
Proximal and Federated Random Reshuffling

Anonymous Author(s)

Affiliation

Address

email

Abstract

Random Reshuffling (RR), also known as Stochastic Gradient Descent (SGD) without replacement, is a popular and theoretically grounded method for finite-sum minimization. We propose two new algorithms: Proximal and Federated Random Reshuffling (ProxRR and FedRR). The first algorithm, ProxRR, solves composite finite-sum minimization problems in which the objective is the sum of a (potentially non-smooth) convex regularizer and an average of n smooth objectives. ProxRR evaluates the proximal operator once per epoch only. When the proximal operator is expensive to compute, this small difference makes ProxRR up to n times faster than algorithms that evaluate the proximal operator in every iteration, such as proximal (stochastic) gradient descent. We give examples of practical optimization tasks where the proximal operator is difficult to compute and ProxRR has a clear advantage. One such task is federated or distributed optimization, where the evaluation of the proximal operator corresponds to communication across the network. We obtain our second algorithm, FedRR, as a special case of ProxRR applied to federated optimization, and prove it has a smaller communication footprint than either distributed gradient descent or Local SGD. Our theory covers both constant and decreasing stepsizes, and allows for importance resampling schemes that can improve conditioning, which may be of independent interest. Our theory covers both convex and nonconvex regimes. Finally, we corroborate our results with experiments on real data sets.

1 Introduction

Modern theory and practice of training supervised machine learning models is based on the paradigm of regularized empirical risk minimization (ERM) [Shalev-Shwartz and Ben-David, 2014]. While the ultimate goal of supervised learning is to train models that generalize well to unseen data, in practice only a finite data set is available during training. Settling for a model merely minimizing the average loss on this training set—the empirical risk—is insufficient, as this often leads to over-fitting and poor generalization performance in practice. Due to this reason, empirical risk is virtually always amended with a suitably chosen regularizer whose role is to encode prior knowledge about the learning task at hand, thus biasing the training algorithm towards better performing models.

The regularization framework is quite general and perhaps surprisingly it also allows us to consider methods for federated learning (FL)—a paradigm in which we aim at training model for a number of clients that do not want to reveal their data [Konečný et al., 2016, McMahan et al., 2017, Kairouz et al., 2019]. The training in FL usually happens on devices with only a small number of model updates being shared with a global host. To this end, Federated Averaging algorithm has emerged that performs Local SGD updates on the clients’ devices and periodically aggregates their average. Its analysis usually requires special techniques and deliberately constructed sequences hindering the research in this direction. We shall see, however, that the convergence of our FedRR follows from merely applying our algorithm for regularized problems to a carefully chosen reformulation.

39 Formally, regularized ERM problems are optimization problems of the form

$$\min_{x \in \mathbb{R}^d} [P(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) + \psi(x)], \quad (1)$$

40 where $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$ is the loss of model parameterized by vector $x \in \mathbb{R}^d$ on the i -th training data
 41 point, and $\psi: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is a regularizer. Let $[n] := \{1, 2, \dots, n\}$. We shall make the
 42 following assumption throughout the paper without explicitly mentioning it:

43 **Assumption 1.** The functions f_i are L_i -smooth, and the regularizer ψ is proper, closed and convex.
 44 Let $L_{\max} := \max_{i \in [n]} L_i$.

45 In some results we will additionally assume that either the individual functions f_i , or their average
 46 $f := \frac{1}{n} \sum_i f_i$, or the regularizer ψ are μ -strongly convex. Whenever we need such additional
 47 assumptions, we will make this explicitly clear. While all these concepts are standard, we review
 48 them briefly in Section A.

49 **Proximal SGD.** When the number n of training data points is huge, as is increasingly common
 50 in practice, the most efficient algorithms for solving (1) are stochastic first-order methods, such
 51 as stochastic gradient descent (SGD) [Bordes et al., 2009], in one or another of its many variants
 52 proposed in the last decade [Shang et al., 2018, Pham et al., 2020]. These method almost invariably
 53 rely on alternating stochastic gradient steps with the evaluation of the proximal operator

$$\text{prox}_{\gamma\psi}(x) := \operatorname{argmin}_{z \in \mathbb{R}^d} \left\{ \gamma\psi(z) + \frac{1}{2} \|z - x\|^2 \right\}.$$

54 The simplest of these has the form

$$x_{k+1}^{\text{SGD}} = \text{prox}_{\gamma_k\psi}(x_k^{\text{SGD}} - \gamma_k \nabla f_{i_k}(x_k^{\text{SGD}})), \quad (2)$$

55 where i_k is an index from $\{1, 2, \dots, n\}$ chosen uniformly at random, and $\gamma_k > 0$ is a properly
 56 chosen learning rate. Our understanding of (2) is quite mature; see [Gorbunov et al., 2020] for a
 57 general treatment which considers methods of this form in conjunction with more advanced stochastic
 58 gradient estimators in place of ∇f_{i_k} .

59 Applications such as training sparse linear models [Tibshirani, 1996], nonnegative matrix factoriza-
 60 tion [Lee and Seung, 1999], image deblurring [Rudin et al., 1992, Bredies et al., 2010], and training
 61 with group selection [Yuan and Lin, 2006] all rely on the use of hand-crafted regularizers. For most of
 62 them, the proximal operator can be evaluated efficiently, and SGD is near or at the top of the list of
 63 efficient training algorithms.

64 **Random reshuffling.** A particularly successful variant of SGD is based on the idea of random
 65 shuffling (permutation) of the training data followed by n iterations of the form (2), with the index
 66 i_k following the pre-selected permutation [Bottou, 2012]. This process is repeated several times,
 67 each time using a new freshly sampled random permutation of the data, and the resulting method is
 68 known under the name *Random Reshuffling* (RR). When the same permutation is used throughout,
 69 the technique is known under the name *Shuffle-Once* (SO).

70 One of the main advantages of this approach is rooted in its intrinsic ability to avoid cache misses when
 71 reading the data from memory, which enables a significantly faster implementation. Furthermore,
 72 RR is often observed to converge in fewer iterations than SGD in practice. This can intuitively be
 73 ascribed to the fact that while due to its sampling-with-replacement approach SGD can miss to learn
 74 from some data points in any given epoch, RR will learn from each data point in each epoch.

75 Understanding the random reshuffling trick, and why it works, has been a non-trivial open problem
 76 for a long time [Bottou, 2009, Recht and Ré, 2012, Gürbüzbalaban et al., 2019, Haochen and Sra,
 77 2019]. Until recent development which lead to a significant simplification of the convergence
 78 analysis technique and proofs [Mishchenko et al., 2020], prior state of the art relied on long and
 79 elaborate proofs requiring sophisticated arguments and tools, such as analysis via the Wasserstein
 80 distance [Nagaraj et al., 2019], and relied on a significant number of strong assumptions about
 81 the objective [Shamir, 2016, Haochen and Sra, 2019]. In alternative recent development, Ahn et al.
 82 [2020] also develop new tools for analyzing the convergence of random reshuffling, in particular using
 83 decreasing stepsizes and for objectives satisfying the Polyak-Łojasiewicz condition, a generalization
 84 of strong convexity [Polyak, 1963, Łojasiewicz, 1963].

85 The difficulty of analyzing RR has been the main obstacle in the development of even some of the
 86 most seemingly benign extensions of the method. Indeed, while all these are well understood in

Algorithm 1 Proximal Random Reshuffling (ProxRR) and Shuffle-Once (ProxSO)

Require: Stepsizes $\gamma_t > 0$, initial vector $x_0 \in \mathbb{R}^d$, number of epochs T

- 1: Sample a permutation $\pi = (\pi_{0u}, \pi_1, \dots, \pi_{n-1})$ of $[n]$ (Do step 1 only for ProxSO)
- 2: **for** epochs $t = 0, 1, \dots, T - 1$ **do**
- 3: Sample a permutation $\pi = (\pi_0, \pi_1, \dots, \pi_{n-1})$ of $[n]$ (Do step 3 only for ProxRR)
- 4: $x_t^0 = x_t$
- 5: **for** $i = 0, 1, \dots, n - 1$ **do**
- 6: $x_t^{i+1} = x_t^i - \gamma_t \nabla f_{\pi_i}(x_t^i)$
- 7: $x_{t+1} = \text{prox}_{\gamma_t n \psi}(x_t^n)$

combination with its much simpler-to-analyze cousin SGD, *to the best of our knowledge, there exists no theoretical analysis of proximal, parallel, and importance sampling variants of RR with both constant and decreasing stepsizes, and in most cases it is not even clear how should such methods be constructed.* Empowered by and building on the recent advances of [Mishchenko et al. \[2020\]](#), in this paper we address all these challenges.

2 Contributions

In this section we outline the key contributions of our work, and also offer a few intuitive explanations motivating some of the development.

• **New algorithm: ProxRR.** Despite rich literature on Proximal SGD [\[Gorbunov et al., 2020\]](#), it is not obvious how one should extend RR to solve problem (1) when a regularizer ψ is present. Indeed, the standard practice for SGD is to apply the proximal operator after each stochastic step [\[Duchi and Singer, 2009\]](#), i.e., in analogy with (2). On the other hand, RR is motivated by the fact that a data pass better approximates the full gradient step. If we applied the proximal operator after each step of RR, we would no longer approximate the full gradient after an epoch, as we illustrate next.

Example 1. Let $n = 2$, $\psi(x) = \frac{1}{2}\|x\|^2$, $f_1(x) = \langle c_1, x \rangle$, $f_2(x) = \langle c_2, x \rangle$ with some $c_1, c_2 \in \mathbb{R}^d$, $c_1 \neq c_2$. Let $x_0 \in \mathbb{R}^d$, $\gamma > 0$ and define $x_1 = x_0 - \gamma \nabla f_1(x_0)$, $x_2 = x_1 - \gamma \nabla f_2(x_1)$. Then, we have $\text{prox}_{2\gamma\psi}(x_2) = \text{prox}_{2\gamma\psi}(x_0 - 2\gamma \nabla f(x_0))$. However, if $\tilde{x}_1 = \text{prox}_{\gamma\psi}(x_0 - \gamma \nabla f_1(x_0))$ and $\tilde{x}_2 = \text{prox}_{\gamma\psi}(x_1 - \gamma \nabla f_2(\tilde{x}_1))$, then $\tilde{x}_2 \neq \text{prox}_{2\gamma\psi}(x_0 - 2\gamma \nabla f(x_0))$.

Motivated by this observation, we propose ProxRR (Algorithm 1), in which the proximal operator is applied at the end of each epoch of RR, i.e., after each pass through all randomly reshuffled data. A notable property of Algorithm 1 is that *only a single proximal operator evaluation is needed during each data pass*. This is in sharp contrast with the way Proximal SGD works, and offers significant advantages in regimes where the evaluation of the proximal mapping is expensive (e.g., comparable to the evaluation of n gradients $\nabla f_1, \dots, \nabla f_n$).

• **Convergence of ProxRR (for strongly convex functions or regularizer).** We establish several convergence results for ProxRR, of which we highlight two here. Both offer a linear convergence rate with a fixed stepsize to a neighborhood of the solution. In both we reply on Assumption 1. Firstly, in the case when in addition, each f_i is μ -strongly convex, we prove the rate (see Theorem 2)

$$\mathbb{E} \left[\|x_T - x_*\|^2 \right] \leq (1 - \gamma\mu)^{nT} \|x_0 - x_*\|^2 + \frac{2\gamma^2 \sigma_{\text{rad}}^2}{\mu},$$

where $\gamma_t = \gamma \leq 1/L_{\max}$ is the stepsize, and σ_{rad}^2 is a *shuffling radius* constant (for precise definition, see (4)). In Theorem 1 we bound the shuffling radius in terms of $\|\nabla f(x_*)\|^2$, n , L_{\max} and the more common quantity $\sigma_*^2 := \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_*) - \nabla f(x_*)\|^2$. Secondly, if ψ is μ -strongly convex, and we choose the stepsize $\gamma_t = \gamma \leq 1/L_{\max}$, we prove the rate (see Theorem 3)

$$\mathbb{E} \left[\|x_T - x_*\|^2 \right] \leq (1 + 2\gamma\mu n)^{-T} \|x_0 - x_*\|^2 + \frac{\gamma^2 \sigma_{\text{rad}}^2}{\mu}.$$

Both mentioned rates show exponential (linear in logarithmic scale) convergence to a neighborhood whose size is proportional to $\gamma^2 \sigma_{\text{rad}}^2$. Since we can choose γ to be arbitrarily small or periodically

decrease it, this implies that the iterates converge to x_* in the limit. Moreover, we show in Section 4 that when $\gamma = \mathcal{O}(\frac{1}{T})$ the error is $\mathcal{O}(\frac{1}{T^2})$, which is superior to the $\mathcal{O}(\frac{1}{T})$ error of SGD.

• **Results for SO.** All of our results apply to the Shuffle-Once algorithm as well. For simplicity, we center the discussion around RR, whose current theoretical guarantees in the nonconvex case are better than that of SO. Nevertheless, the other results are the same for both methods, and ProxRR is identical to ProxSO in terms of our theory too. A study of the empirical differences between RR and SO can be found in [Mishchenko et al., 2020].

• **Application to Federated Learning.** In Section 6 we describe an application of our results to federated learning [Konečný et al., 2016, McMahan et al., 2017, Kairouz et al., 2019]. In this way we obtain the FedRR method, which is similar to Local SGD, except the local solver is a single pass of RR over the local data. Empirically, FedRR can be vastly superior to Local SGD (see Figure 2). Remarkably, we also show that the rate of FedRR *beats the best known lower bound for Local SGD* due to [Woodworth et al., 2020] (we needed to adapt it from the original online to the finite-sum setting we consider in this paper) for large enough n . See Section F for more details.

• **Nonconvex analysis.** In the nonconvex regime, and under suitable assumptions, we establish (see Theorems 5 and 8) an $\mathcal{O}(\frac{1}{\gamma T})$ rate up to a neighborhood of size $\mathcal{O}(\gamma^2)$. For a certain stepsize it yields an $\mathcal{O}(\frac{1}{\varepsilon^3})$ convergence rate.

Besides the above results, we describe several extensions in the appendix, which we now outline.

• **Extension 1: Decreasing stepsizes.** The convergence of RR is not always exact and depends on the parameters of the objective. Similarly, if the shuffling radius σ_{rad}^2 is positive, and we wish to find an ε -approximate solution, the optimal choice of a fixed stepsize for ProxRR will depend on ε . This deficiency can be fixed by using decreasing stepsizes in both vanilla RR [Ahn et al., 2020] and in SGD [Stich, 2019]. We adopt the same technique to our setting. However, we depart from [Ahn et al., 2020] by only adjusting the stepsize *once per epoch* rather than at every iteration, similarly to the concurrent work of Tran et al. [2020] on RR with momentum. For details, see Section I.

• **Extension 2: Importance resampling for Proximal RR.** While importance sampling is a well established technique for speeding up the convergence of SGD [Zhao and Zhang, 2015, Khaled and Richtárik, 2020], no importance sampling variant of RR has been proposed nor analyzed. This is not surprising since the key property of importance sampling in SGD—unbiasedness—does not hold for RR. Our approach to equip ProxRR with importance sampling is via a reformulation of problem (1) into a similar problem with a larger number of summands. In particular, for each $i \in [n]$ we include n_i copies of the function $\frac{1}{n_i} f_i$, and then take average of all $N = \sum_i n_i$ functions constructed this way. The value of n_i depends on the “importance” of f_i , described below. We then apply ProxRR to this reformulation. If f_i is L_i -smooth for all $i \in [n]$ and we let $\bar{L} := \frac{1}{n} \sum_i L_i$, then we choose $n_i = \lceil L_i / \bar{L} \rceil$. It is easy to show that $N \leq 2n$, and hence our reformulation leads to at most a doubling of the number of functions forming the finite sum. However, the overall complexity of ProxRR applied to this reformulation will depend on \bar{L} instead of $\max_i L_i$ (see Theorem 10), which can lead to a significant improvement. For details of the construction and our complexity results, see Section J.

3 Preliminaries

In our analysis, we build upon the notions of *limit points* and *shuffling variance* introduced by Mishchenko et al. [2020] for vanilla (i.e., non-proximal) RR. Given a stepsize $\gamma > 0$ (held constant during each epoch) and a permutation π of $\{1, 2, \dots, n\}$, the inner loop iterates of RR/SO converge to a neighborhood of intermediate limit points $x_*^1, x_*^2, \dots, x_*^n$ defined by

$$x_*^i := x_* - \gamma \sum_{j=0}^{i-1} \nabla f_{\pi_j}(x_*), \quad i = 1, \dots, n. \quad (3)$$

The intuition behind this definition is fairly simple: if we performed i steps starting at x_* , we would end up close to x_*^i . To quantify the closeness, we define the *shuffling radius*.

Definition 1 (Shuffling radius). Given a stepsize $\gamma > 0$ and a random permutation π of $\{1, 2, \dots, n\}$ used in Algorithm 1, define $x_*^i = x_*^i(\gamma, \pi)$ as in (3). Then, the shuffling radius is defined by

$$\sigma_{\text{rad}}^2(\gamma) := \max_{i=0, \dots, n-1} \left[\frac{1}{\gamma^2} \mathbb{E}_\pi [D_{f_{\pi_i}}(x_*^i, x_*)] \right], \quad (4)$$

where the expectation is taken with respect to the randomness in the permutation π . If there are multiple stepsizes $\gamma_1, \gamma_2, \dots$ used in Algorithm 1, we take the maximum of all of them as the shuffling radius, i.e., $\sigma_{\text{rad}}^2 := \max_{t \geq 1} \sigma_{\text{rad}}^2(\gamma_t)$.

