \tau^2 \varepsilon^{-4})$ . Therefore, our gradient complexity for both altSGDA-RR and simSGDA-RR is better when  $\varepsilon$  is small as  $\varepsilon \leq \mathcal{O}(\kappa_2 \tau/\sqrt{n})$ .

For PL( $\Phi$ )-PL objectives, Yang et al. (2020, Theorem 3.3) obtain a convergence rate for with-replacement altSGDA with  $r = \Theta(\kappa_2^2)$ .<sup>4</sup> They apply diminishing step sizes ( $\mathcal{O}(1/t)$ ,  $t \in \mathbb{N}$ ) to derive a gradient complexity bound  $\mathcal{O}\left(\frac{\kappa_1 \kappa_2^4 \tau^2}{\mu_1 \varepsilon^2}\right)$  to achieve  $\mathbb{E}[V_\lambda(\mathbf{z}_T)] \leq \varepsilon^2$ . One can apply the constant step sizes depending on the total number  $T$  of iterations to their analysis and derive a similar complexity with only deterioration in a logarithmic factor. In contrast, our gradient complexity for both sim/altSGDA-RR using constant step sizes can be written as, for small enough  $\varepsilon$ ,

$$\tilde{\mathcal{O}}\left(\frac{\kappa_1 \kappa_2^3 \tau \sqrt{n}}{\varepsilon \sqrt{\mu_1}}\right). \quad (\text{Ours, from Theorem 2})$$

This is a better complexity in  $\varepsilon$  and  $\kappa_2$ , especially when  $\varepsilon \leq \tilde{\mathcal{O}}(\kappa_2 \tau/\sqrt{n \mu_1})$ , even without the requirement of diminishing step size.

## 4.2 COMPARISON WITH OTHER WORKS ON STOCHASTIC WITHOUT-REPLACEMENT SETTING

One of the most relevant works to this paper is Das et al. (2022, Theorem 3). The authors obtain a similar convergence rate to us for the two-sided PL objective, based on linearization of gradients, but for a dissimilar algorithm which they refer to as *AGDA-RR*. The algorithm can be also thought of as *epoch-wise-alternating* SGDA-RR, whereas our algorithm (altSGDA-RR) can be called as *step-wise-alternating* SGDA-RR. In epoch  $k$ , their algorithm (i) performs updates only on  $\mathbf{x}(\mathbf{x}_0^k, \dots, \mathbf{x}_n^k)$  while fixing  $\mathbf{y}$  to  $\mathbf{y}_0^k$ , and then (ii) performs updates only on  $\mathbf{y}(\mathbf{y}_0^k, \dots, \mathbf{y}_n^k)$  while fixing  $\mathbf{x}$  to  $\mathbf{x}_0^{k+1} = \mathbf{x}_n^k$ . We believe that our step-wise algorithm is closer to practice, especially when  $n$  is large. Because of the distinction between algorithms, the proof techniques are also different.

Xie et al. (2021, Theorem 3) present a convergence rate of *CD-MA*, an extension of simSGDA to the cross-device federated learning setup, on nonconvex-PL setting. Their convergence result for *CD-MA* also assumes mini-batch sampling by random reshuffling. As a consequence, they yield a rate analogous to our Theorem 1 if we reduce their result to the single-machine setup. Nevertheless, our convergence bound contains a term that shrinks with the number of components or mini-batches, whereas theirs does not. For a more detailed comparison, please refer to Appendix H.

There are also some works on RR-based (constrained) minimax optimization algorithms other than SGDA, but for convex-concave problems. Maheshwari et al. (2022) present *OGDA-RR*, a gradient-free RR-based optimistic GDA algorithm. Yu et al. (2022) study *stochastic proximal point with RR*, consisting of double-loop epochs. Their analyses exploit convex-concavity and Lipschitz continuity of their objective, based on the arguments by Nagaraj et al. (2019). This enables a direct usage of the duality gap, the difference between primal function  $\Phi(\cdot)$  and dual function  $\Psi(\cdot) = \min_{\mathbf{x}} f(\mathbf{x}; \cdot)$ , as a criterion for optimality. On the contrary, our work relies on a different structure of the functions, which in turn differentiates the constructions of convergence rates.

<sup>4</sup>Although they consider two-sided PL problems, their analysis applies to PL( $\Phi$ )-PL problems as well.

### 4.3 COMPARISON WITH DETERMINISTIC SETTING

Here, we compare our rates with (full-batch) *gradient descent-ascent* (GDA):

$$\begin{cases} \mathbf{x}_k = \mathbf{x}_{k-1} - \alpha \nabla_1 f(\mathbf{x}_{k-1}; \mathbf{y}_{k-1}), \\ \mathbf{y}_k = \mathbf{y}_{k-1} + \beta \nabla_2 f(\mathbf{x}'_k; \mathbf{y}_{k-1}), \end{cases} \quad \text{where } \mathbf{x}'_k = \begin{cases} \mathbf{x}_{k-1}, & (\text{simGDA}), \\ \mathbf{x}_k, & (\text{altGDA}). \end{cases}$$

It uses the whole information of the objective  $f$  at every iteration without any noise. For comparison with GDA, we utilize our extended theorems for arbitrary mini-batch size  $b$  (Theorems 4 and 5 in Appendix A). By letting  $b = n$  and matching our iterate  $\mathbf{z}_0^k = (\mathbf{x}_0^k; \mathbf{y}_0^k)$  to a GDA iterate  $\mathbf{z}_k = (\mathbf{x}_k; \mathbf{y}_k)$ , our results reduce to upper convergence bounds for simGDA and altGDA.

For nonconvex-PŁ problems (Theorems 1 & 4), the convergence rate and *iteration* complexity (*i.e.*, sufficient number of iterations  $K_\varepsilon$ ) become

$$\min_{k \in [K]} \|\nabla \Phi(\mathbf{x}_k)\|^2 \leq \mathcal{O}\left(\frac{\kappa_2^2 L V_\lambda(\mathbf{z}_1)}{K}\right); \quad \text{i.e., } K_\varepsilon = \mathcal{O}\left(\frac{\kappa_2^2 L V_\lambda(\mathbf{z}_1)}{\varepsilon^2}\right), \quad (3)$$

when  $r = \Theta(\kappa_2^2)$ . This is similar to a known rate of simGDA with  $r = \Theta(\kappa_2^2)$  for nonconvex-strongly-concave problems by Lin et al. (2020, Theorem 4.4) as a special case. Their iteration complexity is written as  $\mathcal{O}((\kappa_2^2 L \Delta_\Phi + \kappa_2 L^2 D^2)/\varepsilon^2)$ , where the symbols are already defined in Section 4.1. To see how the two bounds compare in terms of the factors other than  $\varepsilon$ , notice that we have  $\Phi(\mathbf{x}) - f(\mathbf{x}; \mathbf{y}) \leq \frac{L}{2} \|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2$  for any  $(\mathbf{x}; \mathbf{y})$ , due to the  $L$ -smoothness of  $-f$ . Here,  $\mathbf{y}^*(\mathbf{x})$  is an element of  $\mathcal{Y}_\mathbf{x}^* = \arg \max_{\mathbf{y}} f(\mathbf{x}; \mathbf{y})$ . Thus, we have  $V_\lambda(\mathbf{z}_1) = \lambda[\Phi(\mathbf{x}_1) - \Phi^*] + [\Phi(\mathbf{x}_1) - f(\mathbf{z}_1)] \leq \lambda \Delta_\Phi + L D^2/2$ . As a result, we could *loosely* translate our iteration complexity (3) to  $\mathcal{O}((\kappa_2^2 L \Delta_\Phi + \kappa_2^2 L^2 D^2)/\varepsilon^2)$ . We suspect that the discrepancy in terms of  $\kappa_2$  comes from the fact that our analysis does not require the (strong) concavity in terms of  $\mathbf{y}$  or a bounded constraint  $\mathcal{Y}$ : these requirements made a considerable difference in proofs.

For PŁ( $\Phi$ )-PŁ problems (Theorems 2 & 5), the rate and iteration complexity ( $K'_\varepsilon$ ) become

$$V_\lambda(\mathbf{z}_{K+1}) \leq V_\lambda(\mathbf{z}_1) \cdot \exp\left(-\frac{K}{C \kappa_1 \kappa_2^2}\right); \quad \text{i.e., } K'_\varepsilon = \mathcal{O}(\kappa_1 \kappa_2^2 \log(1/\varepsilon)) \quad (4)$$

where  $r = \Theta(\kappa_2^2)$  and  $C$  is a numerical constant. This recovers the linear convergence by Yang et al. (2020, Theorem 3.2) as a special case, where they prove convergence of altGDA with step size ratio  $r = \Theta(\kappa_2^2)$  for two-sided PŁ problem. Following the proof of (Yang et al., 2020, Theorem 3.2), one can show that the bound (4) indeed implies the actual convergence to a global minimax point  $\mathbf{z}^*$ , in the sense that we can achieve  $\|\mathbf{z}_k - \mathbf{z}^*\| \leq \varepsilon$  in  $\mathcal{O}(\kappa_1 \kappa_2^2 \log(1/\varepsilon))$  iterations.

## 5 LOWER BOUND FOR (FULL-BATCH) SIMGDA USING SEPARATE STEP SIZES

As an extension of the discussion from Section 4.3, we characterize a lower complexity bound of deterministic simGDA with separate step sizes  $(\alpha, \beta)$  of arbitrary ratio  $r = \beta/\alpha$ , for smooth strongly-convex-strongly-concave (SC-SC) cases. Surprisingly, at least for  $r \gtrsim \kappa_2^2$ , our lower bound matches the upper complexity bound of GDA for a much wider class of smooth PŁ( $\Phi$ )-PŁ problems,<sup>5</sup> which is quite surprising.

For a smooth PŁ( $\Phi$ )-PŁ problems, simGDA with at least  $r = \Omega(\kappa_2^2)$  has an upper complexity bound  $K = \mathcal{O}(\kappa_1 r \log(1/\varepsilon))$  for a *global*  $\varepsilon$ -convergence  $V_\lambda(\mathbf{z}_K) \leq \varepsilon^2$  in terms of potential function. This means that the lowest complexity is  $\mathcal{O}(\kappa_1 \kappa_2^2 \log(1/\varepsilon))$  achieved when  $r = \Theta(\kappa_2^2)$ . On the other hand, for a  $L$ -smooth  $\mu$ -SC-SC problem with saddle point  $\mathbf{z}^*$ , it is well-known that the simGDA with a single step-size ( $\alpha = \beta$ ) has a tight upper/lower complexity  $K = \Theta(\kappa^2 \log(1/\varepsilon))$  to achieve  $\|\mathbf{z}_K - \mathbf{z}^*\|^2 \leq \varepsilon^2$ , where  $\kappa = L/\mu$  (*e.g.*, Das et al. (2022, Theorem C.1)). The difference of complexity bounds in condition number ( $\kappa_1 \kappa_2^2$  v.s.  $\kappa^2$ ) is somewhat questionable because, at least in smooth minimization problems, strongly convex problems and PŁ problems have identical gradient descent (GD) iteration complexity  $\mathcal{O}(\kappa \log(1/\varepsilon))$  (Karimi et al., 2016, Theorem 1).

One could ask where the discrepancy in terms of  $\kappa$  comes from: is it due to (i) the criteria ( $V_\lambda(\mathbf{z}_K)$  v.s.  $\|\mathbf{z}_K - \mathbf{z}^*\|^2$ ) for  $\varepsilon$ -accuracy, (ii) the function classes (PŁ( $\Phi$ )-PŁ v.s. SC-SC), or (iii) the step size

<sup>5</sup>strongly-convex-strongly-concave (SC-SC)  $\subset$  two-sided PŁ  $\subset$  PŁ( $\Phi$ )-PŁ  $\subset$  nonconvex-PŁ.

ratios ( $\Omega(\kappa_2^2)$  v.s. 1)? We answer the question by showing the following theorem: the discrepancy in  $\kappa$  comes from the step size ratio difference. We defer the proof to Appendix E.

**Theorem 3** (Lower bound, ratio-specific). *Consider a class  $\mathcal{F}(L, \mu_1, \mu_2)$  of functions  $f(\mathbf{x}; \mathbf{y})$  with two arguments  $\mathbf{x}$  and  $\mathbf{y}$ , which is  $L$ -smooth,  $\mu_1$ -strongly-convex in  $\mathbf{x}$ , and  $\mu_2$ -strongly-concave in  $\mathbf{y}$ . Suppose  $\kappa_1 = L/\mu_1 \geq c$  and  $\kappa_2 = L/\mu_2 \geq c$  for some constant  $c > 1$ . Then, for any step size ratio  $r = \beta/\alpha > 0$ , there exists a function  $f \in \mathcal{F}(L, \mu_1, \mu_2)$  with a unique saddle point  $\mathbf{z}^*$ , for which simGDA with any step sizes  $(\alpha, \beta) = (\beta/r, \beta)$  requires at least*

$$K = \begin{cases} \Omega(\kappa_1 r \log(1/\varepsilon)), & \text{if } r \geq \kappa_2/c, \\ \Omega(\kappa_1 \kappa_2 \log(1/\varepsilon)), & \text{if } c/\kappa_1 \leq r \leq \kappa_2/c, \\ \Omega((\kappa_2/r) \log(1/\varepsilon)), & \text{if } 0 < r \leq c/\kappa_1 \end{cases}$$

iterations to achieve either  $\|\mathbf{z}_K - \mathbf{z}^*\|^2 \leq \varepsilon^2$  or  $V_\lambda(\mathbf{z}_K) \leq \varepsilon^2$ .

Thanks to Theorem 3, we can say from Theorem 5 that for any step size ratio  $r \gtrsim \kappa_2^2$ , we have a *tight* upper bound on the iteration complexity  $K = \mathcal{O}(\kappa_1 r \log(1/\varepsilon))$  of simGDA for general PL( $\Phi$ )-PL problems. Note that Theorem 3 also subsumes the existing lower bound of the equal-step-size ( $r = 1$ ) simGDA for  $\mu$ -SC-SC problems.

Given the tightness of bounds for  $r \gtrsim \kappa_2^2$ , a natural next step is to discuss  $1 \lesssim r \lesssim \kappa_2^2$ . Recent work by Li et al. (2022) also discusses the step size ratio of simGDA. In Li et al. (2022, Theorem 4.1), the authors construct a  $y$ -side strongly-concave function<sup>6</sup> and show that simGDA with a step size ratio  $r \leq \kappa_2$  is *impossible* to converge. The *necessity* of  $r \gtrsim \kappa_2$  implied by this theorem also applies to the PL( $\Phi$ )-PL case. Thus, there is no hope for showing an upper convergence bound of simGDA with  $1 \lesssim r \lesssim \kappa_2$  for general nonconvex-PL problems. We remark that their theorem does not contradict nor subsume Theorem 3 because we consider a much smaller function class (SC-SC) to construct the lower bounds.

On the *sufficiency* of  $r \gtrsim \kappa_2$  for convergence, Li et al. (2022, Theorem 4.2) show that simGDA with  $r \geq c\kappa$  (for some  $c > 1$ ) can *locally* converge at the iteration complexity  $\mathcal{O}(\kappa_1 r \log(1/\varepsilon))$  for some nonconvex-strongly-concave problems, which matches the bound in Theorem 3. Our upper bounds (Theorems 4 and 5) do require  $r \gtrsim \kappa_2^2$ , which may look suboptimal, but we claim that our results are not necessarily weaker. One reason is that our convergence guarantee is *global*, i.e., independent of the initialization. Another reason is that their analysis is only valid when a *differential Stackelberg equilibrium*<sup>7</sup> exists, whereas a general PL( $\Phi$ )-PL function may not have such an equilibrium (for an example, see Proposition 13 in Appendix B).

As far as we know, it is still an open problem whether a *global* convergence bound for simGDA on nonconvex-PL problems can be shown when the step size ratio  $r$  is between  $\Omega(\kappa_2)$  and  $\mathcal{O}(\kappa_2^2)$ .

## 6 EXPERIMENTS

To validate our main theoretical findings, here we present some numerical results. We focus on the primal-PL-strongly-concave (or PL( $\Phi$ )-SC, which is PL( $\Phi$ )-PL as well) quadratic games of the form

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^d} f(\mathbf{x}; \mathbf{y}) &= \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{x}^\top \mathbf{B} \mathbf{y} - \frac{1}{2} \mathbf{y}^\top \mathbf{C} \mathbf{y} = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}; \mathbf{y}), \\ \text{where } f_i(\mathbf{x}; \mathbf{y}) &= \frac{1}{2} \mathbf{x}^\top \mathbf{A}_i \mathbf{x} + \mathbf{x}^\top \mathbf{B}_i \mathbf{y} - \frac{1}{2} \mathbf{y}^\top \mathbf{C}_i \mathbf{y} + \mathbf{u}_i^\top \mathbf{x} - \mathbf{v}_i^\top \mathbf{y}. \end{aligned} \quad (5)$$

This toy example is often used to numerically evaluate the minimax algorithms (Yang et al., 2020; Loizou et al., 2021; Das et al., 2022) and appears in various domains such as AUC maximization (Ying et al., 2016), policy evaluation (Du et al., 2017), and imitation learning (Cai et al., 2019)

To make the game in Equation (5) satisfy PL( $\Phi$ )-SC and component  $L$ -smoothness, we should sample the coefficient matrices and vectors carefully. First, they need to be  $\|\mathbf{A}_i\|_2, \|\mathbf{B}_i\|_2, \|\mathbf{C}_i\|_2 \leq L$  and  $\sum_{i=1}^n \mathbf{u}_i = \sum_{i=1}^n \mathbf{v}_i = \mathbf{0}$ . To make the primal function  $\Phi$  a well-defined real-valued function

<sup>6</sup>  $g(x; y) = -\frac{1}{2}x^2 + Lxy - \frac{\mu}{2}y^2$ , where  $L/\mu > 1$ : its primal function is strongly convex.

<sup>7</sup> Loosely speaking, a differential Stackelberg equilibrium is a stationary point  $(\mathbf{x}^*; \mathbf{y}^*)$  where  $f(\mathbf{x}^*; \cdot)$  is locally strongly concave near  $\mathbf{y}^*$  and  $\Phi(\cdot)$  is locally strongly convex near  $\mathbf{x}^*$ .



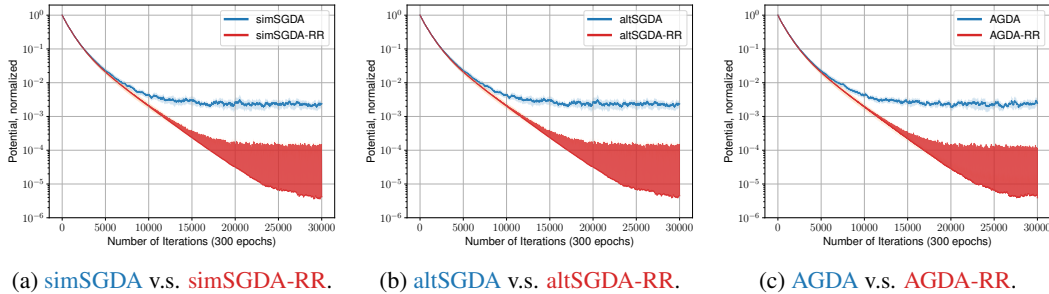


Figure 1: Experimental results on quadratic games (5). Solid lines: average across 10 different runs. Shaded regions: 95% confidence intervals ( $\pm 1.96$  std). Dots: start/end of epochs. The vertical axes are on a *logarithmic scale*.

for any  $\mathbf{x} \in \mathbb{R}^d$ , we choose  $\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i$  to be positive definite, *i.e.*,  $\mu \mathbf{I} \preceq \mathbf{C}$  for an identity matrix  $\mathbf{I}$  and  $\mu > 0$ . Then, the primal function can be explicitly written as

$$\Phi(\mathbf{x}) = \max_{\mathbf{y} \in \mathbb{R}^d} f(\mathbf{x}; \mathbf{y}) = \frac{1}{2} \mathbf{x}^\top (\mathbf{A} + \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^\top) \mathbf{x} := \frac{1}{2} \mathbf{x}^\top \mathbf{M} \mathbf{x}.$$

We construct a matrix  $\mathbf{M} := \mathbf{A} + \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^\top$  to be rank-deficient positive semi-definite. Letting the smallest nonzero eigenvalue of  $\mathbf{M}$  by  $\mu$ , we ensure that  $\Phi$  is  $\mu$ -PL but not strongly convex. We emphasize that the objective function  $f$  is not even (strongly-)convex in  $\mathbf{x}$  in general.

We compare six algorithms in total: simSGDA-RR, altSGDA-RR, AGDA-RR (as defined in Das et al. (2022)), and the with-replacement counterparts of these three algorithms. To this end, on 5 different randomly-generated quadratic games and under 2 random seeds per game (*i.e.*, 10 runs per algorithm), we run each algorithm for the same number of epochs using constant step sizes of ratio  $\beta/\alpha = c\kappa_2^2$  for some constant  $c$  and  $\kappa_2 = L/\mu$ .

We report the potential function values ( $V_\lambda$ , defined in Equation (2)) at every iteration.<sup>8</sup> Results are presented in Figure 1: the values are normalized by dividing them by the initial value. As we discussed in Section 4.1, we observe that the random reshuffling considerably accelerates the convergence of the algorithms. Furthermore, all three algorithms with random reshuffling show more or less the same performance. Specifically, the plots for simSGDA (resp. simSGDA-RR) and altSGDA (resp. altSGDA-RR) are almost identical. We believe this is because we choose a random seed for each of the 10 different runs and share it across different algorithms.

Please refer to Appendix G for more detailed construction, discussion, and comparative study of the experimental results.

## 7 CONCLUSION

We investigated stochastic algorithms based on without-replacement component sampling, called simSGDA-RR and altSGDA-RR, for solving smooth nonconvex finite-sum minimax optimization problems. We established convergence rates under the  $\mathbf{y}$ -side PL condition (nonconvex-PL) and, additionally, the primal PL condition (PL( $\Phi$ )-PL). We ascertain that the SGDA-RR can achieve a faster rate than its with-replacement counterpart, which agrees with the existing theory on without-replacement SGD for minimization. Lastly, we provided complexity lower bounds for simGDA with an arbitrarily fixed step size ratio  $r$ , demonstrating that the full-batch upper bound with  $r \gtrsim \kappa_2^2$  for PL( $\Phi$ )-PL functions is tight.

Possible future directions include widening our results beyond sim/altSGDA (*e.g.*, extra-gradient or optimistic GDA) and beyond RR (*e.g.*, single/adversarial shuffling). As also discussed in Section 5, an interesting open question remains open: can we identify tight convergence rates for stochastic (with-/without-replacement) and/or deterministic GDA with step size ratio  $r$  satisfying  $\kappa_2 \lesssim r \lesssim \kappa_2^2$ , for general nonconvex-PL problems?

<sup>8</sup>As described in Section 4.2, AGDA-RR uses only one-side gradient ( $\nabla_1$  or  $\nabla_2$ ) at each iteration: given a fixed budget of gradient computations, it should access components twice as many times as SGDA-RR. Hence, we report the values at every other iteration of AGDA & AGDA-RR, for a fair comparison.

## ACKNOWLEDGMENTS

This work was supported by Institute of Information & communications Technology Planning & evaluation (IITP) grant (No. 2019-0-00075, Artificial Intelligence Graduate School Program (KAIST)) funded by the Korea government (MSIT). The work was also supported by the National Research Foundation of Korea (NRF) grant (No. NRF-2019R1A5A1028324) funded by the Korea government (MSIT). CY acknowledges support from a grant funded by Samsung Electronics Co., Ltd.

