

---

# On the Convergence to a Global Solution of Shuffling-Type Gradient Algorithms

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Stochastic gradient descent (SGD) algorithm is the method of choice in many  
2 machine learning tasks thanks to its scalability and efficiency in dealing with  
3 large-scale problems. In this paper, we focus on the shuffling version of SGD  
4 which matches the mainstream practical heuristics. We show the convergence  
5 to a global solution of shuffling SGD for a class of non-convex functions under  
6 over-parameterized settings. Our analysis employs more relaxed non-convex  
7 assumptions than previous literature. Nevertheless, we maintain the desired computational  
8 complexity as shuffling SGD has achieved in the general convex setting.

## 9 1 Introduction

10 In the last decade, neural network-based models have shown great success in many machine learning  
11 applications such as natural language processing [Collobert and Weston, 2008, Goldberg et al., 2018],  
12 computer vision and pattern recognition [Goodfellow et al., 2014, He and Sun, 2015]. The training  
13 task of many learning models boils down to the following finite-sum minimization problem:

$$\min_{w \in \mathbb{R}^d} \left\{ F(w) := \frac{1}{n} \sum_{i=1}^n f(w; i) \right\}, \quad (1)$$

14 where  $f(\cdot; i) : \mathbb{R}^d \rightarrow \mathbb{R}$  is smooth and possibly non-convex for  $i \in [n] := \{1, \dots, n\}$ . Solving the  
15 empirical risk minimization (1) had been a difficult task for a long time due to the non-convexity  
16 and the complicated learning models. Later progress with stochastic gradient descent (SGD) and its  
17 variants [Robbins and Monro, 1951, Duchi et al., 2011, Kingma and Ba, 2014] have shown great  
18 performance in training deep neural networks. These stochastic first-order methods are favorable  
19 thanks to its scalability and efficiency in dealing with large-scale problems. At each iteration SGD  
20 samples an index  $i$  uniformly from the set  $\{1, \dots, n\}$ , and uses the individual gradient  $\nabla f(\cdot; i)$  to  
21 update the weight.

22 While there has been much attention on the theoretical aspect of the traditional i.i.d. (independently  
23 identically distributed) version of SGD [Nemirovski et al., 2009, Ghadimi and Lan, 2013, Bottou  
24 et al., 2018], practical heuristics often use without-replacement data sampling schemes. Also known  
25 as shuffling sampling schemes, these methods generate some random or deterministic permutations  
26 of the index set  $\{1, 2, \dots, n\}$  and apply gradient updates using these permutation orders. Intuitively,  
27 a collection of such  $n$  individual updates is a pass over all the data, or an epoch. One may choose to  
28 create a new random permutation at the beginning of each epoch (in Random Reshuffling scheme) or  
29 use a random permutation for every epoch (in Single Shuffling scheme). Alternatively, one may use a  
30 Incremental Gradient scheme with a fixed deterministic order of indices. In this paper, we use the  
31 term unified shuffling SGD for SGD method using *any* data permutations, which includes the three  
32 special schemes described above.

33 Although shuffling sampling schemes usually show a better empirical performance than SGD [Bottou,  
 34 2009], the theoretical guarantees for these schemes are often more limited than vanilla SGD version,  
 35 due to the lack of statistical independence. Recent works have shown improvement in computational  
 36 complexity for shuffling schemes over SGD in various settings [Gürbüzbalaban et al., 2019, Haochen  
 37 and Sra, 2019, Safran and Shamir, 2020, Nagaraj et al., 2019, Nguyen et al., 2021, Mishchenko et al.,  
 38 2020, Ahn et al., 2020]. In particular, in a general non-convex setting, shuffling sampling schemes  
 39 improve the computational complexity in terms of  $\hat{\epsilon}$  for SGD from  $\mathcal{O}(\sigma^2/\hat{\epsilon}^2)$  to  $\mathcal{O}(n\sigma/\hat{\epsilon}^{3/2})$ ,  
 40 where  $\sigma$  is the bounded variance constant [Ghadimi and Lan, 2013, Nguyen et al., 2021, Mishchenko  
 41 et al., 2020]<sup>1</sup>. We summarize the detailed literature for multiple settings later in Table 1.

42 While global convergence is a desirable property for neural network training, the non-convexity  
 43 landscape of complex learning models leads to difficulties in finding the global minimizer. In  
 44 addition, there is little to no work studying the convergence to a global solution of shuffling-type  
 45 SGD algorithms for a general non-convex setting. The closest line of research investigates the Polyak-  
 46 Lojasiewicz (PL) condition (a generalization of strong-convexity), which demonstrates similar  
 47 convergence rates as the strongly convex rates for shuffling SGD methods [Haochen and Sra, 2019,  
 48 Ahn et al., 2020, Nguyen et al., 2021]. In another direction, Gower et al. [2021] and Khaled and  
 49 Richtárik [2020] investigates the global convergence for some class of non-convex functions, however  
 50 for vanilla SGD method. Beznosikov and Takáč [2021] investigate a random shuffle version of  
 51 variance reduction methods (e.g. SARAH algorithm Nguyen et al. [2017]), but this approach only  
 52 can show convergence to stationary points. With a target on shuffling SGD methods and specific  
 53 learning architectures, we come up with the central question of this paper:

54 *How can we establish the convergence to global solution for a class of non-convex functions using*  
 55 *shuffling-type SGD algorithms? Can we exploit the structure of neural networks to achieve this goal?*

56 We answer this question affirmatively, and our contributions are summarized below. **Contributions.**

- 57 • We investigate a new framework for the convergence of a shuffling-type gradient algorithm  
 58 to a global solution. We consider a relaxed set of assumptions and discuss their relations  
 59 with previous settings. We show that our average-PL inequality (Assumption 3) holds for a  
 60 wide range of neural networks equipped with squared loss function.
- 61 • Our analysis generalizes the class function called star- $M$ -smooth-convex. This class contains  
 62 non-convex functions and is more general than the class of star-convex smooth functions  
 63 with respect to the minimizer (in the over-parameterized settings). In addition, our analysis  
 64 does not use any bounded gradient or bounded weight assumptions.
- 65 • We show the total complexity of  $\mathcal{O}(\frac{n}{\hat{\epsilon}^{3/2}})$  for a class of non-convex functions to reach an  
 66  $\hat{\epsilon}$ -accurate global solution. This result matches the same gradient complexity to a stationary  
 67 point for unified shuffling methods in non-convex settings, however, we are able to show the  
 68 convergence to a global minimizer.

## 69 1.1 Related Work

70 In recent years, there have been different approaches to investigate the global convergence for machine  
 71 learning optimization. This includes a popular line of research that studies some specific neural  
 72 networks and utilizes their architectures. The most early works show the global convergence of  
 73 Gradient Descent (GD) for simple linear networks and two-layer networks [Brutzkus et al., 2018,  
 74 Soudry et al., 2018, Arora et al., 2019, Du et al., 2019b]. These results are further extended to deep  
 75 learning architectures [Allen-Zhu et al., 2019, Du et al., 2019a, Zou and Gu, 2019]. This line of  
 76 research continues with Stochastic Gradient Descent (SGD) algorithm, which proves the global  
 77 convergence of SGD for deep neural networks for some probability depending on the initialization  
 78 process and the number of input data [Brutzkus et al., 2018, Allen-Zhu et al., 2019, Zou et al.,  
 79 2018, Zou and Gu, 2019]. The common theme that appeared in most of these references is the over-  
 80 parameterized setting, which means that the number of parameters in the network are excessively  
 81 large [Brutzkus et al., 2018, Soudry et al., 2018, Allen-Zhu et al., 2019, Du et al., 2019a, Zou and Gu,  
 82 2019]. This fact is closely related to our setting, and we will discuss it throughout our paper.

---

<sup>1</sup>The computational complexity is the number of (individual) gradient computations needed to reach an  
 $\hat{\epsilon}$ -accurate stationary point (i.e. a point  $\hat{w} \in \mathbb{R}^d$  that satisfies  $\|\nabla F(\hat{w})\|^2 \leq \hat{\epsilon}$ .)

83 **Polyak-Lojasiewicz (PL) condition and related assumptions.** An alternative approach is to  
 84 investigate some conditions on the optimization problem that may guarantee global convergence. A  
 85 popular assumption is the Polyak-Lojasiewicz (PL) inequality, a generalization of strong-convexity  
 86 [Polyak, 1964, Karimi et al., 2016, Nesterov and Polyak, 2006]. Using this PL assumption, it can  
 87 be shown that (stochastic) gradient descent achieves the same theoretical rate as in the strongly  
 88 convex setting (i.e linear convergence for GD and sublinear convergence for SGD) [Karimi et al.,  
 89 2016, De et al., 2017, Gower et al., 2021]. Recent works demonstrate similar results for shuffling  
 90 type SGD [Haochen and Sra, 2019, Ahn et al., 2020, Nguyen et al., 2021], both for unified and  
 91 randomized shuffling schemes. On the other hand, [Schmidt and Roux, 2013, Vaswani et al., 2019]  
 92 propose to use a new assumption called the Strong Growth Condition (SGC) that controls the rates  
 93 at which the stochastic gradients decay comparing to the full gradient. This condition implies that  
 94 the stochastic gradients and their variances converge to zero at the optimum solution [Schmidt and  
 95 Roux, 2013, Vaswani et al., 2019]. While the PL condition for  $F$  implies that every stationary point  
 96 of  $F$  is also a global solution, the SGC implies that such a point is also a stationary point of every  
 97 individual function. However, complicated models as deep feed-forward neural networks generally  
 98 have non-optimal stationary points [Karimi et al., 2016]. Thus, these assumptions are somewhat  
 99 strong for non-convex settings.

100 Although there are plenty of works investigating the PL condition for the objective function  $F$   
 101 [De et al., 2017, Vaswani et al., 2019, Gower et al., 2021], not many materials devoted to study  
 102 the PL inequality for the individual functions  $f(\cdot; i)$ . A recent work [Sankararaman et al., 2020]  
 103 analyzes SGD with the specific notion of gradient confusion for over-parameterized settings where the  
 104 individual functions satisfy PL condition. They show that the neighborhood where SGD converges  
 105 linearly depends on the level of gradient confusion (i.e. how much the individual gradients are  
 106 negatively correlated). Taking a different approach, we investigate the PL property for individual  
 107 functions and further show that our condition holds for a general class of neural networks with  
 108 quadratic loss.

109 **Over-paramaterized settings for neural networks.** Most of the modern learning architectures  
 110 contain deep and large networks, where the number of parameters are often far more than the number  
 111 of input data. This leads to the fact that the objective loss function is trained closer and closer to zero.  
 112 Understandably, in such settings all the individual functions  $f(\cdot; i)$  are minimized simultaneously at  
 113 0 and they share a common minimizer. This condition is called the interpolation property (see e.g.  
 114 [Schmidt and Roux, 2013, Ma et al., 2018, Meng et al., 2020, Loizou et al., 2021]) and is studied  
 115 well in the literature (see e.g. [Zhou et al., 2019, Gower et al., 2021]). For a comparison, functions  
 116 satisfying the strong growth condition necessarily satisfy the interpolation property. This property  
 117 implies zero variance of individual gradients at the global minimizer, which allows good behavior  
 118 for SGD near the solution. In this work, we slightly change this assumption which requires a small  
 119 variance up to some level of the threshold  $\varepsilon$ . Note that when letting  $\varepsilon \rightarrow 0$ , our assumption exactly  
 120 recovers the interpolation property.