The shuffling radius is related by a multiplicative factor in the stepsize to the shuffling variance introduced by Mishchenko et al. [2020]. When the stepsize is held fixed, the difference between the two notions is minimal. When the stepsize is decreasing, however, the shuffling radius is easier to work with, since it can be upper bounded by problem constants independent of the stepsizes.

Armed with a special lemma for sampling without replacement, we can upper bound the shuffling radius using the smoothness constant L_{\max} , size of the vector $\nabla f(x_*)$, and the variance σ_*^2 of the gradient vectors $\nabla f_1(x_*), \dots, \nabla f_n(x_*)$.

Theorem 1 (Bounding the shuffling radius). For any stepsize $\gamma > 0$ and any random permutation π of $\{1, 2, \dots, n\}$ we have $\sigma_{\text{rad}}^2 \leq \frac{L_{\max}}{2} n(n \|\nabla f(x_*)\|^2 + \frac{1}{2} \sigma_*^2)$, where x_* is a solution of Problem (1) and σ_*^2 is the population variance at the optimum

$$\sigma_*^2 := \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_*) - \nabla f(x_*)\|^2. \quad (5)$$

All proofs are relegated to the supplementary material. In order to better understand the bound given by Theorem 1, note that if there is no proximal operator (i.e., $\psi = 0$) then $\nabla f(x_*) = 0$ and we get that $\sigma_{\text{rad}}^2 \leq \frac{L_{\max} n \sigma_*^2}{4}$. This recovers the existing upper bound on the shuffling variance of Mishchenko et al. [2020] for vanilla RR. On the other hand, if $\nabla f(x_*) \neq 0$ then we get an additive term of size proportional to the squared norm of $\nabla f(x_*)$.

4 Theory for strongly convex losses f_1, \dots, f_n

Our first theorem establishes a convergence rate for Algorithm 1 applied with a constant stepsize to Problem (1) when each objective f_i is strongly convex. This assumption is commonly satisfied in machine learning applications where each f_i represents a regularized loss on some data points, as in ℓ_2 regularized linear regression and ℓ_2 regularized logistic regression.

Theorem 2. Let Assumption 1 be satisfied. Further, assume that each f_i is μ -strongly convex. If Algorithm 1 is run with constant stepsize $\gamma_t = \gamma \leq 1/L_{\max}$, then its iterates satisfy

$$\mathbb{E} \left[\|x_T - x_*\|^2 \right] \leq (1 - \gamma\mu)^{nT} \|x_0 - x_*\|^2 + \frac{2\gamma^2 \sigma_{\text{rad}}^2}{\mu}.$$

We can convert the guarantee of Theorem 2 to a convergence rate by properly tuning the stepsize and using the upper bound of Theorem 1 on the shuffling radius. In particular, if we choose the stepsize as $\gamma = \min \left\{ \frac{1}{L_{\max}}, \frac{\sqrt{\varepsilon\mu}}{\sqrt{2}\sigma_{\text{rad}}} \right\}$, and let $\kappa := L_{\max}/\mu$ and $r_0 := \|x_0 - x_*\|^2$, then we obtain

$\mathbb{E} \left[\|x_T - x_*\|^2 \right] = \mathcal{O}(\varepsilon)$ provided that the total number of iterations $K_{\text{RR}} = nT$ is at least

$$K_{\text{RR}} \geq \left[\left(\kappa + \frac{\sqrt{\kappa n}}{\sqrt{\varepsilon\mu}} (\sqrt{n} \|\nabla f(x_*)\| + \sigma_*) \right) \log \left(\frac{2r_0}{\varepsilon} \right) \right]. \quad (6)$$

Comparison with vanilla RR. If there is no proximal operator, then $\|\nabla f(x_*)\| = 0$ and we recover the earlier result of Mishchenko et al. [2020] on the convergence of RR without proximal, which is optimal in ε up to logarithmic factors. On the other hand, when the proximal operator is nonzero, we get an extra term in the complexity proportional to $\|\nabla f(x_*)\|$: thus, even when all the functions are the same (i.e., $\sigma_* = 0$), we do not recover the linear convergence of Proximal Gradient Descent [Karimi et al., 2016, Beck, 2017]. This can be easily explained by the fact that Algorithm 1 performs n gradient steps per one proximal step. Hence, even if $f_1 = \dots = f_n$, Algorithm 1 does not reduce to Proximal Gradient Descent. We note that other algorithms for composite optimization which may not take a proximal step at every iteration (for example, using stochastic projection steps) also suffer from the same dependence [Patrascu and Irofti, 2021].

Comparison with proximal SGD. In order to compare (6) against the complexity of Proximal SGD (Algorithm 2), we recall that Proximal SGD achieves $\mathbb{E} \left[\|x_K - x_*\|^2 \right] = \mathcal{O}(\varepsilon)$ if either f or ψ is μ -strongly convex and

$$K_{\text{SGD}} \geq \left(\kappa + \frac{\sigma_*^2}{\varepsilon\mu^2} \right) \log \left(\frac{2r_0}{\varepsilon} \right). \quad (7)$$

Algorithm 2 Proximal SGD

Require: Stepsizes $\gamma_k > 0$, initial vector $x_0 \in \mathbb{R}^d$, number of steps K

- 1: **for** steps $k = 0, 1, \dots, K - 1$ **do**
 - 2: Sample i_k uniformly at random from $[n]$
 - 3: $x_{k+1} = \text{prox}_{\gamma_k \psi}(x_k - \gamma_k \nabla f_{i_k}(x_k))$
-

This result is standard [Needell et al., 2016, Gower et al., 2019], with the exception that we do not know any proof in the literature for the case when ψ is strongly convex. For completeness, we prove it in Appendix C, but since our proof is a minor modification of that in [Gower et al., 2019], we do not provide it here.

By comparing K_{SGD} (given by (7)) and K_{RR} (given by (6)), we see that ProxRR has milder dependence on ε than Proximal SGD. In particular, ProxRR converges faster whenever the target accuracy ε is small enough to satisfy $\varepsilon \leq \frac{1}{L_{\max} n \mu} \left(\frac{\sigma_*^4}{n \|\nabla f(x_*)\|^2 + \sigma_*^2} \right)$. Furthermore, ProxRR is much better when we consider *proximal iteration complexity* (# of proximal operator access), in which case the complexity of ProxRR (6) is reduced by a factor of n (because we take one proximal step every n iterations), while the proximal iteration complexity of Proximal SGD remains the same as (7). In this case, ProxRR is better whenever the accuracy ε satisfies

$$\varepsilon \geq \frac{n}{L_{\max} \mu} \left[n \|\nabla f(x_*)\|^2 + \sigma_*^2 \right] \quad \text{or} \quad \varepsilon \leq \frac{n}{L_{\max} \mu} \left[\frac{\sigma_*^4}{n \|\nabla f(x_*)\|^2 + \sigma_*^2} \right].$$

We can see that if the target accuracy is large enough or small enough, and if the cost of proximal operators dominates the computation, ProxRR is much quicker to converge than Proximal SGD.

5 Theory for strongly convex regularizer ψ

In Theorem 2, we assume that each f_i is μ -strongly convex. This is motivated by the common practice of using ℓ_2 regularization in machine learning. However, applying ℓ_2 regularization in every step of Algorithm 1 can be expensive when the data are sparse and the iterates x_t^i are dense, because it requires accessing each coordinate of x_t^i which can be much more expensive than computing sparse gradients $\nabla f_i(x_t^i)$. Alternatively, we may instead choose to put the ℓ_2 regularization inside ψ and only ask that ψ be strongly convex—this way, we can save a lot of time as we need to access each coordinate of the dense iterates x_t^i only once per epoch rather than every iteration. Theorem 3 gives a convergence guarantee in this setting.

Theorem 3. Let Assumption 1 hold and f_1, \dots, f_n be convex. Further, assume that ψ is μ -strongly convex. If Algorithm 1 is run with constant stepsize $\gamma_t = \gamma \leq 1/L_{\max}$, where $L_{\max} = \max_i L_i$, then its iterates satisfy

$$\mathbb{E} \left[\|x_T - x_*\|^2 \right] \leq (1 + 2\gamma\mu n)^{-T} \|x_0 - x_*\|^2 + \frac{\gamma^2 \sigma_{\text{rad}}^2}{\mu}.$$

Using Theorem 3 and choosing the stepsize as

$$\gamma = \min \left\{ \frac{1}{L_{\max}}, \frac{\sqrt{\varepsilon\mu}}{\sigma_{\text{rad}}} \right\}, \quad (8)$$

we get $\mathbb{E} \left[\|x_T - x_*\|^2 \right] = \mathcal{O}(\varepsilon)$ provided that the total number of iterations satisfies

$$K \geq \left(\kappa + \frac{\sigma_{\text{rad}}/\mu}{\sqrt{\varepsilon\mu}} + n \right) \log \left(\frac{2r_0}{\varepsilon} \right). \quad (9)$$

This can be converted to a bound similar to (6) by using Theorem 1, in which case the only difference between the two cases is an extra $n \log \left(\frac{1}{\varepsilon} \right)$ term when only the regularizer ψ is μ -strongly convex. Since for small enough accuracies the $1/\sqrt{\varepsilon}$ term dominates, this difference is minimal.

6 FedRR: application of ProxRR to federated learning

Let us consider now the problem of minimizing the average of $N = \sum_{m=1}^M N_m$ functions that are stored on M devices, which have N_1, \dots, N_M samples correspondingly,

$$\min_{x \in \mathbb{R}^d} F(x) + R(x), \quad F(x) = \frac{1}{N} \sum_{m=1}^M F_m(x), \quad F_m(x) = \sum_{j=1}^{N_m} f_{mj}(x). \quad (10)$$

Algorithm 3 Federated Random Reshuffling (FedRR)

Require: Stepsize $\gamma > 0$, initial vector $x_0 = x_0^0 \in \mathbb{R}^d$, number of epochs T

```

1: for epochs  $t = 0, 1, \dots, T - 1$  do
2:   for  $m = 1, \dots, M$  locally in parallel do
3:      $x_{t,m}^0 = x_t$ 
4:     Sample permutation  $\pi_{0,m}, \pi_{1,m}, \dots, \pi_{N_m-1,m}$  of  $\{1, 2, \dots, N_m\}$ 
5:     for  $i = 0, 1, \dots, N_m - 1$  do
6:        $x_{t,m}^{i+1} = x_{t,m}^i - \gamma \nabla f_{\pi_{i,m}}(x_{t,m}^i)$ 
7:      $x_{t,m}^n = x_{t,m}^{N_m}$ 
8:    $x_{t+1} = \frac{1}{M} \sum_{m=1}^M x_{t,m}^n$ 

```

For example, $f_{mj}(x)$ can be the loss associated with a single sample (X_{mj}, y_{mj}) , where pairs (X_{mj}, y_{mj}) follow a distribution D_m that is specific to device m . An important instance of such formulation is federated learning, where M devices train a shared model by communicating periodically with a server. We normalize the objective in (10) by N as this is the total number of functions after we expand each F_m into a sum. We denote the solution of (10) by x_* .

Extending the space. To rewrite the problem as an instance of (1), we are going to consider a bigger product space, which is sometimes used in distributed optimization [Bianchi et al., 2015]. Let us define $n := \max\{N_1, \dots, N_m\}$ and introduce ψ_C , the *consensus* constraint, defined via

$$\psi_C(x_1, \dots, x_M) := \begin{cases} 0, & x_1 = \dots = x_M \\ +\infty, & \text{otherwise} \end{cases}.$$

By introducing dummy variables x_1, \dots, x_M and adding the constraint $x_1 = \dots = x_M$, we arrive at the intermediate problem

$$\min_{x_1, \dots, x_M \in \mathbb{R}^p} \frac{1}{N} \sum_{m=1}^M F_m(x_m) + (R + \psi_C)(x_1, \dots, x_M),$$

where $R + \psi_C$ is defined, with a slight abuse of notation, as $(R + \psi_C)(x_1, \dots, x_M) = R(x_1)$ if $x_1 = \dots = x_M$, and $(R + \psi_C)(x_1, \dots, x_M) = +\infty$ otherwise.

Since we have replaced R with a more complicated regularizer $R + \psi_C$, we need to understand how to compute the proximal operator of the latter. We show (Lemma 7 in the supplementary) that the proximal operator of $(R + \psi_C)$ is merely the projection onto $\{(x_1, \dots, x_M) \mid x_1 = \dots = x_M\}$ followed by the proximal operator of R with a smaller stepsize.

Reformulation. To have n functions in every F_m , we write F_m as a sum with extra $n - N_m$ zero functions, $f_{mj}(x) \equiv 0$ for any $j > N_m$, so that $F_m(x_m) = \sum_{j=1}^n f_{mj}(x_m) = \sum_{j=1}^{N_m} f_{mj}(x_m) + \sum_{j=N_m+1}^n 0$. We can now stick the vectors together into $\mathbf{x} = (x_1, \dots, x_M) \in \mathbb{R}^{M \cdot d}$ and multiply the objective by $\frac{N}{n}$, which gives the following reformulation:

$$\min_{\mathbf{x} \in \mathbb{R}^{M \cdot d}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) + \psi(\mathbf{x}), \quad (11)$$

where $\psi(\mathbf{x}) := \frac{N}{n}(R + \psi_C)$ and

$$f_i(\mathbf{x}) = f_i(x_1, \dots, x_M) := \sum_{m=1}^M f_{mi}(x_m).$$

In other words, function $f_i(\mathbf{x})$ includes i -th data sample from each device and contains at most one loss from every device, while $F_m(x)$ combines all data losses on device m . Note that the solution of (11) is $\mathbf{x}_* := (x_*^\top, \dots, x_*^\top)^\top$ and the gradient of the extended function $f_i(\mathbf{x})$ is given by $\nabla f_i(\mathbf{x}) = (\nabla f_{1i}(x_1)^\top, \dots, \nabla f_{Mi}(x_M)^\top)^\top$. Therefore, a stochastic gradient step that uses $\nabla f_i(\mathbf{x})$ corresponds to updating all local models with the gradient of i -th data sample, without any communication.

Algorithm 1 for this specific problem can be written in terms of x_1, \dots, x_M , which results in Algorithm 3. Note that since $f_{mi}(x_i)$ depends only on x_i , computing its gradient does not require communication. Only once the local epochs are finished, the vectors are averaged as the result of projecting onto the set $\{(x_1, \dots, x_M) \mid x_1 = \dots = x_M\}$.

Reformulation properties. To analyze FedRR, the only thing that we need to do is understand the properties of the reformulation (11) and then apply Theorem 2 or Theorem 3. The following lemma gives us the smoothness and strong convexity properties of (11).

Lemma 1. Let function f_{mi} be L_i -smooth and μ -strongly convex for every m . Then, f_i from reformulation (11) is L_i -smooth and μ -strongly convex.

The previous lemma shows that the conditioning of the reformulation is $\kappa = \frac{L_{\max}}{\mu}$ just as we would expect. Moreover, it implies that the requirement on the stepsize remains exactly the same: $\gamma \leq 1/L_{\max}$. What remains unknown is the value of σ_{rad}^2 , which plays a key role in the convergence bounds for ProxRR and ProxSO. To find an upper bound on σ_{rad}^2 , let us define

$$\sigma_{m,*}^2 := \frac{1}{N_m} \sum_{j=1}^n \left\| \nabla f_{mj}(x_*) - \frac{1}{N_m} \nabla F_m(x_*) \right\|^2,$$

which is the variance of local gradients on device m . This quantity characterizes the convergence rate of local SGD [Yuan et al., 2020], so we should expect it to appear in our bounds too. The next lemma explains how to use it to upper bound σ_{rad}^2 .

Lemma 2. The shuffling radius σ_{rad}^2 of the reformulation (11) is upper bounded by

$$\sigma_{\text{rad}}^2 \leq L_{\max} \cdot \sum_{m=1}^M \left(\left\| \nabla F_m(x_*) \right\|^2 + \frac{n}{4} \sigma_{m,*}^2 \right).$$

The lemma shows that the upper bound on σ_{rad}^2 depends on the sum of local variances $\sum_{m=1}^M \sigma_{m,*}^2$ as well as on the local gradient norms $\sum_{m=1}^M \left\| \nabla F_m(x_*) \right\|^2$. Both of these sums appear in the existing literature on convergence of Local GD/SGD [Khaled et al., 2019, Woodworth et al., 2020, Yuan et al., 2020]. We are now ready to present formal convergence results. For simplicity, we will consider heterogeneous and homogeneous cases separately and assume that $N_1 = \dots = N_M = n$. To further illustrate generality of our results, we will present the heterogeneous assuming strong convexity R and the homogeneous under strong convexity of functions f_{mi} .

Heterogeneous data. In the case when the data are heterogeneous, we provide the first local RR method. We can apply either Theorem 2 or Theorem 3, but for brevity, we give only the corollary obtained from Theorem 3.

Theorem 4. Assume that functions f_{mi} are convex and L_i -smooth for each m and i . If R is μ -strongly convex and $\gamma \leq 1/L_{\max}$, then we have for the iterates produced by Algorithm 3

$$\mathbb{E} \left[\|x_T - x_*\|^2 \right] \leq (1 + 2\gamma\mu n)^{-T} \|x_0 - x_*\|^2 + \frac{\gamma^2 L_{\max}}{M\mu} \sum_{m=1}^M \left(\left\| \nabla F_m(x_*) \right\|^2 + \frac{N}{4M} \sigma_{m,*}^2 \right).$$

For nonconvex analysis, we consider $R \equiv 0$ and require the following standard assumption.

