

FINITE-SUM SMOOTH OPTIMIZATION WITH SARAH

LAM M. NGUYEN*, MARTEN VAN DIJK†, DZUNG T. PHAN‡, PHUONG HA NGUYEN§, TSUI-WEI WENG¶, AND JAYANT R. KALAGNANAM||

Abstract. We introduce NC-SARAH for non-convex optimization as a *practical* slightly modified version of the original SARAH algorithm that was developed for convex optimization. NC-SARAH is the first to achieve two crucial performance properties at the same time – allowing flexible minibatch sizes and large step sizes to achieve fast convergence in practice as verified by experiments. NC-SARAH has a close to optimal asymptotic convergence rate equal to existing prior variants of SARAH called SPIDER and SpiderBoost that either use an order of magnitude smaller step size or a fixed minibatch size. For convex optimization, we propose SARAH++ with sublinear convergence for general convex and linear convergence for strongly convex problems; and we provide a practical version for which numerical experiments on various datasets show an improved performance.

Key words. Finite-sum, smooth, non-convex, convex, stochastic algorithm, variance reduction

AMS subject classifications. 90C25, 90C26, 90C30

1. Introduction. We are interested in solving the *finite-sum* minimization problem

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

where each $f_i, i \in [n] \stackrel{\text{def}}{=} \{1, \dots, n\}$, has a Lipschitz continuous gradient. Throughout the paper, we consider the case where F has a finite lower bound F^* . We would like to attain an ϵ -accurate solution satisfying $\mathbb{E}[\|\nabla F(\tilde{w})\|^2] \leq \epsilon$ for the outputted approximation \tilde{w} .

Problems of form (1.1) cover a wide range of convex and nonconvex problems in machine learning applications including but not limited to logistic regression, neural networks, multi-kernel learning, etc. In many of these applications, the number of component functions n is very large, which makes the classical Gradient Descent (GD) method less efficient since it requires to compute a full gradient many times. In recent years, a large number of improved variants of stochastic gradient algorithms called variance reduction methods have been proposed to obtain better computational cost compared to GD, in particular, SAG/SAGA [16, 4], SDCA [17], MISO [9], SVRG/S2GD [6, 7], SARAH [11], etc. These methods were first analyzed for strongly convex problems of form (1.1). Due to recent interest in deep neural networks, *nonconvex* problems of form (1.1) have been studied and analyzed by considering a number of different approaches including many variants of variance reduction techniques (see e.g. [13, 8, 1, 2, 5], etc.)

SARAH is the variance reduction algorithm which was originally proposed in [11] in the *convex case*. We introduce a slight modification to SARAH in Algorithm 1.1, called NC-SARAH, for the non-convex case. SARAH's as well as NC-SARAH's iterations are divided into an outer loop where a full gradient is computed and an inner loop where only one stochastic gradient is computed. We use upper index (s) to indicate the s -th outer loop and lower index t to indicate the t -th iteration in the inner loop. The key update rule, which

*IBM T.J. Watson Research Center, Yorktown Heights, NY (LamNguyen.MLTD@ibm.com).

[†]University of Connecticut, Storrs, CT (marten.van_dijk@uconn.edu).

[‡]IBM T.J. Watson Research Center, Yorktown Heights, NY (phandu@us.ibm.com).
[§]

³University of Connecticut, Storrs, CT (phuongha.ntu@gmail.com).

[¶]Massachusetts Institute of Technology, Cambridge, MA (twweng@mit.edu).

^{||}IBM T.J. Watson Research Center, Yorktown Heights, NY (jayant@us.ibm.com).