121 **Star-convexity and related conditions.** There have been many attentions to a class of structured  
 122 non-convex functions called star-convex [Nesterov and Polyak, 2006, Lee and Valiant, 2016, Bjorck  
 123 et al., 2021]. Star-convexity can be understood as convexity between an arbitrary point  $w$  and the  
 124 global minimizer  $w_*$ . The name star-convex comes from the fact that each sublevel set is star-shaped  
 125 [Nesterov and Polyak, 2006, Lee and Valiant, 2016]. Zhou et al. [2019] shows that if SGD follows a  
 126 star-convex path and there exists a common global minimizer for all component functions, then SGD  
 127 converges to a global minimum.

128 In recent progress, Hinder et al. [2020] considers the class of quasar-convex functions, which further  
 129 generalizes star-convexity. This property was introduced originally in [Hardt et al., 2018] under  
 130 the name ‘weakly quasi-convex’, and investigated recently in literature [Hinder et al., 2020, Jin,  
 131 2020, Gower et al., 2021]. This class uses a parameter  $\zeta \in (0, 1]$  to control the non-convexity of the  
 132 function, where  $\zeta = 1$  yields the star-convexity and  $\zeta$  approaches 0 indicates more non-convexity  
 133 [Hinder et al., 2020]. Intuitively, quasar-convex functions are unimodal on all lines that pass through  
 134 a global minimizer. Gower et al. [2021] investigates the performance of SGD for smooth and quasar-  
 135 convex functions using an expected residual assumption (which is comparable to the interpolation  
 136 property). They show a convergence rate of  $\mathcal{O}(1/\sqrt{K})$  for i.i.d. sampling SGD with the number  
 137 of total iterations  $K$ , which translates to the computational complexity of  $\mathcal{O}(1/\varepsilon^2)$ . To the best  
 138 of our knowledge, this paper is the first work studying the relaxation of star-convexity and global  
 139 convergence for SGD with shuffling sample schemes, not for the i.i.d. version.

140 **2 Theoretical Setting**

141 We first present the shuffling-type gradient algorithm below. Our convergence results hold for any  
 142 permutation of the training data  $\{1, 2, \dots, n\}$ , including deterministic and random ones. Thus, our  
 143 theoretical framework is general and applicable for any shuffling strategy, including Incremental  
 144 Gradient, Single Shuffling, and Random Reshuffling.

---

**Algorithm 1** (Shuffling-Type Gradient Algorithm for Solving (1))

---

```

1: Initialization: Choose an initial point  $\tilde{w}_0 \in \text{dom}(F)$ .
2: for  $t = 1, 2, \dots, T$  do
3:   Set  $w_0^{(t)} := \tilde{w}_{t-1}$ ;
4:   Generate any permutation  $\pi^{(t)}$  of  $[n]$  (either deterministic or random);
5:   for  $i = 1, \dots, n$  do
6:     Update  $w_i^{(t)} := w_{i-1}^{(t)} - \eta_i^{(t)} \nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))$ ;
7:   end for
8:   Set  $\tilde{w}_t := w_n^{(t)}$ ;
9: end for

```

---

145 We further specify the choice of learning rate  $\eta_i^{(t)}$  in the detailed analysis. Now we proceed to  
 146 describe the set of assumptions used in our paper.

147 **Assumption 1.** Suppose that  $f_i^* := \min_{w \in \mathbb{R}^d} f(w; i) > -\infty$ ,  $i \in \{1, \dots, n\}$ .

148 **Assumption 2.** Suppose that  $f(\cdot; i)$  is  $L$ -smooth for all  $i \in \{1, \dots, n\}$ , i.e. there exists a constant  
 149  $L \in (0, +\infty)$  such that:

$$\|\nabla f(w; i) - \nabla f(w'; i)\| \leq L\|w - w'\|, \quad \forall w, w' \in \mathbb{R}^d. \quad (2)$$

150 Assumption 1 is required in any algorithm to guarantee the well-definedness of (1). In most applica-  
 151 tions, the component losses are bounded from below. By Assumption 2, the objective function  $F$  is  
 152 also  $L$ -smooth. This Lipschitz smoothness Assumption is widely used for gradient-type methods. In  
 153 addition, we denote the minimum value of the objective function  $F_* = \min_{w \in \mathbb{R}^d} F(w)$ . It is worthy  
 154 to note the following relationship between  $F_*$  and the component minimum values:

$$F_* = \min_{w \in \mathbb{R}^d} F(w) = \frac{1}{n} \min_{w \in \mathbb{R}^d} \left( \sum_{i=1}^n f(w; i) \right) \geq \frac{1}{n} \sum_{i=1}^n \min_{w \in \mathbb{R}^d} (f(w; i)) = \frac{1}{n} \sum_{i=1}^n f_i^*. \quad (3)$$

155 We are interested in the case where the set of minimizers of  $F$  is not empty. The equality  $F_* =$   
 156  $\frac{1}{n} \sum_{i=1}^n f_i^*$  attains if and only if a minimizer of  $F$  is also the common minimizer for all component  
 157 functions. This condition implies that the variance of individual functions is 0 at the common  
 158 minimizer.

159 **2.1 PL Condition for Component Functions**

160 Now we are ready to discuss the Polyak-Lojasiewicz condition as follows.

161 **Definition 1** (Polyak-Lojasiewicz condition). We say that  $f$  satisfies Polyak-Lojasiewicz (PL) in-  
 162 equality for some constant  $\mu > 0$  if

$$\|\nabla f(w)\|^2 \geq 2\mu[f(w) - f^*], \quad \forall w \in \mathbb{R}^d, \quad (4)$$

163 where  $f^* := \min_{w \in \mathbb{R}^d} f(w)$ .

164 The PL condition for the objective function  $F$  is sufficient to show a global convergence for (stochas-  
 165 tic) gradient descent [Karimi et al., 2016, Nesterov and Polyak, 2006, Polyak, 1964]. It is well known  
 166 that a function satisfying the PL condition is not necessarily convex [Karimi et al., 2016]. However,  
 167 this assumption on  $F$  is somewhat strong because it implies that every stationary point of  $F$  is also a  
 168 global minimizer. Our goal is to consider a class of non-convex function which is more relaxed than  
 169 the PL condition on  $F$ , while still having the good global convergence properties. In this paper, we  
 170 formulated an assumption called ‘‘average PL inequality’’, specifically for the finite sum setting:

171 **Assumption 3.** Suppose that  $f(\cdot; i)$  satisfies average PL inequality for some constant  $\mu > 0$  such  
 172 that

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f(w; i)\|^2 \geq 2\mu \frac{1}{n} \sum_{i=1}^n [f(w; i) - f_i^*], \quad \forall w \in \mathbb{R}^d. \quad (5)$$

173 where  $f_i^* := \min_{w \in \mathbb{R}^d} f(w; i)$ .

174 **Comparisons.** Assumption 3 is weaker than assuming the PL inequality for every component function  
 175  $f(\cdot; i)$ . In general setting, Assumption 3 is not comparable to assuming the PL inequality for  $F$ .  
 176 Formally, if  $F$  satisfies PL the condition for some parameter  $\tau > 0$ , then we have:

$$2\tau[F(w) - F_*] \leq \|\nabla F(w)\|^2 \leq \frac{1}{n} \sum_{i=1}^n \|\nabla f(w; i)\|^2. \quad (6)$$

177 However, by equation (3) we have that  $[F(w) - F_*] \leq \frac{1}{n} \sum_{i=1}^n [f(w; i) - f_i^*]$ . Therefore, the PL  
 178 inequality for each function  $f(\cdot; i)$ , cannot directly imply the PL condition on  $F$  and vice versa.

179 In the interpolated setting where there is a common minimizer for all component function  $f(\cdot; i)$ , it  
 180 can be shown that the PL condition on  $F$  is stronger than our average PL assumption:

$$2\tau \frac{1}{n} \sum_{i=1}^n [f(w; i) - f_i^*] = 2\tau[F(w) - F_*] \leq \|\nabla F(w)\|^2 \leq \frac{1}{n} \sum_{i=1}^n \|\nabla f(w; i)\|^2.$$

181 On the other hand, our assumption cannot imply the PL inequality on  $F$  unless we impose a strong  
 182 relationship that upper bound the sum of individual squared gradients  $\frac{1}{n} \sum_{i=1}^n \|\nabla f(w; i)\|^2$  in terms  
 183 of the full squared gradient  $\|\nabla F(w)\|^2$ , for every  $w \in \mathbb{R}^d$ . For these reasons, the average PL  
 184 Assumption 3 is arguably more reasonable than assuming the PL inequality for the objective function  
 185  $F$ . Moreover, we show that Assumption 3 holds for a general class of neural networks with a final  
 186 bias layer and squared loss function. We have the following theorem.

187 **Theorem 1.** Let  $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$  is a training data set where  $x^{(i)} \in \mathbb{R}^m$  is the input data and  
 188  $y^{(i)} \in \mathbb{R}^c$  is the output data for  $i = 1, \dots, n$ . We consider an architecture  $h(w; i)$  with  $w$  be the  
 189 vectorized weight and  $h$  consists of a final bias layer  $b$ :

$$h(w; i) = W^T z(\theta; i) + b,$$

190 where  $w = \text{vec}(\{\theta, W, b\})$  and  $z(\theta; i)$  are some inner architectures, which can be chosen arbitrarily.  
 191 Next, we consider the squared loss  $f(w; i) = \frac{1}{2} \|h(w; i) - y^{(i)}\|^2$ . Then

$$\|\nabla f(w; i)\|^2 \geq 2[f(w; i) - f_i^*], \quad \forall w \in \mathbb{R}^d, \quad (7)$$

192 where  $f_i^* := \min_{w \in \mathbb{R}^d} f(w; i)$ .

193 Therefore, for this application, Assumption 3 holds with  $\mu = 1$ .

## 194 2.2 Small Variance at Global Solutions

195 In this section, we change the interpolation property in previous literature [Ma et al., 2018, Meng  
 196 et al., 2020, Loizou et al., 2021] by a small threshold. For any global solution  $w_*$  of  $F$ , let us define

$$\sigma_*^2 := \inf_{w_* \in \mathcal{W}_*} \left( \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_*; i)\|^2 \right). \quad (8)$$

197 We can show that when there is a common minimizer for all component functions (i.e. when the  
 198 equality  $F_* = \frac{1}{n} \sum_{i=1}^n f_i^*$  holds), the best variance  $\sigma_*^2$  is 0. It is sufficient for our Theorem to impose  
 199 a  $\mathcal{O}(\varepsilon)$ -level upper bound on the variance  $\sigma_*^2$ :

200 **Assumption 4.** Suppose that the best variance at  $w_*$  is small, that is, for  $\varepsilon > 0$

$$\sigma_*^2 \leq P\varepsilon, \quad (9)$$

201 for some  $P > 0$ .

202 It is important to note that in current literature, Assumption 4 alone (or, assuming  $\sigma_*^2 = 0$  alone)  
 203 is not sufficient enough to guarantee a global convergence property for SGD. Typically, some  
 204 other conditions on the good landscape of the loss function are needed to complement the over-  
 205 parameterized setting. Thus, we have motivation to introduce our next assumption.