## REFERENCES

- Kwangjun Ahn, Chulhee Yun, and Suvrit Sra. SGD with shuffling: optimal rates without component convexity and large epoch requirements. *Advances in Neural Information Processing Systems*, 33: 17526–17535, 2020. [2](#), [22](#)
- Aleksandr Beznosikov, Eduard Gorbunov, Hugo Berard, and Nicolas Loizou. Stochastic gradient descent-ascent: Unified theory and new efficient methods. *arXiv preprint arXiv:2202.07262*, 2022. [2](#)
- Qi Cai, Mingyi Hong, Yongxin Chen, and Zhaoran Wang. On the global convergence of imitation learning: A case for linear quadratic regulator. *arXiv preprint arXiv:1901.03674*, 2019. [8](#)
- Aniket Das, Bernhard Schölkopf, and Michael Muehlebach. Sampling without replacement leads to faster rates in finite-sum minimax optimization. *arXiv preprint arXiv:2206.02953*, 2022. [2](#), [4](#), [6](#), [7](#), [8](#), [9](#), [22](#)
- Simon S Du, Jianshu Chen, Lihong Li, Lin Xiao, and Dengyong Zhou. Stochastic variance reduction methods for policy evaluation. In *International Conference on Machine Learning*, pp. 1049–1058. PMLR, 2017. [8](#)
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. [1](#)
- Zhishuai Guo, Zhuoning Yuan, Yan Yan, and Tianbao Yang. Fast objective & duality gap convergence for nonconvex-strongly-concave min-max problems. *arXiv preprint arXiv:2006.06889*, 2020. [2](#), [4](#), [19](#)
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. [3](#)
- Roger A Horn and Charles R Johnson. *Matrix Analysis*. Cambridge University Press, Cambridge, England, 2 edition, October 2012. [37](#)
- Chi Jin, Praneeth Netrapalli, and Michael Jordan. What is local optimality in nonconvex-nonconcave minimax optimization? In *International conference on machine learning*, pp. 4880–4889. PMLR, 2020. [5](#), [19](#)
- Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In *Joint European conference on machine learning and knowledge discovery in databases*, pp. 795–811. Springer, 2016. [4](#), [7](#), [16](#), [17](#), [19](#)
- Haochuan Li, Farzan Farnia, Subhro Das, and Ali Jadbabaie. On convergence of gradient descent ascent: A tight local analysis. In *International Conference on Machine Learning*, pp. 12717–12740. PMLR, 2022. [3](#), [8](#), [19](#), [42](#)
- Shihui Li, Yi Wu, Xinyue Cui, Honghua Dong, Fei Fang, and Stuart Russell. Robust multi-agent reinforcement learning via minimax deep deterministic policy gradient. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 4213–4220, 2019. [1](#)

- Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pp. 6083–6093. PMLR, 2020. [2](#), [3](#), [5](#), [6](#), [7](#)
- Mingrui Liu, Zhuoning Yuan, Yiming Ying, and Tianbao Yang. Stochastic AUC maximization with deep neural networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJepXaVYDr>. [1](#), [4](#)
- Nicolas Loizou, Hugo Berard, Gauthier Gidel, Ioannis Mitliagkas, and Simon Lacoste-Julien. Stochastic gradient descent-ascent and consensus optimization for smooth games: Convergence analysis under expected co-coercivity. *Advances in Neural Information Processing Systems*, 34: 19095–19108, 2021. [2](#), [8](#), [15](#), [44](#)
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*. OpenReview.net, 2018. [1](#)
- Chinmay Maheshwari, Chih-Yuan Chiu, Eric Mazumdar, Shankar Sastry, and Lillian Ratliff. Zeroth-order methods for convex-concave min-max problems: Applications to decision-dependent risk minimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 6702–6734. PMLR, 2022. [2](#), [6](#)
- Konstantin Mishchenko, Ahmed Khaled, and Peter Richtárik. Random reshuffling: Simple analysis with vast improvements. *Advances in Neural Information Processing Systems*, 33:17309–17320, 2020. [2](#), [21](#), [22](#), [44](#)
- Dheeraj Nagaraj, Prateek Jain, and Praneeth Netrapalli. SGD without replacement: Sharper rates for general smooth convex functions. In *International Conference on Machine Learning (ICML)*, pp. 4703–4711. PMLR, 2019. [2](#), [6](#)
- Lam M Nguyen, Quoc Tran-Dinh, Dzung T Phan, Phuong Ha Nguyen, and Marten Van Dijk. A unified convergence analysis for shuffling-type gradient methods. *The Journal of Machine Learning Research*, 22(1):9397–9440, 2021. [2](#), [3](#), [22](#), [25](#)
- Maher Nouiehed, Maziar Sanjabi, Tianjian Huang, Jason D Lee, and Meisam Razaviyayn. Solving a class of non-convex min-max games using iterative first order methods. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019. [2](#), [4](#), [17](#), [18](#), [42](#)
- Shashank Rajput, Anant Gupta, and Dimitris Papailiopoulos. Closing the convergence gap of SGD without replacement. In *International Conference on Machine Learning*, pp. 7964–7973. PMLR, 2020. [2](#)
- Aman Sinha, Hongseok Namkoong, and John Duchi. Certifiable distributional robustness with principled adversarial training. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=Hk6kPgZA->. [1](#)
- Jiahao Xie, Chao Zhang, Zebang Shen, Weijie Liu, and Hui Qian. Efficient cross-device federated learning algorithms for minimax problems. *arXiv e-prints*, pp. arXiv-2105, 2021. [6](#), [14](#), [45](#), [46](#)
- Yan Yan, Yi Xu, Qihang Lin, Wei Liu, and Tianbao Yang. Optimal epoch stochastic gradient descent ascent methods for min-max optimization. *Advances in Neural Information Processing Systems*, 33:5789–5800, 2020. [2](#)
- Junchi Yang, Negar Kiyavash, and Niao He. Global convergence and variance reduction for a class of nonconvex-nonconcave minimax problems. *Advances in Neural Information Processing Systems*, 33:1153–1165, 2020. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [18](#), [40](#)
- Junchi Yang, Antonio Orvieto, Aurelien Lucchi, and Niao He. Faster single-loop algorithms for minimax optimization without strong concavity. In *International Conference on Artificial Intelligence and Statistics*, pp. 5485–5517. PMLR, 2022. [6](#)
- Yiming Ying, Longyin Wen, and Siwei Lyu. Stochastic online AUC maximization. *Advances in neural information processing systems*, 29, 2016. [1](#), [8](#)

Yaodong Yu, Tianyi Lin, Eric V Mazumdar, and Michael Jordan. Fast distributionally robust learning with variance-reduced min-max optimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 1219–1250. PMLR, 2022. [2](#), [6](#)

Zhuoning Yuan, Yan Yan, Milan Sonka, and Tianbao Yang. Large-scale robust deep AUC maximization: A new surrogate loss and empirical studies on medical image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3040–3049, 2021. [1](#)

Chulhee Yun, Suvrit Sra, and Ali Jadbabaie. Open problem: Can single-shuffle SGD be better than reshuffling SGD and GD? In *Conference on Learning Theory*, pp. 4653–4658. PMLR, 2021. [2](#)

Chulhee Yun, Shashank Rajput, and Suvrit Sra. Minibatch vs local SGD with shuffling: Tight convergence bounds and beyond. In *International Conference on Learning Representations*. OpenReview.net, 2022. [2](#)

## CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Summary of our contributions . . . . .	2
<b>2</b>	<b>Problem setup</b>	<b>2</b>
2.1	Notation . . . . .	2
2.2	Algorithms: simSGDA-RR & altSGDA-RR . . . . .	3
2.3	Assumptions and definitions . . . . .	3
<b>3</b>	<b>Main results</b>	<b>4</b>
3.1	Potential function $V_\lambda$ . . . . .	4
3.2	Main theorems: upper bounds of convergence rates . . . . .	4
<b>4</b>	<b>Comparison with related works</b>	<b>5</b>
4.1	Comparison with stochastic with-replacement setting . . . . .	5
4.2	Comparison with other works on stochastic without-replacement setting . . . . .	6
4.3	Comparison with deterministic setting . . . . .	7
<b>5</b>	<b>Lower bound for (full-batch) simGDA using separate step sizes</b>	<b>7</b>
<b>6</b>	<b>Experiments</b>	<b>8</b>
<b>7</b>	<b>Conclusion</b>	<b>9</b>
<b>A</b>	<b>Mini-batch SGDA-RR and convergence rates</b>	<b>15</b>
<b>B</b>	<b>Technical propositions</b>	<b>16</b>
B.1	Function classes: PL condition, smoothness, and more . . . . .	16
B.2	Without-replacement sampling . . . . .	20
B.3	Basic recurrence inequality . . . . .	22
<b>C</b>	<b>Proofs for (mini-batch) simultaneous SGDA-RR</b>	<b>22</b>
C.1	Warm-up: proof sketch for $b = 1$ . . . . .	22
C.2	Epoch-wise representations and bounding noise terms . . . . .	23
C.3	Recurrence inequalities for general smooth nonconvex-PL objective . . . . .	27
C.4	Convergence rates for smooth nonconvex-PL problem . . . . .	30
C.5	Convergence rates for smooth primal-PL-PL problem . . . . .	30
<b>D</b>	<b>Proofs for (mini-batch) alternating SGDA-RR: focusing on changes in the proof</b>	<b>32</b>
D.1	Epoch-wise representations and bounding noise terms . . . . .	32
D.2	Bounding noise terms: a bit different proof of Lemma 18 . . . . .	33
D.3	Recurrence inequalities for general smooth nonconvex-PL objective . . . . .	34

D.4	Small step size assumptions . . . . .	34
<b>E</b>	<b>Proofs for lower bound of deterministic full-batch simGDA</b>	<b>35</b>
<b>F</b>	<b>Remark on smoothness assumptions and Lower bound of with-replacement SGD(A)</b>	<b>39</b>
<b>G</b>	<b>Experiments: quadratic games</b>	<b>41</b>
G.1	Parameter choices . . . . .	41
G.2	Comparison: the effect of component discrepancy . . . . .	43
G.3	Comparison: the effect of condition number . . . . .	43
G.4	Comparison: the effect of batch size . . . . .	44
<b>H</b>	<b>Omitted comparison with related works</b>	<b>45</b>
H.1	Comparison with <a href="#">Xie et al. (2021)</a> . . . . .	45

## A MINI-BATCH SGDA-RR AND CONVERGENCE RATES

In this appendix, we present an algorithm that extends simSGDA-RR and altSGDA-RR by using mini-batches of size  $b \geq 1$ . For simplicity, we assume that the number of components  $n$  is an integer multiple of the mini-batch size  $b$  in our analysis; *i.e.*,  $n = bq$  for some integer  $q \geq 1$ . One can extend this to the case when  $n$  is not necessarily a multiple of  $b$  (*e.g.*,  $n = b(q-1) + s$ , where  $q \geq 1$ ,  $s \in [b]$ ) so that there are  $q-1$  mini-batches of size  $b$  and one more mini-batch of size  $s \leq b$ .

---

### Algorithm 2 Mini-batch simSGDA/altSGDA-RR

---

- 1: **Given:** The number of components  $n = b(q-1) + s$  ( $q$ : number of iterations per epoch); **mini-batch size**  $b$ ; the number of epochs  $K$ ; step sizes  $\alpha, \beta > 0$
  - 2: **Initialize:**  $(\mathbf{x}_0^1; \mathbf{y}_0^1) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$
  - 3: **for**  $k \in [K]$  **do**
  - 4:   Sample  $\sigma_k \sim \text{Unif}(\mathbb{S}_n)$                     $\triangleright$  RR: uniformly randomly shuffle the indices every epoch
  - 5:   **for**  $t \in [q]$  **do**
  - 6:      $\mathcal{B}_t^k := \{\sigma_k(j) : b(t-1) < j \leq bt, j \in [n]\}$     $\triangleright$  Mini-batch : a set of component indices
  - 7:      $\mathbf{x}_t^k = \mathbf{x}_{t-1}^k - \frac{\alpha}{b} \sum_{i \in \mathcal{B}_t^k} \nabla_1 f_i(\mathbf{x}_{t-1}^k; \mathbf{y}_{t-1}^k)$
  - 8:     **if** simSGDA-RR **then**
  - 9:        $\mathbf{y}_t^k = \mathbf{y}_{t-1}^k + \frac{\beta}{b} \sum_{i \in \mathcal{B}_t^k} \nabla_2 f_i(\mathbf{x}_{t-1}^k; \mathbf{y}_{t-1}^k)$                     $\triangleright$  simultaneous update:  $\mathbf{x}$  &  $\mathbf{y}$
  - 10:    **else if** altSGDA-RR **then**
  - 11:       $\mathbf{y}_t^k = \mathbf{y}_{t-1}^k + \frac{\beta}{b} \sum_{i \in \mathcal{B}_t^k} \nabla_2 f_i(\mathbf{x}_t^k; \mathbf{y}_{t-1}^k)$                     $\triangleright$  alternating update:  $\mathbf{x} \rightarrow \mathbf{y}$
  - 12:     $(\mathbf{x}_0^{k+1}; \mathbf{y}_0^{k+1}) = (\mathbf{x}_{n/b}^k; \mathbf{y}_{n/b}^k)$
- 

Next, we illustrate the generalized versions of our main results (Theorems 1 and 2) for Algorithm 2 with mini-batches of size  $b \geq 1$ . Let us assume  $n \geq 2$  because the case  $n = 1$  trivially boils down to simGDA or altGDA. We defer the proofs for simultaneous updates to Appendix C. We present the parts that change in the proof for alternating updates in Appendix D.

**Theorem 4** (Nonconvex-PL, mini-batch SGDA-RR). *Suppose  $f$  satisfies Assumptions 1, 2, 3, and 4. Let  $\lambda = 4$ . Choose the step sizes  $\alpha$  and  $\beta$  by  $\alpha = \beta/r$  for some  $r \geq 14\kappa_2^2$  and*

$$\beta = b \cdot \min \left\{ \frac{1}{6Ln \sqrt{1 + \frac{n-b}{n-1} \cdot \frac{A}{n}}}, \frac{1}{c} \left( \frac{V_\lambda(\mathbf{z}_0^1)}{Ln^2 \left(\frac{n-b}{n-1}\right) BK} \right)^{\frac{1}{3}} \right\},$$

for some numerical constant  $c > 0$ . Then, mini-batch simSGDA-RR and altSGDA-RR with mini-batch size  $b$  (a divisor of  $n$ ) satisfy

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \|\nabla \Phi(\mathbf{x}_0^k)\|^2 \leq \frac{6rLV_\lambda(\mathbf{z}_0^1)}{K} \sqrt{1 + \left(\frac{n-b}{n-1}\right) \frac{A}{n}} + 2cr \left( \frac{L^2 B V_\lambda(\mathbf{z}_0^1)^2}{nK^2} \cdot \frac{n-b}{n-1} \right)^{1/3}.$$

**Theorem 5** (PL( $\Phi$ )-PL, mini-batch SGDA-RR). *Suppose  $f$  satisfies Assumptions 1, 2, 3, 4, and 5. Let  $\lambda = 4$ . Choose the step sizes  $\alpha$  and  $\beta$  by  $\alpha = \beta/r$  for some  $r \geq 14\kappa_2^2$  and*

$$\beta = b \cdot \min \left\{ \frac{1}{6Ln \sqrt{1 + \frac{n-b}{n-1} \cdot \frac{A}{n}}}, \frac{2r}{\mu_1 n K} \max \left\{ 1, \log \left( \frac{V_\lambda(\mathbf{z}_0^1) \mu_1 n K^2}{8c^3 \kappa_1^2 r^3 \left(\frac{n-b}{n-1}\right) B} \right) \right\} \right\},$$

for some numerical constant  $c > 0$ . Then, mini-batch simSGDA-RR and altSGDA-RR with mini-batch size  $b$  (a divisor of  $n$ ) satisfy

$$\mathbb{E}[V_\lambda(\mathbf{z}_0^{K+1})] \leq \mathcal{O} \left( V_\lambda(\mathbf{z}_0^1) \cdot \exp \left( -\frac{K}{12\kappa_1 r \sqrt{1 + \frac{n-b}{n-1} \cdot \frac{A}{n}}} \right) \right) + \tilde{\mathcal{O}} \left( \frac{\kappa_1^2 r^3 B}{\mu_1 n K^2} \right) \cdot \frac{n-b}{n-1}.$$

As a side remark, some works consider a sampling method called *b-minibatch sampling* where all the elements in each mini-batch are distinct (*i.e.*, without-replacement component sampling per mini-batch), *e.g.*, Loizou et al. (2021, Definition 2.1). However, there is a significant gap between this method and ours: any two distinct mini-batches sampled by the *b*-minibatch sampling can intersect with each other (*i.e.*, mini-batches are sampled with replacement), whereas, in each epoch of our Algorithm 2, all the mini-batches are mutually disjoint.

## B TECHNICAL PROPOSITIONS

**Notation.** Throughout this appendix, we use  $\mathcal{X} = \mathbb{R}^{d_x}$  and  $\mathcal{Y} = \mathbb{R}^{d_y}$ . Given a closed set  $\mathcal{S} \subset \mathbb{R}^d$ , we denote the set of all projection(s) of  $\mathbf{v} \in \mathbb{R}^d$  onto  $\mathcal{S}$ , *i.e.*, the nearest point(s) in  $\mathcal{S}$  from  $\mathbf{v}$ , by  $\Pi_{\mathcal{S}}(\mathbf{v}) := \arg \min_{\mathbf{w} \in \mathcal{S}} \|\mathbf{v} - \mathbf{w}\|$ .

### B.1 FUNCTION CLASSES: PŁ CONDITION, SMOOTHNESS, AND MORE

**Proposition 6** ( $\kappa \geq 1$ ). *Let  $g$  be an  $L$ -smooth function which is bounded below by  $g^*$ . Then, for any  $\mathbf{x}$ ,*

$$\|\nabla g(\mathbf{x})\|^2 \leq 2L [g(\mathbf{x}) - g^*].$$

*If  $g$  is  $\mu$ -PŁ as well, then  $\mu \leq L$ . Consequently, the condition number  $\kappa := L/\mu$  of  $g$  is  $\geq 1$ .*

*Proof.* Since  $g$  is  $L$ -smooth, for any  $\mathbf{x}$  and  $\mathbf{y}$ ,

$$g^* \leq g(\mathbf{y}) \leq g(\mathbf{x}) + \langle \nabla g(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad (6)$$

Now define a convex quadratic function  $h_x(\mathbf{y})$  of  $\mathbf{y}$  as

$$h_x(\mathbf{y}) := g(\mathbf{x}) + \langle \nabla g(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Since its gradient is

$$\nabla h_x(\mathbf{y}) = \nabla g(\mathbf{x}) + L(\mathbf{y} - \mathbf{x}),$$

$\mathbf{y}^* := \mathbf{x} - \frac{1}{L} \nabla g(\mathbf{x})$  is a minimum of  $h_x$ . Plugging  $\mathbf{y} = \mathbf{y}^*$  to the equation (6), we get

$$g^* \leq g(\mathbf{x}) + \left\langle \nabla g(\mathbf{x}), -\frac{1}{L} \nabla g(\mathbf{x}) \right\rangle + \frac{L}{2} \left\| -\frac{1}{L} \nabla g(\mathbf{x}) \right\|^2 = g(\mathbf{x}) - \frac{1}{2L} \|\nabla g(\mathbf{x})\|^2.$$

Rearranging the terms,

$$\|\nabla g(\mathbf{x})\|^2 \leq 2L [g(\mathbf{x}) - g^*].$$

If we additionally utilize PŁ inequality with  $g^* := \min g(\mathbf{x})$ ,

$$\|\nabla g(\mathbf{x})\|^2 \geq 2\mu [g(\mathbf{x}) - g^*],$$

we directly yield  $\mu \leq L$  and thus  $\kappa = L/\mu \geq 1$ .  $\square$

**Definition 1** (Karimi et al. (2016)). *Consider  $g : \mathcal{X} \rightarrow \mathbb{R}$ . Let  $\mathbf{x}_p \in \Pi_{\mathcal{X}^*}(\mathbf{x})$  be a projection of  $\mathbf{x}$  onto the optimal set  $\mathcal{X}^* = \arg \min_{\mathbf{x} \in \mathcal{X}} g(\mathbf{x})$ .*

- (1) *We say  $g$  satisfies  $\mu$ -strong convexity (SC) if  $g(\mathbf{x}') \geq g(\mathbf{x}) + \langle \nabla g(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{x}' - \mathbf{x}\|^2$  for any  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ .*
- (2) *We say  $g$  satisfies  $\mu$ -restricted secant inequality (RSI) if  $\langle \nabla g(\mathbf{x}), \mathbf{x} - \mathbf{x}_p \rangle \geq \mu \|\mathbf{x}_p - \mathbf{x}\|^2$  for any  $\mathbf{x} \in \mathcal{X}$ .*
- (3) *We say  $g$  satisfies  $\mu$ -error bound (EB) condition if  $\|\nabla g(\mathbf{x})\| \geq \mu \|\mathbf{x}_p - \mathbf{x}\|$  for any  $\mathbf{x} \in \mathcal{X}$ .*
- (4) *We say  $g$  satisfies  $\mu$ -quadratic growth (QG) condition if  $g(\mathbf{x}) - \min_{\mathbf{x}'} g(\mathbf{x}') \geq \frac{\mu}{2} \|\mathbf{x}_p - \mathbf{x}\|^2$  for any  $\mathbf{x} \in \mathcal{X}$ .*

**Proposition 7.** *From Definition 1, The following implications are true.*

- $\mu$ -SC implies  $\mu$ -PŁ and  $\mu$ -RSI.
- $\mu$ -PŁ implies  $\mu$ -QG and  $\mu$ -EB.
- $\mu$ -RSI implies  $\mu$ -EB.
- $\mu$ -EB and  $L$ -smoothness together imply  $(\mu^2/L)$ -PŁ.



*Proof.* Most of the proofs originated from Karimi et al. (2016, Theorem 2).

(SC  $\Rightarrow$  PŁ) Substitute  $\mathbf{x}$  to  $\mathbf{x}_p$  and  $\mathbf{x}'$  to  $\mathbf{x}$ , respectively, from Definition 1.(1).

(PŁ  $\Rightarrow$  QG & EB) See the proof in Karimi et al. (2016, Theorem 2)

(SC  $\Rightarrow$  RSI) We know  $\mu$ -SC  $\Rightarrow$   $\mu$ -PŁ  $\Rightarrow$   $\mu$ -QG. From Definition 1.(1) & 1.(4),

$$\begin{aligned} \langle \nabla g(\mathbf{x}), \mathbf{x} - \mathbf{x}_p \rangle &\stackrel{\text{SC}}{\geq} g(\mathbf{x}) - g(\mathbf{x}_p) + \frac{\mu}{2} \|\mathbf{x}_p - \mathbf{x}\|^2 \\ &\stackrel{\text{QG}}{\geq} \frac{\mu}{2} \|\mathbf{x}_p - \mathbf{x}\|^2 + \frac{\mu}{2} \|\mathbf{x}_p - \mathbf{x}\|^2 = \mu \|\mathbf{x}_p - \mathbf{x}\|^2. \end{aligned}$$

This implies  $\mu$ -RSI.

(RSI  $\Rightarrow$  EB) See the proof in Karimi et al. (2016, Theorem 2).

(EB & smooth  $\Rightarrow$  PŁ) We use  $\nabla g(\mathbf{x}_p) = \mathbf{0}$ . By  $L$ -smoothness and  $\mu$ -EB condition,

$$\begin{aligned} g(\mathbf{x}) - g(\mathbf{x}_p) &\stackrel{\text{smooth}}{\leq} \langle \nabla g(\mathbf{x}_p), \mathbf{x} - \mathbf{x}_p \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_p\|^2 = \frac{L}{2} \|\mathbf{x} - \mathbf{x}_p\|^2 \\ &\stackrel{\text{EB}}{\leq} \frac{L}{2\mu^2} \|\nabla g(\mathbf{x})\|^2. \end{aligned}$$

This implies  $(\mu^2/L)$ -PŁ condition on  $g$ .  $\square$

**Proposition 8** (Lipschitz continuity-like property of  $\mathbf{y}^*(\mathbf{x})$ ). *For an  $L$ -smooth function  $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ , suppose  $-g(\mathbf{x}; \cdot)$  is  $\mu_2$ -PŁ. Let  $\kappa_2 = L/\mu_2$ .*

*Consider any  $\mathbf{x}_0, \mathbf{x}_1 \in \mathcal{X}$ . For any  $\mathbf{y}_0^* \in \mathcal{Y}_{\mathbf{x}_0}^* = \arg \max_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}_0; \mathbf{y})$ , there exists a  $\mathbf{y}_1^* \in \mathcal{Y}_{\mathbf{x}_1}^* = \arg \max_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}_1; \mathbf{y})$  such that  $\|\mathbf{y}_0^* - \mathbf{y}_1^*\| \leq \kappa_2 \|\mathbf{x}_0 - \mathbf{x}_1\|$ .*

*In fact, it is enough to choose  $\mathbf{y}_1^*$  as a projection of  $\mathbf{y}_0^*$  onto the set  $\mathcal{Y}_{\mathbf{x}_1}^*$ , namely,  $\mathbf{y}_1^* \in \Pi_{\mathcal{Y}_{\mathbf{x}_1}^*}(\mathbf{y}_0^*)$ .*

*Proof.* We borrow the proof from Nouiheed et al. (2019, Lemma A.3).

Recall  $\Phi(\mathbf{x}) := \max_{\mathbf{y}' \in \mathcal{Y}} g(\mathbf{x}; \mathbf{y}')$ . By PŁ inequality and smoothness of  $g$ ,

$$\begin{aligned} 2\mu_2 (\Phi(\mathbf{x}_1) - g(\mathbf{x}_1; \mathbf{y}_0^*)) &\leq \|\nabla_2 g(\mathbf{x}_1; \mathbf{y}_0^*)\|^2 \\ &= \|\nabla_2 g(\mathbf{x}_1; \mathbf{y}_0^*) - \nabla_2 g(\mathbf{x}_0; \mathbf{y}_0^*)\|^2 \leq L^2 \|\mathbf{x}_1 - \mathbf{x}_0\|^2. \end{aligned}$$

The second equality applies  $\nabla_2 g(\mathbf{x}_0; \mathbf{y}_0^*) = \mathbf{0}$ , since  $\mathbf{y}_0^* \in \arg \max_{\mathbf{y}} g(\mathbf{x}_0; \mathbf{y})$ .

Moreover, note that  $-g(\mathbf{x}_1; \cdot)$  satisfies  $\mu_2$ -QG condition ( $\cdot$ : Proposition 7). To apply this, we utilize our choice of  $\mathbf{y}_1^*$ :

$$\Phi(\mathbf{x}_1) - g(\mathbf{x}_1; \mathbf{y}_0^*) \geq \frac{\mu_2}{2} \|\mathbf{y}_1^* - \mathbf{y}_0^*\|^2.$$

As a result, we have  $\mu_2^2 \|\mathbf{y}_0^* - \mathbf{y}_1^*\|^2 \leq L^2 \|\mathbf{x}_0 - \mathbf{x}_1\|^2$ . This completes the proof.  $\square$

**Proposition 9** (Smoothness of primal function). *Consider the same function  $g$  as Proposition 8. Then, the function  $\Phi(\mathbf{x}) := \max_{\mathbf{y}' \in \mathcal{Y}} g(\mathbf{x}; \mathbf{y}')$  is differentiable with*

$$\nabla \Phi(\mathbf{x}) = \nabla_1 g(\mathbf{x}; \mathbf{y}^*(\mathbf{x})), \quad \text{regardless of the choice of } \mathbf{y}^*(\mathbf{x}) \in \arg \max_{\mathbf{y}' \in \mathcal{Y}} g(\mathbf{x}; \mathbf{y}').$$

*Moreover,  $\Phi$  is  $L(\kappa_2 + 1)$ -smooth, where  $\kappa_2 = L/\mu_2$ .*

*Proof.* This is already proved in Lemma A.5 of Nouiheed et al. (2019). However, we present a bit different proof without using second-order Taylor expansion. To start, recall  $\mathcal{Y}_{\mathbf{x}}^* := \arg \max_{\mathbf{y} \in \mathcal{Y}} g(\mathbf{x}; \mathbf{y})$ . That is, we could choose any  $\mathbf{y}^*(\mathbf{x}) \in \mathcal{Y}_{\mathbf{x}}^*$ .

We first show the differentiability of  $\Phi$ . Fix a unit vector  $\mathbf{u} \in \mathcal{X} = \mathbb{R}^{d_x}$ :  $\|\mathbf{u}\| = 1$ . Let any  $h > 0$ . We first claim that there exists a path  $\mathbf{p} : (-h, h) \rightarrow \mathcal{Y} = \mathbb{R}^{d_y}$  which is continuous at  $t = 0$  and  $\mathbf{p}(t) \in \mathcal{Y}_{(\mathbf{x}+t\mathbf{u})}^*$ . In fact, let  $\mathbf{p}(t)$  be a projection of  $\mathbf{y}^*(\mathbf{x})$  (that we chose) onto the set  $\mathcal{Y}_{(\mathbf{x}+t\mathbf{u})}^*$ .

