

---

# Proximal and Federated Random Reshuffling

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

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

## 21 1 Introduction

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

30 The regularization framework is quite general and perhaps surprisingly it also allows us to consider  
31 methods for federated learning (FL)—a paradigm in which we aim at training model for a number of  
32 clients that do not want to reveal their data [Konečný et al., 2016, McMahan et al., 2017, Kairouz  
33 et al., 2019]. The training in FL usually happens on devices with only a small number of model  
34 updates being shared with a global host. To this end, Federated Averaging algorithm has emerged  
35 that performs Local SGD updates on the clients’ devices and periodically aggregates their average.  
36 Its analysis usually requires special techniques and deliberately constructed sequences hindering the  
37 research in this direction. We shall see, however, that the convergence of our FedRR follows from  
38 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  $\psi$  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  $\psi$  are  $\mu$ -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  $\nabla f_{i_k}$ .

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:    $x_{t+1} = \text{prox}_{\gamma_t n \psi}(x_t^n)$

---

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

## 92 2 Contributions

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

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

101 **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  $c_1, c_2 \in \mathbb{R}^d$ ,  
 102  $c_1 \neq c_2$ . 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  
 103 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  
 104  $\tilde{x}_2 = \text{prox}_{\gamma\psi}(x_1 - \gamma \nabla f_2(\tilde{x}_1))$ , then  $\tilde{x}_2 \neq \text{prox}_{2\gamma\psi}(x_0 - 2\gamma \nabla f(x_0))$ .

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

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

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

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

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

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

121 decrease it, this implies that the iterates converge to  $x_*$  in the limit. Moreover, we show in Section 4  
 122 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.

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

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

135 • **Nonconvex analysis.** In the nonconvex regime, and under suitable assumptions, we establish (see  
 136 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  
 137 an  $\mathcal{O}(\frac{1}{\varepsilon^3})$  convergence rate.

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

139 • **Extension 1: Decreasing stepsizes.** The convergence of RR is not always exact and depends on  
 140 the parameters of the objective. Similarly, if the shuffling radius  $\sigma_{\text{rad}}^2$  is positive, and we wish to find  
 141 an  $\varepsilon$ -approximate solution, the optimal choice of a fixed stepsize for ProxRR will depend on  $\varepsilon$ . This  
 142 deficiency can be fixed by using decreasing stepsizes in both vanilla RR [Ahn et al., 2020] and in  
 143 SGD [Stich, 2019]. We adopt the same technique to our setting. However, we depart from [Ahn et al.,  
 144 2020] by only adjusting the stepsize *once per epoch* rather than at every iteration, similarly to the  
 145 concurrent work of Tran et al. [2020] on RR with momentum. For details, see Section I.

146 • **Extension 2: Importance resampling for Proximal RR.** While importance sampling is a well  
 147 established technique for speeding up the convergence of SGD [Zhao and Zhang, 2015, Khaled and  
 148 Richtárik, 2020], no importance sampling variant of RR has been proposed nor analyzed. This is not  
 149 surprising since the key property of importance sampling in SGD—unbiasedness—does not hold for  
 150 RR. Our approach to equip ProxRR with importance sampling is via a reformulation of problem (1)  
 151 into a similar problem with a larger number of summands. In particular, for each  $i \in [n]$  we include  
 152  $n_i$  copies of the function  $\frac{1}{n_i} f_i$ , and then take average of all  $N = \sum_i n_i$  functions constructed this  
 153 way. The value of  $n_i$  depends on the “importance” of  $f_i$ , described below. We then apply ProxRR  
 154 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  
 155  $n_i = \lceil L_i / \bar{L} \rceil$ . It is easy to show that  $N \leq 2n$ , and hence our reformulation leads to at most a doubling  
 156 of the number of functions forming the finite sum. However, the overall complexity of ProxRR  
 157 applied to this reformulation will depend on  $\bar{L}$  instead of  $\max_i L_i$  (see Theorem 10), which can lead  
 158 to a significant improvement. For details of the construction and our complexity results, see Section J.

### 159 3 Preliminaries

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

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

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

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

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

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

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

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

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

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

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

## 186 4 Theory for strongly convex losses $f_1, \dots, f_n$

187 Our first theorem establishes a convergence rate for Algorithm 1 applied with a constant stepsize to  
 188 Problem (1) when each objective  $f_i$  is strongly convex. This assumption is commonly satisfied in  
 189 machine learning applications where each  $f_i$  represents a regularized loss on some data points, as in  
 190  $\ell_2$  regularized linear regression and  $\ell_2$  regularized logistic regression.