206 **2.3 Generalized Star-Smooth-Convex Condition for Shuffling Type Algorithm**

207 We introduce the definition of star-smooth-convex function as follows.

208 **Definition 2.** *The function  $g$  is star- $M$ -smooth-convex with respect to a reference point  $\hat{w} \in \mathbb{R}^d$  if*

$$\|\nabla g(w) - \nabla g(\hat{w})\|^2 \leq M \langle \nabla g(w) - \nabla g(\hat{w}), w - \hat{w} \rangle, \quad \forall w \in \mathbb{R}^d. \quad (10)$$

209 It is well known that when  $g$  is  $L$ -smooth and convex [Nesterov, 2004], we have the following general  
210 inequality for every  $w, w' \in \mathbb{R}^d$ :

$$\|\nabla g(w) - \nabla g(w')\|^2 \leq L \langle \nabla g(w) - \nabla g(w'), w - w' \rangle \quad (11)$$

211 Our class of star-smooth-convex function requires a similar inequality to hold only for the special  
212 point  $w' = \hat{w}$ . Interestingly, this is related to a class of star-convex functions, which satisfies the  
213 convex inequality for the minimizer  $\hat{w}$ :

$$(\text{star-convexity w.r.t } \hat{w}) \quad g(w) - g(\hat{w}) \leq \langle \nabla g(w), w - \hat{w} \rangle, \quad \forall w \in \mathbb{R}^d, \quad (12)$$

214 This class of functions contains non-convex objective losses and is well studied in the literature (see  
215 e.g. [Zhou et al., 2019]). Our Lemma 2 in the Appendix shows that the class of star-smooth-convex  
216 function is broader than the class of  $L$ -smooth and star-convex functions. Therefore, our problem of  
217 interest is non-convex in general.

218 For the analysis of shuffling type algorithm in this paper, we consider the general assumption called  
219 the *generalized star-smooth-convex condition for shuffling algorithms*:

220 **Assumption 5.** *Using Algorithm 1, let us assume that there exist some constants  $M > 0$  and  $N > 0$   
221 such that at each epoch  $t = 1, \dots, T$ , we have for  $i = 1, \dots, n$ :*

$$\begin{aligned} \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i)) - \nabla f(w_*; \pi^{(t)}(i))\|^2 &\leq M \langle \nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i)) - \nabla f(w_*; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_* \rangle \\ &\quad + N \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2, \end{aligned} \quad (13)$$

222 where  $w_*$  is a global solution of  $F$ .

223 We note that when the individual function  $f(\cdot; i)$  is star- $M$ -smooth-convex with respect to  $w_*$  for  
224 every  $i = 1, \dots, n$ , Assumption 5 holds for the case  $N = 0$ . Assumption 5 is more flexible than the  
225 one in (10) because the right-hand side term  $\langle \nabla f(w; i) - \nabla f(w_*; i), w - w_* \rangle$  could be negative for  
226 some  $w \in \mathbb{R}^d$ . An additional term  $N \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2$  for some constant  $N > 0$  will allow  
227 for extra flexibility in our setting. Note that we do not impose any assumptions on bounded weights  
228 or bounded gradients. Therefore, the term  $\frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2$  cannot be uniformly bounded by  
229 any universal constant.

230 **3 New Framework for Convergence to a Global Solution**

231 In this section, we present our theoretical results. Our Lemma 1 first provides a recursion to bound  
232 the squared distance term  $\|\tilde{w}_t - w_*\|^2$ :

233 **Lemma 1.** *Assume that Assumptions 1, 2, 3, and 5 hold. Let  $\{\tilde{w}_t\}_{t=1}^T$  be the sequence generated by  
234 Algorithm 1 with  $0 < \eta_t \leq \min\{\frac{n}{2M}, \frac{1}{2L}\}$ . For every  $\gamma > 0$  we have*

$$\|\tilde{w}_t - w_*\|^2 \leq (1 + C_1 \eta_t^3) \|\tilde{w}_{t-1} - w_*\|^2 + C_2 \eta_t \sigma_*^2 - C_3 \eta_t [F(\tilde{w}_{t-1}) - F_*]. \quad (14)$$

235 where  $w_*$  is a global solution of  $F$ ,  $F_* = \min_{w \in \mathbb{R}^d} F(w)$ , and

$$\begin{cases} C_1 = \frac{8L^2}{3} + \frac{14NL^2}{M} + \frac{4\gamma L^4}{6M}, \\ C_2 = \frac{2}{M} + 1 + \frac{5}{6L^2} + \frac{8N}{3ML^2} + \frac{5\gamma}{12M}, \\ C_3 = \frac{\gamma}{\gamma+1} \frac{\mu}{M}. \end{cases} \quad (15)$$

236 Rearranging the results of Lemma 1, we have

$$F(\tilde{w}_{t-1}) - F_* \leq \frac{1}{C_3} \left( \frac{1}{\eta_t} + C_1 \eta_t^2 \right) \|\tilde{w}_{t-1} - w_*\|^2 - \frac{1}{C_3 \eta_t} \|\tilde{w}_t - w_*\|^2 + \frac{C_2}{C_3} \sigma_*^2. \quad (16)$$

237 Therefore, with an appropriate choice of learning rate that guarantee  $(1/\eta_t + C_1 \eta_t^2) \leq 1/\eta_{t-1}$ , we  
238 can unroll the recursion from Lemma 1. Thus we have our main result in the next Theorem.

239 **Theorem 2.** Assume that Assumptions 1, 2, 3, and 5 hold. Let  $\{\tilde{w}_t\}_{t=1}^T$  be the sequence generated  
 240 by Algorithm 1 with the learning rate  $\eta_t^{(t)} = \frac{\eta_t}{n}$  where  $0 < \eta_t \leq \min\{\frac{n}{2M}, \frac{1}{2L}\}$ . Let the number of  
 241 iterations  $T = \frac{\lambda}{\varepsilon^{3/2}}$  for some  $\lambda > 0$  and  $\varepsilon > 0$ . Constants  $C_1$ ,  $C_2$ , and  $C_3$  are defined in (15) for any  
 242  $\gamma > 0$ . We further define  $K = 1 + C_1 D^3 \varepsilon^{3/2}$  and specify the learning rate  $\eta_t = K \eta_{t-1} = K^t \eta_0$   
 243 and  $\eta_0 = \frac{D\sqrt{\varepsilon}}{K \exp(\lambda C_1 D^3)}$  such that  $\frac{D\sqrt{\varepsilon}}{K} \leq \min\{\frac{n}{2M}, \frac{1}{2L}\}$  for some constant  $D > 0$ . Then we have

$$\frac{1}{T} \sum_{t=1}^T [F(\tilde{w}_{t-1}) - F_*] \leq \frac{K \exp(\lambda C_1 D^3)}{C_3 D \lambda} \|\tilde{w}_0 - w_*\|^2 \cdot \varepsilon + \frac{C_2}{C_3} \sigma_*^2, \quad (17)$$

244 where  $F_* = \min_{w \in \mathbb{R}^d} F(w)$  and  $\sigma_*^2$  is defined in (8).

245 Our analysis holds for arbitrarily constant values of the parameters  $\gamma$ ,  $\lambda$  and  $D$ . In addition, we show  
 246 our current analysis for every shuffling scheme. An interesting research question arises: whether the  
 247 convergence results can be improved if one chooses to analyze a randomized shuffling scheme in this  
 248 framework. However, we leave that question to future works.

249 Using Assumption 4, we can show the total complexity of Algorithm 1 for our setting.

250 **Corollary 1.** Suppose that the conditions in Theorem 2 and Assumption 4 hold. Choose  $C_1 D \lambda = 1$   
 251 and  $\varepsilon = \hat{\varepsilon}/G$  such that  $0 < \hat{\varepsilon} \leq G$  with the constants

$$G = \frac{2C_1 D^2 e}{C_3} \|\tilde{w}_0 - w_*\|^2 + \frac{C_2 P}{C_3}, \text{ where}$$

$$\begin{cases} C_1 = \frac{8L^2}{3} + \frac{14NL^2}{M} + \frac{4L^2}{3M}, \\ C_2 = \frac{2}{M} + 1 + \frac{5}{6L^2} + \frac{8N}{3ML^2} + \frac{5}{12ML}, \\ C_3 = \frac{1}{L^2+1} \frac{\mu}{M}. \end{cases}$$

252 Then, the we need  $T = \frac{\lambda G^{3/2}}{\varepsilon^{3/2}}$  epochs to guarantee

$$\min_{1 \leq t \leq T} [F(\tilde{w}_{t-1}) - F_*] \leq \frac{1}{T} \sum_{t=1}^T [F(\tilde{w}_{t-1}) - F_*] \leq \hat{\varepsilon}.$$

253 **Computational complexity.** Our global convergence result in this Corollary holds for a fixed value of  
 254  $\hat{\varepsilon}$  in Assumption 4. Thus, when  $\varepsilon \rightarrow 0$ , this assumption is equivalent to assuming  $\sigma_*^2 = 0$ . The total  
 255 complexity of Corollary 1 is  $\mathcal{O}\left(\frac{n}{\varepsilon^{3/2}}\right)$ . This rate matches the best known rate for unified sampling  
 256 schemes for SGD in convex setting [Mishchenko et al., 2020, Nguyen et al., 2021]. However, our  
 257 result holds for a broader class of functions that are possibly non-convex. Comparing to the non-  
 258 convex setting, current literature [Mishchenko et al., 2020, Nguyen et al., 2021] also matches our  
 259 rate to the order of  $\hat{\varepsilon}$ , however, we can only prove that SGD converges to a stationary point with a  
 260 weaker criteria  $\|\nabla F(w)\|^2 \leq \hat{\varepsilon}$  for general non-convex funtions. Table 1 shows these comparisons in  
 261 various settings. Note that when using a randomized shuffling scheme, SGD often performs a better  
 262 rate in terms of the data  $n$  in various settings with and without (strongly) convexity. For example, in  
 263 strongly convex and/or PL setting, the convergence rate of RR is  $\tilde{\mathcal{O}}(\sqrt{n}/\sqrt{\hat{\varepsilon}})$ , which is better than  
 264 unified schemes with  $\tilde{\mathcal{O}}(n/\sqrt{\hat{\varepsilon}})$  [Ahn et al., 2020]. However, for a fair comparison, we do not report  
 265 these results in Table 1 as our theoretical analysis is derived for unified shuffling scheme.

If we further assume that  $L, M, N > 1$ , the detailed complexity with respect to these constants is

$$\mathcal{O}\left(\frac{L^4(M+N)^{3/2}}{\mu^{3/2}} \cdot \frac{n}{\varepsilon^{3/2}}\right).$$

266 We present all the detailed proofs in the Appendix. Our theoretical framework is new and adapted to  
 267 the finite-sum minimization problem. Moreover, it utilizes the choice of shuffling sample schemes to  
 268 show a better complexity in terms of  $\hat{\varepsilon}$  than the complexity of vanilla i.i.d. sampling scheme.

Table 1: Comparisons of computational complexity (the number of individual gradient evaluations) needed by SGD algorithm to reach an  $\hat{\epsilon}$ -accurate solution  $w$  that satisfies  $F(w) - F(w_*) \leq \hat{\epsilon}$  (or  $\|\nabla F(w)\|^2 \leq \hat{\epsilon}$  in the non-convex case).