38 is called the SARAH update [11] for the inner loop, is

39 (1.2)
$$v_t^{(s)} = \nabla f_{i_t}(w_t^{(s)}) - \nabla f_{i_t}(w_{t-1}^{(s)}) + v_{t-1}^{(s)},$$

41 where i_t is chosen uniformly at random in $[n]$. The computed $v_t^{(s)}$ is used to update $w_{t+1}^{(s)} =$
 42 $w_t^{(s)} - \eta v_t^{(s)}$. In NC-SARAH, after m iterations in the inner loop, the outer loop remembers
 43 the last computed $w_{m+1}^{(s)}$ and starts its loop anew – first with a full gradient computation
 44 before again entering the inner loop with updates (1.2). Instead of remembering $\tilde{w}_s = w_{m+1}^{(s)}$
 45 for the next outer loop, the original SARAH algorithm in [11] uses $\tilde{w}_s = w_t^{(s)}$ with t chosen
 46 uniformly at random from $\{0, 1, \dots, m\}$; the authors of [11] chose to do this in order to being
 47 able to analyze the convergence rate for a single outer loop.

Algorithm 1.1 NC-SARAH

Parameters: the learning rate $\eta > 0$, the inner loop size m , and the outer loop size S

Initialize: \tilde{w}_0

Iterate:

for $s = 1, 2, \dots, S$ **do**

$$w_0^{(s)} = \tilde{w}_{s-1}$$

$$v_0^{(s)} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_0^{(s)})$$

$$w_1^{(s)} = w_0^{(s)} - \eta v_0^{(s)}$$

Iterate:

for $t = 1, \dots, m$ **do**

Sample i_t uniformly at random from $[n]$

$$v_t^{(s)} = \nabla f_{i_t}(w_t^{(s)}) - \nabla f_{i_t}(w_{t-1}^{(s)}) + v_{t-1}^{(s)}$$

$$w_{t+1}^{(s)} = w_t^{(s)} - \eta v_t^{(s)}$$

end for

$$\text{Set } \tilde{w}_s = w_{m+1}^{(s)}$$

end for

48 We notice that in [12] SARAH was extended to deal with *mini-batch* updates by, instead
 49 of choosing a single sample i_t in (1.2), we choose b samples uniformly at random from $[n]$
 50 for updating v_t in the inner loop. This gives a SARAH update rule for mini-batches:

51 (1.3)
$$v_t^{(s)} = \frac{1}{b} \sum_{i \in I_t} [\nabla f_i(w_t^{(s)}) - \nabla f_i(w_{t-1}^{(s)})] + v_{t-1}^{(s)},$$

 52

53 where we choose a mini-batch $I_t \subseteq [n]$ of size b uniformly at random at each iteration of the
 54 inner loop. NC-SARAH in Algorithm 1.1 is for the single batch case and if we replace the
 55 update rule in the inner loop by (1.3) we get NC-SARAH for the mini-batch case. SARAH
 56 in [12] for mini-batches was analyzed for the non-convex case for only a single outer loop
 57 giving a total complexity of $\mathcal{O}(n + \frac{L^2}{\epsilon^2})$. With our modification to \tilde{w}_s in NC-SARAH we are
 58 able to provide a ‘multiple outer loop’ analysis for the non-convex case for single batches and
 59 mini-batches.

60 SPIDER [5], a recent variant of SARAH for the non-convex case, is the first work that
 61 achieves the best known total¹ complexity of $\mathcal{O}(n + L\sqrt{n}/\epsilon)$ for the non-convex case, where
 62 L is the Lipschitz constant of the gradients. Its complexity matches the lower-bound worst

¹Measured as the total number of gradient computations needed to achieve an ϵ -accurate solution.

63 case complexity of $\mathcal{O}(\sqrt{n}/\epsilon)$ in [5] up to a constant factor when $n \leq \mathcal{O}(\epsilon^{-2})$. Another
 64 variant of SARAH [18] provides an improved version of SPIDER called SpiderBoost which
 65 allows a larger learning rate but restricting on the choice of the mini-batch size. Both SPIDER
 66 and SpiderBoost use the SARAH update rule (1.2) as originally proposed in [11] and use
 67 the mini-batch version of the update rule (1.3) in [12]. SPIDER and SpiderBoost do not
 68 divide into an outer loop and inner loop like SARAH, although SPIDER and SpiderBoost do
 69 similarly perform a full gradient update after a certain fixed number of iterations.