Assumption 2 (Bounded variance and dissimilarity). There exist constants $\sigma, \zeta > 0$ such that for any $x \in \mathbb{R}^d$ and

$$\frac{1}{n} \sum_{i=1}^n \left\| \nabla f_{mi} - \frac{1}{n} \nabla F_m(x) \right\|^2 \leq \sigma^2 \quad \text{and} \quad \frac{1}{M} \sum_{m=1}^M \left\| \frac{1}{n} \nabla F_m(x) - \nabla F(x) \right\|^2 \leq \zeta^2.$$

Note that above $\frac{1}{n} \nabla F_m(x) = \frac{1}{N_m} \nabla F_m(x)$ is the gradient of a local dataset and $\nabla F(x) = \frac{1}{N} \sum_{l=1}^M \nabla F_l(x)$ is the full gradient on all data.

Theorem 5 (Nonconvex convergence). Let Assumptions 1 and 2 be satisfied, and $R \equiv 0$ (no prox).

Then, the communication complexity to achieve $\mathbb{E} \left[\left\| \nabla F(x_T) \right\|^2 \right] \leq \varepsilon^2$ is

$$T = \mathcal{O} \left(\left(\frac{1}{\varepsilon^2} + \frac{\sigma}{\sqrt{n\varepsilon^3}} + \frac{\zeta}{\varepsilon^3} \right) (F(x_0) - F_*) \right).$$

Notice that by replicating the data locally on each device and thereby increasing the value of n without changing the objective, we can improve the second term in the communication complexity. In particular, if the data are not too dissimilar ($\sigma \gg \zeta$) and ε is small ($\frac{1}{\varepsilon^3} \gg \frac{1}{\varepsilon^2}$), the second term in the complexity dominates, and it helps to have more local steps. However, if the data are less similar, the nodes have to communicate more frequently to get more information about other objectives.

Homogeneous data. For simplicity, in the homogeneous (i.e., i.i.d.) data case we provide guarantees without the proximal operator. Since then we have $F_1(x) = \dots = F_M(x)$, for any m it holds $\nabla F_m(x_*) = 0$, and thus $\sigma_{m,*}^2 = \frac{1}{n} \sum_{j=1}^n \left\| \nabla f_{mj}(x_*) \right\|^2$. The full variance is then given by

$$\sum_{m=1}^M \sigma_{m,*}^2 = \frac{1}{n} \sum_{m=1}^M \sum_{i=1}^n \left\| \nabla f_{mi}(x_*) \right\|^2 = \frac{N}{n} \sigma_*^2 = M \sigma_*^2,$$

where $\sigma_*^2 := \frac{1}{N} \sum_{i=1}^n \sum_{m=1}^M \left\| \nabla f_{mi}(x_*) \right\|^2$ is the variance of the gradients over all data.

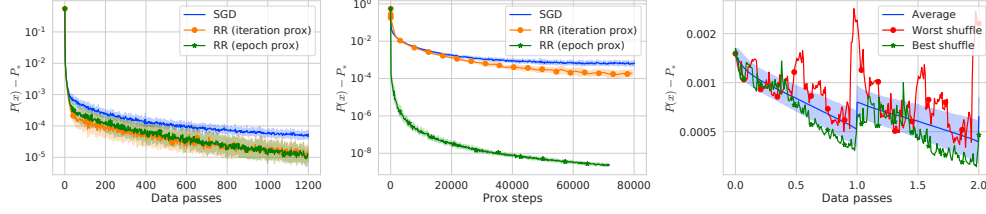


Figure 1: Experimental results for problem (12). The first two plots show with average and confidence intervals estimated on 20 random seeds and clearly demonstrate that one can save a lot of proximal operator computations with our method. The right plot shows the best/worst convergence of ProxSO over 20,000 sampled permutations.

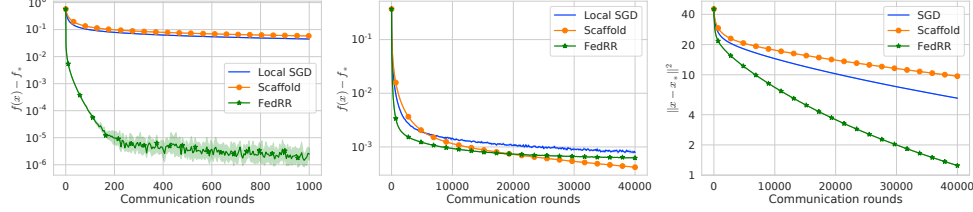


Figure 2: FedRR vs Local-SGD and Scaffold: i.i.d. data (left) and heterogeneous data (middle and right). We set $\lambda_1 = 0$ and estimate the averages and standard deviations by running 10 random seeds for each method.

Theorem 6. Let $R(x) \equiv 0$ (no prox) and the data be i.i.d., that is $\nabla F_m(x_*) = 0$ for any m , where x_* is the solution of (10). Let $\sigma_*^2 := \frac{1}{N} \sum_{i=1}^n \sum_{m=1}^M \|\nabla f_{mi}(x_*)\|^2$. If each f_{mj} is L_{\max} -smooth and μ -strongly convex, then the iterates of Algorithm 3 satisfy

$$\mathbb{E} [\|x_T - x_*\|^2] \leq (1 - \gamma\mu)^{nT} \|x_0 - x_*\|^2 + \frac{\gamma^2 L_{\max} N \sigma_*^2}{M\mu}.$$

The most important part of this result is that the last term in Theorem 6 has a factor of M in the denominator, meaning that the convergence bound improves with the number of devices involved.

7 Experiments¹

ProxRR vs SGD. In Figure 1, we look at the logistic regression loss with the elastic net regularization,

$$\frac{1}{N} \sum_{i=1}^N f_i(x) + \lambda_1 \|x\|_1 + \frac{\lambda_2}{2} \|x\|^2, \quad (12)$$

where each $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as $f_i(x) := -(b_i \log(h(a_i^\top x)) + (1 - b_i) \log(1 - h(a_i^\top x)))$, and where $(a_i, b_i) \in \mathbb{R}^d \times \{0, 1\}$, $i = 1, \dots, N$ are the data samples, $h : t \rightarrow 1/(1 + e^{-t})$ is the sigmoid function, and $\lambda_1, \lambda_2 \geq 0$ are parameters. We set minibatch sizes to 32 for all methods and use theoretical stepsizes, without any tuning. We denote the heuristic version of RR that performs proximal operator step after each iteration as ‘RR (iteration prox)’. From the experiments, we can see that all methods behave more or less the same way. However, the algorithm that we propose needs only a small fraction of proximal operator evaluations, which gives it a huge advantage whenever the operator takes more time to compute than stochastic gradients.

FedRR vs Local SGD and Scaffold. We also compare the performance of FedRR, Local SGD and Scaffold Karimireddy et al. [2020] on homogeneous (i.e., i.i.d.) and heterogeneous data. Since Local SGD and Scaffold require smaller stepsizes to converge, they are significantly slower in the i.i.d. regime, as can be seen in Figure 2. FedRR, however, does not need small initial stepsize and very quickly converges to a noisy neighborhood of the solution. We obtain heterogeneous regime by sorting data with respect to the labels and mixing the sorted dataset with the unsorted one. In this scenario, we also use the same small stepsize for every method to address the data heterogeneity. Clearly, Scaffold is the best in terms of functional values because it does variance reduction with respect to the data. Extending FedRR in the same way might be useful too, but this goes beyond the scope of our paper and we leave it for future work. We also note that in terms of distances from the optimum, FedRR still performs much better than Local SGD and Scaffold.

¹Our code is provided in the supplementary. More experimental details are in the appendix.

References

- Kwangjun Ahn, Chulhee Yun, and Suvrit Sra. SGD with shuffling: optimal rates without component convexity and large epoch requirements. *arXiv preprint arXiv:2006.06946. Neural Information Processing Systems (NeurIPS) 2020*, 2020. (Cited on pages 2, 4, and 31)
- Amir Beck. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017. doi: 10.1137/1.9781611974997. (Cited on page 5)
- Pascal Bianchi, Walid Hachem, and Franck Iutzeler. A coordinate descent primal-dual algorithm and application to distributed asynchronous optimization. *IEEE Transactions on Automatic Control*, 61(10):2947–2957, 2015. (Cited on page 7)
- Antoine Bordes, Léon Bottou, and Patrick Gallinari. SGD-QN: Careful quasi-Newton stochastic gradient descent. 2009. (Cited on page 2)
- Léon Bottou. Curiously fast convergence of some stochastic gradient descent algorithms. Unpublished open problem offered to the attendance of the SLDS 2009 conference, 2009. URL <http://leon.bottou.org/papers/bottou-slds-open-problem-2009>. (Cited on page 2)
- Léon Bottou. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade*, pages 421–436. Springer, 2012. (Cited on page 2)
- Kristian Bredies, Karl Kunisch, and Thomas Pock. Total generalized variation. *SIAM Journal on Imaging Sciences*, 3(3):492–526, 2010. (Cited on page 2)
- Gong Chen and Marc Teboulle. Convergence Analysis of a Proximal-Like Minimization Algorithm Using Bregman Functions. *SIAM Journal on Optimization*, 3(3):538–543, 1993. doi: 10.1137/0803026. (Cited on page 19)
- John Duchi and Yoram Singer. Efficient online and batch learning using forward backward splitting. *Journal of Machine Learning Research*, 10(Dec):2899–2934, 2009. (Cited on page 3)
- Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. A Unified Theory of SGD: Variance Reduction, Sampling, Quantization and Coordinate Descent. volume 108 of *Proceedings of Machine Learning Research*, pages 680–690, Online, 26–28 Aug 2020. PMLR. (Cited on pages 2, 3, 18, and 34)
- Robert M. Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. SGD: General Analysis and Improved Rates. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5200–5209, Long Beach, California, USA, 09–15 Jun 2019. PMLR. (Cited on page 6)
- Robert M. Gower, Peter Richtárik, and Francis Bach. Stochastic quasi-gradient methods: variance reduction via Jacobian sketching. *Mathematical Programming*, pages 1–58, 2020. ISSN 0025-5610. doi: 10.1007/s10107-020-01506-0. (Cited on page 34)
- Mert Gürbüzbalaban, Asuman Özdağlar, and Pablo A. Parrilo. Why random reshuffling beats stochastic gradient descent. *Mathematical Programming*, Oct 2019. ISSN 1436-4646. doi: 10.1007/s10107-019-01440-w. (Cited on page 2)
- Jeff Haochen and Suvrit Sra. Random Shuffling Beats SGD after Finite Epochs. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2624–2633, Long Beach, California, USA, 09–15 Jun 2019. PMLR. (Cited on page 2)
- Peter Kairouz et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019. (Cited on pages 1 and 4)
- Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear Convergence of Gradient and Proximal-Gradient Methods Under the Polyak-Łojasiewicz Condition. In *European Conference on Machine Learning and Knowledge Discovery in Databases - Volume 9851*, ECML PKDD 2016, page 795–811, Berlin, Heidelberg, 2016. Springer-Verlag. (Cited on page 5)

389 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian U. Stich, and
390 Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learning. In
391 *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020. (Cited on pages 9
392 and 30)

393 Ahmed Khaled and Peter Richtárik. Better theory for SGD in the nonconvex world. *arXiv Preprint*
394 *arXiv:2002.03329*, 2020. (Cited on pages 4 and 31)

395 Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. First Analysis of Local GD on
396 Heterogeneous Data. *arXiv preprint arXiv:1909.04715*, 2019. (Cited on page 8)

397 Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for Local SGD on
398 identical and heterogeneous data. In *International Conference on Artificial Intelligence and*
399 *Statistics*, pages 4519–4529. PMLR, 2020. (Cited on page 30)

400 Jakub Konečný, H. Brendan McMahan, Felix Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave
401 Bacon. Federated learning: strategies for improving communication efficiency. In *NIPS Private*
402 *Multi-Party Machine Learning Workshop*, 2016. (Cited on pages 1 and 4)

403 Daniel D. Lee and H. Sebastian Seung. Learning the parts of objects by non-negative matrix
404 factorization. *Nature*, 401(6755):788–791, 1999. (Cited on page 2)

405 Stanislaw Łojasiewicz. A topological property of real analytic subsets. *Coll. du CNRS, Les équations*
406 *aux dérivées partielles*, 117:87–89, 1963. (Cited on page 2)

407 H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas.
408 Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the*
409 *20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017. (Cited on
410 pages 1 and 4)

411 Konstantin Mishchenko, Ahmed Khaled, and Peter Richtárik. Random Reshuffling: Simple Analysis
412 with Vast Improvements. *arXiv preprint arXiv:2006.05988. Neural Information Processing Systems*
413 *(NeurIPS) 2020*, 2020. (Cited on pages 2, 3, 4, 5, 16, 19, 20, 25, and 26)

414 Dheeraj Nagaraj, Prateek Jain, and Praneeth Netrapalli. SGD without Replacement: Sharper Rates
415 for General Smooth Convex Functions. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors,
416 *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings*
417 *of Machine Learning Research*, pages 4703–4711, Long Beach, California, USA, 09–15 Jun 2019.
418 PMLR. (Cited on page 2)

419 Deanna Needell, Nathan Srebro, and Rachel Ward. Stochastic gradient descent, weighted sampling,
420 and the randomized Kaczmarz algorithm. *Mathematical Programming*, 155(1):549–573, Jan 2016.
421 ISSN 1436-4646. doi: 10.1007/s10107-015-0864-7. (Cited on pages 6 and 34)

422 Neal Parikh and Stephen Boyd. Proximal Algorithms. *Foundations and Trends in Optimization*, 1(3):
423 127–239, January 2014. ISSN 2167-3888. doi: 10.1561/24000000003. (Cited on pages 16 and 29)

424 Andrei Patrascu and Paul Irofti. Stochastic proximal splitting algorithm for composite minimization.
425 *Optimization Letters*, pages 1–19, 2021. (Cited on page 5)

426 Nhan H. Pham, Lam M. Nguyen, Dzung T. Phan, and Quoc Tran-Dinh. ProxSARAH: An efficient
427 algorithmic framework for stochastic composite nonconvex optimization. *Journal of Machine*
428 *Learning Research*, 21(110):1–48, 2020. (Cited on page 2)

429 Boris T. Polyak. Gradient methods for minimizing functionals. *Zhurnal Vychislitel’noi Matematiki i*
430 *Matematicheskoi Fiziki*, 3(4):643–653, 1963. (Cited on page 2)

431 Benjamin Recht and Christopher Ré. Toward a noncommutative arithmetic-geometric mean in-
432 equality: Conjectures, case-studies, and consequences. In S. Mannor, N. Srebro, and R. C.
433 Williamson, editors, *Proceedings of the 25th Annual Conference on Learning Theory*, volume 23,
434 page 11.1–11.24, 2012. Edinburgh, Scotland. (Cited on page 2)

435 Leonid I. Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal
436 algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992. (Cited on page 2)

437 Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: from theory to algo-*
438 *rithms*. Cambridge University Press, 2014. (Cited on page 1)

439 Ohad Shamir. Without-replacement sampling for stochastic gradient methods. In *Advances in neural*
440 *information processing systems*, pages 46–54, 2016. (Cited on page 2)

441 Fanhua Shang, Licheng Jiao, Kaiwen Zhou, James Cheng, Yan Ren, and Yufei Jin. ASVRG:
442 Accelerated Proximal SVRG. In Jun Zhu and Ichiro Takeuchi, editors, *Proceedings of Machine*
443 *Learning Research*, volume 95, pages 815–830. PMLR, 14–16 Nov 2018. (Cited on page 2)

444 Sebastian U. Stich. Unified Optimal Analysis of the (Stochastic) Gradient Method. *arXiv preprint*
445 *arXiv:1907.04232*, 2019. (Cited on pages 4 and 31)

446 Ruo-Yu Sun. Optimization for Deep Learning: An Overview. *Journal of the Operations Research*
447 *Society of China*, 8(2):249–294, Jun 2020. ISSN 2194-6698. doi: 10.1007/s40305-020-00309-6.
448 (Cited on page 31)

449 Junqi Tang, Karen Egiazarian, Mohammad Golbabaee, and Mike Davies. The practicality of stochastic
450 optimization in imaging inverse problems. *IEEE Transactions on Computational Imaging*, 6:1471–
451 1485, 2020. (Cited on page 34)

452 Robert Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical*
453 *Society: Series B (Methodological)*, 58(1):267–288, 1996. (Cited on page 2)

454 Trang H. Tran, Lam M. Nguyen, and Quoc Tran-Dinh. Shuffling gradient-based methods with
455 momentum. *arXiv preprint arXiv:2011.11884*, 2020. (Cited on pages 4 and 31)

456 Blake Woodworth, Kumar Kshitij Patel, and Nathan Srebro. Minibatch vs Local SGD for Hetero-
457 geneous Distributed Learning. *arXiv preprint arXiv:2006.04735*. *Neural Information Processing*
458 *Systems (NeurIPS) 2020*, 2020. (Cited on pages 4, 8, and 24)

459 Honglin Yuan, Manzil Zaheer, and Sashank Reddi. Federated composite optimization. *arXiv preprint*
460 *arXiv:2011.08474*, 2020. (Cited on page 8)

461 Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. *Journal*
462 *of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006. (Cited on
463 page 2)

464 Peilin Zhao and Tong Zhang. Stochastic optimization with importance sampling for regularized loss
465 minimization. In *Proceedings of the 32nd International Conference on Machine Learning, PMLR*,
466 volume 37, pages 1–9, 2015. (Cited on page 4)