Settings	References	Complexity	Shuffling Schemes	Global Solution
Convex	Nemirovski et al. [2009], Shamir and Zhang [2013] <sup>(1)</sup>	$\mathcal{O}\left(\frac{\Delta_0^2 + G^2}{\hat{\epsilon}^2}\right)$	✗	✓
	Mishchenko et al. [2020], Nguyen et al. [2021] <sup>(2)</sup>	$\mathcal{O}\left(\frac{n}{\hat{\epsilon}^{3/2}}\right)$	✓	✓
PL condition	Nguyen et al. [2021]	$\tilde{\mathcal{O}}\left(\frac{n\sigma^2}{\hat{\epsilon}^{1/2}}\right)$	✓	✓
Star-convex related	Gower et al. [2021] <sup>(3)</sup>	$\mathcal{O}\left(\frac{1}{\hat{\epsilon}^2}\right)$	✗	✓
Non-convex	Ghadimi and Lan [2013] <sup>(5)</sup>	$\mathcal{O}\left(\frac{\sigma^2}{\hat{\epsilon}^2}\right)$	✗	✗
	Nguyen et al. [2021], Mishchenko et al. [2020] <sup>(5)</sup>	$\mathcal{O}\left(\frac{n\sigma}{\hat{\epsilon}^{3/2}}\right)$	✓	✗
<b><i>Our setting (non-convex)</i></b>	<b>This paper, Corollary 1</b> <sup>(4)</sup>	$\mathcal{O}\left(\frac{n(N \vee 1)^{3/2}}{\hat{\epsilon}^{3/2}}\right)$	✓	✓

<sup>(0)</sup> We note that the assumptions in this table are not comparable and we only show the roughly complexity in terms of  $\hat{\epsilon}$ . In addition, to make fair comparisons, we only report the complexity of unified shuffling schemes.

<sup>(1)</sup> Standard results for SGD in convex literature often use a different set of assumptions from the one in this paper (e.g. bounded domain that  $\|w - w_*\|^2 \leq \Delta_0$  for each iterate  $w$  and/or bounded gradient that  $\mathbb{E}[\|\nabla f(w; i)\|] \leq G$ ). We report the corresponding complexity for a rough comparison.

<sup>(2)</sup> [Mishchenko et al., 2020] shows a bound for Incremental Gradient while [Nguyen et al., 2021] has a unified setting. We translate these results for unified shuffling schemes from these references to the convex setting.

<sup>(3)</sup> Since we cannot find a reference containing the convergence rate for vanilla SGD and star-convex functions, we adapt the reference Gower et al. [2021] here. This paper shows a result for  $L$ -smooth and quasar convex function with an additional Expected Residual (ER) assumption, which is weaker than assuming smoothness for  $f(\cdot; i)$  and interpolation property. The star-convex results hold when the quasar-convex parameter is 1.

<sup>(4)</sup> Since we use a different set of assumptions than the other references, we only report the rough comparison in  $n$ ,  $N$  and  $\hat{\epsilon}$ , where  $N$  is the constant from Assumption 5 and  $N \vee 1 = \max(N, 1)$ . Note that  $N = 0$  in the framework of star-smooth-convex function. In addition, we need  $\sigma_*^2 = 0$  so that the complexity holds with arbitrary  $\hat{\epsilon}$ . We explain the detailed complexity below and in the Appendix.

<sup>(5)</sup> Standard literature for SGD in non-convex setting assumes a bounded variance that  $\mathbb{E}_i[\|f(w; i) - \nabla F(w)\|^2] \leq \sigma^2$ , we report the rough comparison.

## 269 4 Numerical Experiments

270 In this section, we show some experiments for shuffling-type SGD algorithms to demonstrate our  
 271 theoretical results of convergence to a global solution. Following the setting of Theorem 1, we  
 272 consider the non-convex regression problem with squared loss function. We choose fully connected  
 273 neural networks in our implementation. We experiment with different regression datasets: the  
 274 Diabetes dataset from sklearn library [Efron et al., 2004, Pedregosa et al., 2011] with 442 samples in  
 275 dimension 10; the Life expectancy dataset from WHO [Repository, 2016] with 1649 trainable data  
 276 points and 19 features. In addition, we test with the California Housing data from StatLib repository  
 277 [Repository, 1997, Pedregosa et al., 2011] with a training set of 16514 samples and 8 features.

278 For the small Diabetes dataset, we use the classic LeNet-300-100 model [LeCun et al., 1998]. For  
 279 other larger datasets, we use similar fully connected neural networks with an additional starting layer  
 280 of 900 neurons. We apply the randomized reshuffling scheme using PyTorch framework [Paszke et al.,  
 281 2019]. This shuffling scheme is the common heuristic in training neural networks and is implemented  
 282 in many deep learning platforms (e.g. TensorFlow, PyTorch, and Keras [Abadi et al., 2015, Paszke  
 283 et al., 2019, Chollet et al., 2015]).

284 For each dataset  $\{x_i, y_i\}$ , we preprocess and modify the initial data lightly to guarantee the over-  
 285 parameterized setting in our experiment. We do this by using a pre-training stage: firstly we use  
 286 GD/SGD algorithm to find a weight  $w$  that yields a sufficiently small value for the loss function  
 287 (for Diabetes dataset we train to  $10^{-8}$  and for other datasets we train to  $10^{-2}$ ). Letting the input  
 288 data  $x_i$  be fixed, we change the label data to  $\hat{y}_i$  such that the weight  $w$  yields a small loss function  
 289  $\mathcal{O}(\epsilon)$  for the optimization associated with data  $\{x_i, \hat{y}_i\}$ , and the distance between  $\hat{y}_i$  and  $y_i$  is small.  
 290 Then the modified data is ready for the next stage. We summarize the data (after modification) in our  
 291 experiments in Table 2.

Table 2: Datasets used in our experiments

Data name	# Samples	# Features	Networks layers	Sources
Diabetes	442	10	300-100	Efron et al. [2004]
Life Expectancy	1649	19	900-300-100	Repository [2016]
California Housing	16514	8	900-300-100	Repository [1997]

292 For each dataset, we first tune the step size using a coarse grid search  $[0.0001, 0.001, 0.01, 0.1, 1]$   
 293 for 100 epochs. Then, for example, if 0.01 performs the best, we test the second grid search  
 294  $[0.002, 0.005, 0.01, 0.02, 0.05]$  for 5000 epochs. Finally, we progress to the training stage with  $10^6$   
 295 epochs and repeat that experiment for 5 random seeds. We report the average results with confidence  
 296 intervals in Figure 1.

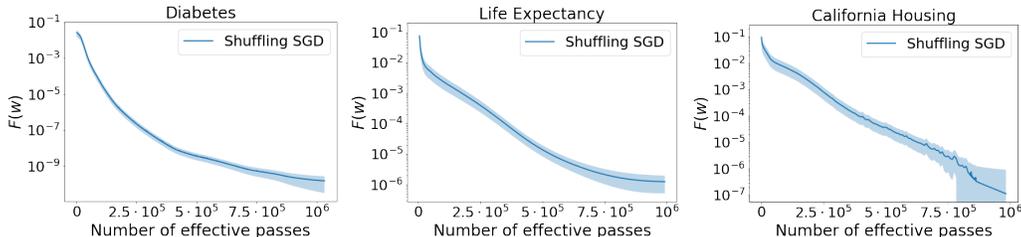


Figure 1: The train loss produced by Shuffling SGD algorithm for three datasets: Diabetes, Life Expectancy and California Housing.

297 For California Housing data, Shuffling SGD fluctuates toward the end of the training process.  
 298 Nevertheless, for all three datasets it converges steadily to a small value of loss function. In summary,  
 299 this experiment confirms our theoretical guarantee that demonstrates a convergence to global solution  
 300 for shuffling-type SGD algorithm in neural network settings.

## 301 5 Conclusion

302 In this paper, we study the global convergence property for shuffling-type SGD methods. We  
 303 consider a relaxed set of assumptions in the framework of star-smooth-convex functions and show  
 304 the total complexity of  $\mathcal{O}(\frac{n}{\epsilon^{3/2}})$  to reach an  $\hat{\epsilon}$ -accurate global solution. This result matches the  
 305 previous computational complexity of unified shuffling methods in convex settings. Our theoretical  
 306 framework utilizes the choice of shuffling sample schemes for finite-sum minimization problems in  
 307 machine learning. We provide discussions on the relations of our framework and the well-known  
 308 over-parameterized settings, as well as current literature on the star-convexity class of functions. In  
 309 addition, we show the connections to neural network architectures and discuss how these learning  
 310 models fit into our optimization frameworks. Potential research questions arising from our paper  
 311 include practical network designs and relaxed theoretical settings that support the global convergence  
 312 of Shuffling SGD methods. Moreover, the global convergence framework for other stochastic gradient  
 313 methods [Duchi et al., 2011, Kingma and Ba, 2014] and variance reduction methods [Nguyen et al.,  
 314 2017, Beznosikov and Takáč, 2021] with shuffling sampling schemes is also an interesting direction.

## 315 References

- 316 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S.  
317 Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew  
318 Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath  
319 Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah,  
320 Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent  
321 Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg,  
322 Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on  
323 heterogeneous systems, 2015. URL <https://www.tensorflow.org/>. Software available from  
324 tensorflow.org.
- 325 Kwangjun Ahn, Chulhee Yun, and Suvrit Sra. Sgd with shuffling: optimal rates without compo-  
326 nent convexity and large epoch requirements. In H. Larochelle, M. Ranzato, R. Hadsell, M. F.  
327 Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,  
328 pages 17526–17535. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/  
329 paper/2020/file/cb8acb1dc9821bf74e6ca9068032d623-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/cb8acb1dc9821bf74e6ca9068032d623-Paper.pdf).
- 330 Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-  
331 parameterization. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of  
332 the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine  
333 Learning Research*, pages 242–252. PMLR, 09–15 Jun 2019. URL [http://proceedings.mlr.  
334 press/v97/allen-zhu19a.html](http://proceedings.mlr.press/v97/allen-zhu19a.html).
- 335 Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of  
336 optimization and generalization for overparameterized two-layer neural networks. In Kamalika  
337 Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference  
338 on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 322–332.  
339 PMLR, 09–15 Jun 2019. URL <http://proceedings.mlr.press/v97/arora19a.html>.
- 340 Aleksandr Beznosikov and Martin Takáč. Random-reshuffled sarah does not need a full gradient  
341 computations. *arXiv preprint arXiv:2111.13322*, 2021.
- 342 Johan Bjorck, Anmol Kabra, Kilian Q. Weinberger, and Carla Gomes. Characterizing the loss  
343 landscape in non-negative matrix factorization. *Proceedings of the AAAI Conference on Artificial  
344 Intelligence*, 35(8):6768–6776, May 2021. URL [https://ojs.aaai.org/index.php/AAAI/  
345 article/view/16836](https://ojs.aaai.org/index.php/AAAI/article/view/16836).
- 346 L. Bottou, F. E. Curtis, and J. Nocedal. Optimization Methods for Large-Scale Machine Learning.  
347 *SIAM Rev.*, 60(2):223–311, 2018.
- 348 Léon Bottou. Curiously fast convergence of some stochastic gradient descent algorithms, 2009.
- 349 Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. SGD learns over-  
350 parameterized networks that provably generalize on linearly separable data. In *International  
351 Conference on Learning Representations*, 2018. URL [https://openreview.net/forum?id=  
352 rJ33wvxRb](https://openreview.net/forum?id=rJ33wvxRb).
- 353 Francois Chollet et al. Keras. *GitHub*, 2015. URL <https://github.com/fchollet/keras>.
- 354 Ronan Collobert and Jason Weston. A unified architecture for natural language processing: deep  
355 neural networks with multitask learning. In William W. Cohen, Andrew McCallum, and Sam T.  
356 Roweis, editors, *Machine Learning, Proceedings of the Twenty-Fifth International Conference  
357 (ICML 2008), Helsinki, Finland, June 5-9, 2008*, volume 307 of *ACM International Conference  
358 Proceeding Series*, pages 160–167. ACM, 2008. doi: 10.1145/1390156.1390177. URL [https:  
359 //doi.org/10.1145/1390156.1390177](https://doi.org/10.1145/1390156.1390177).
- 360 Soham De, Abhay Yadav, David Jacobs, and Tom Goldstein. Automated Inference with Adaptive  
361 Batches. In Aarti Singh and Jerry Zhu, editors, *Proceedings of the 20th International Conference  
362 on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*,  
363 pages 1504–1513. PMLR, 20–22 Apr 2017. URL [https://proceedings.mlr.press/v54/  
364 de17a.html](https://proceedings.mlr.press/v54/de17a.html).