Then,  $\mathbf{p}(0) = \mathbf{y}^*(\mathbf{x})$ , and by Proposition 8, we have  $\|\mathbf{p}(0) - \mathbf{p}(t)\| \leq \kappa_2 \|\mathbf{x} - (\mathbf{x} + t\mathbf{u})\| = \kappa_2 t$ . This shows the continuity of  $\mathbf{p}(t)$  at  $t = 0$ . Now, note that there exists a  $t_1 \in (0, h)$  such that,

$$\begin{aligned} & \Phi(\mathbf{x} + h\mathbf{u}) - \Phi(\mathbf{x}) \\ &= g(\mathbf{x} + h\mathbf{u}; \mathbf{p}(h)) - g(\mathbf{x}; \mathbf{p}(0)) \\ &= \{g(\mathbf{x} + h\mathbf{u}; \mathbf{p}(h)) - g(\mathbf{x} + h\mathbf{u}; \mathbf{p}(0))\} + \{g(\mathbf{x} + h\mathbf{u}; \mathbf{p}(0)) - g(\mathbf{x}; \mathbf{p}(0))\} \\ &\geq 0 + \langle \nabla_1 g(\mathbf{x} + t_1\mathbf{u}; \mathbf{p}(0)), h\mathbf{u} \rangle, \end{aligned}$$

by mean value theorem (applied to the first argument). We have the inequality at the last line because  $g(\mathbf{x} + h\mathbf{u}; \mathbf{p}(h)) \geq g(\mathbf{x} + h\mathbf{u}; \mathbf{p}(0))$ , since  $\mathbf{p}(h) \in \mathcal{Y}_{(\mathbf{x}+h\mathbf{u})}^*$ . With a similar logic, there exists a  $t_2 \in (0, h)$  such that,

$$\begin{aligned} & \Phi(\mathbf{x} + h\mathbf{u}) - \Phi(\mathbf{x}) \\ &= g(\mathbf{x} + h\mathbf{u}; \mathbf{p}(h)) - g(\mathbf{x}; \mathbf{p}(0)) \\ &= \{g(\mathbf{x} + h\mathbf{u}; \mathbf{p}(h)) - g(\mathbf{x}; \mathbf{p}(h))\} + \{g(\mathbf{x}; \mathbf{p}(h)) - g(\mathbf{x}; \mathbf{p}(0))\} \\ &\leq \langle \nabla_1 g(\mathbf{x} + t_2\mathbf{u}; \mathbf{p}(h)), h\mathbf{u} \rangle + 0. \end{aligned}$$

To combine these two inequalities into a single line,

$$\langle \nabla_1 g(\mathbf{x} + t_1\mathbf{u}; \mathbf{p}(0)), \mathbf{u} \rangle \leq \frac{\Phi(\mathbf{x} + h\mathbf{u}) - \Phi(\mathbf{x})}{h} \leq \langle \nabla_1 g(\mathbf{x} + t_2\mathbf{u}; \mathbf{p}(h)), \mathbf{u} \rangle.$$

Using the continuity of  $\mathbf{p}(\cdot)$  and  $\nabla_1 g(\cdot; \cdot)$  ( $\cdot$ :  $g$  has Lipschitz continuous gradient), we can deduce that the directional derivative of  $\Phi$  in a direction  $\mathbf{u}$  (denoted by  $D_{\mathbf{u}}\Phi$ ) is in fact

$$D_{\mathbf{u}}\Phi(\mathbf{x}) = \langle \nabla_1 g(\mathbf{x}; \mathbf{y}^*(\mathbf{x})), \mathbf{u} \rangle,$$

by taking the limit  $h \rightarrow 0+$ . Since  $\mathbf{u}$  is arbitrary, we can conclude that  $\nabla\Phi(\mathbf{x}) = \nabla_1 g(\mathbf{x}; \mathbf{y}^*(\mathbf{x}))$ .

The proof of Lipschitz smoothness of  $\Phi$  exactly follows the proof by Nouiehed et al. (2019). Consider any  $\mathbf{x}_0, \mathbf{x}_1 \in \mathcal{X}$ . As in Proposition 8, choose any  $\mathbf{y}_0^* \in \mathcal{Y}_{\mathbf{x}_0}^*$  and  $\mathbf{y}_1^* \in \Pi_{\mathcal{Y}_{\mathbf{x}_1}^*}(\mathbf{y}_0^*)$ . Then,

$$\begin{aligned} & \|\nabla\Phi(\mathbf{x}_0) - \nabla\Phi(\mathbf{x}_1)\| \\ &= \|\nabla_1 g(\mathbf{x}_0; \mathbf{y}_0^*) - \nabla_1 g(\mathbf{x}_1; \mathbf{y}_1^*)\| \\ &\leq \|\nabla_1 g(\mathbf{x}_0; \mathbf{y}_0^*) - \nabla_1 g(\mathbf{x}_1; \mathbf{y}_0^*)\| + \|\nabla_1 g(\mathbf{x}_1; \mathbf{y}_0^*) - \nabla_1 g(\mathbf{x}_1; \mathbf{y}_1^*)\| \\ &\leq L \{\|\mathbf{x}_0 - \mathbf{x}_1\| + \|\mathbf{y}_0^* - \mathbf{y}_1^*\|\} \\ &\leq L(1 + \kappa_2) \|\mathbf{x}_0 - \mathbf{x}_1\|. \end{aligned}$$

The last inequality holds because of Proposition 8.  $\square$

**Proposition 10** ( $x$ -side PŁ  $\Rightarrow$  primal PŁ). *Suppose  $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is  $L$ -smooth and two-sided PŁ with constants  $\mu_1$  and  $\mu_2$ . Then,  $g$  satisfies primal PŁ condition: the function  $\Phi(\mathbf{x}) := \max_{\mathbf{y}' \in \mathcal{Y}} g(\mathbf{x}; \mathbf{y}')$  is  $\mu_1$ -PŁ. As a result, a smooth two-sided PŁ function is PŁ( $\Phi$ )-PŁ.*

*Proof.* See Lemma A.3 of Yang et al. (2020).  $\square$

**Definition 2.** Consider  $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ . Then, the point  $(\mathbf{x}^*; \mathbf{y}^*) \in \mathcal{X} \times \mathcal{Y}$  is called

- (i) a stationary point of  $g$  if  $\nabla_1 g(\mathbf{x}^*; \mathbf{y}^*) = \nabla_2 g(\mathbf{x}^*; \mathbf{y}^*) = 0$ .
- (ii) a saddle point of  $g$  if  $g(\mathbf{x}^*; \mathbf{y}) \leq g(\mathbf{x}^*; \mathbf{y}^*) \leq g(\mathbf{x}; \mathbf{y}^*)$  for all  $\mathbf{x}, \mathbf{y}$ .
- (iii) a global minimax point of  $g$  if  $g(\mathbf{x}^*; \mathbf{y}) \leq g(\mathbf{x}^*; \mathbf{y}^*) \leq \max_{\mathbf{y}'} g(\mathbf{x}; \mathbf{y}')$  for all  $\mathbf{x}, \mathbf{y}$ .
- (iv) a global maximin point of  $g$  if  $\min_{\mathbf{x}'} g(\mathbf{x}'; \mathbf{y}) \leq g(\mathbf{x}^*; \mathbf{y}^*) \leq g(\mathbf{x}; \mathbf{y}^*)$  for all  $\mathbf{x}, \mathbf{y}$ .

**Proposition 11.** Consider a function  $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ .

(1) In general, a saddle point of  $g$  is a global minimax/maximin point.

(2) Let  $\Phi(\mathbf{x}) := \max_{\mathbf{y}} g(\mathbf{x}; \mathbf{y})$  and  $\Phi^* := \min_{\mathbf{x}} \Phi(\mathbf{x})$  be well-defined. Let  $\lambda > 0$  be a constant. In general, a point  $(\mathbf{x}^*; \mathbf{y}^*)$  is a global minimax point of  $g$  if and only if

$$V_{\lambda}(\mathbf{x}^*; \mathbf{y}^*) := \lambda[\Phi(\mathbf{x}^*) - \Phi^*] + [\Phi(\mathbf{x}^*) - g(\mathbf{x}^*; \mathbf{y}^*)] = 0.$$

- (3) If  $g$  is smooth nonconvex-PŁ, then a global minimax point is a stationary point.
- (4) If  $g$  is PŁ( $\Phi$ )-PŁ, then there exists a global minimax point  $(\mathbf{x}^*; \mathbf{y}^*)$  of  $g$ . As a result, if  $g$  is also smooth, then the point  $(\mathbf{x}^*; \mathbf{y}^*)$  is a stationary point.
- (5) If  $g$  is smooth two-sided PŁ, every stationary point is a saddle point. As a result, there exists a saddle point  $(\mathbf{x}^*; \mathbf{y}^*)$  of  $g$ .

In particular, smooth two-sided PŁ functions enjoy the “minimax theorem,” which establishes “minimax = maximin.”

*Proof.* (1) (saddle point  $\Rightarrow$  global minimax & global maximin) This is straightforward by the definitions: for any  $\mathbf{x}$  and  $\mathbf{y}$ ,

$$\min_{\mathbf{x}'} g(\mathbf{x}'; \mathbf{y}) \leq g(\mathbf{x}^*; \mathbf{y}) \leq g(\mathbf{x}^*; \mathbf{y}^*) \leq g(\mathbf{x}; \mathbf{y}^*) \leq \max_{\mathbf{y}'} g(\mathbf{x}; \mathbf{y}').$$

(2) (global minimax  $\iff V_\lambda = 0$ ) The terms  $\Phi(\mathbf{x}) - \Phi^*$  and  $\Phi(\mathbf{x}) - g(\mathbf{x}; \mathbf{y})$  are non-negative. Hence,  $V_\lambda(\mathbf{x}; \mathbf{y})$  is non-negative, and  $V_\lambda(\mathbf{x}^*; \mathbf{y}^*) = 0$  if and only if  $\Phi^* = \Phi(\mathbf{x}^*) = g(\mathbf{x}^*; \mathbf{y}^*)$ , which is equivalent to the global minimax point condition.

(3) (smooth nonconvex-PŁ: global minimax  $\Rightarrow$  stationary) Suppose  $(\mathbf{x}^*; \mathbf{y}^*)$  is a global minimax point. Since  $g(\mathbf{x}^*; \mathbf{y}) \leq g(\mathbf{x}^*; \mathbf{y}^*)$  for any  $\mathbf{y}$ ,  $\Phi(\mathbf{x}^*) = \max_{\mathbf{y}} g(\mathbf{x}^*; \mathbf{y}) = g(\mathbf{x}^*; \mathbf{y}^*)$ . Thus,  $\Phi$  has a minimum  $g(\mathbf{x}^*; \mathbf{y}^*)$  at  $\mathbf{x} = \mathbf{x}^*$ . By Proposition 9,  $\Phi(\cdot)$  is a differentiable function and we have

$$\nabla_1 g(\mathbf{x}^*; \mathbf{y}^*) = \nabla \Phi(\mathbf{x}^*) = 0.$$

Also, since a differentiable function  $g(\mathbf{x}^*; \mathbf{y})$  has a maximum at  $\mathbf{y} = \mathbf{y}^*$ , we also have  $\nabla_2 g(\mathbf{x}^*; \mathbf{y}^*) = 0$ . Therefore,  $(\mathbf{x}^*; \mathbf{y}^*)$  is a stationary point.

(4) (PŁ( $\Phi$ )-PŁ:  $\exists$  global minimax) Let  $\mathbf{x}^* \in \arg \min_{\mathbf{x}} \Phi(\mathbf{x})$  and  $\mathbf{y}^* \in \arg \max_{\mathbf{y}} f(\mathbf{x}^*; \mathbf{y})$ . Then,  $f(\mathbf{x}^*, \mathbf{y}^*) = \Phi(\mathbf{x}^*) = \Phi^*$ . as noted in (2),  $(\mathbf{x}^*, \mathbf{y}^*)$  is a global minimax point. By (3), it is in fact a stationary point, when  $g$  is smooth as well.

(5) (smooth two-sided PŁ: stationary  $\Rightarrow$  saddle) Let  $(\mathbf{x}^*; \mathbf{y}^*)$  be a stationary point. By PŁ inequalities, for any  $\mathbf{x}$  and  $\mathbf{y}$ ,

$$\begin{aligned} 0 &= \|\nabla_2 g(\mathbf{x}^*; \mathbf{y}^*)\|^2 \geq 2\mu_2(\max_{\mathbf{y}} g(\mathbf{x}^*; \mathbf{y}) - g(\mathbf{x}^*; \mathbf{y}^*)) \geq 0, \\ 0 &= \|\nabla_1 g(\mathbf{x}^*; \mathbf{y}^*)\|^2 \geq 2\mu_1(g(\mathbf{x}^*; \mathbf{y}^*) - \min_{\mathbf{x}} g(\mathbf{x}; \mathbf{y}^*)) \geq 0. \end{aligned}$$

Since  $\mu_1, \mu_2 > 0$ , these imply  $\max_{\mathbf{y}} g(\mathbf{x}^*; \mathbf{y}) = g(\mathbf{x}^*; \mathbf{y}^*) = \min_{\mathbf{x}} g(\mathbf{x}; \mathbf{y}^*)$ . Thus,  $(\mathbf{x}^*; \mathbf{y}^*)$  is a saddle point. Note that (4) and Proposition 10 together proves the existence of a stationary point of  $g$ . Therefore, there must exist a saddle point, which is also pointed out by Guo et al. (2020, Lemma 8). This concludes the proof.  $\square$

We remark that, in the proof above, (3) is false for general (nonconvex-nonconcave) functions. Only local minimax point can ensure stationarity (Jin et al., 2020). As remarked by Jin et al. (2020) (Figure 2 of their paper), the function  $xy - \cos(y)$  has non-stationary global minimax points  $(0, \pm\pi)$ .

The following two propositions are for showing that general two-sided PŁ function may not have a differential Stackelberg equilibrium defined as Li et al. (2022, Definition 3.1).

**Proposition 12.** Let  $g$  be a  $\mu$ -strongly convex function on  $\mathbb{R}^n$ . Consider any matrix  $\mathbf{M} \in \mathbb{R}^{n \times m}$  with a positive rank. Suppose that  $\theta$  is the smallest nonzero singular value of  $\mathbf{M}$ . Then  $g(\mathbf{M}\mathbf{y})$  is a  $\mu\theta^2$ -PŁ function of  $\mathbf{y} \in \mathbb{R}^m$ .

*Proof.* See Karimi et al. (2016, Appendix B) for the proof.  $\square$

**Proposition 13.** Consider a twice continuously differentiable strongly-convex-strongly-concave function  $h : \mathbb{R}^r \times \mathbb{R}^s \rightarrow \mathbb{R}$ . That is, for some constants  $\mu_1, \mu_2 > 0$ ,  $h(\mathbf{x}; \mathbf{y})$  is  $\mu_1$ -strongly-convex in  $\mathbf{x}$  and  $-h(\mathbf{x}; \mathbf{y})$  is  $\mu_2$ -strongly-convex in  $\mathbf{y}$ . Let  $(\mathbf{x}^*; \mathbf{y}^*)$  be the unique stationary point of  $h$ . Of

course, it is a **differential Stackelberg equilibrium** of  $h$ . That is, if the hessian matrix  $\nabla^2 h(\mathbf{x}^*; \mathbf{y}^*)$  at that point is written as

$$\nabla^2 h(\mathbf{x}^*; \mathbf{y}^*) = \begin{bmatrix} \nabla_{1,1}^2 h(\mathbf{x}^*; \mathbf{y}^*) & \nabla_{1,2}^2 h(\mathbf{x}^*; \mathbf{y}^*) \\ \nabla_{2,1}^2 h(\mathbf{x}^*; \mathbf{y}^*) & \nabla_{2,2}^2 h(\mathbf{x}^*; \mathbf{y}^*) \end{bmatrix} = \begin{bmatrix} \mathbf{C} & \mathbf{B} \\ \mathbf{B}^\top & -\mathbf{A} \end{bmatrix},$$

then  $\mathbf{A}$  and  $\mathbf{C} - \mathbf{B}\mathbf{A}^{-1}\mathbf{B}^\top$  are both positive definite matrices. Consider a function  $g : \mathbb{R}^p \times \mathbb{R}^q \rightarrow \mathbb{R}$  defined by  $g(\mathbf{x}; \mathbf{y}) = h(\mathbf{M}\mathbf{x}; \mathbf{N}\mathbf{y})$  for some matrices  $\mathbf{M} \in \mathbb{R}^{r \times p}$ ,  $\mathbf{N} \in \mathbb{R}^{s \times q}$ . Then,  $g$  is two-sided PL. Moreover, each stationary point of  $g$  may not be a differential Stackelberg equilibrium in general, for example, when  $s < q$ .

*Proof.* Because of Proposition 12,  $g$  is clearly a two-sided PL function.

If  $(\mathbf{x}; \mathbf{y})$  is a stationary point of  $g$ , then it must be an element of an affine set  $\{(\mathbf{x}; \mathbf{y}) \in \mathbb{R}^p \times \mathbb{R}^q : \mathbf{M}\mathbf{x} = \mathbf{x}^*; \mathbf{N}\mathbf{y} = \mathbf{y}^*\}$ . This is because

$$\nabla g(\mathbf{x}; \mathbf{y}) = \begin{bmatrix} \nabla_1 g(\mathbf{x}; \mathbf{y}) \\ \nabla_2 g(\mathbf{x}; \mathbf{y}) \end{bmatrix} = \begin{bmatrix} \mathbf{M}^\top \nabla_1 h(\mathbf{M}\mathbf{x}; \mathbf{N}\mathbf{y}) \\ \mathbf{N}^\top \nabla_2 h(\mathbf{M}\mathbf{x}; \mathbf{N}\mathbf{y}) \end{bmatrix} = \mathbf{0}$$

if and only if  $\nabla_1 h(\mathbf{M}\mathbf{x}; \mathbf{N}\mathbf{y}) = \mathbf{0}$  and  $\nabla_2 h(\mathbf{M}\mathbf{x}; \mathbf{N}\mathbf{y}) = \mathbf{0}$ , being equivalent to  $\mathbf{M}\mathbf{x} = \mathbf{x}^*$  and  $\mathbf{N}\mathbf{y} = \mathbf{y}^*$ . Furthermore, the hessian of  $g$  at  $(\mathbf{x}; \mathbf{y})$  is

$$\begin{aligned} \nabla^2 g(\mathbf{x}; \mathbf{y}) &= \begin{bmatrix} \mathbf{M}^\top \nabla_{1,1}^2 h(\mathbf{M}\mathbf{x}; \mathbf{N}\mathbf{y}) \mathbf{M} & \mathbf{M}^\top \nabla_{1,2}^2 h(\mathbf{M}\mathbf{x}; \mathbf{N}\mathbf{y}) \mathbf{N} \\ \mathbf{N}^\top \nabla_{2,1}^2 h(\mathbf{M}\mathbf{x}; \mathbf{N}\mathbf{y}) \mathbf{M} & \mathbf{N}^\top \nabla_{2,2}^2 h(\mathbf{M}\mathbf{x}; \mathbf{N}\mathbf{y}) \mathbf{N} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{M}^\top \mathbf{C} \mathbf{M} & \mathbf{M}^\top \mathbf{B} \mathbf{N} \\ (\mathbf{M}^\top \mathbf{B} \mathbf{N})^\top & -\mathbf{N}^\top \mathbf{A} \mathbf{N} \end{bmatrix}. \end{aligned}$$

If  $s < q$ , the  $q \times q$  matrix  $\mathbf{N}^\top \mathbf{A} \mathbf{N}$  cannot have a full rank, thereby it cannot be even invertible. This implies the stationary point  $(\mathbf{x}; \mathbf{y})$  cannot be a differential Stackelberg equilibrium.  $\square$

## B.2 WITHOUT-REPLACEMENT SAMPLING

In this subsection, we provide a useful proposition for analysis of mini-batching approach under without-replacement sampling. We consider the case of mutually disjoint mini-batches in a whole epoch, not only applying without-replacement sampling to each individual mini-batch.

Consider a collection of  $n$  vectors  $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^d$ . Suppose we uniformly randomly sample a permutation  $\sigma : [n] \rightarrow [n]$ ; i.e.,  $\sigma \sim \text{Unif}(\mathbb{S}_n)$ . Define

$$\mathbf{m} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \quad (\text{sample mean}) \quad \text{and} \quad \tau^2 = \frac{1}{n} \sum_{i=1}^n \|\mathbf{v}_i - \mathbf{m}\|^2 \quad (\text{sample variance}).$$

Fix any  $b \in [n]$  and let  $n = b(q-1) + s$  for some integers  $q \geq 1$  and  $s \in [b]$ . Now, divide the indices  $[n]$  into  $q$  batches, with exactly  $b$  items per batch (except for the last batch when  $s < b$ ), as follows:

$$\mathcal{W}_t = \{\sigma(j) : b(t-1) < j \leq bt, j \in [n]\} \quad (t \in [q]).$$

For each batch  $\mathcal{W}_t$ , define

$$\mathbf{w}_t = \frac{1}{|\mathcal{W}_t|} \sum_{i \in \mathcal{W}_t} \mathbf{v}_i \quad (\text{batch mean}).$$

For any  $k \in [q-1]$ , define

$$\mathbf{m}_k := \frac{1}{k} \sum_{t=1}^k \mathbf{w}_t \quad (\text{accumulative average of batch means over } 1 \leq t \leq k).$$

Of course, we may simply take  $\mathbf{m}_q = \mathbf{m}$  (deterministically) for  $k = q$ . Thus, because of the randomness of  $\sigma$ , we can obtain the mean (vector) and the variance (scalar) of  $\mathbf{m}_k$  as follows.

**Proposition 14** (Without-replacement sampling). *Given the setup above, for any  $k < q$  and  $n > 1$ ,*

$$\mathbb{E}[\mathbf{m}_k] = \mathbf{m} \quad \text{and} \quad \mathbb{E}[\|\mathbf{m}_k - \mathbf{m}\|^2] = \frac{(n-bk)}{bk(n-1)} \tau^2.$$

(Of course, if  $k = q$  or  $n = 1 = q$ ,  $\mathbb{E}[\|\mathbf{m}_q - \mathbf{m}\|^2] = 0$  since  $\mathbf{m}_q = \mathbf{m}$ .)

**Remark.** As a special case, if  $n = bq$  (namely,  $b$  divides  $n$  and  $s = b$ ), then for any  $k \leq q$ ,

$$\mathbb{E} \left[ \|\mathbf{m}_k - \mathbf{m}\|^2 \right] = \frac{(q-k)}{k(n-1)} \tau^2.$$

If we further assume  $b = s = 1$  and  $q = n$ , this proposition recovers Lemma 1 of [Mishchenko et al. \(2020\)](#).

*Proof of Proposition 14.* Since  $\sigma$  is a uniformly randomly sampled permutation, it is easy to obtain that

$$\mathbb{E}[\mathbf{v}_{\sigma(i)}] = \mathbb{E}[\mathbf{w}_t] = \mathbb{E}[\mathbf{m}_k] = \mathbf{m},$$

for any  $i \in [n]$ ,  $t \in [q]$ , and  $k \in [q]$ .

The covariances between  $\mathbf{v}_{\sigma(i)}$ 's can be deduced from the proof by [Mishchenko et al. \(2020, Lemma 1\)](#) as follows:

$$\text{Cov}(\mathbf{v}_{\sigma(i)}, \mathbf{v}_{\sigma(j)}) := \mathbb{E} \left[ \langle \mathbf{v}_{\sigma(i)} - \mathbf{m}, \mathbf{v}_{\sigma(j)} - \mathbf{m} \rangle \right] = \begin{cases} -\frac{\tau^2}{n-1}, & \text{if } i \neq j, \\ \tau^2 & \text{if } i = j. \end{cases}$$

Thus, for each  $t \in [q]$ , the variance of  $\mathbf{w}_t$  is obtained as

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{w}_t - \mathbf{m}\|^2 \right] &= \mathbb{E} \left[ \left\| \frac{1}{|\mathcal{W}_t|} \sum_{i \in \mathcal{W}_t} (\mathbf{v}_i - \mathbf{m}) \right\|^2 \right] \\ &= \frac{1}{|\mathcal{W}_t|^2} \left\{ \sum_{i \in \mathcal{W}_t} \mathbb{E} \left[ \|\mathbf{v}_i - \mathbf{m}\|^2 \right] + \sum_{\substack{i, j \in \mathcal{W}_t \\ i \neq j}} \text{Cov}(\mathbf{v}_i, \mathbf{v}_j) \right\} \\ &= \frac{1}{|\mathcal{W}_t|^2} \left\{ |\mathcal{W}_t| \tau^2 + |\mathcal{W}_t| (|\mathcal{W}_t| - 1) \left( -\frac{\tau^2}{n-1} \right) \right\} = \frac{n - |\mathcal{W}_t|}{|\mathcal{W}_t|(n-1)} \tau^2, \end{aligned}$$

which can also be directly deduced by Lemma 1 of [Mishchenko et al. \(2020\)](#). We notice that this does not depend on the size of the batch  $\mathcal{W}_t$ .

Next, we look at the covariances between distinct  $\mathbf{w}_t$ 's. For a pair of distinct integers  $t, u \in [q]$ , by the bi-linearity of covariance,

$$\begin{aligned} \text{Cov}(\mathbf{w}_t, \mathbf{w}_u) &= \frac{1}{|\mathcal{W}_t| \cdot |\mathcal{W}_u|} \sum_{(i, j) \in \mathcal{W}_t \times \mathcal{W}_u} \text{Cov}(\mathbf{v}_i, \mathbf{v}_j) \\ &= \frac{1}{|\mathcal{W}_t| \cdot |\mathcal{W}_u|} \sum_{(i, j) \in \mathcal{W}_t \times \mathcal{W}_u} \left( -\frac{\tau^2}{n-1} \right) = -\frac{\tau^2}{n-1}. \end{aligned}$$

The second equality holds because  $\mathcal{W}_t$  and  $\mathcal{W}_u$  are a disjoint set of integers whenever  $t \neq u$ .