Checklist

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
- (b) Did you describe the limitations of your work? [Yes]
- (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]

2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- (b) Did you include complete proofs of all theoretical results? [Yes]

3. If you ran experiments...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? [Yes]
- (b) Did you mention the license of the assets? [N/A]
- (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]

5. If you used crowdsourcing or conducted research with human subjects...

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

Appendix

Contents

505	1 Introduction	1
506	2 Contributions	3
507	3 Preliminaries	4
508	4 Theory for strongly convex losses f_1, \dots, f_n	5
509	5 Theory for strongly convex regularizer ψ	6
510	6 FedRR: application of ProxRR to federated learning	6
511	7 Experiments	9
512	I Proofs	16
513	A Basic notions and preliminaries	16
514	A.1 Bregman divergence	16
515	A.2 Properties of the proximal operator	16
516	B Proof of Theorem 1 (Bounding the shuffling radius)	17
517	C Proof of Convergence of Proximal SGD	17
518	D Proofs of Theorem 2 and Theorem 3 (Main convergence results)	19
519	D.1 A key lemma for shuffling-based methods	19
520	D.2 Proof of Theorem 2	20
521	D.3 Proof of Theorem 3	21
522	E Proofs for federated learning	21
523	E.1 Lemma for the extended proximal operator	21
524	E.2 Proof of Lemma 1	22
525	E.3 Proof of Lemma 2	22
526	E.4 Proof of Theorem 4	23
527	E.5 Proof of Theorem 6	23
528	F FedRR beats distributed GD and Local SGD	23
529	F.1 Heterogeneous Data	23
530	F.1.1 Distributed gradient descent	24
531	F.1.2 Local SGD	24

532	G Nonconvex analysis	24
533	G.1 A key lemma	26
534	G.2 Main theorem	27
535	G.3 Proof of Theorem 5	29
536	H Further experimental details	29
537	II Extensions	31
538	I Extension: Decreasing stepsizes	31
539	I.1 A recursion Lemma	31
540	I.2 Proof of Theorem 9	33
541	J Extension: Importance resampling	34

Part I

Proofs

A Basic notions and preliminaries

We say that an extended real-valued function $\phi: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is proper if its domain, $\text{dom } \phi := \{x : \phi(x) < +\infty\}$, is nonempty. We say that it is convex (resp. closed) if its epigraph, $\text{epi } \phi := \{(x, t) \in \mathbb{R}^d \times \mathbb{R} : \phi(x) \leq t\}$, is a convex (resp. closed) set. Equivalently, ϕ is convex if $\text{dom } \phi$ is a convex set and $\phi(\alpha x + (1 - \alpha)y) \leq \alpha\phi(x) + (1 - \alpha)\phi(y)$ for all $x, y \in \text{dom } \phi$ and $\alpha \in (0, 1)$. Finally, ϕ is μ -strongly convex if $\phi(x) - \frac{\mu}{2} \|x\|^2$ is convex, and L -smooth if $\frac{L}{2} \|x\|^2 - \phi(x)$ is convex.

One useful fact that we will need is that for any vectors $a_1, \dots, a_M \in \mathbb{R}^d$ we have

$$\sum_{m=1}^m \|a_i\|^2 = \frac{1}{M} \left\| \sum_{m=1}^M a_m \right\|^2 + \sum_{m=1}^m \left\| a_m - \frac{1}{M} \sum_{l=1}^M a_l \right\|^2. \quad (13)$$

The identity above is sometimes called bias-variance decomposition.

To prove the upper bound in Theorem 1, we rely on a lemma due to [Mishchenko et al. \[2020\]](#) that bounds the variance when sampling without replacement.

Lemma 3 (Lemma 1 in [\[Mishchenko et al., 2020\]](#)). Let $X_1, \dots, X_n \in \mathbb{R}^d$ be fixed vectors, let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ be their mean, and let $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \|X_i - \bar{X}\|^2$ be their variance. Fix any $i \in [n]$ and let $X_{\pi_0}, \dots, X_{\pi_{i-1}}$ be sampled uniformly without replacement from $\{X_1, \dots, X_n\}$ and $\bar{X}_\pi = \frac{1}{i} \sum_{j=0}^{i-1} X_{\pi_j}$ be their average. Then, the sample average and variance are given by

$$\mathbb{E} [\bar{X}_\pi] = \bar{X}, \quad \mathbb{E} [\|\bar{X}_\pi - \bar{X}\|^2] = \frac{n-i}{i(n-1)} \sigma^2. \quad (14)$$

Finally, we define $[n] := \{1, 2, \dots, n\}$.

A.1 Bregman divergence

These notions have a more useful characterization in the case of real valued and continuously differentiable functions $\phi: \mathbb{R}^d \rightarrow \mathbb{R}$. The Bregman divergence of such ϕ is defined by $D_\phi(x, y) := \phi(x) - \phi(y) - \langle \nabla \phi(y), x - y \rangle$. A continuously differentiable function ϕ is called μ -strongly convex if

$$\frac{\mu}{2} \|x - y\|^2 \leq D_\phi(x, y), \quad \forall x, y \in \mathbb{R}^d.$$

It is convex if this holds with $\mu = 0$. Moreover, a continuously differentiable function ϕ is called L -smooth if

$$-\frac{L}{2} \|x - y\|^2 \leq D_\phi(x, y) \leq \frac{L}{2} \|x - y\|^2, \quad \forall x, y \in \mathbb{R}^d. \quad (15)$$

Note that the first inequality is redundant for convex ϕ because convexity implies $0 \leq D_\phi(x, y)$.

A.2 Properties of the proximal operator

Before we proceed to the proofs of convergence, we should state some basic and well-known properties of the regularized objectives. The following lemma explains why the solution of (1) is a fixed point of the proximal-gradient step for any stepsize.

Lemma 4. Let Assumption 1 be satisfied.² Then point x_* is a minimizer of $P(x) = f(x) + \psi(x)$ if and only if for any $\gamma, b > 0$ we have

$$x_* = \text{prox}_{\gamma b \psi}(x_* - \gamma b \nabla f(x_*)).$$

Proof. This follows by writing the first-order optimality conditions for problem (1), see [\[Parikh and Boyd, 2014, p.32\]](#) for a full proof. ■

²We only need the part about ψ .

575 The lemma above only shows that proximal-gradient step does not hurt if we are at the solution. In
 576 addition, we will rely on the following a bit stronger result which postulates that the proximal operator
 577 is a contraction (resp. strong contraction) if the regularizer ψ is convex (resp. strongly convex).

578 **Lemma 5.** Let Assumption 1 be satisfied.³ If ψ is μ -strongly convex with $\mu \geq 0$, then for any $\gamma > 0$
 579 we have

$$\|\text{prox}_{\gamma n \psi}(x) - \text{prox}_{\gamma n \psi}(y)\|^2 \leq \frac{1}{1 + 2\gamma \mu n} \|x - y\|^2, \quad (16)$$

580 for all $x, y \in \mathbb{R}^d$.

581 *Proof.* Let $u := \text{prox}_{\gamma n \psi}(x)$ and $v := \text{prox}_{\gamma n \psi}(y)$. By definition, $u = \text{argmin}_w \{\psi(w) + \frac{1}{2\gamma n} \|w - x\|^2\}$.
 582 By first-order optimality, we have $0 \in \partial\psi(u) + \frac{1}{\gamma n}(u - x)$ or simply $x - u \in \gamma n \partial\psi(u)$. Using
 583 a similar argument for v , we get $x - u - (y - v) \in \gamma n (\partial\psi(u) - \partial\psi(v))$. Thus, by strong convexity
 584 of ψ , we get

$$\langle x - u - (y - v), u - v \rangle \geq \gamma \mu n \|u - v\|^2.$$

585 Hence,

$$\begin{aligned} \|x - y\|^2 &= \|u - v + (x - u - (y - v))\|^2 \\ &= \|u - v\|^2 + 2\langle x - u - (y - v), u - v \rangle + \|x - u - (y - v)\|^2 \\ &\geq \|u - v\|^2 + 2\langle x - u - (y - v), u - v \rangle \\ &\geq (1 + 2\gamma \mu n) \|u - v\|^2. \end{aligned} \quad \blacksquare$$

586 B Proof of Theorem 1 (Bounding the shuffling radius)

587 *Proof.* By the L_i -smoothness of f_i and the definition of x_*^i , we can replace the Bregman divergence
 588 in (4) with the bound

$$\begin{aligned} \mathbb{E} [D_{f_{\pi_i}}(x_*^i, x_*)] &\stackrel{(15)}{\leq} \mathbb{E} \left[\frac{L_{\pi_i}}{2} \|x_*^i - x_*\|^2 \right] \leq \frac{L_{\max}}{2} \mathbb{E} [\|x_*^i - x_*\|^2] \\ &\stackrel{(3)}{=} \frac{\gamma^2 L_{\max}}{2} \mathbb{E} \left[\left\| \sum_{j=0}^{i-1} \nabla f_{\pi_j}(x_*) \right\|^2 \right] \\ &= \frac{\gamma^2 L_{\max} i^2}{2} \mathbb{E} \left[\left\| \frac{1}{i} \sum_{j=0}^{i-1} \nabla f_{\pi_j}(x_*) \right\|^2 \right] \\ &= \frac{\gamma^2 L_{\max} i^2}{2} \mathbb{E} [\|\bar{X}_\pi\|^2], \end{aligned} \quad (17)$$

589 where $\bar{X}_\pi = \frac{1}{j} \sum_{j=0}^{i-1} X_{\pi_j}$ with $X_j := \nabla f_j(x_*)$ for $j = 1, 2, \dots, n$. Since $\bar{X} = \nabla f(x_*)$, by
 590 applying Lemma 3 we get

$$\mathbb{E} [\|\bar{X}_\pi\|^2] = \|\bar{X}\|^2 + \mathbb{E} [\|\bar{X}_\pi - \bar{X}\|^2] \stackrel{(14)+(5)}{=} \|\nabla f(x_*)\|^2 + \frac{n-i}{i(n-1)} \sigma_*^2. \quad (18)$$

591 It remains to combine (17) and (18), use the bounds $i^2 \leq n^2$ and $i(n-i) \leq \frac{n(n-1)}{2}$, which holds for
 592 all $i \in \{0, 1, \dots, n-1\}$, and divide both sides of the resulting inequality by γ^2 . \blacksquare

593 C Proof of Convergence of Proximal SGD

594 **Theorem 7** (Proximal SGD). Let Assumption 1 hold. Further, suppose that either $f := \frac{1}{n} \sum_{i=1}^n f_i$
 595 is μ -strongly convex or that ψ is μ -strongly convex. If Algorithm 2 is run with a constant stepsize
 596 $\gamma_k = \gamma > 0$ satisfying $\gamma \leq \frac{1}{2L_{\max}}$, then the final iterate after K steps satisfies

$$\mathbb{E} [\|x_K - x_*\|^2] \leq (1 - \gamma \mu)^K \|x_0 - x_*\|^2 + \frac{2\gamma \sigma_*^2}{\mu}.$$

³We only need the part about ψ .

597 *Proof.* We will prove the case when ψ is μ -strongly convex. The other result follows as a straightfor-
 598 ward special case of [Gorbunov et al., 2020, Theorem 4.1]. We start by analyzing one step of SGD
 599 with stepsize $\gamma_k = \gamma$ and using Lemma 4

$$\begin{aligned} \|x_{k+1} - x_*\|^2 &= \|\text{prox}_{\gamma\psi}(x_k - \gamma\nabla f_\xi(x_k)) - \text{prox}_{\gamma\psi}(x_* - \gamma\nabla f(x_*))\|^2 \\ &\leq \frac{1}{1 + 2\gamma\mu} \|x_k - \gamma\nabla f_\xi(x_k) - (x_* - \gamma\nabla f(x_*))\|^2. \end{aligned} \quad (19)$$

600 We may write the squared norm term in (19) as

$$\begin{aligned} \|x_k - \gamma\nabla f_\xi(x_k) - (x_* - \gamma\nabla f(x_*))\|^2 &= \|x_k - x_*\|^2 - 2\gamma \langle x_k - x_*, \nabla f_\xi(x_k) - \nabla f(x_*) \rangle \\ &\quad + \gamma^2 \|\nabla f_\xi(x_k) - \nabla f(x_*)\|^2. \end{aligned} \quad (20)$$

601 We denote by $\mathbb{E}_k[\cdot]$ expectation conditional on x_k . Note that the gradient estimate is condition-
 602 ally unbiased, i.e., that $\mathbb{E}_k[\nabla f_\xi(x_k)] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_k) = \nabla f(x_k)$. Hence, taking conditional
 603 expectation in (20) and using unbiasedness we have

$$\begin{aligned} \mathbb{E}_k \left[\|x_k - \gamma\nabla f_\xi(x_k) - (x_* - \gamma\nabla f(x_*))\|^2 \right] &= \|x_k - x_*\|^2 - 2\gamma \langle x_k - x_*, \nabla f(x_k) - \nabla f(x_*) \rangle \\ &\quad + \gamma^2 \mathbb{E}_k \left[\|\nabla f_\xi(x_k) - \nabla f(x_*)\|^2 \right]. \end{aligned} \quad (21)$$

604 By the convexity of f we have

$$\langle x_k - x_*, \nabla f(x_k) - \nabla f(x_*) \rangle \geq D_f(x_k, x_*).$$

605 Furthermore, we may estimate the third term in (21) by first using the fact that $\|x + y\|^2 \leq 2\|x\|^2 +$
 606 $2\|y\|^2$ for any two vectors $x, y \in \mathbb{R}^d$

$$\begin{aligned} \mathbb{E}_k \left[\|\nabla f_\xi(x_k) - \nabla f(x_*)\|^2 \right] &\leq 2\mathbb{E}_k \left[\|\nabla f_\xi(x_k) - \nabla f_\xi(x_*)\|^2 \right] + 2\mathbb{E}_k \left[\|\nabla f_\xi(x_*) - \nabla f(x_*)\|^2 \right] \\ &= 2\mathbb{E}_k \left[\|\nabla f_\xi(x_k) - \nabla f_\xi(x_*)\|^2 \right] + 2\sigma_*^2. \end{aligned}$$

607 We now use that by the L_{\max} -smoothness of f_i we have that

$$\|\nabla f_i(x_k) - \nabla f_i(x_*)\|^2 \leq 2L_{\max} \cdot D_{f_i}(x_k, x_*).$$

608 Hence

$$\begin{aligned} \mathbb{E}_k \left[\|\nabla f_\xi(x_k) - \nabla f_\xi(x_*)\|^2 \right] &= \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f_i(x_*)\|^2 \\ &\leq \frac{2L_{\max}}{n} \sum_{i=1}^n [f_i(x_k) - f_i(x_*) - \langle \nabla f_i(x_*), x_k - x_* \rangle] \\ &= 2L_{\max} [f(x_k) - f(x_*) - \langle \nabla f(x_*), x_k - x_* \rangle] \\ &= 2L_{\max} D_f(x_k, x_*). \end{aligned} \quad (22)$$

609 Combining equations (21)–(22) we obtain

$$\mathbb{E}_k \left[\|x_k - \gamma\nabla f_\xi(x_k) - (x_* - \gamma\nabla f(x_*))\|^2 \right] \leq \|x_k - x_*\|^2 - 2\gamma(1 - 2\gamma L_{\max}) D_f(x_k, x_*) + 2\gamma^2 \sigma_*^2.$$

610 Since $\gamma \leq \frac{1}{2L_{\max}}$ by assumption we have that $1 - 2\gamma L_{\max} \geq 0$. Since $D_f(x_k, x_*) \geq 0$ by the
 611 convexity of f we arrive at

$$\mathbb{E}_k \left[\|x_k - \gamma\nabla f_\xi(x_k) - (x_* - \gamma\nabla f(x_*))\|^2 \right] \leq \|x_k - x_*\|^2 + 2\gamma^2 \sigma_*^2.$$

612 Taking unconditional expectation and combining (41) with the last equation we have

$$\begin{aligned} \mathbb{E} \left[\|x_{k+1} - x_*\|^2 \right] &\leq \frac{1}{1 + 2\gamma\mu} \left(\mathbb{E} \left[\|x_k - x_*\|^2 \right] + 2\gamma^2 \sigma_*^2 \right) \\ &= \frac{1}{1 + 2\gamma\mu} \mathbb{E} \left[\|x_k - x_*\|^2 \right] + \frac{2\gamma^2 \sigma_*^2}{1 + 2\gamma\mu} \\ &\leq \frac{1}{1 + 2\gamma\mu} \mathbb{E} \left[\|x_k - x_*\|^2 \right] + 2\gamma^2 \sigma_*^2. \end{aligned}$$

613 To simplify this further, we use that for any $x \leq \frac{1}{2}$ we have that $\frac{1}{1+2x} \leq 1 - x$ and that $\gamma\mu \leq$
614 $\frac{\mu}{2L_{\max}} \leq \frac{1}{2}$, hence

$$\mathbb{E} [\|x_{k+1} - x_*\|^2] \leq (1 - \gamma\mu) \mathbb{E} [\|x_k - x_*\|^2] + 2\gamma^2\sigma_*^2.$$

615 Recursing the above inequality for K steps yields

$$\begin{aligned} \mathbb{E} [\|x_K - x_*\|^2] &\leq (1 - \gamma\mu)^K \|x_0 - x_*\|^2 + 2\gamma^2\sigma_*^2 \left(\sum_{k=0}^{K-1} (1 - \gamma\mu)^k \right) \\ &\leq (1 - \gamma\mu)^K \|x_0 - x_*\|^2 + 2\gamma^2\sigma_*^2 \left(\sum_{k=0}^{\infty} (1 - \gamma\mu)^k \right) \\ &= (1 - \gamma\mu)^K \|x_0 - x_*\|^2 + \frac{2\gamma\sigma_*^2}{\mu}. \end{aligned} \quad \blacksquare$$

616 Furthermore, by choosing the stepsize γ as $\gamma = \min \left\{ \frac{1}{2L_{\max}}, \frac{\varepsilon\mu}{4\sigma_*} \right\}$, we get that $\mathbb{E} [\|x_K - x_*\|^2] =$
617 $\mathcal{O}(\varepsilon)$ provided that the number of iterations is at least

$$K_{\text{SGD}} \geq \left(\kappa + \frac{\sigma_*^2}{\varepsilon\mu^2} \right) \log \left(\frac{2r_0}{\varepsilon} \right),$$

618 which we previously stated in (7).