- 365 Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global  
366 minima of deep neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors,  
367 *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings*  
368 *of Machine Learning Research*, pages 1675–1685. PMLR, 09–15 Jun 2019a. URL [http://](http://proceedings.mlr.press/v97/du19c.html)  
369 [proceedings.mlr.press/v97/du19c.html](http://proceedings.mlr.press/v97/du19c.html).
- 370 Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes  
371 over-parameterized neural networks. In *International Conference on Learning Representations*,  
372 2019b. URL <https://openreview.net/forum?id=S1eK3i09YQ>.
- 373 John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and  
374 stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- 375 Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Diabetes dataset, 2004. URL  
376 <https://www4.stat.ncsu.edu/~boos/var.select/diabetes.html>.
- 377 S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic program-  
378 ming. *SIAM J. Optim.*, 23(4):2341–2368, 2013.
- 379 Yoav Goldberg, Graeme Hirst, Yang Liu, and Meng Zhang. Neural network methods for natural  
380 language processing. *Comput. Linguistics*, 44(1), 2018. doi: 10.1162/COLI\_r\_00312. URL  
381 [https://doi.org/10.1162/COLI\\_r\\_00312](https://doi.org/10.1162/COLI_r_00312).
- 382 Ian J. Goodfellow, Yaroslav Bulatov, Julian Ibarz, Sacha Arnoud, and Vinay D. Snet. Multi-digit  
383 number recognition from street view imagery using deep convolutional neural networks. In Yoshua  
384 Bengio and Yann LeCun, editors, *2nd International Conference on Learning Representations*,  
385 *ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. URL  
386 <http://arxiv.org/abs/1312.6082>.
- 387 Robert Gower, Othmane Sebbouh, and Nicolas Loizou. Sgd for structured nonconvex functions:  
388 Learning rates, minibatching and interpolation. In *International Conference on Artificial Intelli-*  
389 *gence and Statistics*, pages 1315–1323. PMLR, 2021.
- 390 M. Gürbüzbalaban, A. Ozdaglar, and P. A. Parrilo. Why random reshuffling beats stochastic  
391 gradient descent. *Mathematical Programming*, Oct 2019. ISSN 1436-4646. doi: 10.1007/  
392 [s10107-019-01440-w](http://dx.doi.org/10.1007/s10107-019-01440-w). URL <http://dx.doi.org/10.1007/s10107-019-01440-w>.
- 393 Jeff Haochen and Suvrit Sra. Random shuffling beats sgd after finite epochs. In *International*  
394 *Conference on Machine Learning*, pages 2624–2633. PMLR, 2019.
- 395 Moritz Hardt, Tengyu Ma, and Benjamin Recht. Gradient descent learns linear dynamical systems.  
396 *Journal of Machine Learning Research*, 19(29):1–44, 2018. URL [http://jmlr.org/papers/](http://jmlr.org/papers/v19/16-465.html)  
397 [v19/16-465.html](http://jmlr.org/papers/v19/16-465.html).
- 398 Kaiming He and Jian Sun. Convolutional neural networks at constrained time cost. In *IEEE*  
399 *Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June*  
400 *7-12, 2015*, pages 5353–5360. IEEE Computer Society, 2015. doi: 10.1109/CVPR.2015.7299173.  
401 URL <https://doi.org/10.1109/CVPR.2015.7299173>.
- 402 Oliver Hinder, Aaron Sidford, and Nimit Sohoni. Near-optimal methods for minimizing star-convex  
403 functions and beyond. In *Conference on Learning Theory*, pages 1894–1938. PMLR, 2020.
- 404 Jikai Jin. On the convergence of first order methods for quasars-convex optimization. *arXiv preprint*  
405 *arXiv:2010.04937*, 2020.
- 406 Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-  
407 gradient methods under the Polyak-Łojasiewicz condition. In Paolo Frasconi, Niels Landwehr,  
408 Giuseppe Manco, and Jilles Vreeken, editors, *Machine Learning and Knowledge Discovery in*  
409 *Databases*, pages 795–811, Cham, 2016. Springer International Publishing.
- 410 Ahmed Khaled and Peter Richtárik. Better theory for sgd in the nonconvex world. *arXiv preprint*  
411 *arXiv:2002.03329*, 2020.

- 412 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*,  
413 abs/1412.6980, 2014.
- 414 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to  
415 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 416 Jasper C.H. Lee and Paul Valiant. Optimizing star-convex functions. In *2016 IEEE 57th Annual*  
417 *Symposium on Foundations of Computer Science (FOCS)*, pages 603–614, 2016. doi: 10.1109/  
418 FOCS.2016.71.
- 419 Nicolas Loizou, Sharan Vaswani, Issam Hadj Laradji, and Simon Lacoste-Julien. Stochastic polyak  
420 step-size for sgd: An adaptive learning rate for fast convergence. In Arindam Banerjee and Kenji  
421 Fukumizu, editors, *Proceedings of The 24th International Conference on Artificial Intelligence and*  
422 *Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 1306–1314. PMLR,  
423 13–15 Apr 2021. URL <http://proceedings.mlr.press/v130/loizou21a.html>.
- 424 Siyuan Ma, Raef Bassily, and Mikhail Belkin. The power of interpolation: Understanding the  
425 effectiveness of SGD in modern over-parametrized learning. In Jennifer Dy and Andreas Krause,  
426 editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of  
427 *Proceedings of Machine Learning Research*, pages 3325–3334. PMLR, 10–15 Jul 2018. URL  
428 <http://proceedings.mlr.press/v80/ma18a.html>.
- 429 Si Yi Meng, Sharan Vaswani, Issam Hadj Laradji, Mark Schmidt, and Simon Lacoste-Julien. Fast  
430 and furious convergence: Stochastic second order methods under interpolation. In Silvia Chiappa  
431 and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on*  
432 *Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*,  
433 pages 1375–1386. PMLR, 26–28 Aug 2020. URL [http://proceedings.mlr.press/v108/  
434 meng20a.html](http://proceedings.mlr.press/v108/meng20a.html).
- 435 Konstantin Mishchenko, Ahmed Khaled Ragab Bayoumi, and Peter Richtárik. Random reshuffling:  
436 Simple analysis with vast improvements. *Advances in Neural Information Processing Systems*, 33,  
437 2020.
- 438 Dheeraj Nagaraj, Prateek Jain, and Praneeth Netrapalli. Sgd without replacement: Sharper rates  
439 for general smooth convex functions. In *International Conference on Machine Learning*, pages  
440 4703–4711, 2019.
- 441 A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to  
442 stochastic programming. *SIAM J. on Optimization*, 19(4):1574–1609, 2009.
- 443 Yurii Nesterov. *Introductory lectures on convex optimization : a basic course*. Applied optimization.  
444 Kluwer Academic Publ., Boston, Dordrecht, London, 2004. ISBN 1-4020-7553-7.
- 445 Yurii Nesterov and Boris T Polyak. Cubic regularization of Newton method and its global performance.  
446 *Mathematical Programming*, 108(1):177–205, 2006.
- 447 Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method for machine  
448 learning problems using stochastic recursive gradient. In *Proceedings of the 34th International*  
449 *Conference on Machine Learning-Volume 70*, pages 2613–2621. JMLR. org, 2017.
- 450 Lam M. Nguyen, Quoc Tran-Dinh, Dzung T. Phan, Phuong Ha Nguyen, and Marten van Dijk. A  
451 unified convergence analysis for shuffling-type gradient methods. *Journal of Machine Learning*  
452 *Research*, 22(207):1–44, 2021. URL <http://jmlr.org/papers/v22/20-1238.html>.
- 453 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
454 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward  
455 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,  
456 Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep  
457 learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035.  
458 Curran Associates, Inc., 2019.

- 459 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-  
460 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and  
461 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,  
462 12:2825–2830, 2011.
- 463 Boris T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR*  
464 *Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- 465 Global Health Observatory Data Repository. Life expectancy and healthy life expectancy, 2016. URL  
466 <https://apps.who.int/gho/data/view.main.SDG2016LEXREGv?lang=en>.
- 467 StatLib Repository. California housing, 1997. URL [https://www.dcc.fc.up.pt/~ltorgo/](https://www.dcc.fc.up.pt/~ltorgo/Regression/cal_housing.html)  
468 [Regression/cal\\_housing.html](https://www.dcc.fc.up.pt/~ltorgo/Regression/cal_housing.html).
- 469 Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical*  
470 *Statistics*, 22(3):400–407, 1951.
- 471 Itay Safran and Ohad Shamir. How good is sgd with random shuffling? In *Conference on Learning*  
472 *Theory*, pages 3250–3284. PMLR, 2020.
- 473 Karthik Abinav Sankararaman, Soham De, Zheng Xu, W Ronny Huang, and Tom Goldstein. The  
474 impact of neural network overparameterization on gradient confusion and stochastic gradient  
475 descent. In *International conference on machine learning*, pages 8469–8479. PMLR, 2020.
- 476 Mark Schmidt and Nicolas Le Roux. Fast convergence of stochastic gradient descent under a strong  
477 growth condition, 2013.
- 478 Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization: Convergence  
479 results and optimal averaging schemes. In Sanjoy Dasgupta and David McAllester, editors,  
480 *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings*  
481 *of Machine Learning Research*, pages 71–79, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.  
482 URL <https://proceedings.mlr.press/v28/shamir13.html>.
- 483 Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit  
484 bias of gradient descent on separable data. *J. Mach. Learn. Res.*, 19(1):2822–2878, January 2018.  
485 ISSN 1532-4435.
- 486 Sharan Vaswani, Francis Bach, and Mark Schmidt. Fast and faster convergence of sgd for over-  
487 parameterized models and an accelerated perceptron. In Kamalika Chaudhuri and Masashi  
488 Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial*  
489 *Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 1195–  
490 1204. PMLR, 16–18 Apr 2019. URL [https://proceedings.mlr.press/v89/vaswani19a.](https://proceedings.mlr.press/v89/vaswani19a.html)  
491 [html](https://proceedings.mlr.press/v89/vaswani19a.html).
- 492 Yi Zhou, Junjie Yang, Huishuai Zhang, Yingbin Liang, and Vahid Tarokh. Sgd converges to global  
493 minimum in deep learning via star-convex path. *arXiv preprint arXiv:1901.00451*, 2019.
- 494 Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural  
495 networks. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc,  
496 Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems*  
497 *32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December*  
498 *8-14, 2019, Vancouver, BC, Canada*, pages 2053–2062, 2019.
- 499 Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes  
500 over-parameterized deep relu networks, 2018.