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

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

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

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

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

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

207 **Comparison with proximal SGD.** In order to compare (6) against the complexity of Proximal SGD  
 208 (Algorithm 2), we recall that Proximal SGD achieves  $\mathbb{E} \left[ \|x_K - x_*\|^2 \right] = \mathcal{O}(\varepsilon)$  if either  $f$  or  $\psi$  is  
 209  $\mu$ -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:      $x_{k+1} = \text{prox}_{\gamma_k \psi}(x_k - \gamma_k \nabla f_{i_k}(x_k))$
- 

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

214 By comparing  $K_{\text{SGD}}$  (given by (7)) and  $K_{\text{RR}}$  (given by (6)), we see that ProxRR has milder  
 215 dependence on  $\varepsilon$  than Proximal SGD. In particular, ProxRR converges faster whenever the target  
 216 accuracy  $\varepsilon$  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  
 217 better when we consider *proximal iteration complexity* (# of proximal operator access), in which case  
 218 the complexity of ProxRR (6) is reduced by a factor of  $n$  (because we take one proximal step every  $n$   
 219 iterations), while the proximal iteration complexity of Proximal SGD remains the same as (7). In this  
 220 case, ProxRR is better whenever the accuracy  $\varepsilon$  satisfies

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

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

## 223 5 Theory for strongly convex regularizer $\psi$

224 In Theorem 2, we assume that each  $f_i$  is  $\mu$ -strongly convex. This is motivated by the common practice  
 225 of using  $\ell_2$  regularization in machine learning. However, applying  $\ell_2$  regularization in every step  
 226 of Algorithm 1 can be expensive when the data are sparse and the iterates  $x_t^i$  are dense, because it  
 227 requires accessing each coordinate of  $x_t^i$  which can be much more expensive than computing sparse  
 228 gradients  $\nabla f_i(x_t^i)$ . Alternatively, we may instead choose to put the  $\ell_2$  regularization inside  $\psi$  and  
 229 only ask that  $\psi$  be strongly convex—this way, we can save a lot of time as we need to access each  
 230 coordinate of the dense iterates  $x_t^i$  only once per epoch rather than every iteration. Theorem 3 gives a  
 231 convergence guarantee in this setting.