70 The drawback of SPIDER is the utilization of a small learning rate which depends on ϵ ,
 71 but it offers flexibility in the range of mini-batch sizes for the inner loop $[1, \sqrt{n}]$. SpiderBoost
 72 has a larger stepsize independent on ϵ , which gives a big practical improvement for solving
 73 real applications. However, one needs to fix the mini-batch size of \sqrt{n} for SpiderBoost,
 74 which limits the design space to such mini-batch sizes. Besides achieving the state-of-the-
 75 art asymptotic total complexity $\mathcal{O}(n + L\sqrt{n}/\epsilon)$ like SPIDER and SpiderBoost, our proposed
 76 variant of NC-SARAH mitigates at the same time both the learning rate limitation of SPIDER
 77 and the mini-batch limitation of SpiderBoost: In fact our learning rate is higher than those
 78 of SPIDER as well as SpiderBoost, and our mini-batch size for the inner loop can be freely
 79 selected from $[1, \sqrt{n}]$ (see Section 3 for more detail).

80 **Contributions:** We summarize our key contributions as follows.

81 1. **Smooth Non-Convex.**

- 82 • We provide a new convergence analysis for a new variant of the SARAH al-
 83 gorithm (NC-SARAH) for non-convex problems. We show that NC-SARAH
 84 achieves the state-of-the-art total complexity² for finding a first-order station-
 85 ary point in the non-convex case based on **only** the average smooth assumption;
 86 see Theorem 1 and Corollary 2 for the single batch case and Theorem 2 and
 87 Corollary 3 for the mini-batch case. We notice that our convergence analysis
 88 framework is simple and intuitive (Lemma 1 and Theorem 1).
- 89 • We rigorously show that, given a fixed mini-batch size for the inner loop and
 90 given a fixed number of inner loop iterations, NC-SARAH can adopt an *or-
 91 ders of magnitude larger learning rate* compared to SPIDER (Corollary 5) and
 92 a *strictly larger learning rate* compared to SpiderBoost (Corollary 6). Nu-
 93 matical experiments show how NC-SARAH outperforms both SPIDER and
 94 SpiderBoost (Section 3.3).
- 95 • NC-SARAH allows *a range of mini-batch sizes* for the inner loop similar to
 96 SPIDER (Section 3). In this sense NC-SARAH adopts the advantage of SPI-
 97 DER and, unlike SpiderBoost, does not need to give up on the flexibility of
 98 choosing mini-batch sizes in order to achieve practical large learning rates.
 99 Numerical experiments show that a flexible mini-batch size improves perfor-
 100 mance – a mini-batch size of about $n^{0.1} - n^{0.2}$ rather than the fixed mini-batch
 101 size of $n^{0.5}$ in SpiderBoost leads to best performance in our case study (Sec-
 102 tion 3.3).

103 2. **Smooth Convex.** In order to complete the picture, we study SARAH+ [11] which
 104 was designed as a variant of SARAH for convex optimization. SARAH+ provides a
 105 stopping criteria for the inner loop and shows the efficiency over SARAH. SARAH+
 106 suggests to empirically choose parameter without theoretical guarantee. We propose
 107 a novel variant of SARAH+ called SARAH++. Here, we study the *iteration com-
 108 plexity* measured by the total number of iterations (which counts one full gradient
 109 computation as adding one iteration to the complexity) – and leave an analysis of the

²State-of-the-art complexity matches the lower-bound worst case complexity of $\mathcal{O}(\sqrt{n}/\epsilon)$ in [5] up to a con-
 stant factor when $n \leq \mathcal{O}(\epsilon^{-2})$.

total complexity as an open problem. For SARAH++ we show a sublinear convergence rate in the general convex case (Theorem 3) and a linear convergence rate in the strongly convex case (Theorem 4). SARAH itself may already lead to good convergence and there may no need to introduce SARAH++; in numerical experiments we show the advantage of SARAH++ over SARAH. We further propose a practical version called *SARAH Adaptive* which improves the performance of SARAH and SARAH++ for convex problems – numerical experiments on various data sets show good overall performance.