619 D Proofs of Theorem 2 and Theorem 3 (Main convergence results)

620 D.1 A key lemma for shuffling-based methods

621 The intermediate limit points x_*^i are extremely important for showing tight convergence guarantees
622 for Random Reshuffling even without proximal operator. The following lemma illustrates that by
623 giving a simple recursion, whose derivation follows [Mishchenko et al., 2020, Proof of Theorem 1].
624 The proof is included for completeness.

625 **Lemma 6** (Theorem 1 in [Mishchenko et al., 2020]). Suppose that each f_i is L_i -smooth and λ -
626 strongly convex (where $\lambda = 0$ means each f_i is just convex). Then the inner iterates generated by
627 Algorithm 1 satisfy

$$\mathbb{E} [\|x_t^{i+1} - x_*^{i+1}\|^2] \leq (1 - \gamma\lambda) \mathbb{E} [\|x_t^i - x_*^i\|^2] - 2\gamma(1 - \gamma L_{\max}) \mathbb{E} [D_{f_{\pi_i}}(x_t^i, x_*)] + 2\gamma^3\sigma_{\text{rad}}^2, \quad (23)$$

628 where x_*^i is as in (3), $i = 0, 1, \dots, n-1$, and x_* is any minimizer of P .

629 *Proof.* By definition of x_t^{i+1} and x_*^{i+1} , we have

$$\begin{aligned} \mathbb{E} [\|x_t^{i+1} - x_*^{i+1}\|^2] &= \mathbb{E} [\|x_t^i - x_*^i\|^2] - 2\gamma \mathbb{E} [\langle \nabla f_{\pi_i}(x_t^i) - \nabla f_{\pi_i}(x_*) , x_t^i - x_*^i \rangle] \\ &\quad + \gamma^2 \mathbb{E} [\|\nabla f_{\pi_i}(x_t^i) - \nabla f_{\pi_i}(x_*)\|^2]. \end{aligned} \quad (24)$$

630 Note that the third term in (24) can be bounded as

$$\|\nabla f_{\pi_i}(x_t^i) - \nabla f_{\pi_i}(x_*)\|^2 \leq 2L_{\max} \cdot D_{f_{\pi_i}}(x_t^i, x_*). \quad (25)$$

631 We may rewrite the second term in (24) using the three-point identity [Chen and Teboulle, 1993,
632 Lemma 3.1] as

$$\langle \nabla f_{\pi_i}(x_t^i) - \nabla f_{\pi_i}(x_*), x_t^i - x_*^i \rangle = D_{f_{\pi_i}}(x_*^i, x_t^i) + D_{f_{\pi_i}}(x_t^i, x_*) - D_{f_{\pi_i}}(x_*^i, x_*). \quad (26)$$

633 Combining (24), (25), and (26) we obtain

$$\begin{aligned} \mathbb{E} [\|x_t^{i+1} - x_*^{i+1}\|^2] &\leq \mathbb{E} [\|x_t^i - x_*^i\|^2] - 2\gamma \cdot \mathbb{E} [D_{f_{\pi_i}}(x_*^i, x_t^i)] + 2\gamma \cdot \mathbb{E} [D_{f_{\pi_i}}(x_t^i, x_*)] \\ &\quad - 2\gamma(1 - \gamma L_{\max}) \mathbb{E} [D_{f_{\pi_i}}(x_t^i, x_*)]. \end{aligned} \quad (27)$$

634 Using λ -strong convexity of f_{π_i} , we derive

$$\frac{\lambda}{2} \|x_t^i - x_*^i\|^2 \leq D_{f_{\pi_i}}(x_*^i, x_t^i). \quad (28)$$

635 Furthermore, by the definition of shuffling radius (Definition 1), we have

$$\mathbb{E} [D_{f_{\pi_i}}(x_*^i, x_*)] \leq \max_{i=0, \dots, n-1} \mathbb{E} [D_{f_{\pi_i}}(x_*^i, x_*)] = \gamma^2 \sigma_{\text{rad}}^2. \quad (29)$$

636 Using (28) and (29) in (27) yields (23). ■

637 D.2 Proof of Theorem 2

638 *Proof.* Starting with Lemma 6 with $\lambda = \mu$, we have

$$\mathbb{E} [\|x_t^{i+1} - x_*^{i+1}\|^2] \leq (1 - \gamma\mu) \mathbb{E} [\|x_t^i - x_*^i\|^2] - 2\gamma(1 - \gamma L_{\max}) \mathbb{E} [D_{f_{\pi_i}}(x_t^i, x_*)] + 2\gamma^3 \sigma_{\text{rad}}^2.$$

639 Since $D_{f_{\pi_i}}(x_t^i, x_*)$ is a Bregman divergence of a convex function, it is nonnegative. Combining this
640 with the fact that the stepsize satisfies $\gamma \leq 1/L_{\max}$, we have

$$\mathbb{E} [\|x_t^{i+1} - x_*^{i+1}\|^2] \leq (1 - \gamma\mu) \mathbb{E} [\|x_t^i - x_*^i\|^2] + 2\gamma^3 \sigma_{\text{rad}}^2.$$

641 Unrolling this recursion for n steps, we get

$$\begin{aligned} \mathbb{E} [\|x_t^n - x_*^n\|^2] &\leq (1 - \gamma\mu)^n \mathbb{E} [\|x_t^0 - x_*^0\|^2] + 2\gamma^3 \sigma_{\text{rad}}^2 \left(\sum_{j=0}^{n-1} (1 - \gamma\mu)^j \right) \\ &= (1 - \gamma\mu)^n \mathbb{E} [\|x_t - x_*\|^2] + 2\gamma^3 \sigma_{\text{rad}}^2 \left(\sum_{j=0}^{n-1} (1 - \gamma\mu)^j \right), \end{aligned} \quad (30)$$

642 where we used the fact that $x_t^0 - x_*^0 = x_t - x_*$. Since x_* minimizes P , we have by Lemma 4 that

$$x_* = \text{prox}_{\gamma n \psi} \left(x_* - \gamma \sum_{i=0}^{n-1} \nabla f_{\pi_i}(x_*) \right) = \text{prox}_{\gamma n \psi}(x_*^n).$$

643 Moreover, by Lemma 5 we obtain that

$$\|x_{t+1} - x_*\|^2 = \|\text{prox}_{\gamma n \psi}(x_t^n) - \text{prox}_{\gamma n \psi}(x_*^n)\|^2 \leq \|x_t^n - x_*^n\|^2.$$

644 Using this in (30) yields

$$\mathbb{E} [\|x_{t+1} - x_*\|^2] \leq (1 - \gamma\mu)^n \mathbb{E} [\|x_t - x_*\|^2] + 2\gamma^3 \sigma_{\text{rad}}^2 \left(\sum_{j=0}^{n-1} (1 - \gamma\mu)^j \right).$$

645 We now unroll this recursion again for T steps

$$\mathbb{E} [\|x_T - x_*\|^2] \leq (1 - \gamma\mu)^{nT} \mathbb{E} [\|x_0 - x_*\|^2] + 2\gamma^3 \sigma_{\text{rad}}^2 \left(\sum_{j=0}^{n-1} (1 - \gamma\mu)^j \right) \left(\sum_{i=0}^{T-1} (1 - \gamma\mu)^{ni} \right). \quad (31)$$

646 Following Mishchenko et al. [2020], we rewrite and bound the product in the last term as

$$\begin{aligned} \left(\sum_{j=0}^{n-1} (1 - \gamma\mu)^j \right) \left(\sum_{i=0}^{T-1} (1 - \gamma\mu)^{ni} \right) &= \sum_{j=0}^{n-1} \sum_{i=0}^{T-1} (1 - \gamma\mu)^{ni+j} \\ &= \sum_{k=0}^{nT-1} (1 - \gamma\mu)^k \\ &\leq \sum_{k=0}^{\infty} (1 - \gamma\mu)^k = \frac{1}{\gamma\mu}. \end{aligned}$$

647 It remains to plug this bound into (31). ■

648 D.3 Proof of Theorem 3

649 *Proof.* Starting with Lemma 6 with $\lambda = 0$, we have

$$\mathbb{E} \left[\|x_t^{i+1} - x_*^{i+1}\|^2 \right] \leq \mathbb{E} \left[\|x_t^i - x_*^i\|^2 \right] - 2\gamma(1 - \gamma L_{\max}) \mathbb{E} [D_{f_{\pi_i}}(x_t^i, x_*)] + 2\gamma^3 \sigma_{\text{rad}}^2.$$

650 Since $\gamma \leq 1/L_{\max}$ and $D_{f_{\pi_i}}(x_t^i, x_*)$ is nonnegative we may simplify this to

$$\mathbb{E} \left[\|x_t^{i+1} - x_*^{i+1}\|^2 \right] \leq \mathbb{E} \left[\|x_t^i - x_*^i\|^2 \right] + 2\gamma^3 \sigma_{\text{rad}}^2.$$

651 Unrolling this recursion over an epoch we have

$$\mathbb{E} \left[\|x_t^n - x_*^n\|^2 \right] \leq \mathbb{E} \left[\|x_t^0 - x_*^0\|^2 \right] + 2\gamma^3 \sigma_{\text{rad}}^2 n = \mathbb{E} \left[\|x_t - x_*\|^2 \right] + 2\gamma^3 \sigma_{\text{rad}}^2 n. \quad (32)$$

652 Since x_* minimizes P , we have by Lemma 4 that

$$x_* = \text{prox}_{\gamma n \psi} \left(x_* - \gamma \sum_{i=0}^{n-1} \nabla f_{\pi_i}(x_*) \right) = \text{prox}_{\gamma n \psi}(x_*^n).$$

653 Hence, $x_{t+1} - x_* = \text{prox}_{\gamma n \psi}(x_t^n) - \text{prox}_{\gamma n \psi}(x_*^n)$. We may now use Lemma 5 to get

$$(1 + 2\gamma\mu n) \mathbb{E} \left[\|x_{t+1} - x_*\|^2 \right] \leq \mathbb{E} \left[\|x_t^n - x_*^n\|^2 \right].$$

654 Combining this with (32), we obtain

$$\mathbb{E} \left[\|x_{t+1} - x_*\|^2 \right] \leq \frac{1}{1 + 2\gamma\mu n} \mathbb{E} \left[\|x_t - x_*\|^2 \right] + \frac{2\gamma^3 \sigma_{\text{rad}}^2 n}{1 + 2\gamma\mu n}.$$

655 We may unroll this recursion again, this time for T steps, and then use that $\sum_{j=1}^{T-1} (1 + 2\gamma\mu n)^{-j} \leq$

656 $\sum_{j=1}^{\infty} (1 + 2\gamma\mu n)^{-j} = 1/(2\gamma\mu n)$:

$$\begin{aligned} \mathbb{E} \left[\|x_T - x_*\|^2 \right] &\leq (1 + 2\gamma\mu n)^{-T} \mathbb{E} \left[\|x_0 - x_*\|^2 \right] + \frac{2\gamma^3 \sigma_{\text{rad}}^2 n}{1 + 2\gamma\mu n} \left(\sum_{j=0}^{T-1} (1 + 2\gamma\mu n)^{-j} \right) \\ &= (1 + 2\gamma\mu n)^{-T} \mathbb{E} \left[\|x_0 - x_*\|^2 \right] + 2\gamma^3 \sigma_{\text{rad}}^2 n \left(\sum_{j=1}^T (1 + 2\gamma\mu n)^{-j} \right) \\ &\leq (1 + 2\gamma\mu n)^{-T} \mathbb{E} \left[\|x_0 - x_*\|^2 \right] + 2\gamma^3 \sigma_{\text{rad}}^2 n \frac{1}{2\gamma\mu n} \\ &= (1 + 2\gamma\mu n)^{-T} \mathbb{E} \left[\|x_0 - x_*\|^2 \right] + \frac{\gamma^2 \sigma_{\text{rad}}^2}{\mu}. \quad \blacksquare \end{aligned}$$

657 E Proofs for federated learning

658 E.1 Lemma for the extended proximal operator

659 **Lemma 7.** Let ψ_C be the consensus constraint and R be a closed convex proximable function.
660 Suppose that x_1, x_2, \dots, x_M are all in \mathbb{R}^d . Then,

$$\text{prox}_{\gamma(R+\psi_C)}(x_1, \dots, x_M) = \text{prox}_{\frac{\gamma}{M}R}(\bar{x}),$$

661 where $\bar{x} = \frac{1}{M} \sum_{m=1}^M x_m$.

662 *Proof.* We have,

$$\text{prox}_{\gamma(R+\psi_C)}(x_1, \dots, x_M) = \begin{pmatrix} \text{prox}_{\frac{\gamma}{M}R}(\bar{x}) \\ \vdots \\ \text{prox}_{\frac{\gamma}{M}R}(\bar{x}) \end{pmatrix} \quad \text{with} \quad \bar{x} = \frac{1}{M} \sum_{m=1}^M x_m.$$

663 This is a simple consequence of the definition of the proximal operator. Indeed, the result of
 664 $\text{prox}_{\gamma(R+\psi_C)}$ must have blocks equal to some vector z such that

$$\begin{aligned} z &= \underset{x}{\operatorname{argmin}} \left\{ \gamma R(x) + \frac{1}{2} \sum_{m=1}^M \|x - x_m\|^2 \right\} \\ &= \underset{x}{\operatorname{argmin}} \left\{ \gamma R(x) + \frac{1}{2} \sum_{m=1}^M (\|x - \bar{x}\|^2 + 2\langle x - \bar{x}, \bar{x} - x_m \rangle) + \|\bar{x} - x_m\|^2 \right\} \\ &= \underset{x}{\operatorname{argmin}} \left\{ \gamma R(x) + \frac{1}{2} M \|x - \bar{x}\|^2 \right\} = \text{prox}_{\frac{\gamma}{M} R}(\bar{x}). \end{aligned}$$

665

666 E.2 Proof of Lemma 1

667 *Proof.* Given some vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{d \cdot M}$, let us use their block representation $\mathbf{x} = (x_1^\top, \dots, x_M^\top)^\top$,
 668 $\mathbf{y} = (y_1^\top, \dots, y_M^\top)^\top$. Since we use the Euclidean norm, we have

$$\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|^2 = \sum_{m=1}^M \|\nabla f_{mi}(x_m) - \nabla f_{mi}(y_m)\|^2 \leq \sum_{m=1}^M L_i^2 \|x_m - y_m\|^2 = L_i^2 \|\mathbf{x} - \mathbf{y}\|^2.$$

669 We can obtain a lower bound by doing the same derivation and applying strong convexity instead of
 670 smoothness:

$$\sum_{m=1}^M \|\nabla f_{mi}(x_m) - \nabla f_{mi}(y_m)\|^2 \geq \mu^2 \sum_{m=1}^M \|x_m - y_m\|^2 = \mu^2 \|\mathbf{x} - \mathbf{y}\|^2.$$

671 Thus, we have $\mu \|\mathbf{x} - \mathbf{y}\| \leq \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\| \leq L_i \|\mathbf{x} - \mathbf{y}\|$, which is exactly μ -strong convexity
 672 and L_i -smoothness of f_i . ■

673 E.3 Proof of Lemma 2

674 *Proof.* By Theorem 1 we have

$$\sigma_{\text{rad}}^2 \leq \frac{L_{\max}}{2} \left(n^2 \|\nabla f(\mathbf{x}_*)\|^2 + \frac{n}{2} \sigma_*^2 \right).$$

675 Due to the separable structure of f , we have for the variance term

$$n \sigma_*^2 := \sum_{i=1}^n \|\nabla f_i(\mathbf{x}_*) - \nabla f(\mathbf{x}_*)\|^2 = \sum_{i=1}^n \sum_{m=1}^M \left\| \nabla f_{mi}(x_*) - \frac{1}{n} \nabla F_m(x_*) \right\|^2.$$

676 The expression inside the summation is not exactly the variance due to the different normalization: $\frac{1}{n}$
 677 instead of $\frac{1}{N_m}$. Nevertheless, we can expand the norm and try to get the actual variance:

$$\begin{aligned} \sum_{i=1}^n \left\| \nabla f_{mi}(x_*) - \frac{1}{n} \nabla F_m(x_*) \right\|^2 &= \sum_{i=1}^{N_m} \left(\left\| \nabla f_{mi}(x_*) - \frac{1}{N_m} \nabla F_m(x_*) \right\|^2 + \left(\frac{1}{N_m} - \frac{1}{n} \right)^2 \|\nabla F_m(x_*)\|^2 \right) \\ &\quad + 2 \sum_{i=1}^{N_m} \left\langle \nabla f_{mi}(x_*) - \frac{1}{N_m} \nabla F_m(x_*), \left(\frac{1}{N_m} - \frac{1}{n} \right) \nabla F_m(x_*) \right\rangle \\ &= N_m \sigma_{m,*}^2 + N_m \left(\frac{1}{N_m} - \frac{1}{n} \right)^2 \|\nabla F_m(x_*)\|^2 \\ &\leq n \sigma_{m,*}^2 + \|\nabla F_m(x_*)\|^2. \end{aligned}$$