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

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

235 Using Theorem 3 and choosing the stepsize as

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

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

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

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

## 240 6 FedRR: application of ProxRR to federated learning

241 Let us consider now the problem of minimizing the average of  $N = \sum_{m=1}^M N_m$  functions that are  
 242 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:    $x_{t+1} = \frac{1}{M} \sum_{m=1}^M x_{t,m}^n$ 

```

---

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

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

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

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

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

253 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  
 254  $x_1 = \dots = x_M$ , and  $(R + \psi_C)(x_1, \dots, x_M) = +\infty$  otherwise.

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

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

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

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

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

264 In other words, function  $f_i(\mathbf{x})$  includes  $i$ -th data sample from each device and contains at most  
 265 one loss from every device, while  $F_m(x)$  combines all data losses on device  $m$ . Note that the  
 266 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  
 267 by  $\nabla f_i(\mathbf{x}) = (\nabla f_{i1}(x_1)^\top, \dots, \nabla f_{iM}(x_M)^\top)^\top$ . Therefore, a stochastic gradient step that uses  
 268  $\nabla f_i(\mathbf{x})$  corresponds to updating all local models with the gradient of  $i$ -th data sample, without any  
 269 communication.

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

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

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

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

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

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

286 **Lemma 2.** The shuffling radius  $\sigma_{\text{rad}}^2$  of the reformulation (11) is upper bounded by

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

287 The lemma shows that the upper bound on  $\sigma_{\text{rad}}^2$  depends on the sum of local variances  $\sum_{m=1}^M \sigma_{m,*}^2$  as  
 288 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  
 289 literature on convergence of Local GD/SGD [Khaled et al., 2019, Woodworth et al., 2020, Yuan et al.,  
 290 2020]. We are now ready to present formal convergence results. For simplicity, we will consider  
 291 heterogeneous and homogeneous cases separately and assume that  $N_1 = \dots = N_M = n$ . To further  
 292 illustrate generality of our results, we will present the heterogeneous assuming strong convexity  $R$   
 293 and the homogeneous under strong convexity of functions  $f_{mi}$ .

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

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

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

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

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

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

302 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  $\nabla F(x) =$   
 303  $\frac{1}{N} \sum_{l=1}^M \nabla F_l(x)$  is the full gradient on all data.

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

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

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

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

311 **Homogeneous data.** For simplicity, in the homogeneous (i.e., i.i.d.) data case we provide guarantees  
 312 without the proximal operator. Since then we have  $F_1(x) = \dots = F_M(x)$ , for any  $m$  it holds  
 313  $\nabla F_m(x_*) = 0$ , 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

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

314 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.

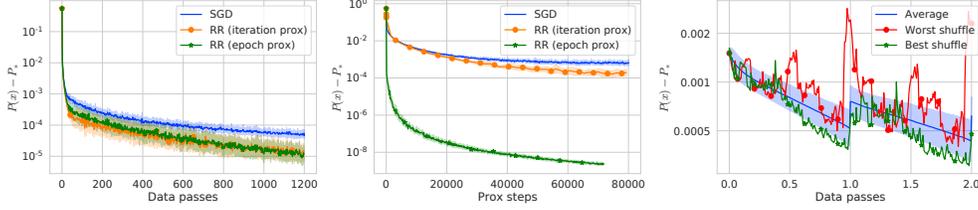


Figure 1: Experimental results for problem (12). The first two plots show with average and confidence intervals estimated on 20 random seeds and clearly demonstrate that one can save a lot of proximal operator computations with our method. The right plot shows the best/worst convergence of ProxSO over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

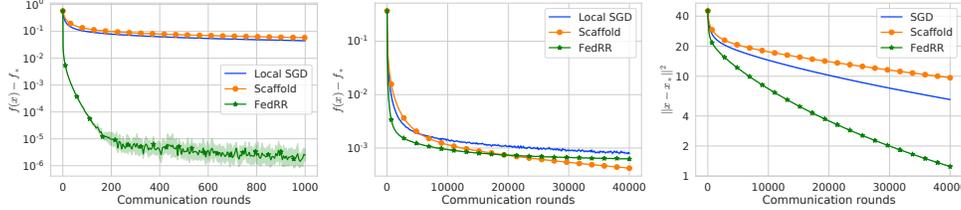


Figure 2: FedRR vs Local-SGD and Scaffold: i.i.d. data (left) and heterogeneous data (middle and right). We set  $\lambda_1 = 0$  and estimate the averages and standard deviations by running 10 random seeds for each method.

315 **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  
 316  $x_*$  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  
 317 and  $\mu$ -strongly convex, then the iterates of Algorithm 3 satisfy

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

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

## 320 7 Experiments<sup>1</sup>

321 **ProxRR vs SGD.** In Figure 1, we look at the logistic regression loss with the elastic net regularization,  
 322

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

323 where each  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as  $f_i(x) := -(b_i \log(h(a_i^\top x)) + (1 - b_i) \log(1 - h(a_i^\top x)))$ ,  
 324 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  
 325 sigmoid function, and  $\lambda_1, \lambda_2 \geq 0$  are parameters. We set minibatch sizes to 32 for all methods and  
 326 use theoretical stepsizes, without any tuning. We denote the heuristic version of RR that performs  
 327 proximal operator step after each iteration as ‘RR (iteration prox)’. From the experiments, we can see  
 328 that all methods behave more or less the same way. However, the algorithm that we propose needs  
 329 only a small fraction of proximal operator evaluations, which gives it a huge advantage whenever the  
 330 operator takes more time to compute than stochastic gradients.