---

# On the Convergence to a Global Solution of Shuffling-Type Gradient Algorithms Supplementary Material, NeurIPS 2023

---

501 **A Theoretical settings: Proof of Theorem 1 and Lemma 2**

502 **A.1 Proof of Theorem 1**

503 *Proof.* Let us use the notation  $f(w; i) = \phi_i(h(w; i)) = \frac{1}{2}\|h(w; i) - y^{(i)}\|^2$ . We consider an  
504 architecture  $h(w; i)$  with  $w$  be the vectorized weight and  $h$  consists of a final bias layer  $b$ :

$$h(w; i) = W^T z(\theta; i) + b \in \mathbb{R}^c,$$

505 where  $w = \mathbf{vec}(\{\theta, W, b\})$  and  $z(\theta; i)$  are some inner architecture, which can be chosen arbitrarily.

506 Firstly, we compute the gradient of  $f(\cdot; i)$  with respect to  $b \in \mathbb{R}^c$ . For  $j = 1, \dots, c$ , we have

$$\frac{\partial f(w; i)}{\partial b_j} = \frac{\partial \phi_i(h(w; i))}{\partial b_j} = \sum_{k=1}^c \frac{\partial h(w; i)_k}{\partial b_j} \cdot \frac{\partial \phi_i(x)}{\partial x_k} \Big|_{x=h(w; i)} = \frac{\partial \phi_i(x)}{\partial x_j} \Big|_{x=h(w; i)}, \quad i = 1, \dots, n. \quad (18)$$

507 The last equality follows since  $\frac{\partial h(w; i)_k}{\partial b_j} = 0$  for every  $k \neq j$  and  $\frac{\partial h(w; i)_k}{\partial b_j} = 1$  for  $k = j$ . In other  
508 words, it is the identity matrix.

509 Let us denote that  $f_i^* = \min_w f(w; i)$  and  $\phi_i^* = \min_x \phi_i(x)$ . We prove the following statement for  
510  $\mu = 1$ :

$$\|\nabla_w f(w; i)\|^2 \geq \|\nabla_x \phi_i(x)|_{x=h(w; i)}\|^2 \geq 2\mu[\phi_i(h(w; i)) - \phi_i^*] \geq 2\mu[f(w; i) - f_i^*],$$

511 for every  $w \in \mathbb{R}^d$ , and  $i = 1, \dots, n$ .

512 We begin with the first inequality:

$$\begin{aligned} \|\nabla_w f(w; i)\|^2 &= \sum_{j=1}^d \left( \frac{\partial f(w; i)}{\partial w_j} \right)^2 \geq \sum_{j=d-c+1}^d \left( \frac{\partial f(w; i)}{\partial w_j} \right)^2 = \sum_{j=1}^c \left( \frac{\partial f(w; i)}{\partial b_j} \right)^2 \\ &\stackrel{(18)}{=} \sum_{j=1}^c \left( \frac{\partial \phi_i(x)}{\partial x_j} \Big|_{x=h(w; i)} \right)^2 = \|\nabla_x \phi_i(x)|_{x=h(w; i)}\|^2. \end{aligned}$$

513 Now let us prove the PL condition for each function  $\phi_i(x)$ , i.e., there exists a constant  $\mu > 0$  such  
514 that:

$$\|\nabla_x \phi_i(x)\|^2 \geq 2\mu[\phi_i(x) - \phi_i^*] \quad \forall x \in \mathbb{R}^c, \quad i = 1, \dots, n.$$

515 Recall the squared loss that  $\phi_i(x) = \frac{1}{2}\|x - y^{(i)}\|^2$  and  $\nabla_x \phi_i(x) = x - y^{(i)}$ . We can see that the  
516 constant  $\mu = 1$  satisfies the following inequality for every  $x \in \mathbb{R}^c$ ,  $i = 1, \dots, n$ :

$$\|\nabla_x \phi_i(x)\|^2 = \|x - y^{(i)}\|^2 = 2 \frac{1}{2} \|x - y^{(i)}\|^2 = 2\mu\phi_i(x) \geq 2\mu[\phi_i(x) - \phi_i^*],$$

517 where the last inequality follows since  $\phi_i^* \geq 0$ .

518 The PL condition for  $\phi_i$  directly implies the second inequality. The last inequality follows from  
519 the facts that  $f(w; i) = \phi_i(h(w; i))$  and  $f_i^* = \min_w f_i \geq \min_x \phi_i(x) = \phi_i^*$ . Hence, Theorem 1 is  
520 proved.  $\square$

521 **A.2 Statement and Proof of Lemma 2**

522 **Lemma 2.** *The function  $g$  is star- $M$ -smooth-convex with respect to  $\hat{w}$  for some constant  $M > 0$  if*  
 523 *one of the two following conditions holds:*

- 524 1.  $g$  is  $L$ -smooth and convex.  
 525 2.  $g$  is  $L$ -smooth and  $g$  is star-convex with respect to  $\hat{w}$ .

526 *Proof.* The first statement of Lemma 2 follows directly from equation (11). We have that  $g$  is  
 527 star- $M$ -smooth-convex with respect to any reference point and  $M = L$ .

528 Now we proceed to the second statement. From the star-convex property of  $g$  with respect to  $\hat{w}$ , we  
 529 have

$$g(w) - g(\hat{w}) \leq \langle \nabla g(w), w - \hat{w} \rangle, \forall w \in \mathbb{R}^d,$$

530 and  $\nabla g(\hat{w}) = 0$  since  $\hat{w}$  is the global minimizer of  $g$ . On the other hand,  $g$  is  $L$ -smooth and we have

$$g(\hat{w}) \leq g\left(w - \frac{1}{L}\nabla g(w)\right) \leq g(w) - \frac{1}{2L}\|\nabla g(w)\|^2,$$

531 which is equivalent to  $\|\nabla g(w)\|^2 \leq 2L[g(w) - g(\hat{w})]$ ,  $i \in [n]$ . Since  $\nabla g(\hat{w}) = 0$ ,  $i \in [n]$ , we have  
 532 for  $\forall w \in \mathbb{R}^d$

$$\|\nabla g(w) - \nabla g(\hat{w})\|^2 \leq 2L[g(w) - g(\hat{w})] \stackrel{(12)}{\leq} 2L\langle \nabla g(w) - \nabla g(\hat{w}), w - w_* \rangle.$$

533 This is a star- $M$ -smooth-convex function as in Definition 2 with  $M = 2L$ . □

534 **B Preliminary results for SGD Shuffling Algorithm**

535 In this section, we present the preliminary results for Algorithm 1. Firstly, from the choice of learning  
 536 rate  $\eta_i^{(t)} := \frac{\eta_t}{n}$  and the update  $w_{i+1}^{(t)} := w_i^{(t)} - \eta_i^{(t)}\nabla f(w_i^{(t)}; \pi^{(t)}(i+1))$  in Algorithm 1, for  $i \in [n]$ ,  
 537 we have

$$w_i^{(t)} = w_{i-1}^{(t)} - \frac{\eta_t}{n}\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i)) = w_0^{(t)} - \frac{\eta_t}{n}\sum_{j=0}^{i-1}\nabla f(w_j^{(t)}; \pi^{(t)}(j+1)). \quad (19)$$

538 Hence,

$$w_0^{(t+1)} = w_n^{(t)} = w_0^{(t)} - \frac{\eta_t}{n}\sum_{j=0}^{n-1}\nabla f(w_j^{(t)}; \pi^{(t)}(j+1)). \quad (20)$$

539 Next, we refer to a Lemma in [Nguyen et al., 2021] to bound the updates of shuffling SGD algorithms.

540 **Lemma 3** (Lemma 5 in Nguyen et al. [2021]). *Suppose that Assumption 2 holds for (1). Let  $\{w_i^{(t)}\}$*   
 541 *be generated by Algorithm 1 with the learning rate  $\eta_i^{(t)} := \frac{\eta_t}{n} > 0$  for a given positive sequence*  
 542  *$\{\eta_t\}$ . If  $0 < \eta_t \leq \frac{1}{2L}$  for all  $t \geq 1$ , we have*

$$\frac{1}{n}\sum_{j=0}^{n-1}\|w_j^{(t)} - w_*\|^2 \leq 4\|w_0^{(t)} - w_*\|^2 + 8\sigma_*^2 \cdot \eta_t^2, \quad (21)$$

$$\frac{1}{n}\sum_{j=0}^{n-1}\|w_j^{(t)} - w_0^{(t)}\|^2 \leq \eta_t^2 \cdot \frac{8L^2}{3}\|w_0^{(t)} - w_*\|^2 + \frac{16L^2\sigma_*^2}{3} \cdot \eta_t^4 + 2\sigma_*^2 \cdot \eta_t^2. \quad (22)$$

543 Now considering the term  $\|w_n^{(t)} - w_0^{(t)}\|^2$ , we get that

$$\|w_n^{(t)} - w_0^{(t)}\|^2 \stackrel{(20)}{\leq} \frac{\eta_t^2}{n} \left\| \frac{1}{n} \sum_{j=0}^{n-1} \nabla f(w_j^{(t)}; \pi^{(t)}(j+1)) \right\|^2$$

$$\begin{aligned}
&= \frac{\eta_t^2}{n} \left\| \frac{1}{n} \sum_{j=0}^{n-1} (\nabla f(w_j^{(t)}; \pi^{(t)}(j+1)) - \nabla f(w_*; \pi^{(t)}(j+1))) \right\|^2 \\
&\leq \frac{\eta_t^2}{n} \frac{1}{n} \sum_{j=0}^{n-1} \left\| \nabla f(w_j^{(t)}; \pi^{(t)}(j+1)) - \nabla f(w_*; \pi^{(t)}(j+1)) \right\|^2 \\
&\stackrel{(2)}{\leq} \frac{L^2 \eta_t^2}{n} \frac{1}{n} \sum_{j=0}^{n-1} \|w_j^{(t)} - w_*\|^2 \\
&\stackrel{(21)}{\leq} \frac{4L^2 \eta_t^2}{n} \|w_0^{(t)} - w_*\|^2 + \frac{8L^2 \eta_t^4}{n} \sigma_*^2.
\end{aligned}$$

544 We further have

$$\begin{aligned}
\frac{1}{n} \sum_{j=0}^n \|w_j^{(t)} - w_0^{(t)}\|^2 &= \frac{1}{n} \sum_{j=0}^{n-1} \|w_j^{(t)} - w_0^{(t)}\|^2 + \frac{1}{n} \|w_n^{(t)} - w_0^{(t)}\|^2 \\
&\leq \eta_t^2 \cdot \frac{8L^2}{3} \|w_0^{(t)} - w_*\|^2 + \frac{16L^2 \sigma_*^2}{3} \cdot \eta_t^4 + 2\sigma_*^2 \cdot \eta_t^2 \\
&\quad + \frac{4L^2 \eta_t^2}{n} \|w_0^{(t)} - w_*\|^2 + \frac{8L^2 \eta_t^4}{n} \sigma_*^2. \tag{23}
\end{aligned}$$

## 545 C Main results: Proofs of Lemma 4, Lemma 1, Theorem 2, and Corollary 1

### 546 C.1 Proof of Lemma 4

547 **Lemma 4.** Let  $\{w_i^{(t)}\}_{t=1}^T$  be the sequence generated by Algorithm 1 with  $\eta_i^{(t)} = \frac{\eta_t}{n}$ , with  $0 < \eta_t \leq$   
548  $\frac{n}{2M}$  for  $\eta_t \leq \frac{1}{2L}$ . Then, under Assumptions 1, 2, and 5, we have

$$\|w_0^{(t+1)} - w_*\|^2 \leq (1 + B_1 \eta_t^3) \|w_0^{(t)} - w_*\|^2 - \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 + B_2 \eta_t \sigma_*^2, \tag{24}$$