678 Moreover, the gradient term has the same block structure, so

$$n^2 \|\nabla f(\mathbf{x}_*)\|^2 = n^2 \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_*) \right\|^2 = \sum_{m=1}^M \left\| \sum_{i=1}^n \nabla f_{mi}(x_*) \right\|^2 = \sum_{m=1}^M \|\nabla F_m(x_*)\|^2.$$

679 Plugging the last two bounds back inside the upper bound on σ_{rad}^2 , we deduce the lemma's statement.
 680 ■

681 E.4 Proof of Theorem 4

682 *Proof.* Since we assume that $N_1 = \dots = N_M = n$, we have $\frac{N}{M} = n$ and the strong convexity
 683 constant of $\psi = \frac{N}{n}(R + \psi_C)$ is equal to $\frac{N}{n} \cdot \frac{\mu}{M} = \mu$. By applying Theorem 3 we obtain

$$\mathbb{E} \left[\|\mathbf{x}_T - \mathbf{x}_*\|^2 \right] \leq (1 + 2\gamma\mu n)^{-T} \|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \frac{\gamma^2 \sigma_{\text{rad}}^2}{\mu}.$$

684 Since $\mathbf{x}_T = \text{prox}_{\gamma N(R + \psi_C)}(\mathbf{x}_{T-1}^n)$, we have $\mathbf{x}_T \in C$, i.e., all of its blocks are equal to each other
 685 and we have $\mathbf{x}_T = (x_T^\top, \dots, x_T^\top)^\top$. Since we use the Euclidean norm, it also implies

$$\mathbb{E} \left[\|\mathbf{x}_T - \mathbf{x}_*\|^2 \right] = M \|\mathbf{x}_T - \mathbf{x}_*\|^2.$$

686 The same is true for \mathbf{x}_0 , so we need to divide both sides of the upper bound on $\|\mathbf{x}_T - \mathbf{x}_*\|^2$ by M .
 687 Doing so together with applying Lemma 2 yields

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{x}_T - \mathbf{x}_*\|^2 \right] &\leq (1 + 2\gamma\mu n)^{-T} \|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \frac{\gamma^2 \sigma_{\text{rad}}^2}{M\mu} \\ &\leq (1 + 2\gamma\mu n)^{-T} \|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \frac{\gamma^2 L_{\max}}{M\mu} \sum_{m=1}^M \left(\|\nabla F_m(x_*)\|^2 + \frac{n}{4} \sigma_{m,*}^2 \right) \\ &= (1 + 2\gamma\mu n)^{-T} \|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \frac{\gamma^2 L_{\max}}{M\mu} \sum_{m=1}^M \left(\|\nabla F_m(x_*)\|^2 + \frac{N}{4M} \sigma_{m,*}^2 \right). \end{aligned}$$

688 ■

689 E.5 Proof of Theorem 6

690 *Proof.* According to Lemma 1, each f_i is μ -strongly convex and L_{\max} -smooth, so we obtain the
 691 result by trivially applying Theorem 2 and upper bounding σ_{rad}^2 the same way as in the proof
 692 of Theorem 4. ■

693 F FedRR beats distributed GD and Local SGD

694 F.1 Heterogeneous Data

695 In this section we compare between FedRR and several known baseline algorithms for Federated
 696 Learning. In particular, we consider the following algorithms:

- 697 1. Distributed gradient descent (DGD)
- 698 2. Local SGD (with M nodes and n local steps per node)

699 To be clear, the problem we are considering is

$$\min_{x \in \mathbb{R}^d} f(x) := \left[\frac{1}{M} \sum_{m=1}^M F_m(x) + R(x) \right],$$

700 where each objective f_m can be written as

$$F_m(x) = \frac{1}{n} \sum_{i=1}^n f_{m,i}(x).$$

701 We further assume that each objective is L -smooth and convex, and that R is μ -strongly convex. This
 702 implies that f is L -smooth and μ -strongly convex. Note that this is a special case of (10) where we
 703 keep $N_1 = N_2 = \dots = n$ for simplicity.

704 **Corollary 1.** Let $c^2 = \zeta_*^2 + \frac{n}{4} \sigma_*^2$, where $\zeta_*^2 := \frac{1}{M} \sum_{m=1}^M \|\nabla F_m(x)\|^2$ and $\sigma_*^2 =$
 705 $\frac{1}{M} \sum_{m=1}^M \|\nabla F(x_*) - \nabla F_m(x_*)\|^2$. Then the communication complexity required by FedRR to
 706 reach an ϵ -accurate solution is

$$T = \Omega \left(\left(\frac{\kappa}{n} + \frac{c}{\mu n} \sqrt{\frac{\kappa}{\epsilon}} \right) \log \left(\frac{r_0}{\epsilon} \right) \right), \quad (33)$$

707 where $r_0 = \|\mathbf{x}_0 - \mathbf{x}_*\|^2$.

708 *Proof.* This is a straightforward consequence of Theorem 4. ■

709 F.1.1 Distributed gradient descent

710 When we compute n gradients on each node per communication round, we are essentially running
 711 distributed gradient descent (DGD). In order to reach an ϵ -accurate solution, DGD requires the
 712 following number of iterations

$$T = \Omega \left(\kappa \log \left(\frac{r_0}{\epsilon} \right) \right).$$

713 Comparing against the result of Corollary 1, we see that FedRR is better whenever the accuracy ϵ
 714 satisfies

$$\frac{1}{\mu L} \left(\frac{\zeta_*^2}{n^2} + \frac{\sigma_*^2}{n} \right) = \frac{c^2}{\mu n^2 L} < \epsilon.$$

715 Note that this guarantee grows more rigorous with increasing levels of heterogeneity—this has been
 716 observed for other local methods as well, such as Local SGD [Woodworth et al., 2020].

717 F.1.2 Local SGD

718 The best current lower bound for Local SGD is given by [Woodworth et al., 2020] in the *stochastic*
 719 case. By stochastic case, we mean that the problem considered is

$$\min_{x \in \mathbb{R}^d} \mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi(x)].$$

720 This is a more general problem than the finite-sum minimization problem (1) and is usually strictly
 721 harder to solve (i.e., requires more iterations to achieve an ϵ -accurate solution). We are not aware of
 722 any analysis of Local SGD specifically for the finite-sum problem, and thus we specialize the result
 723 of Woodworth et al. [2020] anyway. For Local SGD on μ -strongly convex and L smooth functions,
 724 and with n steps of local steps per node, the lower bound they give after T communication rounds is

$$\min \left(\Delta \exp \left(-\frac{\mu T}{L} \right), \frac{L \zeta_*^2}{\mu^2 T^2} \right) + \frac{\sigma^2}{\mu M n T} + \min \left(\Delta, \frac{L \sigma^2}{\mu^2 n^2 T^2} \right), \quad (34)$$

725 where σ^2 is a uniform bound on the variance (i.e., $\mathbb{E}[\|\nabla f_\xi(x) - \nabla f(x)\|^2] \leq \sigma^2$ for all $x \in \mathbb{R}^d$), ζ_* is
 726 defined as in Corollary 1, and Δ is an upper bound on $f(x_0) - f_*$. We note that this lower bound is
 727 *not* actually met by any of the existing analysis for Local SGD. Even ignoring the dependence on
 728 σ (which may not be tight because this is the stochastic case), the first term (i.e., the “optimization
 729 term”) in (34) scales with κ when T is large and $\frac{\sqrt{\kappa} \zeta_*}{\sqrt{\mu \epsilon}}$ when T is small. This is clearly worse than
 730 (33) for large n .

731 G Nonconvex analysis

732 We shall now present our theory for the nonconvex case. To quantify convergence, we define the
 733 proximal-gradient mapping, which was also used in the prior literature to show convergence of
 734 Proximal SGD.

735 **Definition 2.** Given a stepsize $\gamma > 0$, a convex function ψ and arbitrary f , we define the proximal-
 736 gradient mapping as

$$\mathcal{G}_\gamma(x) := \frac{1}{\gamma} [x - \text{prox}_{\gamma\psi}(x - \gamma \nabla f(x))].$$

737 Similarly to Theorem 1, the analysis shows that a gradient term appears in the variance bound.
 738 However, in contrast to the convex settings of Theorem 1, there might not exist an optimum to which
 739 the iterates would converge and we cannot use $\|\nabla f(x_*)\|^2$ in the variance bound. For this reason, we
 740 resort to the following assumption that bounds full gradients in terms of proximal-gradient mapping
 741 and an extra constant.

Assumption 3. There exists a constant $\zeta \geq 0$ such that the full gradient of f is uniformly bounded by the proximal-gradient mapping and ζ

$$\|\nabla f(x)\|^2 \leq \|\mathcal{G}_{\gamma n}(x)\|^2 + \zeta^2$$

for any $x \in \text{dom}(\psi)$ and $\gamma > 0$.

We note that this assumption is trivially satisfied with $\zeta = 0$ if $\psi \equiv 0$ because in that case, $\mathcal{G}_{\gamma}(x) \equiv \nabla f(x)$. Therefore, when there is no proximal term, it is not an extra assumptions compared to the analysis of [Mishchenko et al. \[2020\]](#). We will also rely on the following measure of gradient variance, which we need for the same reason that there might be no optimum x_* to measure the variance the way we did for Theorem 1.

An important property of Assumption 3 is that it is equivalent to the bounded dissimilarity assumption that was previously used for the nonconvex analysis of Local SGD. We formalize this in the following proposition.

Proposition 1. Consider federated learning reformulation (11). If $\psi \equiv \psi_C$, i.e., $R \equiv 0$, then Assumption 3 with constant $\bar{\zeta}^2 := M\zeta^2$ is equivalent to ζ -bounded dissimilarity (Assumption 2):

$$\frac{1}{M} \sum_{m=1}^M \left\| \nabla F_m(x) - \frac{1}{M} \sum_{l=1}^M \nabla F_l(x) \right\|^2 \leq \zeta^2.$$

Proof. First, observe that if $\mathbf{x} \in \text{dom}(\psi)$, then \mathbf{x} has all blocks equal to some $x \in \mathbb{R}^d$, $\mathbf{x} = (x^\top, \dots, x^\top)^\top$. Therefore, for the objective in reformulation (11) and $\mathbf{x} \in \text{dom}(\psi)$, we have

$$\begin{aligned} \nabla f(\mathbf{x}) &= \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \sum_{m=1}^M \nabla f_{mi}(\mathbf{x}) = \frac{1}{n} \sum_{m=1}^M \sum_{i=1}^n \nabla f_{mi}(\mathbf{x}) \\ &= \frac{1}{n} \sum_{m=1}^M F_m(\mathbf{x}) = \begin{pmatrix} \frac{1}{n} \nabla F_1(x_1) \\ \vdots \\ \frac{1}{n} \nabla F_M(x_M) \end{pmatrix} = \begin{pmatrix} \frac{1}{n} \nabla F_1(x) \\ \vdots \\ \frac{1}{n} \nabla F_M(x) \end{pmatrix}. \end{aligned} \quad (35)$$

With the help of bias-variance decomposition, the left-hand side of Assumption 3 can be written as

$$\begin{aligned} \|\nabla f(\mathbf{x})\|^2 &\stackrel{(35)}{=} \frac{1}{n^2} \sum_{m=1}^M \|\nabla F_m(x)\|^2 \\ &\stackrel{(13)}{=} \frac{1}{Mn^2} \left\| \sum_{m=1}^M \nabla F_m(x) \right\|^2 + \frac{1}{n^2} \sum_{m=1}^M \left\| \nabla F_m(x) - \frac{1}{M} \sum_{l=1}^M \nabla F_l(x) \right\|^2. \end{aligned}$$

Let us now work out the proximal-gradient mapping. According to Lemma 7, the proximal operator of ψ is simply the averaging of all blocks, while the full gradient is given in (35), which give when combined

$$\text{prox}_{\gamma n \psi}(\mathbf{x} - \gamma n \nabla f(\mathbf{x})) = \begin{pmatrix} \frac{1}{M} \sum_{m=1}^M (x - \gamma \nabla F_m(x)) \\ \vdots \\ \frac{1}{M} \sum_{m=1}^M (x - \gamma \nabla F_m(x)) \end{pmatrix}. \quad (36)$$

Therefore,

$$\begin{aligned} \|\mathcal{G}_{\gamma n}(\mathbf{x})\|^2 &= \frac{1}{\gamma^2 n^2} \|\mathbf{x} - \text{prox}_{\gamma n \psi}(\mathbf{x} - \gamma n \nabla f(\mathbf{x}))\|^2 \\ &\stackrel{(36)}{=} \frac{1}{\gamma^2 n^2} \sum_{l=1}^M \left\| x - \frac{1}{M} \sum_{m=1}^M (x - \gamma \nabla F_m(x)) \right\|^2 = \frac{M}{n^2} \left\| \frac{1}{M} \sum_{m=1}^M \nabla F_m(x) \right\|^2. \end{aligned} \quad (37)$$

Having the expressions for both sides, we can write

$$\|\nabla f(\mathbf{x})\|^2 = \|\mathcal{G}_{\gamma n}(\mathbf{x})\|^2 + \sum_{m=1}^M \left\| \frac{1}{n} \nabla F_m(x) - \frac{1}{N} \sum_{l=1}^M \nabla F_l(x) \right\|^2.$$

From this expression and the fact $\frac{1}{N} \sum_{l=1}^M \nabla F_l(x) = \nabla F(x)$, it is easy to see the equivalence. ■

To analyze ProxRR, we also need to measure variance differently from how it was done the strongly convex case because we cannot rely on convergence of iterates to x_* . To this end, we introduce the following assumption, which is quite standard in the literature on SGD.

Assumption 4. There exists a constant $\sigma > 0$ such that $\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x)\|^2 \leq \sigma^2$ for any $x \in \mathbb{R}^d$.

This assumption is more restrictive than the one from [Mishchenko et al. \[2020\]](#) and, in fact, we could relax it a little by introducing extra terms in the right-hand side. Nevertheless, for the sake of simplicity and readability, we prefer the stronger version as presented above.

G.1 A key lemma

For notational convenience, we define

$$g_t := \frac{1}{\gamma n} (x_t - x_t^n) = \frac{1}{n} \sum_{i=0}^{n-1} \nabla f_{\pi_i}(x_t^i),$$

which is equivalent to $x_t^n = x_t - \gamma n g_t$.

Lemma 8. Let functions f_1, \dots, f_n be L_{\max} -smooth, Assumptions 3 and 4 be satisfied and $\gamma \leq \frac{1}{2L_{\max}n}$. Then,

$$\mathbb{E}_t \left[\|\nabla f(x_t) - g_t\|^2 \right] \leq \gamma^2 L_{\max}^2 n^2 (\|\mathcal{G}_{\gamma n}(x_t)\|^2 + \zeta^2) + \gamma^2 L_{\max}^2 n \sigma^2. \quad (38)$$

Proof. We start with the observation that gradient Lipschitzness reduces the left-hand side to a difference of iterates:

$$\begin{aligned} \|\nabla f(x_t) - g_t\|^2 &= \left\| \frac{1}{n} \sum_{i=0}^{n-1} [\nabla f_{\pi_i}(x_t) - \nabla f_{\pi_i}(x_t^i)] \right\|^2 \\ &\leq \frac{1}{n} \sum_{i=0}^{n-1} \|\nabla f_{\pi_i}(x_t) - \nabla f_{\pi_i}(x_t^i)\|^2 \\ &\leq \frac{1}{n} \sum_{i=0}^{n-1} L_{\max}^2 \|x_t - x_t^i\|^2. \end{aligned}$$

Define $V_t := \sum_{i=0}^{n-1} \|x_t^i - x_t\|^2$. Clearly, it is sufficient to bound $\mathbb{E}[V_t]$ to finish the proof. Also note that for any intermediate iterate x_t^k within epoch t we do not use proximal step, so the following identity holds:

$$x_t^k = x_t - \gamma \sum_{i=0}^{k-1} \nabla f_{\pi_i}(x_t^i).$$

This identity only includes gradients, so to bound the deviation of x_t^k from x_t we apply Jensen's inequality and gradient Lipschitzness

$$\begin{aligned} \mathbb{E}_t [\|x_t^k - x_t\|^2] &= \gamma^2 \mathbb{E}_t \left[\left\| \sum_{i=0}^{k-1} \nabla f_{\pi_i}(x_t^i) \right\|^2 \right] \\ &\leq 2\gamma^2 \mathbb{E}_t \left[\left\| \sum_{i=0}^{k-1} (\nabla f_{\pi_i}(x_t^i) - \nabla f_{\pi_i}(x_t)) \right\|^2 \right] + 2\gamma^2 \mathbb{E}_t \left[\left\| \sum_{i=0}^{k-1} \nabla f_{\pi_i}(x_t) \right\|^2 \right] \\ &\leq 2\gamma^2 k \sum_{i=0}^{k-1} \mathbb{E}_t [\|\nabla f_{\pi_i}(x_t^i) - \nabla f_{\pi_i}(x_t)\|^2] + 2\gamma^2 \mathbb{E}_t \left[\left\| \sum_{i=0}^{k-1} \nabla f_{\pi_i}(x_t) \right\|^2 \right] \\ &\leq 2\gamma^2 L_{\max}^2 k \sum_{i=0}^{k-1} \mathbb{E}_t [\|x_t^i - x_t\|^2] + 2\gamma^2 \mathbb{E}_t \left[\left\| \sum_{i=0}^{k-1} \nabla f_{\pi_i}(x_t) \right\|^2 \right]. \end{aligned}$$