Now, fix any  $k \in [q-1]$ . Note that, by our mini-batching strategy,  $|\mathcal{W}_t| = b$  for every  $t < q$ . Therefore, by definition of  $\mathbf{m}_k$ ,

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{m}_k - \mathbf{m}\|^2 \right] &= \mathbb{E} \left[ \left\| \frac{1}{k} \sum_{t=1}^k (\mathbf{w}_t - \mathbf{m}) \right\|^2 \right] \\ &= \frac{1}{k^2} \left\{ \sum_{t=1}^k \mathbb{E} \left[ \|\mathbf{w}_t - \mathbf{m}\|^2 \right] + \sum_{\substack{t, u \in [k] \\ t \neq u}} \text{Cov}(\mathbf{w}_t, \mathbf{w}_u) \right\} \\ &= \frac{1}{k^2} \left\{ k \cdot \left( \frac{n-b}{b(n-1)} \tau^2 \right) + k(k-1) \cdot \left( -\frac{\tau^2}{n-1} \right) \right\} = \frac{n-bk}{bk(n-1)} \tau^2. \end{aligned}$$

□

### B.3 BASIC RECURRENCE INEQUALITY

In this subsection, we present a basic result of a recurrence inequality. It serves as a stepping-stone of our convergence bound, particularly at the end of the proof (Appendix C.5).

**Proposition 15.** *Let  $\{a_k\}_{k=1}^{\infty}$  be a sequence of non-negative numbers satisfying the following recurrence inequality:*

$$a_{k+1} \leq (1 - b\eta)a_k + c\eta^{m+1},$$

where  $b, c$ , and  $\eta$  are non-negative real numbers such that  $b\eta \in (0, 1)$ , and  $m$  is a non-negative integer. Then, for any integer  $K \geq 1$ , we have

$$a_{K+1} \leq (1 - b\eta)^K a_1 + c\eta^m/b.$$

*Proof.* We proceed with induction on  $K = 0, 1, 2, \dots$ . Note that

$$a_1 \leq (1 - b\eta)^0 a_1 + c\eta^m/b.$$

This shows the case when  $K = 0$ . On the other hand, if  $K \geq 1$ , by an inductive assumption,

$$\begin{aligned} a_{K+1} &\leq (1 - b\eta)a_K + c\eta^{m+1} \\ &\leq (1 - b\eta) \cdot ((1 - b\eta)^{K-1} a_1 + c\eta^m/b) + c\eta^{m+1} \\ &= (1 - b\eta)^K a_1 + c\eta^m/b. \end{aligned}$$

□

## C PROOFS FOR (MINI-BATCH) SIMULTANEOUS SGDA-RR

In this appendix, we provide a convergence analysis for the mini-batch **simSGDA-RR** (Algorithm 2) on both general nonconvex-PŁ problems and primal-PŁ-PŁ problems. The two cases mostly share the same proof strategies; they only diverge at the end of the proofs. The proof is long; we first provide the sketch of proof in subsection C.1; then, we provide the full proof by dividing it into 4 follow-up subsections of this appendix. The proof for the alternating counterpart (minibatch altSGDA-RR) can be done with some modifications illustrated in Appendix D. All technical propositions required for the proofs can be found in Appendix B.

### C.1 WARM-UP: PROOF SKETCH FOR $b = 1$

Here we simply consider the proofs of Theorem 1 and 2 for **simSGDA-RR**, which is a fully stochastic case (mini-batches of size  $b = 1$ ). The proofs for altSGDA-RR can be done with slight modifications.

We start the proof by aggregating all updates throughout an epoch to obtain an “epoch-wise” update:

$$\begin{aligned} \mathbf{x}_0^{k+1} &= \mathbf{x}_0^k - n\alpha \mathbf{g}^k, \quad \mathbf{g}^k = \frac{1}{n} \sum_{i=1}^n \nabla_1 f_{\sigma_k(i)}(\mathbf{z}_{i-1}^k), \\ \mathbf{y}_0^{k+1} &= \mathbf{y}_0^k + n\beta \mathbf{h}^k, \quad \mathbf{h}^k = \frac{1}{n} \sum_{i=1}^n \nabla_2 f_{\sigma_k(i)}(\mathbf{z}_{i-1}^k). \end{aligned}$$

The reason is that the sampled components in each epoch are dependent to each other so that it is much harder to deal with each iteration individually. The strategy of update-aggregation is quite general for analysis of optimization algorithms involving without-replacement sampling (Ahn et al., 2020; Mishchenko et al., 2020; Nguyen et al., 2021; Das et al., 2022). We assume that the intermediate iterates  $\mathbf{z}_1^k, \dots, \mathbf{z}_n^k$  stay close to the starting iterate  $\mathbf{z}_0^k$  of an epoch  $k$ , which can be ensured by small step sizes. Then, we can approximate the aggregated epoch of SGDA-RR as a step of simGDA applied to  $f = \frac{1}{n} \sum_{i=1}^n f_i$ , with approximations of  $\mathbf{g}^k \approx \nabla_1 f(\mathbf{z}_0^k)$  and  $\mathbf{h}^k \approx \nabla_2 f(\mathbf{z}_0^k)$ .

With Assumptions 1 and 4, note that the primal function  $\Phi(\cdot)$  is  $(L + L^2/\mu_2)$ -smooth (Proposition 9). Applying this and  $L$ -smoothness of  $-f$ , we can have the following inequality (Lemma 16):

$$\begin{aligned} V_\lambda(\mathbf{z}_0^{k+1}) - V_\lambda(\mathbf{z}_0^k) &\leq -((\lambda + 1)/2) n\alpha \|\nabla \Phi(\mathbf{x}_0^k)\|^2 + (\lambda + 1)n\alpha \|\nabla \Phi(\mathbf{x}_0^k) - \nabla_1 f(\mathbf{z}_0^k)\|^2 \\ &\quad + (n\alpha/2) \|\nabla_1 f(\mathbf{z}_0^k)\|^2 - (n\beta/2) \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \\ &\quad + (\lambda + 1/2) n\alpha \|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2 + (n\beta/2) \|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2. \end{aligned}$$

Hence, to guarantee the fast decrease of  $V_\lambda(\mathbf{z}_0^k)$ , it is important to control the “noise” terms for GDA approximations,  $\|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2$  and  $\|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2$ , in the last line of inequality above. By applying the tools for without-replacement sampling (Proposition 14), we can actually upper-bound the conditional expectations of both noise terms by

$$2L^2n(n+A)\left(\alpha^2\|\nabla_1 f(\mathbf{z}_0^k)\|^2 + \beta^2\|\nabla_2 f(\mathbf{z}_0^k)\|^2\right) + 2L^2n(\alpha^2 + \beta^2)B. \quad (\text{Lemma 17 \& 18})$$

Then, by taking advantage of several properties of smooth nonconvex-PŁ functions (e.g., Propositions 7, 8, and 9) and some small-step-size assumptions (e.g.,  $\beta = \mathcal{O}(1/nL)$ ,  $\beta/\alpha = r \gtrsim \kappa_2^2$ ), we eventually have

$$\mathbb{E}[V_\lambda(\mathbf{z}_0^{k+1})] - \mathbb{E}[V_\lambda(\mathbf{z}_0^k)] \leq -n\alpha\mathbb{E}\left[\|\nabla\Phi(\mathbf{x}_0^k)\|^2\right] - (L\kappa_2n\alpha/2)\mathbb{E}[\Phi(\mathbf{x}_0^k) - f(\mathbf{z}_0^k)] + C\alpha^3,$$

where  $C \geq 0$  is a constant (with respect to  $k$ ) depending on  $L$ ,  $n$ ,  $B$ , and  $r = \beta/\alpha$ . (Lemma 20). We note that the step size ratio  $r \gtrsim \kappa_2^2$  is crucial for showing that the coefficient in front of the term  $\mathbb{E}[\Phi(\mathbf{x}_0^k) - f(\mathbf{z}_0^k)]$  is non-positive: even if it is possible with  $r \lesssim \kappa_2^2$ , we must assume that  $\kappa_2$  upper-bounded by a positive numerical constant, which is not desirable for showing convergence bounds. Thus, we expect that a different proof strategy should be applied to avoid the requirement  $r \gtrsim \kappa_2^2$  on the step size ratio.

The proofs of Theorems 1 and 2 diverge from here. The rest of the proof is mostly about choosing appropriate step sizes and solving the recurrence inequalities.

The full proof of Theorems 4 and 5 starts from the following subsection.

## C.2 EPOCH-WISE REPRESENTATIONS AND BOUNDING NOISE TERMS

Before starting the proof, we again remark that we assume that the mini-batch size  $b$  divides the number of components  $n$  (namely,  $q := n/b$  is a positive integer) for simplicity: thus, readers who want to read proofs for fully stochastic case (i.e.,  $b = 1$ ) can substitute  $n$  to every  $q$ . Also, there is no problem in treating any fraction with a positive numerator and a zero denominator as  $+0$ . Moreover, we simply regard  $(q-1)/(n-1) = 1$  when  $n = 1$ .

We start the proof by aggregating all updates throughout an epoch to obtain an “epoch-wise” update equation. The reason is that the sampled components in each epoch depend on each other, so it is much harder to deal with each iteration individually. At iteration  $t \in [n/b] = [q]$  of epoch  $k \in [K]$ , we use a mini-batch

$$\mathcal{B}_t^k := \{\sigma_k(j) : b(t-1) < j \leq bt, j \in [n]\}.$$

To ease the analysis of Algorithm 2, define the following sums associated with (partial) gradient oracles at a point  $\mathbf{z} = (\mathbf{x}; \mathbf{y})$  over the mini-batch:

$$\mathbf{g}_t^k(\mathbf{z}) := \frac{1}{b} \sum_{i \in \mathcal{B}_t^k} \nabla_1 f_i(\mathbf{z}), \quad \mathbf{h}_t^k(\mathbf{z}) := \frac{1}{b} \sum_{i \in \mathcal{B}_t^k} \nabla_2 f_i(\mathbf{z}).$$

By Assumption 1,  $\mathbf{g}_t^k$  and  $\mathbf{h}_t^k$  are  $L$ -Lipschitz continuous. Computing the average of them over a whole epoch  $(\mathbf{z}_0^k, \dots, \mathbf{z}_{q-1}^k)$ , we define

$$\mathbf{g}^k := \frac{1}{q} \sum_{t=1}^q \mathbf{g}_t^k(\mathbf{z}_{t-1}^k), \quad \mathbf{h}^k := \frac{1}{q} \sum_{t=1}^q \mathbf{h}_t^k(\mathbf{z}_{t-1}^k).$$

Then, by summing up the updates in the epoch  $k$ , we can summarize the epoch as follows.

$$\mathbf{x}_0^{k+1} = \mathbf{x}_0^k - q\alpha\mathbf{g}^k, \quad \mathbf{y}_0^{k+1} = \mathbf{y}_0^k + q\beta\mathbf{h}^k. \quad (\text{simSGDA-RR})$$

We may assume that the intermediate iterates  $\mathbf{z}_1^k, \dots, \mathbf{z}_q^k$  stay close to the starting iterate  $\mathbf{z}_0^k$  of an epoch  $k$ , which results from, e.g., small step sizes. Then, we can approximate the aggregated epoch of SGDA-RR as a step of simGDA applied to  $f = \frac{1}{n} \sum_{i=1}^n f_i$ :  $\mathbf{g}^k \approx \nabla_1 f(\mathbf{z}_0^k)$ ,  $\mathbf{h}^k \approx \nabla_2 f(\mathbf{z}_0^k)$ . In other words,

$$\mathbf{x}_0^{k+1} \approx \mathbf{x}_0^k - q\alpha\nabla_1 f(\mathbf{z}_0^k), \quad \mathbf{y}_0^{k+1} \approx \mathbf{y}_0^k + q\beta\nabla_2 f(\mathbf{z}_0^k), \quad (\approx\text{simGDA})$$

With Assumptions 1, 3 and 4, we can yield a naive (but complicated) upper bound of the gap  $V_\lambda(\mathbf{z}_0^{k+1}) - V_\lambda(\mathbf{z}_0^k)$ , only applying the smoothness of  $\Phi$  and  $-f$ , without any assumptions on step sizes.

**Lemma 16.** *Suppose that Assumptions 1, 3 and 4 hold. Let  $\kappa_2 = L/\mu_2$ , where  $\mu_2$  is PL constant of  $-f(\mathbf{x}; \cdot)$ . Then, the mini-batch simSGDA-RR satisfies that*

$$\begin{aligned}
& V_\lambda(\mathbf{z}_0^{k+1}) - V_\lambda(\mathbf{z}_0^k) \\
& \leq -\left(\frac{\lambda+1}{2}\right) q\alpha \|\nabla\Phi(\mathbf{x}_0^k)\|^2 + (\lambda+1)q\alpha \|\nabla\Phi(\mathbf{x}_0^k) - \nabla_1 f(\mathbf{z}_0^k)\|^2 \\
& \quad + \frac{q\alpha}{2} \|\nabla_1 f(\mathbf{z}_0^k)\|^2 - \frac{q\beta}{2} \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \\
& \quad + \left(\lambda + \frac{1}{2}\right) q\alpha \|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2 + \frac{q\beta}{2} \|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2 \\
& \quad - [\lambda - \{(\lambda+1)(\kappa_2+1) + 1\} Lq\alpha] \frac{q\alpha}{2} \|\mathbf{g}^k\|^2 - (1-Lq\beta) \frac{q\beta}{2} \|\mathbf{h}^k\|^2. \tag{7}
\end{aligned}$$

*Proof.* By definition of  $V_\lambda$ , the following equation holds:

$$V_\lambda(\mathbf{z}_0^{k+1}) - V_\lambda(\mathbf{z}_0^k) = (\lambda+1) [\Phi(\mathbf{x}_0^{k+1}) - \Phi(\mathbf{x}_0^k)] + [f(\mathbf{z}_0^k) - f(\mathbf{z}_0^{k+1})]. \tag{8}$$

First, we seek for an upper bound of  $\Phi(\mathbf{x}_0^{k+1}) - \Phi(\mathbf{x}_0^k)$ . By Proposition 9,  $\Phi$  is  $L(\kappa_2+1)$ -smooth. Hence, we have

$$\begin{aligned}
& \Phi(\mathbf{x}_0^{k+1}) - \Phi(\mathbf{x}_0^k) \\
& \leq \langle \nabla\Phi(\mathbf{x}_0^k), \mathbf{x}_0^{k+1} - \mathbf{x}_0^k \rangle + \frac{L(\kappa_2+1)}{2} \|\mathbf{x}_0^{k+1} - \mathbf{x}_0^k\|^2 \\
& = -q\alpha \langle \nabla\Phi(\mathbf{x}_0^k), \mathbf{g}^k \rangle + \frac{L(\kappa_2+1)}{2} q^2 \alpha^2 \|\mathbf{g}^k\|^2 \\
& = -\frac{q\alpha}{2} \left\{ \|\nabla\Phi(\mathbf{x}_0^k)\|^2 + \|\mathbf{g}^k\|^2 - \|\nabla\Phi(\mathbf{x}_0^k) - \mathbf{g}^k\|^2 \right\} + \frac{L(\kappa_2+1)}{2} q^2 \alpha^2 \|\mathbf{g}^k\|^2 \\
& = -\frac{q\alpha}{2} \|\nabla\Phi(\mathbf{x}_0^k)\|^2 + \frac{q\alpha}{2} \|\nabla\Phi(\mathbf{x}_0^k) - \mathbf{g}^k\|^2 - \frac{q\alpha}{2} (1-L(\kappa_2+1)q\alpha) \|\mathbf{g}^k\|^2 \\
& \leq -\frac{q\alpha}{2} \|\nabla\Phi(\mathbf{x}_0^k)\|^2 + q\alpha \|\nabla\Phi(\mathbf{x}_0^k) - \nabla_1 f(\mathbf{z}_0^k)\|^2 + q\alpha \|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2 \\
& \quad - \frac{q\alpha}{2} (1-L(\kappa_2+1)q\alpha) \|\mathbf{g}^k\|^2. \tag{9}
\end{aligned}$$

The third line is due to polarization equality<sup>9</sup> and the last inequality applies Young's inequality.<sup>10</sup>

Next, applying Assumption 1,  $L$ -smoothness of  $-f(\cdot; \cdot)$  yields an upper bound of  $f(\mathbf{z}_0^k) - f(\mathbf{z}_0^{k+1})$ .

$$\begin{aligned}
& f(\mathbf{z}_0^k) - f(\mathbf{z}_0^{k+1}) \\
& \leq -\langle \nabla f(\mathbf{z}_0^k), \mathbf{z}_0^{k+1} - \mathbf{z}_0^k \rangle + \frac{L}{2} \|\mathbf{z}_0^{k+1} - \mathbf{z}_0^k\|^2 \\
& = -\langle \nabla_1 f(\mathbf{z}_0^k), \mathbf{x}_0^{k+1} - \mathbf{x}_0^k \rangle - \langle \nabla_2 f(\mathbf{z}_0^k), \mathbf{y}_0^{k+1} - \mathbf{y}_0^k \rangle + \frac{L}{2} \|\mathbf{x}_0^{k+1} - \mathbf{x}_0^k\|^2 + \frac{L}{2} \|\mathbf{y}_0^{k+1} - \mathbf{y}_0^k\|^2 \\
& = q\alpha \langle \nabla_1 f(\mathbf{z}_0^k), \mathbf{g}^k \rangle - q\beta \langle \nabla_2 f(\mathbf{z}_0^k), \mathbf{h}^k \rangle + \frac{L}{2} q^2 \alpha^2 \|\mathbf{g}^k\|^2 + \frac{L}{2} q^2 \beta^2 \|\mathbf{h}^k\|^2 \\
& = \frac{q\alpha}{2} \|\nabla_1 f(\mathbf{z}_0^k)\|^2 - \frac{q\alpha}{2} \|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2 + \frac{q\alpha}{2} (1+Lq\alpha) \|\mathbf{g}^k\|^2 \\
& \quad - \frac{q\beta}{2} \|\nabla_2 f(\mathbf{z}_0^k)\|^2 + \frac{q\beta}{2} \|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2 - \frac{q\beta}{2} (1-Lq\beta) \|\mathbf{h}^k\|^2. \tag{10}
\end{aligned}$$

The last equality is due to polarization equality. Lastly, substituting (9) and (10) to (8) finishes the proof.  $\square$

We remark that the last two terms of the inequality (7) can be simply ignored by applying small enough step sizes. However, the terms in the third line of (7) are non-negatives terms related to the

<sup>9</sup>For any  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ ,  $2\langle \mathbf{a}, \mathbf{b} \rangle = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - \|\mathbf{a} - \mathbf{b}\|^2$ .

<sup>10</sup>For any  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ ,  $\|\mathbf{a} + \mathbf{b}\|^2 \leq 2\|\mathbf{a}\|^2 + 2\|\mathbf{b}\|^2$ .



“noise” of approximation  $\mathbf{g}^k \approx \nabla_1 f(\mathbf{z}_0^k)$ ,  $\mathbf{h}^k \approx \nabla_2 f(\mathbf{z}_0^k)$ . Hence, it is important to control the noise terms  $\|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2$  and  $\|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2$  to guarantee a fast decrease of  $V_\lambda(\mathbf{z}_0^k)$ .

**Lemma 17.** *For mini-batch simSGDA-RR, define*

$$G_k := \frac{1}{q} \sum_{t=1}^q \|\mathbf{z}_{t-1}^k - \mathbf{z}_0^k\|^2. \quad (11)$$

With Assumption 1, then

$$\|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2 \leq L^2 G_k \quad \text{and} \quad \|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2 \leq L^2 G_k.$$

As a side remark,  $G_k = 0$  when  $q = 1$  and, in particular,  $n = 1$ .

*Proof.* Recall that  $\frac{1}{q} \sum_{t=1}^q \mathbf{g}_t^k(\mathbf{z}) = \nabla_1 f(\mathbf{z})$  and  $\frac{1}{q} \sum_{t=1}^q \mathbf{h}_t^k(\mathbf{z}) = \nabla_2 f(\mathbf{z})$ . By Lipschitz continuity and Jensen’s inequality,<sup>11</sup>

$$\begin{aligned} \|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2 &= \left\| \frac{1}{q} \sum_{t=1}^q [\mathbf{g}_t^k(\mathbf{z}_{t-1}^k) - \mathbf{g}_t^k(\mathbf{z}_0^k)] \right\|^2 \\ &\leq \frac{1}{q} \sum_{t=1}^q \|\mathbf{g}_t^k(\mathbf{z}_{t-1}^k) - \mathbf{g}_t^k(\mathbf{z}_0^k)\|^2 \leq \frac{L^2}{q} \sum_{t=1}^q \|\mathbf{z}_{t-1}^k - \mathbf{z}_0^k\|^2. \end{aligned}$$

Similarly,

$$\|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2 \leq \frac{L^2}{q} \sum_{t=1}^q \|\mathbf{z}_{t-1}^k - \mathbf{z}_0^k\|^2.$$

This concludes the proof.  $\square$

Thanks to the lemma, it suffices to bound the term  $G_k$ . One can notice that it also represents how far the intermediate iterates  $\mathbf{z}_t^k$  are from the pivot  $\mathbf{z}_0^k$  in average. Before moving on, we define an algorithm-specific symbol denoting a conditional expectation.

**Definition 3.** *We denote a conditional expectation of a random variable  $X$  given all iterates of the first  $k - 1$  epochs by  $\mathbb{E}_k[X] = \mathbb{E}[X | \mathbf{z}_0^1, \mathbf{z}_1^1, \dots, \mathbf{z}_n^{k-1}]$ . In particular, if  $k = 1$ , it boils down to a conditional expectation given only the initial iterate  $\mathbf{z}_0^1$ .*

We get an upper bound of a (conditional) expectation  $\mathbb{E}_k[G_k]$  in the following lemma, which extends a lemma of [Nguyen et al. \(2021, Lemma 6\)](#) to our minimax problems.

**Lemma 18.** *Suppose that Assumptions 1 and 2 hold. Assume that the permutation  $\sigma_k$  is sampled uniformly at random from  $\mathcal{S}_n$ . Then, for any step sizes  $\alpha, \beta$  satisfying  $\alpha^2 + \beta^2 \leq \frac{1}{3q(q-1)L^2}$ , the iterates  $\{\mathbf{z}_t^k\}_{t=0}^{q-1}$  of the  $k$ -th epoch of mini-batch simSGDA-RR satisfies (for  $n > 1$ )*

$$\mathbb{E}_k G_k \leq 2 \left( q^2 + \frac{q(q-1)}{n-1} A \right) \left( \alpha^2 \|\nabla_1 f(\mathbf{z}_0^k)\|^2 + \beta^2 \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right) + \frac{2q(q-1)}{n-1} (\alpha^2 + \beta^2) B.$$

*Proof.* Note that  $G_k = 0$  when  $q = 1$  by its definition. From now, we may assume  $q > 1$  and  $n > 1$  in this proof. By summing the first  $t \in [q - 1]$  updates of the  $k$ -th epoch of mini-batch simSGDA-RR, we have

$$\mathbf{x}_t^k = \mathbf{x}_0^k - t\alpha \left( \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_{j-1}^k) \right), \quad \mathbf{y}_t^k = \mathbf{y}_0^k + t\beta \left( \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_{j-1}^k) \right).$$

<sup>11</sup>For any  $n$  vectors  $a_1, \dots, a_n$ ,  $\left\| \frac{1}{n} \sum_{j=1}^n a_j \right\|^2 \leq \frac{1}{n} \sum_{j=1}^n \|a_j\|^2$ .