549 where

$$\begin{cases} B_1 = \frac{8L^2}{3} + \frac{14NL^2}{M}, \\ B_2 = \frac{2}{M} + 1 + \frac{5}{6L^2} + \frac{8N}{3ML^2}. \end{cases} \tag{25}$$

550 *Proof.* We start with Assumption 5. Using the inequality  $\frac{1}{2}\|a\|^2 - \|b\|^2 \leq \|a - b\|^2$ , we have for  
551  $t = 1, \dots, T$  and  $i = 1, \dots, n$ :

$$\begin{aligned}
&\frac{1}{2} \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 - \|\nabla f(w_*; \pi^{(t)}(i))\|^2 \\
&\leq \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i)) - \nabla f(w_*; \pi^{(t)}(i))\|^2 \\
&\stackrel{(13)}{\leq} M \langle \nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i)) - \nabla f(w_*; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_* \rangle + N \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2 \\
&= M \langle \nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_* \rangle - M \langle \nabla f(w_*; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_* \rangle \\
&\quad + N \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2,
\end{aligned}$$

552 This statement is equivalent to

$$\begin{aligned}
-\langle \nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_* \rangle &\leq -\frac{1}{2M} \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 + \frac{1}{M} \|\nabla f(w_*; \pi^{(t)}(i))\|^2 \\
&\quad - \langle \nabla f(w_*; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_* \rangle
\end{aligned}$$

$$+ \frac{N}{M} \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2, \quad (26)$$

553 For any  $w_* \in W^*$ , from the update (19) we have,

$$\begin{aligned} \|w_i^{(t)} - w_*\|^2 &\stackrel{(19)}{=} \|w_{i-1}^{(t)} - w_*\|^2 - \frac{2\eta_t}{n} \langle \nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_* \rangle + \frac{\eta_t^2}{n^2} \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 \\ &\stackrel{(26)}{\leq} \|w_{i-1}^{(t)} - w_*\|^2 - \frac{2\eta_t}{2Mn} \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 + \frac{2\eta_t}{Mn} \|\nabla f(w_*; \pi^{(t)}(i))\|^2 \\ &\quad - \frac{2\eta_t}{n} \langle \nabla f(w_*; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_* \rangle + \frac{2\eta_t N}{Mn} \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2 \\ &\quad + \frac{\eta_t^2}{n^2} \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 \\ &\stackrel{(a)}{\leq} \|w_{i-1}^{(t)} - w_*\|^2 - \frac{\eta_t}{2Mn} \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 + \frac{2\eta_t}{Mn} \|\nabla f(w_*; \pi^{(t)}(i))\|^2 \\ &\quad - \frac{2\eta_t}{n} \langle \nabla f(w_*; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_* \rangle + \frac{2\eta_t N}{Mn} \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2 \\ &= \|w_{i-1}^{(t)} - w_*\|^2 - \frac{\eta_t}{2Mn} \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 + \frac{2\eta_t}{Mn} \|\nabla f(w_*; \pi^{(t)}(i))\|^2 \\ &\quad - \frac{2\eta_t}{n} \langle \nabla f(w_*; \pi^{(t)}(i)), w_{i-1}^{(t)} - w_0^{(t)} \rangle - \frac{2\eta_t}{n} \langle \nabla f(w_*; \pi^{(t)}(i)), w_0^{(t)} - w_* \rangle \\ &\quad + \frac{2\eta_t N}{Mn} \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2 \\ &\stackrel{(b)}{\leq} \|w_{i-1}^{(t)} - w_*\|^2 - \frac{\eta_t}{2Mn} \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 + \frac{2\eta_t}{Mn} \|\nabla f(w_*; \pi^{(t)}(i))\|^2 \\ &\quad + \frac{\eta_t}{n} \|\nabla f(w_*; \pi^{(t)}(i))\|^2 + \frac{\eta_t}{n} \|w_{i-1}^{(t)} - w_0^{(t)}\|^2 \\ &\quad - \frac{2\eta_t}{n} \langle \nabla f(w_*; \pi^{(t)}(i)), w_0^{(t)} - w_* \rangle + \frac{2\eta_t N}{Mn} \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2, \end{aligned}$$

554 where (a) follows since  $\eta_t \leq \frac{n}{2M}$  and (b) follows by the inequality  $2\langle a, b \rangle \leq \|a\|^2 + \|b\|^2$ .

555 Note that  $\frac{1}{n} \sum_{i=1}^n \langle \nabla f(w_*; \pi^{(t)}(i)), w_0^{(t)} - w_* \rangle = \langle \nabla F(w_*), w_0^{(t)} - w_* \rangle = 0$  since  $w_*$  is a global  
556 solution of  $F$ . Now we sum the derived statement for  $i = 1, \dots, n$  and get

$$\begin{aligned} \|w_n^{(t)} - w_*\|^2 &\leq \|w_0^{(t)} - w_*\|^2 - \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 \\ &\quad + \eta_t \left( \frac{2}{M} + 1 \right) \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_*; \pi^{(t)}(i))\|^2 + \frac{\eta_t}{n} \sum_{i=1}^n \|w_{i-1}^{(t)} - w_0^{(t)}\|^2 \\ &\quad + \frac{2N\eta_t}{M} \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2 \\ &\stackrel{(8),(22)}{\leq} \|w_0^{(t)} - w_*\|^2 - \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 \\ &\quad + \left( \frac{2}{M} + 1 \right) \eta_t \sigma_*^2 + \frac{8L^2\eta_t^3}{3} \|w_0^{(t)} - w_*\|^2 + \frac{16L^2\eta_t^5}{3} \sigma_*^2 + 2\eta_t^3 \sigma_*^2 \\ &\quad + \frac{2N\eta_t}{M} \frac{1}{n} \sum_{i=1}^n \|w_i^{(t)} - w_0^{(t)}\|^2 \\ &\stackrel{(23)}{\leq} \|w_0^{(t)} - w_*\|^2 - \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 \end{aligned}$$

$$\begin{aligned}
& + \left( \frac{2}{M} + 1 \right) \eta_t \sigma_*^2 + \frac{8L^2 \eta_t^3}{3} \|w_0^{(t)} - w_*\|^2 + \frac{16L^2 \eta_t^5}{3} \sigma_*^2 + 2\eta_t^3 \sigma_*^2 \\
& + \frac{16NL^2 \eta_t^3}{3M} \|w_0^{(t)} - w_*\|^2 + \frac{32NL^2 \eta_t^5}{3M} \sigma_*^2 + \frac{4N \eta_t^3}{M} \sigma_*^2 \\
& + \frac{8NL^2 \eta_t^3}{Mn} \|w_0^{(t)} - w_*\|^2 + \frac{16NL^2 \eta_t^5}{Mn} \sigma_*^2,
\end{aligned}$$

557 where we apply the derivations from Lemma 3. Now noting that  $\eta_t \leq \frac{1}{2L}$ ,  $n \leq 1$  and rearranging the  
558 terms we get:

$$\begin{aligned}
\|w_n^{(t)} - w_*\|^2 & \leq \|w_0^{(t)} - w_*\|^2 - \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 \\
& + \left( \frac{8L^2}{3} + \frac{16NL^2}{3M} + \frac{8NL^2}{M} \right) \eta_t^3 \|w_0^{(t)} - w_*\|^2 \\
& + \left( \frac{2}{M} + 1 + \frac{1}{3L^2} + \frac{1}{2L^2} + \frac{2N}{3ML^2} + \frac{N}{ML^2} + \frac{N}{ML^2} \right) \eta_t \sigma_*^2
\end{aligned}$$

559 Since  $w_n^{(t)} = w_0^{(t+1)} = \tilde{w}_t$ , we have the desired result in (24).  $\square$

## 560 C.2 Proof of Lemma 1

561 *Proof.* From (24) where  $B_1$  and  $B_2$  are defined in (25), we have

$$\begin{aligned}
& \|w_0^{(t+1)} - w_*\|^2 \\
& \leq (1 + B_1 \eta_t^3) \|w_0^{(t)} - w_*\|^2 - \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i))\|^2 + B_2 \eta_t \sigma_*^2 \\
& \stackrel{(a)}{\leq} (1 + B_1 \eta_t^3) \|w_0^{(t)} - w_*\|^2 - \frac{\gamma}{\gamma + 1} \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_0^{(t)}; \pi^{(t)}(i))\|^2 \\
& \quad + \frac{\eta_t \gamma}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_{i-1}^{(t)}; \pi^{(t)}(i)) - \nabla f(w_0^{(t)}; \pi^{(t)}(i))\|^2 + B_2 \eta_t \sigma_*^2 \\
& \stackrel{(b)}{\leq} (1 + B_1 \eta_t^3) \|w_0^{(t)} - w_*\|^2 - \frac{\gamma}{\gamma + 1} \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_0^{(t)}; \pi^{(t)}(i))\|^2 \\
& \quad + \frac{\eta_t \gamma L^2}{2M} \frac{1}{n} \sum_{i=1}^n \|w_{i-1}^{(t)} - w_0^{(t)}\|^2 + B_2 \sigma_*^2 \\
& \stackrel{(22)}{\leq} (1 + B_1 \eta_t^3) \|w_0^{(t)} - w_*\|^2 - \frac{\gamma}{\gamma + 1} \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_0^{(t)}; \pi^{(t)}(i))\|^2 \\
& \quad + \frac{\eta_t^3 \gamma L^2}{2M} \left( \frac{8L^2}{3} \|w_0^{(t)} - w_*\|^2 + \frac{16L^2 \sigma_*^2}{3} \cdot \eta_t^2 + 2\sigma_*^2 \right) + B_2 \eta_t \sigma_*^2 \\
& = \left( 1 + B_1 \eta_t^3 + \frac{4\eta_t^3 \gamma L^4}{3M} \right) \|w_0^{(t)} - w_*\|^2 + \left( B_2 + \frac{\eta_t^2 \gamma L^2}{M} + \frac{8\eta_t^4 \gamma L^4}{3M} \right) \eta_t \sigma_*^2 \\
& \quad - \frac{\gamma}{\gamma + 1} \frac{\eta_t}{2M} \frac{1}{n} \sum_{i=1}^n \|\nabla f(w_0^{(t)}; \pi^{(t)}(i))\|^2 \\
& \stackrel{(5)}{\leq} \left( 1 + \eta_t^3 \left( B_1 + \frac{4\gamma L^4}{3M} \right) \right) \|w_0^{(t)} - w_*\|^2 + \left( B_2 + \frac{\eta_t^2 \gamma L^2}{M} + \frac{8\eta_t^4 \gamma L^4}{3M} \right) \eta_t \sigma_*^2 \\
& \quad - \frac{\gamma}{\gamma + 1} \frac{2\mu \eta_t}{2M} \frac{1}{n} \sum_{i=1}^n [f(w_0^{(t)}; \pi^{(t)}(i)) - f_i^*] \\
& \stackrel{(3)}{\leq} \left( 1 + \eta_t^3 \left( B_1 + \frac{4\gamma L^4}{3M} \right) \right) \|w_0^{(t)} - w_*\|^2 + \left( B_2 + \frac{\eta_t^2 \gamma L^2}{M} + \frac{8\eta_t^4 \gamma L^4}{3M} \right) \eta_t \sigma_*^2
\end{aligned}$$

$$\begin{aligned}
& -\frac{\gamma}{\gamma+1} \frac{\mu\eta_t}{M} [F(w_0^{(t)}) - F_*] \\
\stackrel{(b)}{\leq} & \left(1 + \eta_t^3 \left(B_1 + \frac{4\gamma L^4}{3M}\right)\right) \|w_0^{(t)} - w_*\|^2 + \left(B_2 + \frac{\gamma}{4M} + \frac{\gamma}{6M}\right) \eta_t \sigma_*^2 \\
& -\frac{\gamma}{\gamma+1} \frac{\mu\eta_t}{M} [F(w_0^{(t)}) - F_*],
\end{aligned}$$