3. **Generalized Gradient Descent.** For the convergence analysis of NC-SARAH for the non-convex case and SARAH++ for the convex case we show that the analysis generalizes the total complexity of Gradient Descent (GD) (Remarks 1, 2, and 3), i.e., the analysis reproduces known total complexity results of GD. Up to the best of our knowledge, this is the first variance reduction method having this property.

Table 1: Comparison of results on the total complexity for smooth nonconvex optimization

| Method | Total Complexity | Additional assumption |
|------------------------------|---|--|
| GD [10] | $\mathcal{O}\left(\frac{n}{\epsilon}\right)$ | None |
| SVRG [13] | $\mathcal{O}\left(n + \frac{n^{2/3}}{\epsilon}\right)$ | None |
| SCSG [8] | $\mathcal{O}\left((\frac{\sigma}{\epsilon} \wedge n) + \frac{1}{\epsilon} (\frac{\sigma}{\epsilon} \wedge n)^{2/3}\right)$ $\mathcal{O}\left(n + \frac{n^{2/3}}{\epsilon}\right)$ | Bounded variance |
| SNVRG [19] | $\mathcal{O}\left(\log^3\left(\frac{\sigma}{\epsilon} \wedge n\right) \left[(\frac{\sigma}{\epsilon} \wedge n) + \frac{1}{\epsilon} (\frac{\sigma}{\epsilon} \wedge n)^{1/2} \right]\right)$ $\mathcal{O}\left(\log^3(n) \left(n + \frac{\sqrt{n}}{\epsilon} \right)\right)$ | Bounded variance None ($\sigma \rightarrow \infty$) |
| SPIDER [5] | $\mathcal{O}\left(n + \frac{\sqrt{n}}{\epsilon}\right)$ | None |
| SpiderBoost [18] | $\mathcal{O}\left(n + \frac{\sqrt{n}}{\epsilon}\right)$ | None |
| NC-SARAH (this paper) | $\mathcal{O}\left(n + \frac{\sqrt{n}}{\epsilon}\right)$ | None |

1.1. **Related Work.** Table 1³ shows the comparison of results on the total complexity for smooth non-convex optimization. (a) Each of the complexities in Table 1 also depends on the Lipschitz constant L , however, since we consider smooth optimization, it is custom to assume/design $L = \mathcal{O}(1)$ and we therefore ignore the dependency on L in the complexity results. (b) Although many algorithms have appeared during the past few years, we only compare algorithms having a convergence result which only supposes the smooth assumption. (c) Among algorithms with convergence results that only suppose the smooth assumption, Table 1 only mentions recent state-of-the-art results. (d) Although the bounded variance assumption $\mathbb{E}[\|\nabla f_i(w) - \nabla F(w)\|^2] \leq \sigma^2$ is acceptable in many existing literature, this additional assumption limits the applicability of these convergence results since it adds dependence on σ which can be arbitrarily large. For fair comparison with convergence analysis without the bounded variance assumption, σ must be set to go to infinity – and this is what is mentioned in Table 1. As an example, from Table 1 we observe that SCSG has an advantage over SVRG only if $\sigma = \mathcal{O}(1)$ but, theoretically, by removing the bounded variance assumption, it has the same total complexity as SVRG if $\sigma \rightarrow \infty$.

From Table 1, we observe that NC-SARAH, SPIDER and SpiderBoost achieve the total complexity of $\mathcal{O}(n + \sqrt{n}/\epsilon)$ and dominate the complexity of all other algorithms. Indeed, its complexity matches the lower-bound worst case complexity of $\mathcal{O}(\sqrt{n}/\epsilon)$ in [5] up to a constant factor when $n \leq \mathcal{O}(\epsilon^{-2})$. We note that SPIDER and SpiderBoost can easily be rewritten by using an inner loop and outer loop algorithm description similar to SARAH (see

³ $a \wedge b$ is defined as $\min\{a, b\}$ and $a \vee b$ is defined as $\max\{a, b\}$