Then we can bound the following squared distance.

$$\begin{aligned}
\|\mathbf{x}_t^k - \mathbf{x}_0^k\|^2 &= \alpha^2 t^2 \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_{j-1}^k) \right\|^2 \\
&\leq 3\alpha^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t [\mathbf{g}_j^k(\mathbf{z}_{j-1}^k) - \mathbf{g}_j^k(\mathbf{z}_0^k)] \right\|^2 + \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_0^k) - \nabla_1 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_1 f(\mathbf{z}_0^k)\|^2 \right] \\
&\leq 3\alpha^2 t \sum_{j=1}^t \|\mathbf{g}_j^k(\mathbf{z}_{j-1}^k) - \mathbf{g}_j^k(\mathbf{z}_0^k)\|^2 + 3\alpha^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_0^k) - \nabla_1 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_1 f(\mathbf{z}_0^k)\|^2 \right] \\
&\leq 3\alpha^2 L^2 t \cdot \sum_{j=1}^t \|\mathbf{z}_{j-1}^k - \mathbf{z}_0^k\|^2 + 3\alpha^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_0^k) - \nabla_1 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_1 f(\mathbf{z}_0^k)\|^2 \right] \\
&\leq 3\alpha^2 L^2 t \cdot qG_k + 3\alpha^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_0^k) - \nabla_1 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_1 f(\mathbf{z}_0^k)\|^2 \right]. \tag{12}
\end{aligned}$$

The second and third lines are due to Jensen's inequality. The fourth line is due to  $L$ -Lipschitz continuity of  $\mathbf{g}_j^k$ . Likewise,

$$\|\mathbf{y}_t^k - \mathbf{y}_0^k\|^2 \leq 3\beta^2 L^2 t \cdot qG_k + 3\beta^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k) - \nabla_2 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right]. \tag{13}$$

Summing up (12) and (13),

$$\begin{aligned}
\|\mathbf{z}_t^k - \mathbf{z}_0^k\|^2 &= \|\mathbf{x}_t^k - \mathbf{x}_0^k\|^2 + \|\mathbf{y}_t^k - \mathbf{y}_0^k\|^2 \\
&\leq 3(\alpha^2 + \beta^2) L^2 t q G_k + 3\alpha^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_0^k) - \nabla_1 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_1 f(\mathbf{z}_0^k)\|^2 \right] \\
&\quad + 3\beta^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k) - \nabla_2 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right]. \tag{14}
\end{aligned}$$

Taking (conditional) expectation  $\mathbb{E}_k$  (given  $\mathbf{z}_0^k$ ) to inequality (14),

$$\begin{aligned}
&\mathbb{E}_k \|\mathbf{z}_t^k - \mathbf{z}_0^k\|^2 \\
&\stackrel{(14)}{\leq} 3(\alpha^2 + \beta^2) L^2 t q \cdot (\mathbb{E}_k[G_k]) + 3\alpha^2 t^2 \|\nabla_1 f(\mathbf{z}_0^k)\|^2 + 3\beta^2 t^2 \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \\
&\quad + 3\alpha^2 t^2 \mathbb{E}_k \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_0^k) - \nabla_1 f(\mathbf{z}_0^k) \right\|^2 \right] + 3\beta^2 t^2 \mathbb{E}_k \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k) - \nabla_2 f(\mathbf{z}_0^k) \right\|^2 \right]. \tag{15}
\end{aligned}$$

Here we take advantage of the without-replacement sampling. Putting  $\nabla_s f_i(\mathbf{z}_0^k) \mapsto \mathbf{v}_i$  ( $s \in \{1, 2\}$ ), one can realize a correspondence between the quantities that arise from our algorithm and the symbols in Appendix B.2: for  $s = 1$  ( $\nabla_1 f_i(\mathbf{z}_0^k) \mapsto \mathbf{v}_i$ ),

$$\mathbf{m} = \nabla_1 f(\mathbf{z}_0^k), \quad \tau^2 \leq A \|\nabla_1 f(\mathbf{z}_0^k)\|^2 + B, \quad \mathbf{w}_t = \mathbf{g}_t^k(\mathbf{z}_0^k), \quad \mathbf{m}_t = \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_0^k),$$

and for  $s = 2$  ( $\nabla_2 f_i(\mathbf{z}_0^k) \mapsto \mathbf{v}_i$ ),

$$\mathbf{m} = \nabla_2 f(\mathbf{z}_0^k), \quad \tau^2 \leq A \|\nabla_2 f(\mathbf{z}_0^k)\|^2 + B, \quad \mathbf{w}_t = \mathbf{h}_t^k(\mathbf{z}_0^k), \quad \mathbf{m}_t = \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k).$$

The upper bounds of  $\tau^2$ 's come from Assumption 2. Then by Proposition 14, for any  $t \leq q$ ,

$$\begin{aligned} t^2 \mathbb{E}_k \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{g}_j^k(\mathbf{z}_0^k) - \nabla_1 f(\mathbf{z}_0^k) \right\|^2 &\leq \frac{t(q-t)}{n-1} \left( A \|\nabla_1 f(\mathbf{z}_0^k)\|^2 + B \right), \\ t^2 \mathbb{E}_k \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k) - \nabla_2 f(\mathbf{z}_0^k) \right\|^2 &\leq \frac{t(q-t)}{n-1} \left( A \|\nabla_2 f(\mathbf{z}_0^k)\|^2 + B \right). \end{aligned}$$

Putting these to the inequality (15),

$$\begin{aligned} \mathbb{E}_k \|\mathbf{z}_t^k - \mathbf{z}_0^k\|^2 &\leq 3(\alpha^2 + \beta^2) \left[ L^2 t q \mathbb{E}_k [G_k] + \frac{t(q-t)}{n-1} B \right] \\ &\quad + 3 \left( \alpha^2 \|\nabla_1 f(\mathbf{z}_0^k)\|^2 + \beta^2 \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right) \left[ t^2 + \frac{t(q-t)}{n-1} A \right]. \end{aligned}$$

Taking an average of the inequality above over  $0 \leq t \leq q-1$ ,

$$\begin{aligned} \mathbb{E}_k G_k &= \frac{1}{q} \sum_{t=0}^{q-1} \mathbb{E}_k \|\mathbf{z}_t^k - \mathbf{z}_0^k\|^2 \\ &\leq \frac{3q(q-1)}{2} (\alpha^2 + \beta^2) L^2 \mathbb{E}_k G_k + (\alpha^2 + \beta^2) \frac{q^2 - 1}{2(n-1)} B \\ &\quad + \left( \alpha^2 \|\nabla_1 f(\mathbf{z}_0^k)\|^2 + \beta^2 \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right) \left( \frac{(q-1)(2q-1)}{2} + \frac{q^2 - 1}{2(n-1)} A \right), \quad (16) \end{aligned}$$

where we used the facts

$$\sum_{t=0}^{q-1} t = \frac{q(q-1)}{2}, \quad \frac{1}{q} \sum_{t=0}^{q-1} t^2 = \frac{(q-1)(2q-1)}{6}, \quad \text{and} \quad \frac{1}{q} \sum_{t=0}^{q-1} \frac{t(q-t)}{n-1} = \frac{q^2 - 1}{6(n-1)}.$$

Since we assumed  $\alpha^2 + \beta^2 \leq \frac{1}{3q(q-1)L^2}$ , we have  $1 \leq 2 \left( 1 - \frac{3q(q-1)L^2}{2} (\alpha^2 + \beta^2) \right)$ . Using this,

$$\begin{aligned} \mathbb{E}_k G_k &\leq 2 \left( 1 - \frac{3q(q-1)L^2}{2} (\alpha^2 + \beta^2) \right) \mathbb{E}_k G_k \\ &\stackrel{(16)}{\leq} \left( (q-1)(2q-1) + \frac{q^2 - 1}{(n-1)} A \right) \left( \alpha^2 \|\nabla_1 f(\mathbf{z}_0^k)\|^2 + \beta^2 \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right) + \frac{q^2 - 1}{n-1} (\alpha^2 + \beta^2) B \\ &\leq 2 \left( q^2 + \frac{q(q-1)}{n-1} A \right) \left( \alpha^2 \|\nabla_1 f(\mathbf{z}_0^k)\|^2 + \beta^2 \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right) + \frac{2q(q-1)}{n-1} (\alpha^2 + \beta^2) B, \end{aligned}$$

where the last inequality used  $(q-1)(2q-1) \leq 2q^2$  and  $q+1 \leq 2q$  for  $q \geq 1$ .  $\square$

### C.3 RECURRENCE INEQUALITIES FOR GENERAL SMOOTH NONCONVEX-PŁ OBJECTIVE

Subsequently, we obtain recurrence inequalities about (expected) potential function  $\mathbb{E}_k [V_\lambda(\mathbf{z}_0^k)]$  for nonconvex-PŁ problem. Since primal-PŁ-PŁ problem is a subclass of nonconvex-PŁ problem, the recurrence relations can serve as stepping-stones of our convergence rates.

We introduce some assumptions on *small* step sizes which enable us to get rid of a few troublesome terms from our bound. On top of that, combining the PŁ condition (Assumption 4) with Lemmas 16, 17, and 18, we eventually obtain a much more concise bound on the expected per-epoch change of  $V_\lambda$ . This simple recurrence inequality becomes the key to proving our convergence bounds.

**Lemma 19.** Suppose that Assumptions 1, 2, 3, and 4 hold. Assume that the step sizes  $\alpha$  and  $\beta$  satisfy

$$\alpha \leq \frac{\lambda}{\{(\lambda+1)(\kappa_2+1)+1\}qL}, \quad \beta \leq \frac{1}{qL}, \quad \alpha^2 + \beta^2 \leq \frac{1}{3q(q-1)L^2}, \quad (17)$$

and the condition

$$C_0 := q\beta - 2L^2q \left( q^2 + \frac{q(q-1)}{n-1}A \right) ((2\lambda+1)\alpha + \beta)\beta^2 \geq 0$$

as well. Then, the iterates of mini-batch simSGDA-RR satisfy

$$\mathbb{E}_k[V_\lambda(\mathbf{z}_0^{k+1})] - V_\lambda(\mathbf{z}_0^k) \leq -C_1 \|\nabla\Phi(\mathbf{x}_0^k)\|^2 - C_2 [\Phi(\mathbf{x}_0^k) - f(\mathbf{z}_0^k)] + C_3$$

where

$$\begin{aligned} C_1 &= \left( \frac{\lambda-1}{2} \right) q\alpha - 2L^2q \left( q^2 + \frac{q(q-1)}{n-1}A \right) ((2\lambda+1)\alpha + \beta)\alpha^2, \\ C_2 &= \mu_2 C_0 - 2(\lambda+2)L\kappa_2q\alpha - 4L^3\kappa_2q \left( q^2 + \frac{q(q-1)}{n-1}A \right) ((2\lambda+1)\alpha + \beta)\alpha^2 \\ &= \mu_2q\beta - 2(\lambda+2)L\kappa_2q\alpha - 2L^2\mu_2q \left( q^2 + \frac{q(q-1)}{n-1}A \right) ((2\lambda+1)\alpha + \beta) (2\kappa_2^2\alpha^2 + \beta^2), \\ C_3 &= \left( \frac{L^2q^2(q-1)}{n-1} \right) ((2\lambda+1)\alpha + \beta) (\alpha^2 + \beta^2)B. \end{aligned}$$

*Proof.* The first two inequalities of (17) eliminate the last two terms on the right-hand side of the inequality in Lemma 16. In addition, applying Lemma 17 to Lemma 16 as well, we have

$$\begin{aligned} V_\lambda(\mathbf{z}_0^{k+1}) - V_\lambda(\mathbf{z}_0^k) &\leq - \left( \frac{\lambda+1}{2} \right) q\alpha \|\nabla\Phi(\mathbf{x}_0^k)\|^2 + (\lambda+1)q\alpha \|\nabla\Phi(\mathbf{x}_0^k) - \nabla_1f(\mathbf{z}_0^k)\|^2 \\ &\quad + \frac{q\alpha}{2} \|\nabla_1f(\mathbf{z}_0^k)\|^2 - \frac{q\beta}{2} \|\nabla_2f(\mathbf{z}_0^k)\|^2 + \frac{(2\lambda+1)\alpha + \beta}{2} qL^2G_k. \end{aligned} \quad (18)$$

If we take the conditional expectation  $\mathbb{E}_k$  and apply Lemma 18 (which requires the third inequality of (17) to hold) to (18)

$$\begin{aligned} &\mathbb{E}_k[V_\lambda(\mathbf{z}_0^{k+1})] - V_\lambda(\mathbf{z}_0^k) \\ &\leq - \left( \frac{\lambda+1}{2} \right) q\alpha \|\nabla\Phi(\mathbf{x}_0^k)\|^2 + (\lambda+1)q\alpha \|\nabla\Phi(\mathbf{x}_0^k) - \nabla_1f(\mathbf{z}_0^k)\|^2 \\ &\quad + \frac{1}{2} \left[ q\alpha + 2L^2q \left( q^2 + \frac{q(q-1)}{n-1}A \right) ((2\lambda+1)\alpha + \beta)\alpha^2 \right] \|\nabla_1f(\mathbf{z}_0^k)\|^2 \\ &\quad - \frac{1}{2} \underbrace{\left[ q\beta - 2L^2q \left( q^2 + \frac{q(q-1)}{n-1}A \right) ((2\lambda+1)\alpha + \beta)\beta^2 \right]}_{C_0} \|\nabla_2f(\mathbf{z}_0^k)\|^2 \\ &\quad + \underbrace{\left( \frac{L^2q^2(q-1)}{n-1} \right) ((2\lambda+1)\alpha + \beta) (\alpha^2 + \beta^2)B}_{C_3}. \end{aligned} \quad (19)$$

It is now left to bound terms in (19) using the tools developed so far. First, recall that  $\Phi(\mathbf{x}) := \max_{\mathbf{y}' \in \mathcal{Y}} f(\mathbf{x}; \mathbf{y}')$ . Since  $-f(\mathbf{x}; \mathbf{y})$  is  $\mu_2$ -PL in  $\mathbf{y}$ , we have

$$- \|\nabla_2f(\mathbf{z}_0^k)\|^2 \leq -2\mu_2(\Phi(\mathbf{x}_0^k) - f(\mathbf{z}_0^k)). \quad (20)$$

Given any  $\mathbf{x}$ ,  $\nabla\Phi(\mathbf{x}) = \nabla_1f(\mathbf{x}; \mathbf{y}^*(\mathbf{x}))$  for any  $\mathbf{y}^*(\mathbf{x}) \in \arg \max_{\mathbf{y}' \in \mathcal{Y}} f(\mathbf{x}; \mathbf{y}')$  by Proposition 9. Besides,  $-f(\mathbf{x}; \cdot)$  satisfies QG condition with constant  $\mu_2$  by Proposition 7. Thus, by choosing  $\mathbf{y}^*(\mathbf{x}_0^k)$  to be the projection of  $\mathbf{y}_0^k$  onto  $\arg \max_{\mathbf{y}' \in \mathcal{Y}} f(\mathbf{x}_0^k; \mathbf{y}')$ ,

$$\|\nabla\Phi(\mathbf{x}_0^k) - \nabla_1f(\mathbf{z}_0^k)\|^2 \leq L^2 \|\mathbf{y}^*(\mathbf{x}_0^k) - \mathbf{y}_0^k\|^2 \leq 2L\kappa_2 [\Phi(\mathbf{x}_0^k) - f(\mathbf{z}_0^k)]. \quad (21)$$

Here, the first inequality applies  $L$ -Lipschitz continuity of  $\nabla_1 f(\mathbf{x}_0^k; \cdot)$ , implied by Assumption 1. On top of that, applying the Young's inequality to the term  $\|\nabla_1 f(\mathbf{z}_0^k)\|^2$ ,

$$\begin{aligned} \|\nabla_1 f(\mathbf{z}_0^k)\|^2 &\leq 2\|\nabla\Phi(\mathbf{x}_0^k)\|^2 + 2\|\nabla\Phi(\mathbf{x}_0^k) - \nabla_1 f(\mathbf{z}_0^k)\|^2 \\ &\stackrel{(21)}{\leq} 2\|\nabla\Phi(\mathbf{x}_0^k)\|^2 + 4L\kappa_2 [\Phi(\mathbf{x}_0^k) - f(\mathbf{z}_0^k)] \end{aligned} \quad (22)$$

By applying inequalities (20), (21), and (22) to the bound (19), we conclude the proof.  $\square$

In Lemma 19, we saw that if step sizes are chosen to satisfy certain conditions, then we can simplify the per-epoch progress a great deal. It is now left to choose appropriate step sizes and parameters (e.g.,  $\lambda$ ) so as to make sure not only that  $\alpha$  and  $\beta$  meet the *small* step size conditions (17) but also that the constants  $C_0$ ,  $C_1$ ,  $C_2$ , and  $C_3$  are positive.

**Lemma 20.** *Suppose that Assumptions 1, 2, 3 and 4 hold. Let  $\lambda = 4$  and assume that*

$$0 < \beta \leq \frac{1}{6L\sqrt{q^2 + \frac{q(q-1)}{n-1}A}}, \quad \alpha = \frac{\beta}{r}, \quad \text{where } r \geq 14\kappa_2^2.$$

*Then these satisfy all the inequalities (17) and the terms defined in Lemma 19 satisfy*

$$C_0 > 0, \quad C_1 > q\alpha, \quad C_2 > L\kappa_2 q\alpha/2, \quad C_3 \geq 0.$$

*Consequently, due to the recurrence inequality in Lemma 19, mini-batch simSGDA-RR satisfies, for some numerical constant  $c > 0$ ,*

$$\begin{aligned} &\mathbb{E}_k[V_\lambda(\mathbf{z}_0^{k+1})] - V_\lambda(\mathbf{z}_0^k) \\ &\leq -q\alpha\|\nabla\Phi(\mathbf{x}_0^k)\|^2 - (L\kappa_2 q\alpha/2) [\Phi(\mathbf{x}_0^k) - f(\mathbf{z}_0^k)] + (cr)^3 L^2 \left(\frac{q^2(q-1)}{n-1}\right) B\alpha^3. \quad (\star) \end{aligned}$$

Please note that we mark the recurrence inequality above with a special symbol  $(\star)$  because this inequality is the exact point where the proofs of Theorems 4 and 5 start to deviate.

*Proof.* Regardless of  $A \geq 0$ , we have

$$\beta \leq \frac{1}{6Lq} \quad \text{and} \quad \alpha \leq \frac{1}{6Lqr} \leq \frac{1}{84L\kappa_2^2 q}. \quad (23)$$

This is enough to guarantee that the inequalities (17) hold with  $\lambda = 4$ . Since  $C_0 > C_2/\mu_2$ , it is enough to show  $C_2 > 0$  to prove that  $C_0 > 0$ . Applying  $\lambda = 4$ ,  $\kappa_2 \geq 1$ , and  $\beta/\alpha = r \geq 14\kappa_2^2$ ,

$$\begin{aligned} \frac{C_1}{q\alpha} &= \frac{3}{2} - 2L^2 \left(q^2 + \frac{q(q-1)}{n-1}A\right) (9+r) \alpha^2 \\ &\geq \frac{3}{2} - \frac{2}{6^2} \cdot \frac{9+r}{r^2} \geq \frac{3}{2} - \frac{2 \cdot 23}{6^2 \cdot 14^2} > 1, \\ \frac{C_2}{\mu_2 q \beta} &= 1 - \frac{12\kappa_2^2}{r} - 2L^2 \left(q^2 + \frac{q(q-1)}{n-1}A\right) \left(\frac{9}{r} + 1\right) \left(\frac{2\kappa_2^2}{r^2} + 1\right) \beta^2 \\ &\geq 1 - \frac{12}{14} - \frac{2}{6^2} \left(\frac{9}{14\kappa_2^2} + 1\right) \left(\frac{2}{14^2\kappa_2^2} + 1\right) \geq \frac{2}{14} - \frac{2 \cdot 23 \cdot 198}{6^2 \cdot 14^3} > \frac{1}{2 \cdot 14}. \end{aligned}$$

Thus,  $C_1 > q\alpha$  and

$$C_2 > \frac{\mu_2 q \beta}{2 \cdot 14} = \frac{\mu_2 q r \alpha}{2 \cdot 14} \geq L\kappa_2 q\alpha/2.$$

Then we conclude the proof by bounding the term  $C_3$ . We can already check from the definition that  $C_3 \geq 0$ . We can upper-bound  $C_3$  by

$$C_3 = \left(\frac{L^2 q^2 (q-1)}{n-1}\right) (9+r) (1+r^2) B\alpha^3 \leq (cr)^3 L^2 \left(\frac{q^2 (q-1)}{n-1}\right) B\alpha^3,$$

for some numerical constant  $c > 0$ .  $\square$

#### C.4 CONVERGENCE RATES FOR SMOOTH NONCONVEX-PŁ PROBLEM

In this subsection, we show the convergence bound of general smooth nonconvex-PŁ problems in terms of  $\min_{k \in [K]} \mathbb{E} [\|\nabla\Phi(\mathbf{x}_0^k)\|^2]$ . From the inequality  $(\star)$  in Lemma 20, we can simply ignore the second term

$$-(L\kappa_2q\alpha/2) [\Phi(\mathbf{x}_0^k) - f(\mathbf{z}_0^k)] \leq 0$$

of the right-hand side because  $\Phi(\mathbf{x}) \geq f(\mathbf{x}; \mathbf{y})$  for any  $(\mathbf{x}; \mathbf{y})$ . In other words, we may deal with the inequality

$$\mathbb{E}_k[V_\lambda(\mathbf{z}_0^{k+1})] - V_\lambda(\mathbf{z}_0^k) \leq -q\alpha \|\nabla\Phi(\mathbf{x}_0^k)\|^2 + (cr)^3 L^2 \left( \frac{q^2(q-1)}{n-1} \right) B\alpha^3. \quad (\text{nc-PŁ})$$

Plugging  $q = n/b$ , we eventually show the convergence rate (Theorem 4). (Recall that  $b$  is the size of mini-batches.)

**Theorem 21** (Equivalent to Theorem 4, for simSGDA-RR). *Suppose that  $f$  satisfies Assumptions 1, 2, 3, and 4 are satisfied. Let  $\lambda = 4$ . Choose the step sizes  $\alpha$  and  $\beta$  by  $\alpha = \beta/r$  for some  $r \geq 14\kappa_2^2$  and*

$$\beta = \min \left\{ \frac{1}{6L\sqrt{q^2 + \frac{q(q-1)}{n-1}}A}, \frac{1}{c} \left( \frac{V_\lambda(\mathbf{z}_0^1)}{L^2q^2\left(\frac{q-1}{n-1}\right)BK} \right)^{\frac{1}{3}} \right\},$$

for some numerical constant  $c > 0$ . Then, mini-batch simSGDA-RR satisfies

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} [\|\nabla\Phi(\mathbf{x}_0^k)\|^2] \leq \frac{6rLV_\lambda(\mathbf{z}_0^1)}{K} \sqrt{1 + \left(\frac{q-1}{n-1}\right) \frac{A}{q}} + 2cr \left( \frac{L^2BV_\lambda(\mathbf{z}_0^1)^2}{qK^2} \cdot \frac{q-1}{n-1} \right)^{1/3}.$$

*Proof.* To replace the conditional expectations with unconditional expectations, we take expectation to both sides of the inequality (nc-PŁ):

$$\mathbb{E}[V_\lambda(\mathbf{z}_0^{k+1}) - V_\lambda(\mathbf{z}_0^k)] \leq -q\alpha \mathbb{E} [\|\nabla\Phi(\mathbf{x}_0^k)\|^2] + (cr)^3 L^2 \left( \frac{q^2(q-1)}{n-1} \right) B\alpha^3.$$

Rearranging the terms and taking a sum from  $k = 1$  to  $k = K$ , we have

$$q\alpha \sum_{k=1}^K \mathbb{E} [\|\nabla\Phi(\mathbf{x}_0^k)\|^2] \leq \mathbb{E}[V_\lambda(\mathbf{z}_0^1) - V_\lambda(\mathbf{z}_0^{K+1})] + (cr)^3 L^2 \left( \frac{q^2(q-1)}{n-1} \right) B\alpha^3 K.$$

Dividing both sides by  $qK\alpha$ , we get the following. Note that  $V_\lambda$  is non-negative.

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} [\|\nabla\Phi(\mathbf{x}_0^k)\|^2] \leq \frac{V_\lambda(\mathbf{z}_0^1)}{qK\alpha} + (cr)^3 L^2 \left( \frac{q(q-1)}{n-1} \right) B\alpha^2$$

Since our choice of step sizes implies

$$\alpha = \min \left\{ \frac{1}{6rL\sqrt{q^2 + \frac{q(q-1)}{n-1}}A}, \frac{1}{cr} \left( \frac{V_\lambda(\mathbf{z}_0^1)}{L^2Bq^2\left(\frac{q-1}{n-1}\right)K} \right)^{\frac{1}{3}} \right\},$$

we eventually prove the theorem by using the inequality  $\max\{a, b\} \leq a + b$  (for  $a, b \geq 0$ ).  $\square$

#### C.5 CONVERGENCE RATES FOR SMOOTH PRIMAL-PŁ-PŁ PROBLEM

In this subsection, we prove the convergence bound of primal-PŁ-PŁ (or, PŁ( $\Phi$ )-PŁ) problems in terms of  $\mathbb{E} [V_\lambda(\mathbf{z}_0^{K+1})]$ .

Unlike the previous subsection, we additionally utilize Assumption 5 stating that  $f(\mathbf{x}; \mathbf{y})$  satisfies primal PŁ condition, namely, the primal function  $\Phi(\mathbf{x}) = \max_{\mathbf{y}'} f(\mathbf{x}; \mathbf{y}')$  is a  $\mu_1$ -PŁ function. With this assumption, we yield another recurrence inequality from the inequality  $(\star)$ . We note that it uses the  $\mu_1$ -PŁ condition for  $\Phi$  ( $\because$  Proposition 10) but not necessarily for  $f(\cdot; \mathbf{y})$ .

**Lemma 22.** *Suppose that  $f$  satisfies Assumptions 1, 2, 3, 4, and 5. Then, with the same choice of  $\lambda = 4$  and the same condition of the step sizes  $\alpha$  and  $\beta$  as in Lemma 20, the mini-batch simSGDA-RR satisfies that, for some numerical constant  $c > 0$ ,*

$$\mathbb{E}_k[V_\lambda(\mathbf{z}_0^{k+1})] \leq (1 - \mu_1 q \alpha / 2) V_\lambda(\mathbf{z}_0^k) + (cr)^3 L^2 \left( \frac{q^2(q-1)}{n-1} \right) B \alpha^3. \quad (\text{PL}(\Phi)\text{-PL})$$

*Proof.* Since the primal function  $\Phi$  is a  $\mu_1$ -PL function,

$$-\|\nabla\Phi(\mathbf{x}_0^k)\|^2 \leq -2\mu_1 [\Phi(\mathbf{x}_0^k) - \Phi^*].$$

Also, since  $\mu_1 \leq L$  and  $\kappa_2 \geq 1$ , we know that  $-L\kappa_2 \leq -\mu_1$ . Applying these to the inequality  $(\star)$ , we have

$$\begin{aligned} & \mathbb{E}_k [V_\lambda(\mathbf{z}_0^{k+1})] - V_\lambda(\mathbf{z}_0^k) \\ & \leq -(2\mu_1 q \alpha / \lambda) \cdot \lambda [\Phi(\mathbf{x}_0^k) - \Phi^*] - (\mu_1 q \alpha / 2) [\Phi(\mathbf{x}_0^k) - f(\mathbf{z}_0^k)] + (cr)^3 L^2 \left( \frac{q^2(q-1)}{n-1} \right) B \alpha^3 \\ & = -(\mu_1 q \alpha / 2) \cdot V_\lambda(\mathbf{z}_0^k) + (cr)^3 L^2 \left( \frac{q^2(q-1)}{n-1} \right) B \alpha^3, \end{aligned}$$

since  $\lambda = 4$ . By re-arranging the terms, we conclude the proof.  $\square$

Of course, the multiplier  $1 - \mu_1 q \alpha / 2$  has a value between 0 and 1. To see why, note that from Equation (23),

$$0 < \mu_1 q \alpha / 2 \leq \mu_1 q \cdot \frac{1}{2 \cdot 84L\kappa_2^2 q} = \frac{1}{168\kappa_1\kappa_2^2} < 1.$$

**Theorem 23** (Equivalent to Theorem 5, for simSGDA-RR). *Assume that  $f$  satisfies Assumptions 1, 2, 3, 4, and 5. Let  $\lambda = 4$ . Choose the step sizes by  $\alpha = \beta/r$  for some  $r \geq 14\kappa_2^2$  and*

$$\beta = \min \left\{ \frac{1}{6L\sqrt{q^2 + \frac{q(q-1)}{n-1}}A}, \frac{2r}{\mu_1 q K} \max \left\{ 1, \log \left( \frac{V_\lambda(\mathbf{z}_0^1) \mu_1 q K^2}{8(cr)^3 \kappa_1^2 \left( \frac{q-1}{n-1} \right) B} \right) \right\} \right\},$$

for some numerical constant  $c > 0$ . Then, mini-batch simSGDA-RR satisfies

$$\mathbb{E}[V_\lambda(\mathbf{z}_n^K)] \leq \mathcal{O} \left( V_\lambda(\mathbf{z}_0^1) \cdot \exp \left( -\frac{K}{12\kappa_1 r \sqrt{1 + \left( \frac{q-1}{n-1} \right) \frac{A}{q}}} \right) \right) + \tilde{\mathcal{O}} \left( \frac{\kappa_1^2 r^3 B}{\mu_1 q K^2} \right) \cdot \frac{q-1}{n-1}.$$

*Proof.* To replace the conditional expectations with unconditional expectations, we take expectation to both sides of the inequality (PL( $\Phi$ )-PL):

$$\mathbb{E} [V_\lambda(\mathbf{z}_0^{k+1})] \leq (1 - \mu_1 q \alpha / 2) \mathbb{E} [V_\lambda(\mathbf{z}_0^k)] + (cr)^3 L^2 \left( \frac{q^2(q-1)}{n-1} \right) B \alpha^3.$$

Unrolling the recurrence inequality (Proposition 15) and using the facts  $\beta = 14\kappa_2^2 \alpha$ , we have

$$\begin{aligned} \mathbb{E}[V_\lambda(\mathbf{z}_n^K)] & \leq (1 - \mu_1 q \alpha / 2)^K V_\lambda(\mathbf{z}_0^1) + \frac{2 \cdot (cr)^3 L^2}{\mu_1 q \alpha} \left( \frac{q^2(q-1)}{n-1} \right) B \alpha^3 \\ & \leq \exp(-\mu_1 q K \alpha / 2) V_\lambda(\mathbf{z}_0^1) + 2(cr)^3 \mu_1 \kappa_1^2 \left( \frac{q(q-1)}{n-1} \right) B \alpha^2. \end{aligned} \quad (24)$$

Note that, in the inequality above, the second term of the right hand side becomes zero when  $q = 1$ . In that case, we can prove exponential decay of  $\mathbb{E}[V_\lambda(\mathbf{z}_0^k)]$ . Thus, we simply assume  $q > 1$  hereafter.