562 where (a) follows since  $-\|b\|^2 \leq \gamma\|a-b\|^2 - \frac{\gamma}{\gamma+1}\|a\|^2$  for any  $\gamma > 0$  and (b) follows since  $\eta_t \leq \frac{1}{2L}$ .  
563 Since  $w_0^{(t+1)} = \tilde{w}_t$ , we obtain the desired result in (14).  $\square$

### 564 C.3 Proof of Theorem 2

565 *Proof.* For  $t = 1, \dots, T = \frac{\lambda}{\varepsilon^{3/2}}$  for some  $\lambda > 0$

$$\begin{aligned}
\eta_t &= (1 + C_1 D^3 \varepsilon^{3/2}) \eta_{t-1} = (1 + C_1 D^3 \varepsilon^{3/2})^t \eta_0 \leq (1 + C_1 D^3 \varepsilon^{3/2})^T \eta_0 \\
&= (1 + C_1 D^3 \varepsilon^{3/2})^{\lambda/\varepsilon^{3/2}} \eta_0 = (1 + C_1 D^3 \varepsilon^{3/2})^{\lambda/\varepsilon^{3/2}} \frac{D\sqrt{\varepsilon}}{(1 + C_1 D^3 \varepsilon^{3/2}) \exp(\lambda C_1 D^3)} \\
&\leq \frac{D\sqrt{\varepsilon}}{(1 + C_1 D^3 \varepsilon^{3/2})} \leq \min \left\{ \frac{n}{2M}, \frac{1}{2L} \right\}, \tag{27}
\end{aligned}$$

566 since  $(1+x)^{1/x} \leq e$ ,  $x > 0$ . From (14), we have

$$[F(\tilde{w}_{t-1}) - F_*] \leq \frac{1}{C_3} \left( \frac{1}{\eta_t} + C_1 \eta_t^2 \right) \|\tilde{w}_{t-1} - w_*\|^2 - \frac{1}{C_3 \eta_t} \|\tilde{w}_t - w_*\|^2 + \frac{C_2}{C_3} \sigma_*^2. \tag{28}$$

567 We proceed to prove the following inequality for  $t = 1, \dots, T$ ,

$$\frac{1}{\eta_t} + C_1 \eta_t^2 \leq \frac{1}{\eta_{t-1}}. \tag{29}$$

568 From (27), and  $\eta_t = K \eta_{t-1}$  where  $K = (1 + C_1 D^3 \varepsilon^{3/2})$ , we have

$$\begin{aligned}
C_1 \eta_t^2 &= C_1 K^2 \eta_{t-1}^2 = C_1 K^2 \frac{\eta_{t-1}^3}{\eta_{t-1}} \\
&\leq C_1 K^2 \frac{D^3 \varepsilon^{3/2}}{K^3 \eta_{t-1}} = \frac{C_1 D^3 \varepsilon^{3/2}}{K \eta_{t-1}} && \text{since } \eta_{t-1} \leq \frac{D\sqrt{\varepsilon}}{K} \\
&= \frac{K-1}{K} \frac{1}{\eta_{t-1}} = \frac{1}{\eta_{t-1}} - \frac{1}{K \eta_{t-1}} && \text{since } K = (1 + C_1 D^3 \varepsilon^{3/2}) \\
&= \frac{1}{\eta_{t-1}} - \frac{1}{K \eta_{t-1}} = \frac{1}{\eta_{t-1}} - \frac{1}{\eta_t}, && \text{since } \eta_t = K \eta_{t-1}.
\end{aligned}$$

569 for  $t = 1, \dots, T$ . Hence, from (28), we have

$$\begin{aligned}
[F(\tilde{w}_{t-1}) - F_*] &\leq \frac{1}{C_3} \left( \frac{1}{\eta_t} + C_1 \eta_t^2 \right) \|\tilde{w}_{t-1} - w_*\|^2 - \frac{1}{C_3 \eta_t} \|\tilde{w}_t - w_*\|^2 + \frac{C_2}{C_3} \sigma_*^2 \\
&\leq \frac{1}{C_3 \eta_{t-1}} \|\tilde{w}_{t-1} - w_*\|^2 - \frac{1}{C_3 \eta_t} \|\tilde{w}_t - w_*\|^2 + \frac{C_2}{C_3} \sigma_*^2.
\end{aligned}$$

570 Averaging the statement above for  $t = 1, \dots, T$ , we have

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T [F(\tilde{w}_{t-1}) - F_*] &\leq \frac{1}{C_3 \eta_0 T} \|\tilde{w}_0 - w_*\|^2 + \frac{C_2}{C_3} \sigma_*^2 \\
&\stackrel{(a)}{=} \frac{K \exp(\lambda C_1 D^3)}{C_3 D \sqrt{\varepsilon}} \frac{\varepsilon^{3/2}}{\lambda} \|\tilde{w}_0 - w_*\|^2 + \frac{C_2}{C_3} \sigma_*^2 \\
&= \frac{K \exp(\lambda C_1 D^3)}{C_3 D \lambda} \|\tilde{w}_0 - w_*\|^2 \cdot \varepsilon + \frac{C_2}{C_3} \sigma_*^2,
\end{aligned}$$

571 where (a) follows since  $\eta_0 = \frac{D\sqrt{\varepsilon}}{K \exp(\lambda C_1 D^3)}$  and  $T = \frac{\lambda}{\varepsilon^{3/2}}$ .  $\square$

572 **C.4 Proof of Corollary 1**

573 *Proof.* Choose  $\gamma = \frac{1}{L^2}$ , we have

$$\begin{cases} C_1 = \frac{8L^2}{3} + \frac{14NL^2}{M} + \frac{4L^2}{3M}, \\ C_2 = \frac{2}{M} + 1 + \frac{5}{6L^2} + \frac{8N}{3ML^2} + \frac{5}{12ML}, \\ C_3 = \frac{1}{L^2+1} \frac{\mu}{M}. \end{cases}$$

574 Note that  $K = 1 + C_1 D^3 \varepsilon^{3/2}$  and  $C_1 D^3 = 1/\lambda$ , we get that  $K = 1 + 1/T \leq 2$ . We continue from  
575 the statement of Theorem 2 and the choice  $C_1 D^3 \lambda = 1$ :

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T [F(\tilde{w}_{t-1}) - F_*] &\leq \frac{K \exp(\lambda C_1 D^3)}{C_3 D \lambda} \|\tilde{w}_0 - w_*\|^2 \cdot \varepsilon + \frac{C_2}{C_3} \sigma_*^2 \\ &\leq \frac{2}{C_1 D^3 \lambda} \cdot \frac{C_1 D^2 \exp(\lambda C_1 D^3)}{C_3} \cdot \|\tilde{w}_0 - w_*\|^2 \cdot \varepsilon + \frac{C_2}{C_3} \sigma_*^2 && \text{since } K \leq 2 \\ &\leq \frac{2C_1 D^2 e}{C_3} \|\tilde{w}_0 - w_*\|^2 \cdot \varepsilon + \frac{C_2}{C_3} \sigma_*^2 && \text{since } C_1 D^3 \lambda = 1 \\ &\leq \frac{2C_1 D^2 e}{C_3} \|\tilde{w}_0 - w_*\|^2 \cdot \varepsilon + \frac{C_2}{C_3} P \varepsilon && \text{equation (9)} \\ &\leq \left( \frac{2C_1 D^2 e}{C_3} \|\tilde{w}_0 - w_*\|^2 + \frac{C_2 P}{C_3} \right) \varepsilon = G \varepsilon \end{aligned}$$

576 with

$$G = \frac{2C_1 D^2 e}{C_3} \|\tilde{w}_0 - w_*\|^2 + \frac{C_2 P}{C_3}.$$

577 Let  $0 < \varepsilon \leq 1$  and choose  $\hat{\varepsilon} = G \varepsilon$ . Then the number of iterations  $T$  is

$$\begin{aligned} T &= \frac{\lambda}{\hat{\varepsilon}^{3/2}} = \frac{\lambda G^{3/2}}{\hat{\varepsilon}^{3/2}} \\ &= \frac{1}{\hat{\varepsilon}^{3/2} C_1 D^3} \left( \frac{2C_1 D^2 e \|\tilde{w}_0 - w_*\|^2 + C_2 P}{C_3} \right)^{3/2} \\ &= \frac{1}{\hat{\varepsilon}^{3/2}} \cdot \frac{1}{C_1 D^3 C_3^{3/2}} (2C_1 D^2 e \|\tilde{w}_0 - w_*\|^2 + C_2 P)^{3/2} \\ &= \frac{1}{\hat{\varepsilon}^{3/2}} \cdot \frac{\left( 2D^2 e \|\tilde{w}_0 - w_*\|^2 \left( \frac{8L^2}{3} + \frac{14NL^2}{M} + \frac{4L^2}{3M} \right) + \left( \frac{2+M}{M} + \frac{5}{6L^2} + \frac{8N}{3ML^2} + \frac{5}{12ML} \right) P \right)^{3/2}}{\left( \frac{8L^2}{3} + \frac{14NL^2}{M} + \frac{4L^2}{3M} \right) D^3 \left( \frac{1}{L^2+1} \frac{\mu}{M} \right)^{3/2}} \end{aligned}$$

578 to guarantee

$$\min_{1 \leq t \leq T} [F(\tilde{w}_{t-1}) - F_*] \leq \frac{1}{T} \sum_{t=1}^T [F(\tilde{w}_{t-1}) - F_*] \leq \hat{\varepsilon}.$$

579 Hence, the total complexity (number of individual gradient computations needed to reach  $\hat{\varepsilon}$  accuracy)  
580 is  $\mathcal{O}\left(\frac{n}{\hat{\varepsilon}^{3/2}}\right)$ .

581 If we further assume that  $L, M, N > 1$ :

$$\begin{aligned} T &= \frac{1}{\hat{\varepsilon}^{3/2}} \cdot \frac{\left( 2D^2 e \|\tilde{w}_0 - w_*\|^2 \left( \frac{8ML^2}{3} + 14NL^2 + \frac{4L^2}{3} \right) + \left( 2 + M + \frac{5M}{6L^2} + \frac{8N}{3L^2} + \frac{5}{12L} \right) P \right)^{3/2}}{\left( \frac{8L^2}{3} + \frac{14NL^2}{M} + \frac{4L^2}{3M} \right) D^3 \left( \frac{\mu}{L^2+1} \right)^{3/2}} \\ &\leq \frac{1}{\hat{\varepsilon}^{3/2}} \cdot (\mathcal{O}((M+N)L^2))^{3/2} \cdot \mathcal{O}(1/L^2) \cdot \left( \frac{L^2+1}{\mu} \right)^{3/2} \end{aligned}$$

$$= \mathcal{O}\left(\frac{L^4(M+N)^{3/2}}{\mu^{3/2}} \cdot \frac{1}{\varepsilon^{3/2}}\right)$$

582 and the complexity is  $\mathcal{O}\left(\frac{L^4(M+N)^{3/2}}{\mu^{3/2}} \cdot \frac{n}{\varepsilon^{3/2}}\right)$ .

□