Table 2: Comparison properties among SPIDER, SpiderBoost, and NC-SARAH

| Method | Complexity | Mini-batch size b and Number of inner loop iterations m | Learning Rate |
|--------------------------|---|--|---------------------------|
| SPIDER [5] | $\mathcal{O}\left(n + \frac{\sqrt{n}}{\epsilon}\right)$ | $b = n^{1/2-\gamma}$ and $m = n^{1/2+\gamma}$ $\gamma \in [0, 1/2]$ | Dependent on ϵ |
| SpiderBoost [18] | $\mathcal{O}\left(n + \frac{\sqrt{n}}{\epsilon}\right)$ | $b = n^{1/2}$ and $m = n^{1/2}$ | Independent on ϵ |
| NC-SARAH (this paper) | $\mathcal{O}\left(n + \frac{\sqrt{n}}{\epsilon}\right)$ | $b = n^{1/2-\gamma}$ and $m = n^{1/2+\gamma}$ $\gamma \in [0, 1/2]$ | Independent on ϵ |

143 also Algorithm 1.1). For consistency, we will use the term “inner loop” to indicate where the
 144 SARAH update rule is used in these three algorithms. The advantages of NC-SARAH over
 145 SPIDER and SpiderBoost, respectively, are shown in Table 2. Both SPIDER and NC-SARAH
 146 allow a mini-batch size for the update rule in the inner loop in $b \in [1, \sqrt{n}]$.⁴ SpiderBoost is
 147 restricted in choosing a mini-batch size $b = \sqrt{n}$ for the inner loop while NC-SARAH like
 148 SPIDER has more choices. Our experimental results confirm that the choice of $b = \sqrt{n}$ and
 149 $m = \sqrt{n}$ of SpiderBoost is not the best choice (see Section 3.3). NC-SARAH outperforms
 150 SPIDER in that it can choose a much larger learning rate for the same mini-batch size b and
 151 number of inner loop iterations m . This is because the choice of NC-SARAH’s learning rate
 152 does not depend on ϵ while SPIDER does; the smaller learning rate of SPIDER makes it
 153 converge slowly to small ϵ -accurate solution. Even though the learning rate of SpiderBoost
 154 also does not depend on ϵ , we show that for the same mini-batch size $b = \sqrt{n}$ a number of
 155 inner loop iterations $m = \sqrt{n}$ NC-SARAH can still choose a larger learning rate.

156 **1.2. Paper Organization.** The rest of the paper is organized as follows. Section 2
 157 gives the convergence analysis of NC-SARAH in the non-convex case for both single batch
 158 and mini-batch cases. Section 3 shows the advantages of NC-SARAH over SPIDER and
 159 SpiderBoost in detail. In Section 4, we provide the convergence analysis of SARAH++ in
 160 the convex case and its iteration complexity. Numerical experiments are also given in this
 161 section to show the good performance of SARAH++ and SARAH Adaptive over the original
 162 SARAH. We conclude the paper and discuss future work in Section 5.

163 **2. Non-Convex Case: Convergence Analysis of NC-SARAH.** We will analyze NC-
 164 SARAH for smooth non-convex optimization, i.e., we study (1.1) with the following *average*
 165 *smooth* assumption

166 **ASSUMPTION 1** (average- L -smooth). *The objective function F is L -average-smooth,
 167 i.e., there exists a constant $L > 0$ such that, $\forall w, w' \in \mathbb{R}^d$,*

$$(2.1) \quad \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(w) - \nabla f_i(w')\|^2 \leq L^2 \|w - w'\|^2.$$

170 We notice that, the above assumption is weaker than the assumption on L -smoothness
 171 of each f_i , $i = 1, \dots, n$. Throughout this paper for non-convex results, we only consider
 172 Assumption 1 and no other assumptions. We stress that our convergence analysis only relies
 173 on the above average smooth assumption without bounded variance assumption (as required
 174 in [8, 19]). We note that Assumption 1 implies that F is L -smooth, that is, there exists a
 175 constant $L > 0$ such that, $\forall w, w' \in \mathbb{R}^d$, $\|\nabla F(w) - \nabla F(w')\| \leq L\|w - w'\|$. By Theorem

⁴According to our analysis, NC-SARAH also has the option to choose a mini-batch size $b \in (\sqrt{n}, n]$, but for such a choice we cannot attain a total complexity of $\mathcal{O}\left(\frac{\sqrt{n}}{\epsilon} \vee n\right)$.