*Case 1:* If  $K$  is as large as

$$K > \frac{\kappa_1 r^{3/2}}{\sqrt{\mu_1}} \cdot \sqrt{\frac{8c^3 e B}{V_\lambda(\mathbf{z}_0^1) q} \left( \frac{q-1}{n-1} \right)}, \quad (e = \exp(1))$$

we have a step size  $\alpha$  as

$$\alpha = \min \left\{ \frac{1}{6Lr\sqrt{q^2 + \frac{q(q-1)}{n-1}A}}, \frac{2}{\mu_1 q K} \log(\clubsuit) \right\}, \quad \text{where } \clubsuit = \frac{V_\lambda(z_0^1)\mu_1 q K^2}{8(cr)^3 \kappa_1^2 \kappa_2^6 \left(\frac{q-1}{n-1}\right) B}.$$

Due to the lower bound of epoch size  $K$ , the fraction  $\clubsuit$  inside the log factor is indeed greater than  $e > 1$ , which guarantees the step size is positive. Putting this to the inequality (24) and using the fact that  $\max\{a, b\} \leq a + b$  (for  $a, b \geq 0$ ), we eventually have

$$\begin{aligned} & \mathbb{E} [V_\lambda(z_n^K)] \\ & \leq V_\lambda(z_0^1) \cdot \exp \left( -\frac{K}{12\kappa_1 r \sqrt{1 + \left(\frac{q-1}{n-1}\right) \frac{A}{q}}} \right) + \frac{2 \cdot 8(cr)^3 \kappa_1^2 B}{\mu_1 q K^2} \left(\frac{q-1}{n-1}\right) [1 + \log^2(\clubsuit)] \\ & = V_\lambda(z_0^1) \cdot \exp \left( -\frac{K}{12\kappa_1 r \sqrt{1 + \left(\frac{q-1}{n-1}\right) \frac{A}{q}}} \right) + \tilde{\mathcal{O}} \left( \frac{\kappa_1^2 r^3 B}{\mu_1 q K^2} \right) \cdot \frac{q-1}{n-1}. \end{aligned}$$

*Case 2:* Otherwise, the log factor might have a negative value when  $K$  is too small. However, in this case, we have

$$V_\lambda(z_0^1) \leq \frac{8(cr)^3 e \kappa_1^2 B}{\mu_1 q K^2} \cdot \frac{q-1}{n-1}; \quad \alpha = \min \left\{ \frac{1}{84L\kappa_2^2 \sqrt{q^2 + \frac{q(q-1)}{n-1}A}}, \frac{2}{\mu_1 q K} \right\}.$$

Putting these to the inequality (24), we have

$$\begin{aligned} \mathbb{E} [V_\lambda(z_n^K)] & \leq \frac{8(cr)^3 e \kappa_1^2 B}{\mu_1 q K^2} \left(\frac{q-1}{n-1}\right) \left[ \exp(-\mu_1 q K \alpha / 2) + \frac{1}{e} \cdot (\mu_1 q K \alpha / 2)^2 \right] \\ & \leq \frac{8(cr)^3 e \kappa_1^2 B}{\mu_1 q K^2} \left(\frac{q-1}{n-1}\right) = \mathcal{O} \left( \frac{\kappa_1^2 r^3 B}{\mu_1 q K^2} \right) \cdot \frac{q-1}{n-1}. \end{aligned}$$

The inequality in the last line is due to the fact that  $e^{-t} + t^2/e \leq 1$  for each  $t \in (0, 1]$ , and that  $\mu_1 q K \alpha / 2 \in (0, 1]$ .

Combining both *Case 1* and *Case 2*, we conclude the proof of the theorem.  $\square$

## D PROOFS FOR (MINI-BATCH) ALTERNATING SGDA-RR: FOCUSING ON CHANGES IN THE PROOF

In this appendix, we prove the same convergence rates for altSGDA-RR as the simultaneous update counterpart. Since most of the steps in the proof are similar to those in Appendix C, we only describe which steps change in the proof.

### D.1 EPOCH-WISE REPRESENTATIONS AND BOUNDING NOISE TERMS

To analyze altSGDA-RR, we modify the notation for epoch-wise updates. The only change is that an update  $\mathbf{y}_t^k \mapsto \mathbf{y}_{t+1}^k$  uses  $\mathbf{x}_{t+1}^k$  instead of  $\mathbf{x}_t^k$ . Hence, the definition of  $\mathbf{h}^k$  should be modified. Recall that

$$\mathbf{g}_t^k(z) := \frac{1}{b} \sum_{i \in \mathcal{B}_t^k} \nabla_1 f_i(z), \quad \mathbf{h}_t^k(z) := \frac{1}{b} \sum_{i \in \mathcal{B}_t^k} \nabla_2 f_i(z),$$

where  $\mathcal{B}_t^k$  is a mini-batch of size  $b$  formed at iteration  $t$  of epoch  $k$ . Then, at epoch  $k$ , by re-definition of  $\mathbf{h}^k$ ,

$$\begin{aligned} \mathbf{g}^k & := \frac{1}{q} \sum_{t=1}^q \mathbf{g}_t^k(\mathbf{x}_{t-1}^k; \mathbf{y}_{t-1}^k), \quad \mathbf{h}^k := \frac{1}{q} \sum_{t=1}^q \mathbf{h}_t^k(\mathbf{x}_t^k; \mathbf{y}_{t-1}^k). \\ \mathbf{x}_0^{k+1} & = \mathbf{x}_0^k - q\alpha \mathbf{g}^k, \quad \mathbf{y}_0^{k+1} = \mathbf{y}_0^k + q\beta \mathbf{h}^k. \end{aligned} \quad (\text{altSGDA-RR})$$



We still approximate this epoch-wise update rule to a full-batch simultaneous GDA update ( $\approx \text{simGDA}$ ) with step sizes  $q\alpha$  and  $q\beta$ . Again, we control the “noise” terms  $\|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2$  and  $\|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2$  not to be large. Because of the modification of  $\mathbf{h}^k$ , we have a different result for  $\|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2$  as follows.

**Lemma 24.** *For mini-batch altSGDA-RR, recall that*

$$G_k := \frac{1}{q} \sum_{t=1}^q \|\mathbf{z}_{t-1}^k - \mathbf{z}_0^k\|^2.$$

If we have Assumption 1, then we have

$$\|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2 \leq L^2 G_k + L^2 q \alpha^2 \|\mathbf{g}^k\|^2, \quad \text{whereas} \quad \|\mathbf{g}^k - \nabla_1 f(\mathbf{z}_0^k)\|^2 \leq L^2 G_k. \quad (25)$$

*Proof.* Because of  $L$ -Lipschitz continuity of  $\mathbf{h}_t^k(\cdot; \cdot)$ ,

$$\begin{aligned} \|\mathbf{h}^k - \nabla_2 f(\mathbf{z}_0^k)\|^2 &= \left\| \frac{1}{q} \sum_{t=1}^q [\mathbf{h}_t^k(\mathbf{x}_t^k; \mathbf{y}_{t-1}^k) - \mathbf{h}_t^k(\mathbf{x}_0^k; \mathbf{y}_0^k)] \right\|^2 \\ &\leq \frac{1}{q} \sum_{t=1}^q \|\mathbf{h}_t^k(\mathbf{x}_t^k; \mathbf{y}_{t-1}^k) - \mathbf{h}_t^k(\mathbf{x}_0^k; \mathbf{y}_0^k)\|^2 \\ &\leq \frac{L^2}{q} \sum_{t=1}^q \|\mathbf{z}_{t-1}^k - \mathbf{z}_0^k\|^2 + \frac{L^2}{q} \|\mathbf{x}_q^k - \mathbf{x}_0^k\|^2 = L^2 G_k + L^2 q \alpha^2 \|\mathbf{g}^k\|^2. \end{aligned}$$

The last inequality holds because  $\mathbf{x}_q^k = \mathbf{x}_0^{k+1}$ .  $\square$

## D.2 BOUNDING NOISE TERMS: A BIT DIFFERENT PROOF OF LEMMA 18

We notice that the same result as Lemma 18 holds not only for simultaneous updates but also alternating updates, even though it is not very straightforward. We need to reflect the changes from the previous subsection. That is, we have to be careful when we expand the term  $\|\mathbf{y}_t^k - \mathbf{y}_0^k\|^2$  ( $0 \leq t \leq q-1$ ). Unlike the inequality (12) (in the original proof), we have

$$\begin{aligned} \|\mathbf{y}_t^k - \mathbf{y}_0^k\|^2 &= \beta^2 t^2 \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{x}_j^k; \mathbf{y}_{j-1}^k) \right\|^2 \\ &\leq 3\beta^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t [\mathbf{h}_j^k(\mathbf{x}_j^k; \mathbf{y}_{j-1}^k) - \mathbf{h}_j^k(\mathbf{z}_0^k)] \right\|^2 + \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k) - \nabla_2 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right] \\ &\leq 3\beta^2 t^2 \left[ \frac{1}{t} \sum_{j=1}^t \|\mathbf{h}_j^k(\mathbf{x}_j^k; \mathbf{y}_{j-1}^k) - \mathbf{h}_j^k(\mathbf{z}_0^k)\|^2 + \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k) - \nabla_2 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right] \\ &\leq 3\beta^2 t^2 \left[ \frac{L^2}{t} \left( \|\mathbf{x}_t^k - \mathbf{x}_0^k\|^2 + \sum_{j=1}^t \|\mathbf{z}_{j-1}^k - \mathbf{z}_0^k\|^2 \right) + \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k) - \nabla_2 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right] \\ &\leq 3\beta^2 L^2 t \sum_{j=1}^t \|\mathbf{z}_j^k - \mathbf{z}_0^k\|^2 + 3\beta^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k) - \nabla_2 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right] \\ &\leq 3\beta^2 L^2 t \cdot q G_k + 3\beta^2 t^2 \left[ \left\| \frac{1}{t} \sum_{j=1}^t \mathbf{h}_j^k(\mathbf{z}_0^k) - \nabla_2 f(\mathbf{z}_0^k) \right\|^2 + \|\nabla_2 f(\mathbf{z}_0^k)\|^2 \right]. \end{aligned}$$

The second and third inequality holds by Jensen’s inequality, and the last inequality holds because  $t \leq q - 1$ . The resulting upper bound is identical to the inequality (13). Proving this inequality above suffices to show that the conclusion of Lemma 18 also holds for altSGDA-RR, because we eventually take an average along  $0 \leq t \leq q - 1$  and the other steps in the proof do not utilize the “order” (either simultaneous or alternating) of updates.

### D.3 RECURRENCE INEQUALITIES FOR GENERAL SMOOTH NONCONVEX-PL OBJECTIVE

In the proof for simSGDA-RR, we applied Lemma 16, Lemma 18, and the “small-step-size” assumptions (three inequalities in (17)) to deduce Lemma 19. However, due to Lemma 24 that we obtained for altSGDA-RR, we need slightly different assumptions on step sizes rather than (17).

Fortunately, we notice that the Lemma 16 also holds for altSGDA-RR, with a modified version of  $\mathbf{h}^k$ . This is because the proof of the lemma does not utilize step-wise updates, while the discrepancy between simultaneous and alternating updates only appears in the step-wise updates. Thus, we have the same result as Lemma 19.

**Lemma 25.** *Suppose that Assumptions 1, 2, 3, and 4 hold. Modify the inequalities (17) (from Lemma 19) by*

$$\lambda - \{(\lambda + 1)(\kappa_2 + 1) + 1\} Lq\alpha - L^2q\alpha\beta \geq 0, \quad \beta \leq \frac{1}{qL}, \quad \alpha^2 + \beta^2 \leq \frac{1}{3q(q-1)L^2}. \quad (26)$$

(In fact, only the first one is different.) Then, the result of Lemma 19 still holds for mini-batch altSGDA-RR.

*Proof.* We first apply Lemma 24 to the general bound resulted from Lemma 16:

$$\begin{aligned} & V_\lambda(\mathbf{z}_0^{k+1}) - V_\lambda(\mathbf{z}_0^k) \\ & \leq - \left( \frac{\lambda + 1}{2} \right) q\alpha \|\nabla\Phi(\mathbf{x}_0^k)\|^2 + (\lambda + 1)q\alpha \|\nabla\Phi(\mathbf{x}_0^k) - \nabla_1 f(\mathbf{z}_0^k)\|^2 \\ & \quad + \frac{q\alpha}{2} \|\nabla_1 f(\mathbf{z}_0^k)\|^2 - \frac{q\beta}{2} \|\nabla_2 f(\mathbf{z}_0^k)\|^2 + \frac{(2\lambda + 1)\alpha + \beta}{2} qL^2 G_k \\ & \quad - [\lambda - \{(\lambda + 1)(\kappa_2 + 1) + 1\} Lq\alpha - L^2q\alpha\beta] \frac{q\alpha}{2} \|\mathbf{g}^k\|^2 - (1 - Lq\beta) \frac{q\beta}{2} \|\mathbf{h}^k\|^2. \end{aligned} \quad (27)$$

Hence, the first two inequalities of (26) eliminate the last two terms on the right side of the inequality (27) above:

$$\begin{aligned} V_\lambda(\mathbf{z}_0^{k+1}) - V_\lambda(\mathbf{z}_0^k) & \leq - \left( \frac{\lambda + 1}{2} \right) q\alpha \|\nabla\Phi(\mathbf{x}_0^k)\|^2 + (\lambda + 1)q\alpha \|\nabla\Phi(\mathbf{x}_0^k) - \nabla_1 f(\mathbf{z}_0^k)\|^2 \\ & \quad + \frac{q\alpha}{2} \|\nabla_1 f(\mathbf{z}_0^k)\|^2 - \frac{q\beta}{2} \|\nabla_2 f(\mathbf{z}_0^k)\|^2 + \frac{(2\lambda + 1)\alpha + \beta}{2} qL^2 G_k. \end{aligned}$$

This is identical to the inequality (18) in the proof of Lemma 19. From this point on, the rest of the proof is exactly identical to Lemma 19.  $\square$

Lemma 25 establishes that altSGDA-RR also satisfies a concise bound on the expected per-epoch change of  $V_\lambda$ , albeit under a slightly different set of assumptions (26) on step sizes. Using this result, we can prove the convergence rates for altSGDA-RR that are exactly the same as simSGDA-RR.

### D.4 SMALL STEP SIZE ASSUMPTIONS

It is left to show an altSGDA-RR counterpart for Lemma 20 which establishes the general recurrence inequality ( $\star$ ). In fact, the same choice of step sizes as simSGDA-RR, namely

$$0 < \beta \leq \frac{1}{6L\sqrt{q^2 + \frac{q(q-1)}{n-1}}A} \quad \text{and} \quad \alpha = \frac{\beta}{r} \quad \text{where} \quad r \geq 14\kappa_2^2,$$

actually meets the newly introduced conditions (26). Among the three inequalities, the only one that needs to be checked is

$$\lambda - \{(\lambda + 1)(\kappa_2 + 1) + 1\} Lq\alpha - L^2q\alpha\beta > 0.$$

Note that, regardless of  $A \geq 0$ ,

$$\beta \leq \frac{1}{6Lq} \quad \text{and} \quad \alpha \leq \frac{1}{6Lqr} \leq \frac{1}{84L\kappa_2^2q}$$

In this case,

$$\begin{aligned} & \lambda - \{(\lambda + 1)(\kappa_2 + 1) + 1\} Lq\alpha - L^2q\alpha\beta \\ & \geq 4 - (11\kappa_2 + L\beta)Lq\alpha \geq 4 - \left(11\kappa_2 + \frac{1}{6}\right) \cdot \frac{1}{84\kappa_2^2} > 0. \end{aligned}$$

Therefore, there is no need to modify our choices of  $\lambda$  and the step sizes  $\alpha, \beta$  for the analysis of altSGDA-RR, and the rest of the proof for simSGDA-RR goes through.

## E PROOFS FOR LOWER BOUND OF DETERMINISTIC FULL-BATCH SIMGDA

In this appendix, we illustrate a comprehensive lower bound for full-batch GDA, which is specific to the choice of step size ratio (Theorem 3). Before we start the proof, we define a class of smooth strongly-convex-strongly concave functions.

**Definition 4.** Let  $\mathcal{F}(L, \mu_1, \mu_2)$  be the class of functions  $f(\mathbf{x}; \mathbf{y})$  with two arguments  $\mathbf{x}$  and  $\mathbf{y}$  of any dimension, which is  $L$ -smooth,  $\mu_1$ -strongly-convex in  $\mathbf{x}$ , and  $\mu_2$ -strongly-concave in  $\mathbf{y}$ . Let  $\kappa_1 = L/\mu_1 \geq 1$  and  $\kappa_2 = L/\mu_2 \geq 1$  be condition numbers of the function class. Denote the (unique) saddle (or, global minimax) point by  $\mathbf{z}^* = (\mathbf{x}^*; \mathbf{y}^*)$ .

We restate and prove the Theorem 3 for reader's convenience.

**Theorem 26** (Restatement of Theorem 3). Suppose  $\kappa_1 \geq c$  and  $\kappa_2 \geq c$  for some constant  $c > 1$ . Then, for each step size ratio  $r > 0$ , there exists a function  $f \in \mathcal{F}(L, \mu_1, \mu_2)$  for which simGDA with any step sizes  $\alpha$  and  $\beta$  of ratio  $r = \beta/\alpha$  requires

$$K = \begin{cases} \Omega(\kappa_1 r \log(1/\varepsilon)), & \text{if } r \geq \kappa_2/c, \\ \Omega(\kappa_1 \kappa_2 \log(1/\varepsilon)), & \text{if } c/\kappa_1 \leq r \leq \kappa_2/c, \\ \Omega((\kappa_2/r) \log(1/\varepsilon)), & \text{if } 0 < r \leq c/\kappa_1 \end{cases}$$

iterations to achieve either  $\|\mathbf{z}_k - \mathbf{z}^*\|^2 \leq \varepsilon^2$  or  $V_\lambda(\mathbf{z}_K) \leq \varepsilon^2$ .

*Proof.* The proof is done in case by case, constructing a worst-case function for each of 4 different regimes of step size ratio  $r$ : (1)  $\mu_1/\mu_2 \leq r \leq \kappa_2/c$ , (2)  $c/\kappa_1 \leq r \leq \mu_1/\mu_2$ , (3)  $r \geq \kappa_2/c$ , and (4)  $0 < r \leq c/\kappa_1$ . Readers might notice the similarities of the proofs for (1) $\leftrightarrow$ (2) and (3) $\leftrightarrow$ (4).

Case 1. ( $\mu_1/\mu_2 \leq r \leq \kappa_2/c$ ). Consider

$$f^{(1)}(v, x; y) := \frac{\mu_1}{2}v^2 + \frac{r\mu_2}{2}x^2 - \frac{\mu_2}{2}y^2 + \ell xy,$$

where  $\ell^2 = L^2 - r\mu_2^2 - L\mu_2|r - 1| \geq 0$ . Applying Proposition 28, it can be shown that  $f^{(1)} \in \mathcal{F}(L, \mu_1, \mu_2)$ . Also,  $\mathbf{z}^* = (0, 0; 0)$  is its unique saddle point. Note that, the GDA on  $f^{(1)}$  can be written as

$$v_{t+1} = \left(1 - \frac{\beta\mu_1}{r}\right)v_t, \quad \underbrace{\begin{bmatrix} x_{t+1} \\ y_{t+1} \end{bmatrix}}_{\mathbf{A}} = \underbrace{\begin{bmatrix} 1 - \beta\mu_2 & -\beta\ell/r \\ \beta\ell & 1 - \beta\mu_2 \end{bmatrix}}_{\mathbf{A}} \begin{bmatrix} x_t \\ y_t \end{bmatrix} = \mathbf{A} \begin{bmatrix} x_t \\ y_t \end{bmatrix}.$$

Also, the eigenvalues  $\tau$  of  $\mathbf{A}$  is

$$\begin{aligned} \tau &= 1 - \beta\mu_2 \pm \sqrt{(1 - \beta\mu_2)^2 - ((1 - \beta\mu_2)^2 + \beta^2\ell^2/r)} \\ &= 1 - \beta\mu_2 \pm \frac{\beta\ell}{\sqrt{r}}\sqrt{-1}. \end{aligned}$$

The spectral radius (*i.e.*, maximum absolute eigenvalue) is

$$\rho(\mathbf{A}) = \sqrt{(1 - \beta\mu_2)^2 + \beta^2\ell^2/r}.$$

Since the eigenvalues are complex conjugates of each other (the magnitudes are the same), both eigenvalues have magnitude  $\rho(\mathbf{A})$ . Then, by Proposition 27,  $\rho(\mathbf{A}) < 1$  is necessary for convergence. To this end, we need  $\beta > 0$  satisfying  $\beta < 2\mu_2 r / (r\mu_2^2 + \ell^2)$ .

To guarantee  $\|(v_k, x_k; y_k) - (0, 0; 0)\|^2 \leq \varepsilon^2$ , we need a large enough  $k$  to have  $v_k^2 \leq \mathcal{O}(\varepsilon^2)$ . Such a  $k$  is required to be at least  $\Omega\left(\frac{r}{\beta\mu_1} \log(1/\varepsilon)\right)$ . Now note that, since  $\mu_1/\mu_2 \leq r \leq \kappa_2/c$  and  $\kappa_2 \geq c$ ,

$$\frac{1}{\beta} > \frac{r\mu_2^2 + \ell^2}{2\mu_2 r} = \frac{L^2 - L\mu_2|r-1|}{2\mu_2 r} = \frac{L^2}{2\mu_2 r} \left(1 - \frac{|r-1|}{\kappa_2}\right) \geq \frac{L^2}{2\mu_2 r} \left(1 - \frac{1}{c}\right).$$

The last inequality is true by minimizing  $\left(1 - \frac{|r-1|}{\kappa_2}\right)$  for  $r \in [\mu_1/\mu_2, \kappa_2/c]$ . If  $r \geq 1$ , it has smaller value when  $r$  is larger: by taking  $r = \kappa_2/c$ , we have  $1 - \frac{\kappa_2/c-1}{\kappa_2} \geq 1 - \frac{1}{c}$ . Otherwise ( $r < 1$ ), which is possible only when  $\mu_1 < \mu_2$ , the term has smaller value when  $r$  is smaller: by taking  $r = \mu_1/\mu_2$ , we have  $1 + \frac{\mu_1/\mu_2-1}{\kappa_2} = 1 + \frac{\mu_1-\mu_2}{L} \geq 1 - \frac{1}{\kappa_2} \geq 1 - \frac{1}{c}$ . Thus, we eventually need  $\Omega\left(\frac{L^2}{\mu_1\mu_2} \log(1/\varepsilon)\right)$  iterations.

Case 2. ( $c/\kappa_1 \leq r \leq \mu_1/\mu_2$ ). Consider

$$f^{(2)}(x; y, w) := \frac{\mu_1}{2}x^2 - \frac{\mu_1}{2r}y^2 + \tilde{\ell}xy - \frac{\mu_2}{2}w^2,$$

where  $\tilde{\ell}^2 = L^2 - \mu_1^2/r - L\mu_1|1 - 1/r| \geq 0$ . Applying Proposition 28, it can be shown that  $f^{(2)} \in \mathcal{F}(L, \mu_1, \mu_2)$ , and  $\mathbf{z}^* = (0; 0; 0)$  is its unique saddle point. Note that, the GDA on  $f^{(2)}$  can be written as

$$\begin{bmatrix} x_{t+1} \\ y_{t+1} \end{bmatrix} = \underbrace{\begin{bmatrix} 1 - \beta\mu_1/r & -\beta\ell/r \\ \beta\ell & 1 - \beta\mu_1/r \end{bmatrix}}_{\mathbf{B}} \begin{bmatrix} x_t \\ y_t \end{bmatrix} = \mathbf{B} \begin{bmatrix} x_t \\ y_t \end{bmatrix}, \quad w_{t+1} = (1 - \beta\mu_2)w_t.$$

Also, the eigenvalues  $\tau$  of  $\mathbf{B}$  is

$$\begin{aligned} \tau &= 1 - \beta\mu_1/r \pm \sqrt{(1 - \beta\mu_1/r)^2 - ((1 - \beta\mu_1/r)^2 + \beta^2\ell^2/r)} \\ &= 1 - \frac{\beta\mu_1}{r} \pm \frac{\beta\ell}{\sqrt{r}}\sqrt{-1}. \end{aligned}$$

The spectral radius is

$$\rho(\mathbf{B}) = \sqrt{(1 - \beta\mu_1/r)^2 + \beta^2\ell^2/r}.$$

Since the eigenvalues are complex conjugates of each other (the magnitudes are the same), both eigenvalues have magnitude  $\rho(\mathbf{B})$ . Then, by Proposition 27,  $\rho(\mathbf{B}) < 1$  is necessary for convergence. To this end, we need  $\beta > 0$  satisfying  $\beta < 2\mu_1 / (\mu_1^2/r + \ell^2)$ .