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

<sup>1</sup>Our code is provided in the supplementary. More experimental details are in the appendix.

## 342 References

- 343 Kwangjun Ahn, Chulhee Yun, and Suvrit Sra. SGD with shuffling: optimal rates without component  
344 convexity and large epoch requirements. *arXiv preprint arXiv:2006.06946. Neural Information*  
345 *Processing Systems (NeurIPS) 2020*, 2020. (Cited on pages 2, 4, and 31)
- 346 Amir Beck. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics,  
347 Philadelphia, PA, 2017. doi: 10.1137/1.9781611974997. (Cited on page 5)
- 348 Pascal Bianchi, Walid Hachem, and Franck Iutzeler. A coordinate descent primal-dual algorithm and  
349 application to distributed asynchronous optimization. *IEEE Transactions on Automatic Control*, 61  
350 (10):2947–2957, 2015. (Cited on page 7)
- 351 Antoine Bordes, Léon Bottou, and Patrick Gallinari. SGD-QN: Careful quasi-Newton stochastic  
352 gradient descent. 2009. (Cited on page 2)
- 353 Léon Bottou. Curiously fast convergence of some stochastic gradient descent algorithms. Unpublished  
354 open problem offered to the attendance of the SLDS 2009 conference, 2009. URL [http://leon.](http://leon.bottou.org/papers/bottou-slds-open-problem-2009)  
355 [bottou.org/papers/bottou-slds-open-problem-2009](http://leon.bottou.org/papers/bottou-slds-open-problem-2009). (Cited on page 2)
- 356 Léon Bottou. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade*, pages  
357 421–436. Springer, 2012. (Cited on page 2)
- 358 Kristian Bredies, Karl Kunisch, and Thomas Pock. Total generalized variation. *SIAM Journal on*  
359 *Imaging Sciences*, 3(3):492–526, 2010. (Cited on page 2)
- 360 Gong Chen and Marc Teboulle. Convergence Analysis of a Proximal-Like Minimization Algorithm  
361 Using Bregman Functions. *SIAM Journal on Optimization*, 3(3):538–543, 1993. doi: 10.1137/  
362 0803026. (Cited on page 19)
- 363 John Duchi and Yoram Singer. Efficient online and batch learning using forward backward splitting.  
364 *Journal of Machine Learning Research*, 10(Dec):2899–2934, 2009. (Cited on page 3)
- 365 Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. A Unified Theory of SGD: Variance Reduction,  
366 Sampling, Quantization and Coordinate Descent. volume 108 of *Proceedings of Machine Learning*  
367 *Research*, pages 680–690, Online, 26–28 Aug 2020. PMLR. (Cited on pages 2, 3, 18, and 34)
- 368 Robert M. Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik.  
369 SGD: General Analysis and Improved Rates. In Kamalika Chaudhuri and Ruslan Salakhutdinov,  
370 editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of  
371 *Proceedings of Machine Learning Research*, pages 5200–5209, Long Beach, California, USA,  
372 09–15 Jun 2019. PMLR. (Cited on page 6)
- 373 Robert M. Gower, Peter Richtárik, and Francis Bach. Stochastic quasi-gradient methods: variance  
374 reduction via Jacobian sketching. *Mathematical Programming*, pages 1–58, 2020. ISSN 0025-5610.  
375 doi: 10.1007/s10107-020-01506-0. (Cited on page 34)
- 376 Mert Gürbüzbalaban, Asuman Özdağlar, and Pablo A. Parrilo. Why random reshuffling beats  
377 stochastic gradient descent. *Mathematical Programming*, Oct 2019. ISSN 1436-4646. doi:  
378 10.1007/s10107-019-01440-w. (Cited on page 2)
- 379 Jeff Haochen and Suvrit Sra. Random Shuffling Beats SGD after Finite Epochs. In Kamalika  
380 Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on*  
381 *Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2624–2633,  
382 Long Beach, California, USA, 09–15 Jun 2019. PMLR. (Cited on page 2)
- 383 Peter Kairouz et al. Advances and open problems in federated learning. *arXiv preprint*  
384 *arXiv:1912.04977*, 2019. (Cited on pages 1 and 4)
- 385 Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear Convergence of Gradient and Proximal-  
386 Gradient Methods Under the Polyak-Łojasiewicz Condition. In *European Conference on Machine*  
387 *Learning and Knowledge Discovery in Databases - Volume 9851*, ECML PKDD 2016, page  
388 795–811, 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 29)
- 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 Lojasiewicz. 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: 10.1561/2400000003. (Cited on pages 16 and 30)
- 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 *ritms*. 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)

467 **Checklist**

- 468 1. For all authors...
- 469 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
470 contributions and scope? [Yes]
- 471 (b) Did you describe the limitations of your work? [Yes]
- 472 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 473 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
474 them? [Yes]
- 475 2. If you are including theoretical results...
- 476 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 477 (b) Did you include complete proofs of all theoretical results? [Yes]
- 478 3. If you ran experiments...
- 479 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
480 mental results (either in the supplemental material or as a URL)? [Yes]
- 481 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
482 were chosen)? [Yes]
- 483 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
484 ments multiple times)? [Yes]
- 485 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
486 of GPUs, internal cluster, or cloud provider)? [Yes]
- 487 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 488 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 489 (b) Did you mention the license of the assets? [N/A]
- 490 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 491
- 492 (d) Did you discuss whether and how consent was obtained from people whose data you're  
493 using/curating? [N/A]
- 494 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
495 information or offensive content? [N/A]
- 496 5. If you used crowdsourcing or conducted research with human subjects...
- 497 (a) Did you include the full text of instructions given to participants and screenshots, if  
498 applicable? [N/A]
- 499 (b) Did you describe any potential participant risks, with links to Institutional Review  
500 Board (IRB) approvals, if applicable? [N/A]
- 501 (c) Did you include the estimated hourly wage paid to participants and the total amount  
502 spent on participant compensation? [N/A]