784 Now we are going to use the fact that for any i in RR we have $\mathbb{E}_t [\nabla f_{\pi_i}(x_t)] = \nabla f(x_t)$. Note that
 785 this property does not hold if x_t is not independent of π_i , which is why the result does not hold for
 786 SO. Let us also define $\sigma_t^2 := \frac{1}{n} \sum_{j=1}^n \|\nabla f_j(x_t) - \nabla f(x_t)\|^2$. By Lemma 3 we have

$$\begin{aligned} \mathbb{E}_t \left[\left\| \sum_{i=0}^{k-1} \nabla f_{\pi_i}(x_t) \right\|^2 \right] &= k^2 \|\nabla f(x_t)\|^2 + k^2 \mathbb{E}_t \left[\left\| \frac{1}{k} \sum_{i=0}^{k-1} (\nabla f_{\pi_i}(x_t) - \nabla f(x_t)) \right\|^2 \right] \\ &\stackrel{(14)}{=} k^2 \|\nabla f(x_t)\|^2 + \frac{k(n-k)}{n-1} \sigma_t^2. \end{aligned}$$

787 Plugging this back and using Assumption 4, we derive

$$\begin{aligned} \mathbb{E}_t \left[\|x_t^k - x_t\|^2 \right] &\leq 2\gamma^2 L_{\max}^2 k \sum_{i=0}^{k-1} \mathbb{E}_t \left[\|x_t^i - x_t\|^2 \right] + 2\gamma^2 k^2 \|\nabla f(x_t)\|^2 + 2\gamma^2 \frac{k(n-k)}{n-1} \sigma^2 \\ &\leq 2\gamma^2 L_{\max}^2 k \mathbb{E} [V_t] + 2\gamma^2 k^2 \|\nabla f(x_t)\|^2 + 2\gamma^2 \frac{k(n-k)}{n-1} \sigma^2. \end{aligned}$$

788 Let us use the obtained bound on a single iterate distance $\mathbb{E}_t \left[\|x_t^k - x_t\|^2 \right]$ to upper bound $\mathbb{E} [V_t]$:

$$\begin{aligned} \mathbb{E}_t [V_t] &= \sum_{i=0}^{n-1} \mathbb{E}_t \left[\|x_t^i - x_t\|^2 \right] \\ &\leq \gamma^2 L_{\max}^2 n(n-1) \mathbb{E}_t [V_t] + \frac{1}{3} \gamma^2 (n-1)n(2n-1) \|\nabla f(x_t)\|^2 + \frac{1}{3} \gamma^2 n(n+1) \sigma^2. \end{aligned}$$

789 This inequality has $\mathbb{E}_t [V_t]$ in both sides, so we can rearrange it and use the assumption $\gamma \leq \frac{1}{2L_{\max}n}$,
 790 which results in

$$\begin{aligned} \mathbb{E}_t [V_t] &\leq \frac{4}{3} (1 - \gamma^2 L_{\max}^2 n(n-1)) \mathbb{E}_t [V_t] \\ &\leq \frac{4}{9} \gamma^2 (n-1)n(2n-1) \|\nabla f(x_t)\|^2 + \frac{4}{9} \gamma^2 n(n+1) \sigma^2 \\ &\leq \gamma^2 n^3 \|\nabla f(x_t)\|^2 + \gamma^2 n^2 \sigma^2. \end{aligned}$$

791 To conclude the proof, apply Assumption 3 to $x_t \in \text{dom}(\psi)$ and plug-in the bound on $\mathbb{E}_t [V_t]$ into
 792 the bound on $\mathbb{E}_t [\|\nabla f(x_t) - g_t\|^2]$. ■

793 G.2 Main theorem

794 **Theorem 8** (Convergence result in the nonconvex case). Let Assumptions 3 and 4 hold and choose
 795 any $\gamma \leq \frac{1}{5L_{\max}n}$. Then,

$$\min_{t=0, \dots, T-1} \mathbb{E} [\|\mathcal{G}_{\gamma n}(x_t)\|^2] \leq \frac{4(P(x_0) - P_*)}{\gamma n T} + 2\gamma^2 L_{\max} n^2 \zeta^2 + 2\gamma^2 L_{\max}^2 n \sigma^2.$$

796 *Proof.* Let us introduce

$$w_t := \text{prox}_{\gamma n \psi}(x_t - \gamma n \nabla f(x_t)).$$

797 The idea of our proof is to first obtain a descent recursion for $P(w_t)$ and then bound $P(x_{t+1}) - P(w_t)$.

798 By convexity of ψ , we have for any $g \in \partial \psi(w_t)$

$$\psi(w_t) \leq \psi(x_t) + \langle g, w_t - x_t \rangle.$$

799 Furthermore, the definition of w_t implies by first-order optimality that $x_t - \gamma n \nabla f(x_t) - w_t \in$
 800 $\gamma n \partial \psi(w_t)$, so we can plug it into the bound above to get

$$\begin{aligned} \psi(w_t) &\leq \psi(x_t) + \frac{1}{\gamma n} \langle x_t - \gamma n \nabla f(x_t) - w_t, w_t - x_t \rangle \\ &= \psi(x_t) - \langle \nabla f(x_t), w_t - x_t \rangle - \frac{1}{\gamma n} \|w_t - x_t\|^2. \end{aligned}$$

801 At the same time, by L_{\max} -smoothness of f we have

$$f(w_t) \leq f(x_t) + \langle \nabla f(x_t), w_t - x_t \rangle + \frac{L_{\max}}{2} \|w_t - x_t\|^2.$$

802 Adding the two recursion together yields

$$P(w_t) = f(x_t) + \psi(w_t) \leq P(x_t) + \left(\frac{L_{\max}}{2} - \frac{1}{\gamma n} \right) \|w_t - x_t\|^2.$$

803 Now we shall upper bound $P(x_{t+1})$. Using the convexity of ψ for $x_t^n - x_{t+1} \in \gamma n \partial \psi(x_{t+1})$, we
804 derive

$$\begin{aligned} \psi(x_{t+1}) &\leq \psi(w_t) + \frac{1}{\gamma n} \langle x_t^n - x_{t+1}, x_{t+1} - w_t \rangle = \psi(w_t) - \langle g_t, x_{t+1} - w_t \rangle + \frac{1}{\gamma n} \langle x_t - x_{t+1}, x_{t+1} - w_t \rangle \\ &= \psi(w_t) - \langle g_t, x_{t+1} - w_t \rangle + \frac{1}{2\gamma n} (\|x_t - w_t\|^2 - \|x_t - x_{t+1}\|^2 - \|x_{t+1} - w_t\|^2). \end{aligned}$$

805 Next, we apply L_{\max} -smoothness of f two times, to upper bound $D_f(x_{t+1}, x_t)$ and to lower bound
806 $D_f(w_t, x_t)$:

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t \rangle + \frac{L_{\max}}{2} \|x_{t+1} - x_t\|^2, \\ \text{and} \quad f(x_t) &\leq f(w_t) + \langle \nabla f(x_t), x_t - w_t \rangle + \frac{L_{\max}}{2} \|x_t - w_t\|^2. \end{aligned}$$

807 Therefore,

$$f(x_{t+1}) \leq f(w_t) + \langle \nabla f(x_t), x_{t+1} - w_t \rangle + \frac{L_{\max}}{2} (\|x_{t+1} - x_t\|^2 + \|w_t - x_t\|^2).$$

808 Combining the inequalities for $\psi(x_{t+1})$ and $f(x_{t+1})$, we obtain

$$\begin{aligned} P(x_{t+1}) &\leq P(w_t) + \langle \nabla f(x_t) - g_t, x_{t+1} - w_t \rangle + \left(\frac{L_{\max}}{2} - \frac{1}{2\gamma n} \right) \|x_{t+1} - x_t\|^2 \\ &\quad + \left(\frac{L_{\max}}{2} + \frac{1}{2\gamma n} \right) \|x_t - w_t\|^2 - \frac{1}{2\gamma n} \|x_{t+1} - w_t\|^2. \end{aligned}$$

809 By Young's inequality and Lemma 8 we have

$$\begin{aligned} &\mathbb{E}_t [\langle \nabla f(x_t) - g_t, x_{t+1} - w_t \rangle] \\ &\leq \mathbb{E}_t \left[\frac{\gamma n}{2} \|\nabla f(x_t) - g_t\|^2 + \frac{2}{\gamma n} \|x_{t+1} - w_t\|^2 \right] \\ &\stackrel{(38)}{\leq} \frac{\gamma^3 L_{\max}^2 n^3}{2} \zeta^2 + \frac{\gamma^3 L_{\max}^2 n^3}{2} \|\mathcal{G}_{\gamma n}(x_t)\|^2 + \frac{\gamma^3 L_{\max}^2 n^2}{2} \sigma^2 + \frac{2}{\gamma n} \mathbb{E}_t [\|x_{t+1} - w_t\|^2]. \end{aligned}$$

810 If we plug this back, the term $\|x_{t+1} - w_t\|^2$ will cancel out, giving us for $\gamma \leq \frac{1}{L_{\max} n}$

$$\begin{aligned} &\mathbb{E}_t [P(x_{t+1})] \\ &\leq P(w_t) + \frac{\gamma^3 L_{\max}^2 n^3}{2} \zeta^2 + \frac{\gamma^3 L_{\max}^2 n^3}{2} \|\mathcal{G}_{\gamma n}(x_t)\|^2 + \frac{\gamma^3 L_{\max}^2 n^2}{2} \sigma^2 + \left(\frac{L_{\max}}{2} + \frac{1}{2\gamma n} \right) \|x_t - w_t\|^2 \\ &\quad + \left(\frac{L_{\max}}{2} - \frac{1}{2\gamma n} \right) \mathbb{E}_t [\|x_{t+1} - x_t\|^2] \\ &\leq P(w_t) + \frac{\gamma^3 L_{\max}^2 n^3}{2} \zeta^2 + \frac{\gamma^3 L_{\max}^2 n^3}{2} \|\mathcal{G}_{\gamma n}(x_t)\|^2 + \frac{\gamma^3 L_{\max}^2 n^2}{2} \sigma^2 + \left(\frac{L_{\max}}{2} + \frac{1}{2\gamma n} \right) \|x_t - w_t\|^2. \end{aligned}$$

811 Using the recursion for $P(w_t)$ and our choice $\gamma \leq \frac{1}{5L_{\max} n}$, we finally obtain, after plugging-in

812 $\|x_t - w_t\|^2 = \gamma^2 n^2 \mathcal{G}_{\gamma n}(x_t)$,

$$\begin{aligned} &\mathbb{E}_t [P(x_{t+1})] \\ &\leq P(x_t) + \frac{\gamma^3 L_{\max}^2 n^3}{2} \zeta^2 + \frac{\gamma^3 L_{\max}^2 n^2}{2} \sigma^2 + \left(\frac{\gamma n L_{\max}^2}{2} + \frac{L_{\max}}{2} + \frac{1}{2\gamma n} + \frac{L_{\max}}{2} - \frac{1}{\gamma n} \right) \gamma^2 n^2 \|\mathcal{G}_{\gamma n}(x_t)\|^2 \\ &\leq P(x_t) + \frac{\gamma^3 L_{\max}^2 n^3}{2} \zeta^2 + \frac{\gamma^3 L_{\max}^2 n^2}{2} \sigma^2 + \left(\frac{L_{\max}}{10} + L_{\max} - \frac{1}{2\gamma n} \right) \gamma^2 n^2 \|\mathcal{G}_{\gamma n}(x_t)\|^2 \\ &\leq P(x_t) + \frac{\gamma^3 L_{\max}^2 n^3}{2} \zeta^2 + \frac{\gamma^3 L_{\max}^2 n^2}{2} \sigma^2 - \frac{1}{4\gamma n} \gamma^2 n^2 \|\mathcal{G}_{\gamma n}(x_t)\|^2. \end{aligned}$$

813 Recursing this to $P(x_0)$ and using $P_* \leq P(x_T)$, we get the Theorem's claim. \blacksquare

814 **Obtaining a complexity.** To make the upper bound equal $\mathcal{O}(\varepsilon^2)$, it is sufficient to ensure that every
815 term is equal $\mathcal{O}(\varepsilon^2)$. Therefore, we can impose the following conditions:

$$\gamma nT \geq \frac{P(x_0) - P_*}{\varepsilon^2} \quad \text{and} \quad \gamma^2(L_{\max} n^2 \zeta^2 + L_{\max}^2 n \sigma^2) \leq \varepsilon^2$$

816 To satisfy these conditions, we can choose γ as

$$\gamma = \min \left\{ \frac{1}{5L_{\max} n}, \frac{\varepsilon}{L_{\max} \sqrt{n} \sigma + \sqrt{L_{\max} n} \zeta} \right\}.$$

817 Then, denoting $\delta_0 := P(x_0) - P_*$, we obtain complexity in terms of full number of stochastic
818 gradients nT equal to

$$nT = \mathcal{O} \left(\frac{\delta_0 L_{\max} n}{\varepsilon^2} + \frac{\delta_0 L_{\max} \sqrt{n} \sigma}{\varepsilon^3} + \frac{\delta_0 \sqrt{L_{\max} n} \zeta}{\varepsilon^3} \right).$$

819 G.3 Proof of Theorem 5

820 The federated learning reformulation (11) has different constant scaling than the finite-sum federated
821 learning problem (10), and the only constant that does not change at all is L_{\max} . For the initial error
822 $\bar{\delta}_0$ of the reformulation we have

$$\bar{\delta}_0 = \frac{N}{n} \delta_0 = M \delta_0,$$

823 where $\delta_0 := \frac{1}{N} \sum_{m=1}^M F_m(x_0) - \min_x \frac{1}{N} \sum_{m=1}^M F_m(x)$ and we use only consider the simplified
824 case $N_1 = \dots = N_M = n$ so $\frac{N}{n} = M$. For the variance, we have

$$\bar{\sigma}^2 = \sup_{\mathbf{x}} \mathbb{E} [\|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2] = \sup_{\mathbf{x}} \sum_{m=1}^M \mathbb{E} \left[\left\| \nabla f_{mi}(x_m) - \frac{1}{n} \nabla F_m(x_m) \right\|^2 \right] = M \sigma^2.$$

825 As we derived in (37), the proximal-gradient mapping norm is equal to

$$\mathbb{E} [\|\mathcal{G}_{\gamma n}(\mathbf{x})\|^2] = \frac{M}{n^2} \left\| \frac{1}{M} \sum_{m=1}^M \nabla F_m(x) \right\|^2 = M \left\| \frac{1}{N} \sum_{m=1}^M \nabla F_m(x) \right\|^2 = M \|\nabla F(x)\|^2,$$

826 so to have $\|\nabla F(x_T)\|^2 \leq \varepsilon^2$, we need $\mathbb{E} [\|\mathcal{G}_{\gamma n}(\mathbf{x}_T)\|^2] \leq \bar{\varepsilon}^2 := M \varepsilon^2$. In addition, notice that by
827 Proposition 1 the constant from Assumption 3 is $\bar{\zeta} = \sqrt{M} \zeta$.

828 Thus, Theorem 8 implies, if we ignore L_{\max} , that we need

$$T = \mathcal{O} \left(\frac{\bar{\delta}_0}{\bar{\varepsilon}^2} + \frac{\bar{\delta}_0 \bar{\sigma}}{\sqrt{n} \bar{\varepsilon}^3} + \frac{\bar{\delta}_0 \bar{\zeta}}{\bar{\varepsilon}^3} \right) = \mathcal{O} \left(\frac{\delta_0}{\varepsilon^2} + \frac{\delta_0 \sigma}{\sqrt{n} \varepsilon^3} + \frac{\delta_0 \zeta}{\varepsilon^3} \right)$$

829 communication rounds to achieve $\min_{t=0, \dots, T-1} \mathbb{E} [\|\nabla F(x_T)\|^2] = \mathcal{O}(\varepsilon^2)$.

830 H Further experimental details

831 **Implementation details.** For each i , we have $L_i = \frac{1}{4} \|a_i\|$. For the ℓ_1 -regularized problem, we set
832 $\lambda_2 = 3 \cdot 10^{-5} \cdot L$ and tune λ_1 to obtain a solution with about 25% zero coordinates, which gives
833 $\lambda_1 = 5 \cdot 10^{-5}$. We use stepsizes decreasing as $\mathcal{O}(\frac{1}{t})$ for all methods. We use the 'w8a' dataset⁴ for
834 the experiment with ℓ_1 regularization.