To guarantee  $\|(x_k; y_k, w_k) - (0; 0; 0)\|^2 \leq \varepsilon^2$ , we need a large enough  $k$  to have  $w_k^2 \leq \mathcal{O}(\varepsilon^2)$ . Such a  $k$  is required to be at least  $\Omega\left(\frac{1}{\beta\mu_2} \log(1/\varepsilon)\right)$ . Now note that, since  $c/\kappa_1 \leq r \leq \mu_1/\mu_2$  and  $\kappa_1 \geq c$ ,

$$\frac{1}{\beta} > \frac{\mu_1^2/r + \ell^2}{2\mu_1} = \frac{L^2 - L\mu_1|1 - 1/r|}{2\mu_1} = \frac{L^2}{2\mu_1} \left(1 - \frac{|1 - 1/r|}{\kappa_1}\right) \geq \frac{L^2}{2\mu_1} \left(1 - \frac{1}{c}\right).$$

The last inequality is true by minimizing  $\left(1 - \frac{|1-1/r|}{\kappa_1}\right)$  for  $r \in [c/\kappa_1, \mu_1/\mu_2]$ . If  $1 > 1/r$ , which is possible only when  $\mu_1 > \mu_2$ , it has smaller value when  $r$  is larger: by taking  $r = \mu_1/\mu_2$ , we have  $1 - \frac{1-\mu_2/\mu_1}{\kappa_1} = 1 - \frac{\mu_1-\mu_2}{L} \geq 1 - \frac{1}{\kappa_1} \geq 1 - \frac{1}{c}$ . Otherwise ( $1 < 1/r$ ), the term has smaller value when  $r$  is smaller: by taking  $r = c/\kappa_1$ , we have  $1 + \frac{1-\kappa_1/c}{\kappa_1} \geq 1 - \frac{1}{c}$ . Thus, we eventually need  $\Omega\left(\frac{L^2}{\mu_1\mu_2} \log(1/\varepsilon)\right)$  iterations.

**Case 3.** ( $r \geq \kappa_2/c$ ). Consider  $f^{(3)}(x; y) = \frac{\mu_1}{2}x^2 - \frac{L}{2}y^2$ . Clearly,  $f^{(3)} \in \mathcal{F}(L, \mu_1, L) \subset \mathcal{F}(L, \mu_1, \mu_2)$  and  $\mathbf{z}^* = (0, 0)$  is its unique saddle point. The GDA on  $f^{(3)}$  can be written as

$$x_{k+1} = \left(1 - \frac{\beta\mu_1}{r}\right)x_k, \quad y_{k+1} = (1 - \beta L)y_k.$$

To guarantee  $\|(x_k; y_k) - (0, 0)\|^2 \leq \varepsilon^2$ , we need a large enough  $k$  to have  $x_k^2 \leq \mathcal{O}(\varepsilon^2)$ . Such a  $k$  is required to be at least  $\Omega\left(\frac{r}{\beta\mu_1} \log(1/\varepsilon)\right)$ . Also, we need  $\beta < 2/L$  to guarantee  $y_k \rightarrow 0$  (i.e., otherwise, it diverges). Combining these facts, we eventually need  $\Omega\left(\frac{Lr}{\mu_1} \log(1/\varepsilon)\right)$  iterations.

**Case 4.** ( $0 < r \leq c/\kappa_1$ ). Consider  $f^{(4)}(x; y) = \frac{L}{2}x^2 - \frac{\mu_2}{2}y^2$ . Clearly,  $f^{(4)} \in \mathcal{F}(L, L, \mu_2) \subset \mathcal{F}(L, \mu_1, \mu_2)$  and  $\mathbf{z}^* = (0, 0)$  is its unique saddle point. The GDA on  $f^{(4)}$  can be written as

$$x_{k+1} = \left(1 - \frac{\beta L}{r}\right)x_k, \quad y_{k+1} = (1 - \beta\mu_2)y_k.$$

To guarantee  $\|(x_k; y_k) - (0, 0)\|^2 \leq \varepsilon^2$ , we need a large enough  $k$  to have  $y_k^2 \leq \mathcal{O}(\varepsilon^2)$ . Such a  $k$  is required to be at least  $\Omega\left(\frac{1}{\beta\mu_2} \log(1/\varepsilon)\right)$ . Also, we need  $\beta < 2r/L$  to guarantee  $x_k \rightarrow 0$  (i.e., otherwise, it diverges). Combining these facts, we eventually need  $\Omega\left(\frac{L}{r\mu_2} \log(1/\varepsilon)\right)$  iterations.

Lastly, we note that the lower iteration complexity bound in terms of the potential function  $V_\lambda$  is equivalent to the complexity in terms of squared distance norm from the (unique) saddle point  $\mathbf{z}^*$ , up to constant factors. This is proved in Lemma 29 that we defer its proof.  $\square$

Here are the postponed/omitted proofs from the proof above.

**Proposition 27.** For a square matrix  $\mathbf{A} \in \mathbb{R}^{m \times m}$  and a sequence of  $m$ -dimensional vectors  $(\mathbf{v}_k)$ , the matrix iteration  $\mathbf{v}_{k+1} = \mathbf{A}\mathbf{v}_k$  converges to  $\mathbf{v}_k \rightarrow \mathbf{0}$  if and only if the spectral radius (i.e., maximum absolute eigenvalue) of  $\rho(\mathbf{A})$  of  $\mathbf{A}$  is less than 1. Furthermore, its convergence speed is characterized by  $\mathcal{O}((\rho(\mathbf{A}) + \varepsilon)^k)$  for any (arbitrarily small)  $\varepsilon > 0$ .

*Proof.* See Horn & Johnson (2012, Theorem 5.6.10-12).  $\square$

**Proposition 28.** Let  $\mu_1, \mu_2$ , and  $L$  be positive numbers such that  $L \geq \max\{\mu_1, \mu_2\}$ . Consider a quadratic function  $f$  on  $\mathbb{R} \times \mathbb{R}$  defined by

$$f(x; y) = \frac{\mu_1}{2}x^2 - \frac{\mu_2}{2}y^2 + \ell xy, \quad \text{where } \ell^2 \leq L^2 - \mu_1\mu_2 - L|\mu_1 - \mu_2|.$$

Then,  $f \in \mathcal{F}(L, \mu_1, \mu_2)$ , and its unique saddle point is  $\mathbf{z}^* = (0, 0)$ .

For example, if  $\mu_1 \geq \mu_2$ ,  $\ell^2 = (L - \mu_1)(L + \mu_2)$  is enough to guarantee  $L$ -smoothness.

*Proof.* The strong-convex-strong-concavity is trivially true. Note that the gradient and hessian of  $f$  is

$$\nabla f(x; y) = \mathbf{H}[x \ y]^\top, \quad \mathbf{H} = \begin{bmatrix} \mu_1 & \ell \\ \ell & -\mu_2 \end{bmatrix}.$$

Since  $\mathbf{H}$  is a non-singular matrix,  $f$  has a unique stationary point at origin ( $x = 0, y = 0$ ). By Proposition 11, it is also a unique saddle & global minimax point.

For any two distinct points  $\mathbf{z}_1 = (x_1; y_1)$  and  $\mathbf{z}_2 = (x_2; y_2)$  in  $\mathbb{R} \times \mathbb{R}$ ,

$$\frac{\|\nabla f(\mathbf{z}_1) - \nabla f(\mathbf{z}_2)\|}{\|\mathbf{z}_1 - \mathbf{z}_2\|} = \frac{\|\mathbf{H}(\mathbf{z}_1 - \mathbf{z}_2)\|}{\|\mathbf{z}_1 - \mathbf{z}_2\|} \leq \|\mathbf{H}\|_2,$$

where  $\|\mathbf{H}\|_2$  is spectral norm (i.e., maximum singular value) of  $\mathbf{H}$ . We would like to show that  $\|\mathbf{H}\|_2 \leq L$ . To this end, it is enough to verify the following two inequalities:

$$\begin{aligned} \det(L^2 \mathbf{I} - \mathbf{H}\mathbf{H}^\top) &= L^4 - (\mu_1^2 + \mu_2^2 + 2\ell^2)L^2 + (\mu_1\mu_2 + \ell^2)^2 \geq 0, \\ \text{trace}(\mathbf{H}\mathbf{H}^\top)/2 &= (\mu_1^2 + \mu_2^2)/2 + \ell^2 \leq L^2. \end{aligned}$$

This is because the characteristic polynomial of  $\mathbf{H}\mathbf{H}^\top$ , or  $\det(\omega\mathbf{I} - \mathbf{H}\mathbf{H}^\top)$ , is a quadratic polynomial of  $\omega$ , and its maximum root should not be greater than  $L^2$ . Let  $\ell^2 = L^2 - \mu_1\mu_2 + a$  for some  $a \in \mathbb{R}$ . Plugging this  $\ell^2$  into both inequalities above, we get

$$a^2 - (\mu_1 - \mu_2)^2 L^2 \geq 0 \quad \text{and} \quad a \leq -(\mu_1 - \mu_2)^2 / 2,$$

respectively. One can check that  $a = -L|\mu_1 - \mu_2|$  is the largest possible  $a$  satisfying both inequalities above. This proves the proposition.  $\square$

Subsequently, we show that if our convergence rate is exponential, then the iteration complexity in terms of  $\|z - z^*\|^2$  is equivalent to that in terms of  $V_\lambda(z) = \lambda[\Phi(\mathbf{x}) - \Phi^*] + [\Phi(\mathbf{x}) - f(\mathbf{x})]$  for  $\text{PL}(\Phi)$ - $\text{PL}$  problem, up to constant factors. This also applies to the function class  $\mathcal{F}(L, \mu_1, \mu_2)$  since it is a subclass of smooth  $\text{PL}(\Phi)$ - $\text{PL}$  functions ( $\cdot$ : Propositions 7 and 10).

**Lemma 29.** *Suppose  $f(\mathbf{x}; \mathbf{y})$  is an  $L$ -smooth function satisfying  $\mathbf{y}$ -side  $\mu_2$ - $\text{PL}$  condition and primal  $\mu_1$ - $\text{PL}$  condition (i.e.,  $\text{PL}(\Phi)$ - $\text{PL}$ ). Suppose  $\mathbf{z}^* = (\mathbf{x}^*; \mathbf{y}^*)$  is a global minimax point of  $f$ . Then, it satisfies*

$$\frac{\lambda\mu_1\mu_2^2}{2(\lambda\mu_1\mu_2 + 2L^2)} \|z - z^*\|^2 \leq V_\lambda(z) \leq \frac{(\lambda + 1)L^3}{\mu_2^2} \|z - z^*\|^2.$$

We remark that the second inequality also holds for general smooth nonconvex- $\text{PL}$  problems.

*Proof.* Let  $\kappa_1 = L/\mu_1$  and  $\kappa_2 = L/\mu_2$  be condition numbers. By the conditions of  $f$  (smoothness and  $\text{PL}$  conditions), for any  $\mathbf{x}$  and  $\mathbf{y}$ ,

$$\begin{aligned} \frac{\mu_1}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 &\stackrel{\text{Prop. 10}}{\leq} \Phi(\mathbf{x}) - \Phi^* \stackrel{\text{Prop. 9}}{\leq} \frac{L(\kappa_2 + 1)}{2} \|\mathbf{x} - \mathbf{x}^*\|^2, \\ \frac{\mu_2}{2} \|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2 &\stackrel{\text{Ass. 4}}{\leq} \Phi(\mathbf{x}) - f(\mathbf{x}; \mathbf{y}) \stackrel{\text{Ass. 1}}{\leq} \frac{L}{2} \|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2, \end{aligned}$$

where  $\mathbf{y}^*(\mathbf{x})$  is a projection of  $\mathbf{y}$  to  $\arg \max_{\mathbf{y}'} f(\mathbf{x}; \mathbf{y}')$ . In particular,  $\mathbf{y}^*(\mathbf{x}^*) = \mathbf{y}^*$ . Since  $\mathbf{y}^*(\mathbf{x})$  is a function of  $\mathbf{x}$  and can differ from  $\mathbf{y}^*$ , we need to bound the term  $\|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2$  using  $\|\mathbf{x} - \mathbf{x}^*\|^2$  and  $\|\mathbf{y} - \mathbf{y}^*\|^2$ . To upper-bound the term  $\|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2$ , note that,

$$\begin{aligned} \|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2 &\leq (\|\mathbf{y} - \mathbf{y}^*(\mathbf{x}^*)\| + \|\mathbf{y}^*(\mathbf{x}) - \mathbf{y}^*(\mathbf{x}^*)\|)^2 \\ &\leq (\|\mathbf{y} - \mathbf{y}^*\| + \kappa_2 \|\mathbf{x} - \mathbf{x}^*\|)^2 \\ &\leq (1 + \kappa_2^2) (\|\mathbf{y} - \mathbf{y}^*\|^2 + \|\mathbf{x} - \mathbf{x}^*\|^2). \end{aligned}$$

The first inequality holds by triangle inequality, the second inequality holds by Proposition 8, and the last inequality holds by Cauchy-Schwarz inequality.<sup>12</sup> To lower-bound in a similar way, note that for any constant  $a > 0$ ,

$$\begin{aligned} \|\mathbf{y} - \mathbf{y}^*\|^2 &\leq (\|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\| + \|\mathbf{y}^*(\mathbf{x}) - \mathbf{y}^*(\mathbf{x}^*)\|)^2 \\ &\leq \left( \|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\| + \frac{\kappa_2}{\sqrt{a}} \cdot \sqrt{a} \|\mathbf{x} - \mathbf{x}^*\| \right)^2 \\ &\leq \left( 1 + \frac{\kappa_2^2}{a} \right) (\|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2 + a \|\mathbf{x} - \mathbf{x}^*\|^2). \\ \therefore \|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2 &\geq \frac{1}{1 + \kappa_2^2/a} \|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2 - a \|\mathbf{x} - \mathbf{x}^*\|^2. \end{aligned}$$

<sup>12</sup> $(ax + by)^2 \leq (a^2 + b^2)(x^2 + y^2)$  for real numbers  $a, b, x, y$ .

Now we can prove the inequalities in the lemma. We first show the second one. Applying  $\kappa_2 \geq 1$  multiple times,

$$\begin{aligned} V_\lambda(\mathbf{x}; \mathbf{y}) &= \lambda[\Phi(\mathbf{x}) - \Phi^*] + [\Phi(\mathbf{x}) - f(\mathbf{z})] \\ &\leq \frac{\lambda L(\kappa_2 + 1)}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{L}{2} \|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2 \\ &\leq \left( \frac{\lambda L(\kappa_2 + 1)}{2} + \frac{L(1 + \kappa_2^2)}{2} \right) \|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{L(1 + \kappa_2^2)}{2} \|\mathbf{y} - \mathbf{y}^*\|^2 \\ &\leq (\lambda + 1)L\kappa_2^2 \left( \|\mathbf{x} - \mathbf{x}^*\|^2 + \|\mathbf{y} - \mathbf{y}^*\|^2 \right) = \frac{(\lambda + 1)L^3}{\mu_2^2} \|\mathbf{z} - \mathbf{z}^*\|^2. \end{aligned}$$

To show the first inequality of the lemma, let  $a = \frac{\lambda\mu_1}{2\mu_2}$ .

$$\begin{aligned} V_\lambda(\mathbf{x}; \mathbf{y}) &\geq \frac{\lambda\mu_1}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{\mu_2}{2} \|\mathbf{y} - \mathbf{y}^*(\mathbf{x})\|^2 \\ &\geq \left( \frac{\lambda\mu_1}{2} - \frac{\mu_2 a}{2} \right) \|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{\mu_2}{2(1 + \kappa_2^2/a)} \|\mathbf{y} - \mathbf{y}^*\|^2 \\ &\geq \frac{\lambda\mu_1}{4} \|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{\lambda\mu_1}{4(a + \kappa_2^2)} \|\mathbf{y} - \mathbf{y}^*\|^2 \\ &\geq \frac{\lambda\mu_1}{4(a + \kappa_2^2)} \left( \|\mathbf{x} - \mathbf{x}^*\|^2 + \|\mathbf{y} - \mathbf{y}^*\|^2 \right) = \frac{\lambda\mu_1\mu_2^2}{2(\lambda\mu_1\mu_2 + 2L^2)} \|\mathbf{z} - \mathbf{z}^*\|^2. \end{aligned}$$

This concludes the proof.  $\square$

The equivalence of iteration complexities for achieving  $\|\mathbf{z}_K - \mathbf{z}^*\|^2 \leq \varepsilon^2$  or  $V_\lambda(\mathbf{z}_K) \leq \varepsilon^2$  is quite straightforward from this lemma, as long as the convergence speed is exponential. For example, suppose we have an upper convergence bound  $\|\mathbf{z}_K - \mathbf{z}^*\|^2 \leq a \exp(-K/r)$  for some constants  $a, r > 0$ . This implies an upper iteration complexity bound  $K = \mathcal{O}(r \log(1/\varepsilon))$  sufficient to achieve  $\|\mathbf{z}_K - \mathbf{z}^*\|^2 \leq \varepsilon^2$ . Then by Lemma 29, we also have  $V_\lambda(\mathbf{z}_K)^2 \leq a' \exp(-K/r)$  where  $a' = a(\lambda + 1)L^3/\mu_2^2$  is also a constant. This implies a lower iteration complexity bound  $K = \mathcal{O}(r \log(1/\varepsilon))$  as well, sufficient to achieve  $V_\lambda(\mathbf{z}_K)^2 \leq \varepsilon^2$ . The other way of complexity translation operates with a similar logic.

## F REMARK ON SMOOTHNESS ASSUMPTIONS AND LOWER BOUND OF WITH-REPLACEMENT SGD(A)

During the discussion phase of the conference, a reviewer raised a question about whether or not the *component smoothness* (Assumption 1) is more crucial than the without-replacement component sampling for faster convergence. However, we would like to claim that the component smoothness alone is not sufficient for improving the convergence rate for with-replacement SGD(A). To this end, we provide some formal results on lower convergence bounds. For simplicity, we use mini-batches of size 1 throughout this appendix.

Firstly, the theorem below provides a lower bound on with-replacement SGD for minimization problems. Readers can also verify that an analogous lower bound holds for SGD with unbiased and independently sampled gradient oracle for more general stochastic minimization problems. The proof will appear later in this appendix.

**Theorem 30.** *For any step size  $\eta > 0$ , there exists a real-valued strongly-convex function  $f(\mathbf{x})$  defined on  $\mathbb{R}^d$  with  $f^* := \min_{\mathbf{x}} f(\mathbf{x})$ , satisfying:*

1.  $f$  consists of  $n > 1$  smooth component functions  $f_i$ :  $f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$ , where each component  $f_i$  is smooth;
2. After running  $T > 1$  iterations of with-replacement SGD (with mini-batch size 1) starting from  $x_0 \in \mathbb{R}^d$ , the last iterate  $x_T$  satisfies  $\mathbb{E}[f(x_T) - f^*] \geq \Omega(1/T)$ , where the expectation is taken with respect to the randomness of i.i.d. index choice at each iteration.

Next, we show this theorem naturally induces a convergence lower bound for the minimax counterpart: *with-replacement SGDA*. Consider a (finite-sum) minimax problem  $\min_x \max_y g(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}) - f(\mathbf{y})$ , where  $f = \frac{1}{n} \sum_{i=1}^n f_i$  is a worst-case function in the proof of Theorem 30. Here, the minimax problem on  $g$  can be solved by minimizing  $f$ . Moreover, since the primal function  $\Phi(\mathbf{x}) := \max_y g(\mathbf{x}, \mathbf{y})$  associated with  $g$  is in fact the same as  $f(\mathbf{x}) - f^*$ , the potential function  $V_\lambda(\mathbf{x}, \mathbf{y}) := \lambda[\Phi(\mathbf{x}) - (\min_x \Phi(\mathbf{x}))] + [\Phi(\mathbf{x}) - g(\mathbf{x}, \mathbf{y})]$  becomes the same as  $\lambda(f(\mathbf{x}) - f^*) + (f(\mathbf{y}) - f^*)$  for a constant  $\lambda > 0$ . Combining these facts, we can immediately obtain the following lower convergence bound of with-replacement SGDA.

**Corollary 1.** *There exists a strongly-convex-strongly-concave function  $g(\mathbf{x}, \mathbf{y}) := \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{x}, \mathbf{y})$  consisting of  $n$  smooth component functions  $g_i$ , where the last iterate  $(\mathbf{x}_T, \mathbf{y}_T)$  of with-replacement SGDA satisfies  $\mathbb{E}[V(\mathbf{x}_T, \mathbf{y}_T)] \geq \Omega(1/T)$ .*

Corollary 1 formally proves that with-replacement SGDA on strongly-convex-strongly-concave minimax problems with smooth components has a worst-case convergence rate  $\Omega(1/T)$ . This in fact matches the  $\mathcal{O}(1/T)$  upper bound obtained for primal-PŁ-PŁ problems by Yang et al. (2020). Considering that strongly-convex-strongly-concave functions form a strict subset of primal-PŁ-PŁ functions, Corollary 1 establishes that adding component smoothness assumption does not provide further speed up for with-replacement SGDA.

In contrast, our theoretical result in Theorem 2 shows that SGDA-RR has a much faster convergence rate  $\mathbb{E}[V_\lambda] \leq \tilde{\mathcal{O}}(\frac{1}{nK^2})$  for primal-PŁ-PŁ minimax problems, where  $K$  is the number of epochs. One can check that our  $\tilde{\mathcal{O}}(\frac{1}{nK^2})$  bound is faster than the tight convergence rate  $\Theta(1/T)$  of with-replacement SGDA by simply plugging in  $T = nK$ . In light of Corollary 1 we proved, we can now claim that the improvement can be solely attributed to RR.

Although we do not provide a lower bound for more general nonconvex-PŁ problems here, we believe the more challenging case of nonconvex-PŁ lower bound is a topic for another separate paper. Nonetheless, we conjecture that the speed up by SGDA-RR in nonconvex-PŁ settings is also due to the effect of RR, not component smoothness.

From now on, we provide the postponed proof of Theorem 30.

*Proof of Theorem 30.* We construct worst-case functions with quadratic functions on  $\mathbb{R}$ , which are clearly  $L$ -smooth for a fixed constant  $L > 0$ . Then, it is easy to extend the logic to the functions with domains of higher dimensions. Let  $x_0 \in \mathbb{R}$  be the initial iterate.

**Case 1** ( $\frac{1}{LT} \leq \eta \leq (\frac{2}{L} - \frac{1}{LT})$ ). Note that the condition on the step size,  $\frac{1}{LT} \leq \eta \leq (\frac{2}{L} - \frac{1}{LT})$ , is equivalent to an inequality  $(1 - \eta L)^2 \leq (1 - 1/T)^2$ .

We first assume  $n$  is an even number. We will encounter the case with an odd  $n > 1$  a bit later. Consider  $f(x) = \frac{L}{2}x^2$  consisting of even number of components  $f_i$ 's defined by

$$f_i(x) = \begin{cases} \frac{L}{2}x^2 + \nu x, & (i \leq \frac{n}{2}), \\ \frac{L}{2}x^2 - \nu x, & (i \geq \frac{n}{2} + 1), \end{cases}$$

for some number  $\nu \in \mathbb{R}$ . At each iteration  $t \geq 1$ , we choose a component index  $i(t) \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([n])$  (with-replacement sampling). Then we can write the chosen component function at iteration  $t$  as  $f_{i(t)} = \frac{L}{2}x^2 - s_t \nu x$  for some i.i.d. random variable  $s_t \sim \text{Unif}(\{\pm 1\})$ . Accordingly, an SGD step can be written as

$$x_t = x_{t-1} - \eta \nabla f_{i(t)}(x_{t-1}) = (1 - \eta L)x_{t-1} + \eta s_t \nu.$$

By applying telescopic sum, we have

$$x_T = (1 - \eta L)^T x_0 + \eta \nu \sum_{t=1}^T (1 - \eta L)^{(T-t)} \cdot s_t.$$

Taking squares and expectations (with respect to the random variables  $s_1, \dots, s_T$ ) to both sides, we have

$$\mathbb{E}[x_T^2] = (1 - \eta L)^{2T} x_0^2 + \eta^2 \nu^2 \sum_{t=1}^T (1 - \eta L)^{2(T-t)},$$



by applying the fact that  $s_t$ 's are zero-mean independent random variables with absolute values 1:

$$\mathbb{E}[s_t \cdot s_{t'}] = \begin{cases} 0, & t \neq t' \quad (\because \text{independent}), \\ 1, & t = t' \quad (\because s_t^2 = 1). \end{cases}$$

We calculate the sum above as follows: since  $(1 - \eta L)^2 \leq (1 - 1/T)^2$  and  $(1 - 1/T)^T \leq e^{-1}$ ,

$$\sum_{t=1}^T (1 - \eta L)^{2(T-t)} = \frac{1 - (1 - \eta L)^{2T}}{1 - (1 - \eta L)^2} \geq \frac{1 - (1 - \frac{1}{T})^{2T}}{2\eta L(1 - \frac{\eta L}{2})} \geq \frac{1 - e^{-2}}{2\eta L}.$$

With this inequality, and since  $(1 - \eta L)^{2T} x_0^2 \geq 0$ , we can lower-bound the expectation  $\mathbb{E}[x_T^2]$ :

$$\mathbb{E}[x_T^2] \geq \eta^2 \nu^2 \cdot \frac{1 - e^{-2}}{2\eta L} = \frac{(1 - e^{-2})\nu^2}{2L} \eta \geq \frac{(1 - e^{-2})\nu^2}{2L^2 T}.$$

Since  $f$  has a minimum  $f^* = 0$  at  $x = 0$ , we eventually have

$$\mathbb{E}[f(x_T) - f^*] = \frac{L}{2} \mathbb{E}[x_T^2] \geq \frac{(1 - e^{-2})\nu^2}{4LT} = \Omega\left(\frac{\nu^2}{LT}\right).$$

Now we consider the case when the number of components  $n > 1$  is odd. Consider  $f_n(x) \equiv 0$  and let the remaining  $n - 1$  components be the same as the case above (with an even number of components). Note that the zero-component  $f_n$  does not affect the trajectory of SGD (*i.e.*, the points visited by SGD) and the optimality of  $f$  ( $f^* = 0$  at  $x = 0$ ), while the whole objective function becomes  $f(x) = \frac{n-1}{n} \cdot \frac{L}{2} x^2$ . Thus, it can be easily shown that the  $\Omega\left(\frac{\nu^2}{LT}\right)$  lower bound also holds.

**Case 2** ( $0 < \eta < \frac{1}{LT}$  or  $\eta > (\frac{2}{L} - \frac{1}{LT})$ ). From the condition on the step size, we have  $(1 - \eta L)^2 > (1 - 1/T)^2$ . Consider  $f_i(x) = \frac{L}{2} x^2$  for every  $i \in [n]$ : every components are the same. In this case, we show that the last iterate of SGD is bounded below by a constant with respect to  $T > 1$ .

At each iteration  $t \geq 1$ , we obtain  $x_t = (1 - \eta L)x_{t-1}$  by a step of SGD. Then, applying  $T \geq 2$ ,

$$x_T^2 = (1 - \eta L)^{2T} \cdot x_0^2 > \left(1 - \frac{1}{T}\right)^{2T} x_0^2 \geq \left(1 - \frac{1}{2}\right)^4 x_0^2 = \frac{x_0^2}{16}.$$

Since  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) = \frac{L}{2} x^2$  has a minimum  $f^* = 0$  at  $x = 0$ , we have

$$f(x_T) - f^* > \frac{Lx_0^2}{32} = \Omega(1) \cdot Lx_0^2.$$

□

## G EXPERIMENTS: QUADRATIC GAMES

In this appendix, we provide a more detailed illustration of our numerical evaluations on quadratic games introduced in Section 6. Recall that the objective function  $f$  and its component functions  $f_i$  are given in Equation (5) as

$$\begin{aligned} f(\mathbf{x}; \mathbf{y}) &= \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{x}^\top \mathbf{B} \mathbf{y} - \frac{1}{2} \mathbf{y}^\top \mathbf{C} \mathbf{y}, \\ f_i(\mathbf{x}; \mathbf{y}) &= \frac{1}{2} \mathbf{x}^\top \mathbf{A}_i \mathbf{x} + \mathbf{x}^\top \mathbf{B}_i \mathbf{y} - \frac{1}{2} \mathbf{y}^\top \mathbf{C}_i \mathbf{y} + \mathbf{u}_i^\top \mathbf{x} - \mathbf{v}_i^\top \mathbf{y}. \end{aligned}$$

We choose the same dimensions for the variables  $\mathbf{x} \in \mathbb{R}^{d_x}$  and  $\mathbf{y} \in \mathbb{R}^{d_y}$ : we set  $d_x = d_y = d$ .

### G.1 PARAMETER CHOICES

To sample the matrix  $\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i \in \mathbb{R}^d$  satisfying that  $\mu_C \mathbf{I}_d \preceq \mathbf{C}$  and  $\|\mathbf{C}_i\|_2 \leq L_C$ , we first randomly generate an orthogonal matrix  $\mathbf{Q}_C \in \mathbb{R}^{d \times d}$  (*i.e.*,  $\mathbf{Q}_C \mathbf{Q}_C^\top = \mathbf{I}_d$ ), by taking advantage of

the QR-decomposition of a random matrix. Then, we generate the eigenvalues of  $C_i$ 's as follows. We sample the entries of  $n$  vectors  $\lambda_i^C \in \mathbb{R}^d$  ( $i \in [n]$ ) uniformly from the interval  $[\mu_C, L_C]$ . We add some level of perturbations to some entries of each  $\lambda_i^C$ ; we replace some entries to the numbers in an interval  $[-L_C, \mu_C]$ , keeping the entries of the vector  $\frac{1}{n} \sum_{i=1}^n \lambda_i^C$  in the interval  $[\mu_C, L_C]$ . Finally, we define  $C_i = Q_C \Lambda_i^C Q_C^\top$  where  $\Lambda_i^C = \text{diag}(\lambda_i^C)$ . Because of the perturbation step, some  $C_i$ 's are not positive definite, thereby some components  $f_i$ 's become non-(strongly-)concave in  $\mathbf{y}$ .

Next, we sample the matrix  $B_i$ 's. There are no requirements for  $B$  but  $\|B_i\|_2 \leq L_B$ ;  $B_i$ 's are even not necessarily symmetric when  $d_x \neq d_y$ . Thus, we first generate the orthogonal matrices  $U_i^B$  and  $V_i^B$  by taking advantage of the singular value decomposition of random matrices. Then, we generate the singular values of  $B_i$ 's by sampling the entries of  $n$  vectors  $\sigma_i^B$  uniformly from the interval  $[0, L_C]$ . After that, we define  $B_i = U_i^B \Sigma_i^B V_i^B$  where  $\Sigma_i^B = \text{diag}(\sigma_i^B)$ . We typically want to take a larger  $L_B$  than  $L_C$  to strengthen the interaction term  $\mathbf{x}^\top B \mathbf{y}$ .

Recall that the primal function  $\Phi$  associated with  $f$  is explicitly written as

$$\Phi(\mathbf{x}) = \max_{\mathbf{y} \in \mathbb{R}^d} f(\mathbf{x}; \mathbf{y}) = \frac{1}{2} \mathbf{x}^\top (\mathbf{A} + \mathbf{B} \mathbf{C}^{-1} \mathbf{B}^\top) \mathbf{x} := \frac{1}{2} \mathbf{x}^\top \mathbf{M} \mathbf{x}. \quad (28)$$

Note that the inverse of  $C$  can be efficiently computed as  $C^{-1} = Q_C (\Lambda^C)^{-1} Q_C^\top$ .

Before generating the matrices  $A_i$ 's, we first generate  $M_i$ 's satisfying that  $\frac{1}{n} \sum_{i=1}^n M_i = M$  and the nonzero eigenvalues of positive semidefinite  $M$  are in the interval  $[\mu_M, L_M]$ . The process of sampling  $M_i$ 's is almost identical to how to sample  $C_i$ 's. One notable difference is,  $M_i$ 's and  $M$  are forced to have  $r (< d)$  zero eigenvalues: this makes  $M$  a positive semidefinite (but not strictly positive definite) matrix of rank  $d - r$ . Moreover, we get the  $\mu_M$ -PL( $\Phi$ ) condition in  $\mathbf{x}$  as follows:

**Proposition 31.** *Consider a positive semidefinite matrix  $M \in \mathbb{R}^d$ . If the smallest nonzero eigenvalue of  $M$  is  $\mu$ , then  $\Phi(\mathbf{x}) := \frac{1}{2} \mathbf{x}^\top M \mathbf{x}$  is  $\mu$ -PL in  $\mathbf{x}$ . Also,  $\Phi^* = \min_{\mathbf{x}} \Phi(\mathbf{x}) = 0$ .*

*Proof.* Apply the eigendecomposition of  $M$ :  $M = Q \Lambda Q^\top$ . Let  $\overline{M} = \Lambda^{1/2} Q^\top$ . Then, we have  $\Phi(\mathbf{x}) = \frac{1}{2} \|\overline{M} \mathbf{x}\|^2$ , which implies that  $\Phi(\mathbf{x}) \geq 0$  ( $\forall \mathbf{x}$ ) and in fact  $\Phi^* = 0$ . Note that  $\frac{1}{2} \|\mathbf{x}\|^2$  is 1-strongly convex. Also, the minimum nonzero singular value of  $\overline{M}$  is  $\sqrt{\mu}$  ( $\because M = \overline{M}^\top \overline{M}$ ). Therefore, by the proof of Proposition 12,  $\Phi(\mathbf{x})$  is a  $\mu$ -PL function of  $\mathbf{x}$ . Lastly, we note that  $\Phi(\mathbf{x})$  is not strongly convex in general, especially when  $M$  is a rank-deficient matrix.  $\square$

Typically, the spectral norm  $\|M\|_2$  is known to be bounded above by  $\|A\|_2 + L_B^2/\mu_C$  in *worst-case* (Nouiehed et al., 2019; Li et al., 2022). However, since we sample  $M$  without knowing the exact form of  $A_i$ 's while we want to control the spectral norm  $\|A_i\|_2$  not too large (for smoothness of  $f_i$ ), we (empirically) decide to choose rather smaller  $L_M$ : simply, we choose  $L_M = L_B$ .

Now we let  $A_i = M_i - \mathbf{B} \mathbf{C}^{-1} \mathbf{B}^\top$  and  $\mathbf{A} = \frac{1}{n} \sum_{i=1}^n A_i$  to satisfy Equation (28). We emphasize that  $\mathbf{A}$  may have negative eigenvalues; the objective is nonconvex in  $\mathbf{x}$  in general. We have checked this is true across the experimental settings. Also, we let  $L := \max\{\|\mathbf{A}\|_2, L_B, L_C\}$  for further parameter selection. (In fact, because of our choice of parameter values,  $L$  was always equal to  $L_B$  in our experiments.)

Furthermore, we generate the vectors  $\mathbf{u}_i$ 's and  $\mathbf{v}_i$ 's satisfying  $\sum_{i=1}^n \mathbf{u}_i = \mathbf{0} = \sum_{i=1}^n \mathbf{v}_i$ . The entries of these vectors are uniformly sampled from an interval  $[-\Delta, \Delta]$ , thereby the average of entries is centered to zero. In addition, to verify our theory, we choose the step-sizes of the form  $\beta = c_1 \cdot b/nL$  and  $\alpha = c_0 \cdot \beta/\kappa_2^2$  for some constants  $c_0$  and  $c_1$  and batch size  $b$ .

Lastly, we specify the values of parameters described above:  $n = 100$ ,  $d = 25$ ,  $\mu_M = \mu_C$ , and  $L_C = 1 < L_M = L_B$ . The constants  $c_0$  and  $c_1$  are tuned among  $10^{\{-2, -1.5, \pm 1, \pm 0.5, 0\}}$ . In the following subsections, we investigate the effects of the change of

- (i)  $\Delta \in \{10, 20, 40\}$ , determining the discrepancy between components,
- (ii) condition number  $\kappa_2 \in \{5, 10, 20\}$ , determined by  $L_B$  and  $\mu_C$ , and
- (iii) batch size  $b \in \{1, 25, 50, 100\}$ ,

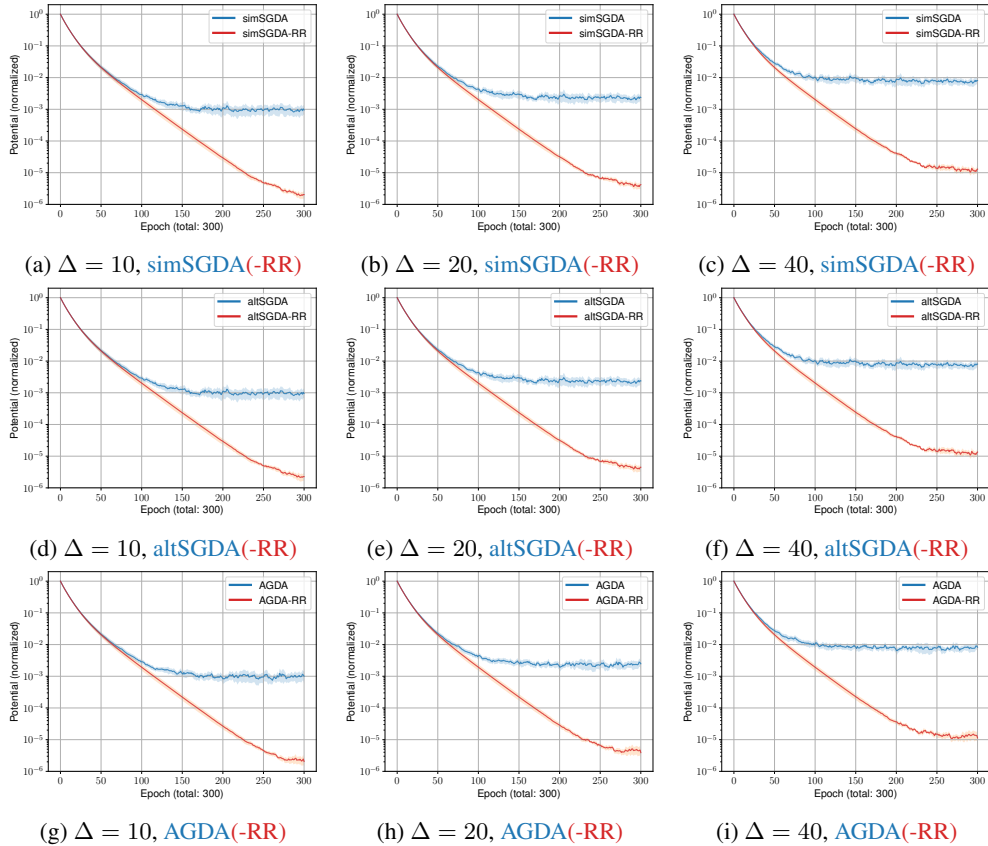


Figure 2: Comparisons by changing the value of  $\Delta \in \{10, 20, 40\}$ . Solid lines: average across 10 different runs. Shaded regions: 95% confidence intervals ( $\pm 1.96$  std). The vertical axes are on a logarithmic scale.

from the plots of the values potential function  $V_\lambda(\mathbf{x}; \mathbf{y}) = (1 + \lambda)\Phi(\mathbf{x}) - f(\mathbf{x}; \mathbf{y})$  over epochs.<sup>13</sup> (Numbers in bold font above are the default values of parameters.)

## G.2 COMPARISON: THE EFFECT OF COMPONENT DISCREPANCY

Notice that the discrepancy between component functions gets larger as  $\Delta$  grows. Technically, one can check that the gradient variance (that we controlled in Assumption 2) is proportional to the norms of the vectors  $\mathbf{u}_i$  and  $\mathbf{v}_i$ . Moreover, we have already discussed that the gap between convergence speeds of SGDA and SGDA-RR becomes larger especially when the gradient variance is large.

Now, we present the results of numerical experiments by varying the values of  $\Delta$  to 10, 20, and 40, while fixing  $L_B = 4$ ,  $\mu_C = 0.4$ ,  $b = 1$ , and other experiment parameters. As shown in Figure 2, we can observe that the difference between the random-reshuffling algorithm and the uniform-sampling algorithm gets larger as  $\Delta$  increases.

## G.3 COMPARISON: THE EFFECT OF CONDITION NUMBER

Here, we present the results of experiments by varying the values of  $\kappa_2$  to 5, 10, and 20, while fixing  $\Delta = 20$ ,  $b = 1$ , and other experiment parameters. To this end, we applied the parameter settings for  $L_B$  and  $\mu_C$  as  $(L_B, \mu_C) = (2.5, 0.5), (4, 0.4), (5, 0.25)$ , respectively.

<sup>13</sup>During and after the discussion phase, we performed some more experiments. As we tried to plot all the results over iterations, the size of the figures in pdf format became too large. Consequently, in this appendix, we only plot the results over epochs to reduce the file size of the figures.

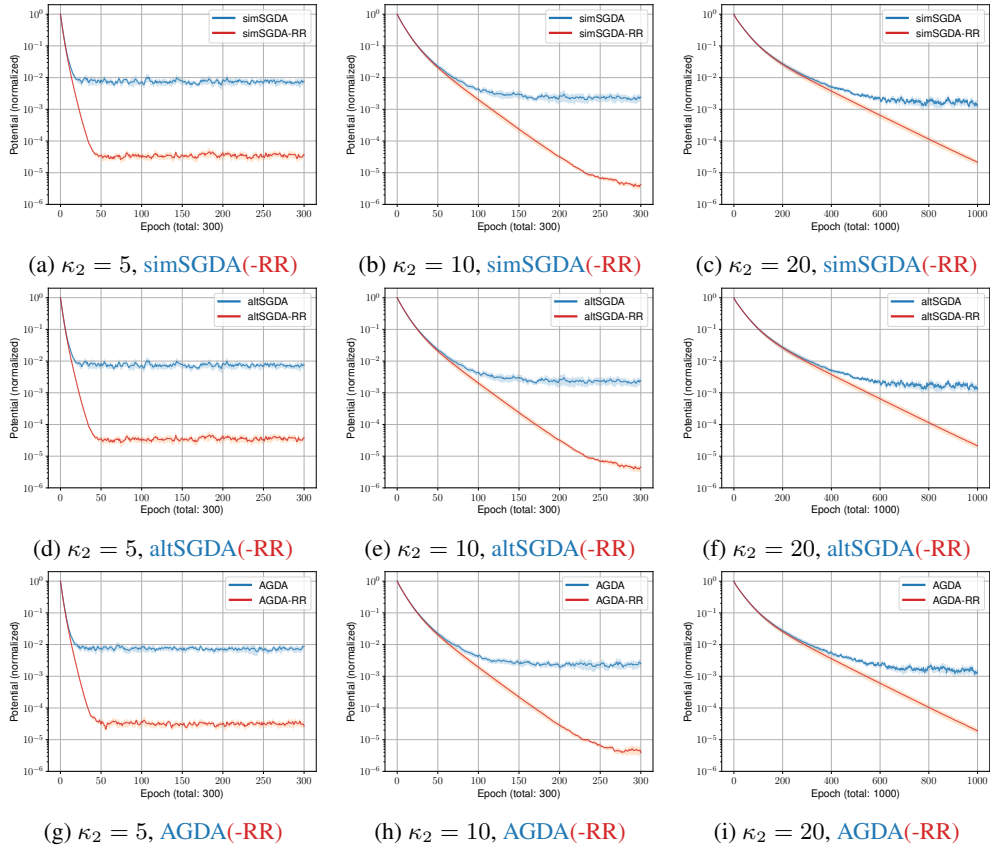


Figure 3: Comparisons by changing the value of  $\kappa_2 = L/\mu_C \in \{5, 10, 20\}$ . Solid lines: average across 10 different runs. Shaded regions: 95% confidence intervals ( $\pm 1.96$  std). The vertical axes are on a *logarithmic scale*. Note: we run **1000 epochs** for  $\kappa_2 = 20$  (see the rightmost column), whereas we run **300 epochs** for the other  $\kappa_2 \in \{5, 10\}$  (see the leftmost & middle columns).

The results are shown in Figure 3. We observe that more epochs are required for convergence when  $\kappa_2$  increases, regardless of the type of algorithm. One may think that the performance gap between RR-based/non-RR-based algorithms is small when  $\kappa_2$  is huge. However, when we run the algorithm for an extended number of epochs, we observe a significant gap in convergence speeds.

#### G.4 COMPARISON: THE EFFECT OF BATCH SIZE

The last comparison is about the effect of batch size  $b \in \{1, 25, 50, 100\}$ . Recall that we linearly scale the step sizes as the batch size changes. However, since the number of epochs is fixed, the number of iterations decreases as  $b$  gets larger.

As the readers can notice, the convergence behavior of SGDA (resp., SGDA-RR) and AGDA (resp., AGDA-RR) are similar in our construction of quadratic games. Thus, in this subsection, we only compare simSGDA and its variants. Rather, we introduce two more methods of component choice other than with-replacement uniform sampling and random reshuffling:

- **WORB**(WithOut-Replacement mini-Batching): every mini-batch is without-replacement & uniformly-randomly sampled, while any pair of mini-batches in an epoch may have some indices in common; the same as *b-minibatch sampling* (Loizou et al., 2021).
- **NS**(No Shuffle): accessing  $1, \dots, n$  in its predefined order to construct mini-batches; without-replacement but deterministic. Remark: for *minimization* problems, SGD with NS is usually referred to as *incremental gradient* (IG) algorithm (Mishchenko et al., 2020).

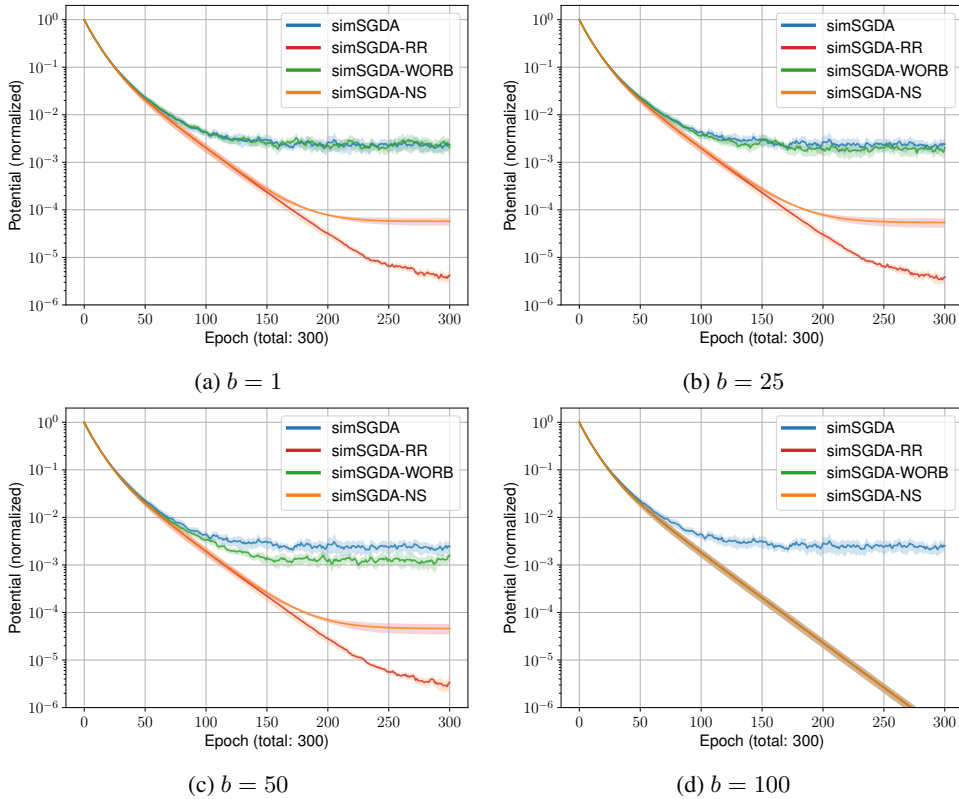


Figure 4: Comparisons of `simSGDA(-RR,-WORB,-NS)` as changing  $b \in \{1, 25, 50, 100\}$ . Solid lines: average across 10 different runs. Shaded regions: 95% confidence intervals ( $\pm 1.96$  std). The vertical axes are on a *logarithmic scale*.

These two methods are somewhat related to without-replacement component sampling, whereas they are both different from RR which uniformly randomly samples a permutation of  $[n]$  every epoch. We call `simSGDA` using mini-batches sampled by WORB and NS as *simSGDA-WORB* and *simSGDA-NS*, respectively. Remarks: If  $b = 1$ , `simSGDA-WORB` becomes the same algorithm as vanilla `simSGDA`. Also, since we choose  $n = 100$ , if  $b = n = 100$ , all three algorithms `simSGDA-RR/-WORB/-NS` become the same as deterministic & full-batch (simultaneous) GDA.

The results are shown in Figure 4. One can notice that the potential plots of `simSGDA`, `simSGDA-RR`, and `simSGDA-NS` are respectively the same even if we change the batch size ( $b < 100$ ). Also, if  $b > 1$ , `simSGDA-WORB` has better performance than vanilla `simSGDA`. These imply that without-replacement mini-batches benefit the convergence speed to some extent in our quadratic game. However, the result of experiments also implies that both (i) without-replacement *per epoch* (*i.e.*, shuffling) and (ii) randomization are indeed essential for fast convergence in our quadratic game experiments. In particular, WORB requires a very large batch size but still has a much slower convergence rate than RR (see Figure 4c which is the case of using half of the total components at each iteration).

## H OMITTED COMPARISON WITH RELATED WORKS

### H.1 COMPARISON WITH XIE ET AL. (2021)

To specialize Xie et al. (2021, Theorem 3) to the single-machine setup and discuss their results in terms of our notation, we need to replace their symbols

$$(T, S, K, \sigma_1^2, \sigma_2^2, G_1^2, G_2^2, L_{12}, L_f, \mu, L_\Phi, \mathcal{L}_0, \eta_t, \gamma_t)$$

with the following symbols from our notation

$$(K, 1, n, 0, 0, B, B, L, L, \mu_2, L(\kappa_2 + 1), V_\lambda(\mathbf{z}_0^1), \alpha, \beta),$$

and also put  $A = 0$  (their analysis only applies uniformly bounded component variance per machine). Then we can *naively* translate the bound of Xie et al. (2021, Theorem 3) to our language as

$$\min_{k \in [K]} \mathbb{E} \left[ \|\Phi(\mathbf{x})\|^2 \right] \stackrel{?)}{\leq} \mathcal{O} \left( \frac{\kappa_2 L V_\lambda(\mathbf{z}_0^1)}{K} + \kappa_2^2 \left( \frac{L^2 B V_\lambda(\mathbf{z}_0^1)^2}{K^2} \right)^{1/3} \right).$$

To the best of our knowledge, however, we believe there may be a mistake in the proof of Xie et al. (2021, Appendix C.4). From the inequalities on the last page of their paper, we notice that the term  $\frac{40L_{12}^2 \mathcal{L}_0}{\mu^2 \gamma K T}$  might be missing in a step, where  $\gamma$  is chosen to be the minimum of several terms including  $\frac{1}{87L_f K}$ . Thus, as far as we can tell, it seems inevitable that this omitted term would lead to an additional term  $\frac{3480L_f L_{12}^2 \mathcal{L}_0}{\mu^2 T}$  in the final bound. By combining this to their bound and re-translating it, we eventually have

$$\min_{k \in [K]} \mathbb{E} \left[ \|\Phi(\mathbf{x})\|^2 \right] \leq \mathcal{O} \left( \frac{\kappa_2^2 L V_\lambda(\mathbf{z}_0^1)}{K} + \kappa_2^2 \left( \frac{L^2 B V_\lambda(\mathbf{z}_0^1)^2}{K^2} \right)^{1/3} \right),$$

since their  $L_{12}^2/\mu^2$  translates to our  $\kappa_2^2$ . Therefore, their result actually shows the same dependency on condition number  $\kappa_2$  as our Theorem 1. Nevertheless, comparing the terms related to the component-wise variance  $B$ , ours is better. In the second term in the bound above does not shrink even when the number of iterations (per machine & per communication) grows. In our case (Theorem 1), however, the dominant term (in  $K$ ) can be briefly written as  $\mathcal{O} \left( \left( \frac{B}{nK^2} \right)^{1/3} \right)$  which can diminish with large  $n$ , *i.e.*, the number of iterations per epoch.