835 **Proximal operator calculation.** It is well-known (see, for instance, [Parikh and Boyd, 2014]) that
836 the proximal operator for $\psi(x) = \lambda_1 \|x\|_1 + \frac{\lambda_2}{2} \|x\|^2$ is given by

$$\text{prox}_{\gamma \psi}(x) = \frac{1}{1 + \gamma \lambda_2} \text{prox}_{\gamma \lambda_1 \|\cdot\|_1}(x),$$

⁴The datasets were downloaded from LibSVM <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html>

837 where the j -th coordinate of $\text{prox}_{\gamma\lambda_1\|\cdot\|_1}(x)$ is

$$[\text{prox}_{\gamma\lambda_1\|\cdot\|_1}(x)]_j = \begin{cases} \text{sign}([x]_j)(|[x]_j| - \gamma\lambda_1), & \text{if } |[x]_j| \geq \gamma\lambda_1, \\ 0, & \text{otherwise.} \end{cases}$$

838 **Federated experiments.** The experiments for the comparison of FedRR, Local SGD and Scaffold
839 use no ℓ_1 regularization and $\lambda_2 = 10^{-5} \cdot L$. To make comparison fair, all methods use n local steps.
840 For FedRR, the initial stepsize was $\frac{1}{L}$ in the i.i.d. regime and $\frac{1}{Ln}$ in the heterogeneous regime. As
841 per Theorem 3 in [Khaled et al., 2020], the stepsizes for Local SGD must satisfy $\gamma_t = \mathcal{O}(1/(LH))$,
842 where H is the number of local steps, a similar result holds for Scaffold Karimireddy et al. [2020].
843 The parallelization of local runs is done using the Ray package⁵. We use the ‘w8a’ dataset for the
844 i.i.d. experiment. For the heterogeneous experiment, we sort ‘a9a’ dataset with respect to the target
845 labels $b \in \{0, 1\}$ and then mix it with the original order in proportion 2:1. For all methods, the local
846 workers used minibatch size 16. Exact implementation can be found in our code.

⁵<https://ray.io/>

Part II

Extensions

Here we discuss two extensions of our theory that significantly matter in practice: using decreasing stepsizes and applying importance resampling.

I Extension: Decreasing stepsizes

Using the theoretical stepsize (8) requires knowing the desired accuracy ε ahead of time as well as estimating σ_{rad} . It also results in extra polylogarithmic factors in the iteration complexity (9), a phenomenon observed and fixed by using decreasing stepsizes in both vanilla RR [Ahn et al., 2020] and in SGD [Stich, 2019].

We show that we can adopt the same technique to our setting. However, we depart from the stepsize scheme of Ahn et al. [2020] by only varying the stepsize *once per epoch* rather than every iteration. This is closer to the common practical heuristic of decreasing the stepsize once every epoch or once every few epochs [Sun, 2020, Tran et al., 2020]. The stepsize scheme we use is inspired by the schemes of [Stich, 2019, Khaled and Richtárik, 2020]: in particular, we fix $T > 0$, let $t_0 = \lceil T/2 \rceil$, and choose the stepsizes $\gamma_t > 0$ by

$$\gamma_t = \begin{cases} \frac{1}{L_{\max}} & \text{if } T \leq \frac{L_{\max}}{2\mu n} \text{ or } t \leq t_0, \\ \frac{7}{\mu n(s+t-t_0)} & \text{if } T > \frac{L_{\max}}{2\mu n} \text{ and } t > t_0, \end{cases} \quad (39)$$

where $s := 7L_{\max}/(4\mu n)$. Hence, we fix the stepsize used in the first $T/2$ iterations and then start decreasing it every epoch afterwards. Using this stepsize schedule, we can obtain the following convergence guarantee when each f_i is smooth and convex and the regularizer ψ is μ -strongly convex.

Theorem 9. Suppose that each f_i is L_{\max} -smooth and convex, and that the regularizer ψ is μ -strongly convex. Fix $T > 0$. Then choosing stepsizes γ_t according to (39) we have that $\gamma_t \leq 1/L_{\max}$ for all t and the final iterate generated by Algorithm 1 satisfies

$$\mathbb{E} [\|x_T - x_*\|^2] = \mathcal{O} \left(\exp \left(-\frac{nT}{\kappa+2n} \right) r_0 + \frac{\sigma_{\text{rad}}^2}{\mu^3 n^2 T^2} \right),$$

where $\kappa := L_{\max}/\mu$, $r_0 := \|x_0 - x_*\|^2$ and $\mathcal{O}(\cdot)$ hides absolute (non-problem-specific) constants.

This guarantee holds for any number of epochs $T > 0$. We believe a similar guarantee can be obtained in the case each f_i is strongly-convex and the regularizer ψ is just convex, but we did not include it as it adds little to the overall message.

In the rest of the section we provide a proof of Theorem 9.

I.1 A recursion Lemma

We first state and prove the following algorithm-independent lemma. This lemma plays a key role in the proof of Theorem 9 and is heavily inspired by the stepsize schemes of Stich [2019] and Khaled and Richtárik [2020] and their proofs.

Lemma 9. Suppose that there exist constants $a, b, c \geq 0$ such that for all $\gamma_t \leq \frac{1}{b}$ we have

$$(1 + \gamma_t a n) r_{t+1} \leq r_t + \gamma_t^3 c. \quad (40)$$

Fix $T > 0$. Let $t_0 = \lceil \frac{T}{2} \rceil$. Then choosing stepsizes $\gamma_t > 0$ by

$$\gamma_t = \begin{cases} \frac{1}{b}, & \text{if } t \leq t_0 \text{ or } T \leq \frac{b}{an}, \\ \frac{7}{an(s+t-t_0)} & \text{if } t > t_0 \text{ and } T > \frac{b}{an}, \end{cases}$$

where $s = \frac{7b}{2an}$. Then

$$r_T \leq \exp \left(-\frac{nT}{2(b/a + n)} \right) r_0 + \frac{1421c}{a^3 n^3 T^2}.$$

880 *Proof.* If $T \leq \frac{7b}{an}$, then we have $\gamma_t = \gamma = \frac{1}{b}$ for all t . Hence recursing we have,

$$r_T \leq (1 + \gamma an)^{-T} r_0 + \frac{\gamma^3 c}{\gamma an} = (1 + \gamma an)^{-T} r_0 + \frac{\gamma^2 c}{an}.$$

881 Note that $\frac{1}{1+x} \leq \exp(-\frac{x}{1+x})$ for all x , hence

$$r_T \leq \exp\left(-\frac{\gamma an T}{1 + \gamma an}\right) r_0 + \frac{\gamma^2 c}{an}$$

882 Substituting for γ yields

$$r_T \leq \exp\left(-\frac{nT}{b/a + n}\right) r_0 + \frac{c}{b^2 an}.$$

883 Note that by assumption we have $\frac{1}{b} \leq \frac{7}{T an}$, hence

$$r_T \leq \exp\left(-\frac{nT}{b/a + n}\right) r_0 + \frac{49c}{T^2 a^3 n^3}. \quad (41)$$

884 If $T > \frac{7b}{an}$, then we have for the first phase when $t \leq t_0$ with stepsize $\gamma_t = \frac{1}{b}$ that

$$r_{t_0} \leq \exp\left(-\frac{nt_0}{b/a + n}\right) r_0 + \frac{c}{b^2 an} \leq \exp\left(-\frac{nT}{2(b/a + n)}\right) r_0 + \frac{c}{b^2 an}. \quad (42)$$

885 Then for $t > t_0$ we have

$$(1 + \gamma_t an) r_{t+1} \leq r_t + \gamma_t^3 c = r_t + \frac{7^3 c}{a^3 n^3 (s + t - t_0)^3}.$$

886 Multiplying both sides by $(s + t - t_0)^3$ yields

$$(s + t - t_0)^3 (1 + \gamma_t an) r_{t+1} \leq (s + t - t_0)^3 r_t + \frac{7^3 c}{a^3 n^3}. \quad (43)$$

887 Note that because t and t_0 are integers and $t > t_0$, we have that $t - t_0 \geq 1$ and therefore $s + t - t_0 \geq 1$.

888 We may use this to lower bound the multiplicative factor in the left hand side of (43) as

$$\begin{aligned} (s + t - t_0)^3 (1 + \gamma_t an) &= (s + t - t_0)^3 \left(1 + \frac{7}{s + t - t_0}\right) \\ &= (s + t - t_0)^3 + 7(s + t - t_0)^2 \\ &= (s + t - t_0)^3 + 3(s + t - t_0)^2 + 3(s + t - t_0)^2 + (s + t - t_0)^2 \\ &\geq (s + t - t_0)^3 + 3(s + t - t_0)^2 + 3(s + t - t_0) + 1 \\ &= (s + t + 1 - t_0)^3. \end{aligned} \quad (44)$$

889 Using (44) in (43) we obtain

$$(s + t + 1 - t_0)^3 r_{t+1} \leq (s + t - t_0)^3 r_t + \frac{7^3 c}{a^3 n^3}.$$

890 Let $w_t = (s + t - t_0)^3$. Then we can rewrite the last inequality as

$$w_{t+1} r_{t+1} - w_t r_t \leq \frac{7^3 c}{a^3 n^3}.$$

891 Summing up and telescoping from $t = t_0$ to T yields

$$w_T r_T \leq w_{t_0} r_{t_0} + \frac{7^3 c}{a^3 n^3} (T - t_0).$$

892 Note that $w_{t_0} = s^3$ and $w_T = (s + T - t_0)^3$. Hence,

$$\begin{aligned} r_T &\leq \frac{s^3}{(s + T - t_0)^3} r_{t_0} + \frac{7^3 c}{a^3 n^3 (s + T - t_0)^2} \frac{T - t_0}{s + T - t_0} \\ &\leq \frac{s^3}{(s + T - t_0)^3} r_{t_0} + \frac{7^3 c}{a^3 n^3 (s + T - t_0)^2}. \end{aligned}$$

893 Since we have $s + T - t_0 \geq T - t_0 \geq T/2$, it holds

$$r_T \leq \frac{8s^3}{T^3} r_{t_0} + \frac{4 \cdot 7^3 c}{a^3 n^3 T^2}. \quad (45)$$

894 The bound in (42) can be rewritten as

$$\frac{s^3}{T^3} r_{t_0} \leq \frac{s^3}{T^3} \exp\left(-\frac{nT}{2(b/a + n)}\right) r_0 + \frac{s^3 c}{b^2 a n T^3}.$$

895 We now rewrite the last inequality, use that $T > 2s$ and further use the fact that $s = \frac{7b}{2an}$:

$$\begin{aligned} \frac{s^3}{T^3} r_{t_0} &\leq \underbrace{\left(\frac{s}{T}\right)^3}_{\leq 1/8} \exp\left(-\frac{nT}{2(b/a + n)}\right) r_0 + \frac{s^2 c}{b^2 a n T^2} \underbrace{\left(\frac{s}{T}\right)}_{\leq 1/2} \\ &\leq \frac{1}{8} \exp\left(-\frac{nT}{2(b/a + n)}\right) r_0 + \frac{s^2 c}{2b^2 a n T^2} \\ &= \frac{1}{8} \exp\left(-\frac{nT}{2(b/a + n)}\right) r_0 + \frac{7^2 c}{8a^3 n^3 T^2}. \end{aligned} \quad (46)$$

896 Plugging in the estimate of (46) into (45) we obtain

$$\begin{aligned} r_T &\leq \exp\left(-\frac{nT}{2(b/a + n)}\right) r_0 + \frac{7^2 c}{a^3 n^3 T^2} + \frac{4 \cdot 7^3 c}{a^3 n^3 T^2} \\ &= \exp\left(-\frac{nT}{2(b/a + n)}\right) r_0 + \frac{1421c}{a^3 n^3 T^2}. \end{aligned} \quad (47)$$

897 Taking the maximum of (41) and (47) we see that for any $T > 0$ we have

$$r_T \leq \exp\left(-\frac{nT}{2(b/a + n)}\right) r_0 + \frac{1421c}{a^3 n^3 T^2}. \quad \blacksquare$$

898 I.2 Proof of Theorem 9

899 *Proof.* Start with Lemma 6 with $\lambda = 0$, $L = L_{\max}$, and $\gamma = \gamma_t$,

$$\mathbb{E} \left[\|x_t^{i+1} - x_*^{i+1}\|^2 \right] \leq \mathbb{E} \left[\|x_t^i - x_*^i\|^2 \right] - 2\gamma(1 - \gamma L_{\max}) \mathbb{E} [D_{f_{\pi_i}}(x_t^i, x_*)] + 2\gamma_t^3 \sigma_{\text{rad}}^2.$$

900 Since $\gamma \leq 1/L_{\max}$ and $D_{f_{\pi}}(x_t^i, x_*)$ is nonnegative we may simplify this to

$$\mathbb{E} \left[\|x_t^{i+1} - x_*^{i+1}\|^2 \right] \leq \mathbb{E} \left[\|x_t^i - x_*^i\|^2 \right] + 2\gamma_t^3 \sigma_{\text{rad}}^2.$$

901 Unrolling this recursion for n steps we get

$$\mathbb{E} \left[\|x_t^n - x_*^n\|^2 \right] \leq \mathbb{E} \left[\|x_t^0 - x_*^0\|^2 \right] + 2n\gamma_t^3 \sigma_{\text{rad}}^2.$$

902 By Lemma 5 and a similar reasoning to Theorem 3 we have

$$(1 + 2\gamma_t \mu n) \mathbb{E} \left[\|x_{t+1} - x_*\|^2 \right] \leq \mathbb{E} \left[\|x_t - x_*\|^2 \right] + 2\gamma_t^3 \sigma_{\text{rad}}^2.$$

903 We may then use Lemma 9 to obtain that

$$\begin{aligned} \mathbb{E} \left[\|x_T - x_*\|^2 \right] &\leq \exp\left(-\frac{nT}{2(L_{\max}/\mu + n)}\right) \|x_0 - x_*\|^2 + \frac{356\sigma_{\text{rad}}^2}{\mu^3 n^2 T^2} \\ &= \mathcal{O}\left(\exp\left(-\frac{nT}{\kappa + 2n}\right) \|x_0 - x_*\|^2 + \frac{\sigma_{\text{rad}}^2}{\mu^3 n^2 T^2}\right). \end{aligned} \quad \blacksquare$$

904 J Extension: Importance resampling

905 Suppose that each f_i is L_i -smooth. Then the iteration complexities of both SGD and RR depend
 906 on L_{\max}/μ , where L_{\max} is the maximum smoothness constant among the smoothness constants
 907 L_1, L_2, \dots, L_n . The maximum smoothness constant can be arbitrarily worse than the average
 908 smoothness constant $\bar{L} = \frac{1}{n} \sum_{i=1}^n L_i$. This situation is in contrast to the complexity of gradient
 909 descent which depends on the smoothness constant L_f of $f = \frac{1}{n} \sum_{i=1}^n f_i$, for which we have
 910 $L_f \leq \bar{L}$. This is a problem commonly encountered with stochastic optimization methods and
 911 may cause significantly degraded performance in practical optimization tasks in comparison with
 912 deterministic methods [Tang et al., 2020].

Importance sampling is a common technique to improve the convergence of SGD (Algorithm 2):
 we sample function $(\bar{L}/L_i)f_i$ with probability p_i proportional to L_i , where $\bar{L} := \frac{1}{n} \sum_{i=1}^n L_i$. In that
 case, the SGD update is still unbiased since

$$\mathbb{E}_i \left[\frac{\bar{L}}{L_i} f_i \right] = \sum_{i=1}^n p_i \frac{\bar{L}}{L_i} f_i = f.$$

913 Moreover, the smoothness of function $(\bar{L}/L_i)f_i$ is \bar{L} for any i , so the guarantees would depend on
 914 \bar{L} instead of $\max_{i=1, \dots, n} L_i$. Importance sampling successfully improves the iteration complexity
 915 of SGD to depend on \bar{L} [Needell et al., 2016], and has been investigated in a wide variety of
 916 settings [Gower et al., 2020, Gorbunov et al., 2020].

917 Importance sampling is a neat technique but it relies heavily on the fact that we use *unbiased* sampling.
 918 How can we obtain a similar result if inside any permutation the sampling is biased? The answer
 919 requires us to think again as to what happens when we replace f_i with $(\bar{L}/L_i)f_i$. To make sure the
 920 problem remains the same, it is sufficient to have $(\bar{L}/L_i)f_i$ inside a permutation exactly L_i/\bar{L} times.
 921 And since L_i/\bar{L} is not necessarily integer, we should use $n_i = \lceil L_i/\bar{L} \rceil$ and solve

$$\min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^n \underbrace{\left(\frac{1}{n_i} f_i(x) + \dots + \frac{1}{n_i} f_i(x) \right)}_{n_i \text{ times}} + \psi(x), \quad (48)$$

where

$$N := n_1 + \dots + n_n = \left\lceil \frac{L_1}{\bar{L}} \right\rceil + \dots + \left\lceil \frac{L_n}{\bar{L}} \right\rceil.$$

922 Clearly, this problem is equivalent to the original formulation in 1. At the same time, we have
 923 improved all smoothness constants to \bar{L} . It might seem that that the new problem has more functions,
 924 but it turns out that the new number of functions satisfies $N \leq 2n$, so any related costs, such as
 925 longer loops or storing duplicates of the data, are negligible, as the next theorem shows.

926 **Theorem 10.** For every i , assume that each f_i is convex and L_i -smooth, and let ψ be μ -strongly
 927 convex. Then, the number of functions N in (48) satisfies $N \leq 2n$, and Algorithm 1 applied to
 928 problem (48) has the same complexity as (9) but proportional to \bar{L} rather than L_{\max} .

929 *Proof.* We show that $N \leq 2n$ as the rest of the theorem's claim trivially follows from Theorem 3.
 930 Firstly, note that for any number $a \in \mathbb{R}$ we have $\lceil a \rceil \leq a + 1$. Therefore,

$$N = \sum_{i=1}^n \left\lceil \frac{L_i}{\bar{L}} \right\rceil \leq \sum_{i=1}^n \left(\frac{L_i}{\bar{L}} + 1 \right) = n + \sum_{i=1}^n \frac{L_i}{\bar{L}} = 2n. \quad \blacksquare$$