176 2.1.5 in [10], we obtain

$$177 \quad (2.2) \quad F(w) \leq F(w') + \nabla F(w')^T(w - w') + \frac{L}{2}\|w - w'\|^2.$$

179 **2.1. Single batch case.** We start analyzing NC-SARAH (Algorithm 1.1) for the case
180 where we choose a single sample i_t uniformly at random from $[n]$ in the inner loop. We then
181 provide the following key lemma.

182 LEMMA 1. Suppose that Assumption 1 holds. Consider a single outer loop iteration in
183 NC-SARAH (Algorithm 1.1) with $\eta \leq \frac{2}{L(\sqrt{1+4m+1})}$. Then, for any $s \geq 1$, we have

$$184 \quad (2.3) \quad \mathbb{E}[F(w_{m+1}^{(s)})] \leq \mathbb{E}[F(w_0^{(s)})] - \frac{\eta}{2} \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2].$$

185 *Proof.* We use some parts of the proof in [12]. By Assumption 1 and $w_{t+1}^{(s)} = w_t^{(s)} -$
186 $\eta v_t^{(s)}$, for any $s \geq 1$, we have

$$\begin{aligned} 187 \quad \mathbb{E}[F(w_{t+1}^{(s)})] &\stackrel{(2.2)}{\leq} \mathbb{E}[F(w_t^{(s)})] - \eta \mathbb{E}[\nabla F(w_t^{(s)})^T v_t^{(s)}] + \frac{L\eta^2}{2} \mathbb{E}[\|v_t^{(s)}\|^2] \\ 188 \quad &= \mathbb{E}[F(w_t^{(s)})] - \frac{\eta}{2} \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] + \frac{\eta}{2} \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] \\ 189 \quad (2.4) \quad &- \left(\frac{\eta}{2} - \frac{L\eta^2}{2} \right) \mathbb{E}[\|v_t^{(s)}\|^2], \end{aligned}$$

191 where the last equality follows from the fact $a^T b = \frac{1}{2} [\|a\|^2 + \|b\|^2 - \|a - b\|^2]$, for any
192 $a, b \in \mathbb{R}^d$. By summing over $t = 0, \dots, m$, we have

$$\begin{aligned} 193 \quad \mathbb{E}[F(w_{m+1}^{(s)})] &\leq \mathbb{E}[F(w_0^{(s)})] - \frac{\eta}{2} \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \\ 194 \quad (2.5) \quad &+ \frac{\eta}{2} \left(\sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] - (1 - L\eta) \sum_{t=0}^m \mathbb{E}[\|v_t^{(s)}\|^2] \right). \end{aligned}$$

196 Now, we would like to determine η such that the expression in (2.5)

$$197 \quad \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] - (1 - L\eta) \sum_{t=0}^m \mathbb{E}[\|v_t^{(s)}\|^2] \leq 0.$$

199 Let $\mathcal{F}_j = \sigma(w_0, w_1, \dots, w_j)$ be the σ -algebra generated by w_0, w_1, \dots, w_j . Note that
200 \mathcal{F}_j also contains all information of v_0, \dots, v_{j-1} . We have

$$\begin{aligned} 201 \quad \mathbb{E}[\|v_j^{(s)} - v_{j-1}^{(s)}\|^2 | \mathcal{F}_j] &\stackrel{(1.2)}{=} \mathbb{E}[\|\nabla f_{i_j}(w_j^{(s)}) - \nabla f_{i_j}(w_{j-1}^{(s)})\|^2 | \mathcal{F}_j] \\ 202 \quad &\stackrel{(2.1)}{\leq} L^2 \|w_j^{(s)} - w_{j-1}^{(s)}\|^2 = L^2 \eta^2 \|v_{j-1}^{(s)}\|^2, \quad j \geq 1. \end{aligned}$$

204 Taking the expectations to both sides yields

$$205 \quad (2.6) \quad \mathbb{E}[\|v_j^{(s)} - v_{j-1}^{(s)}\|^2] \leq L^2 \eta^2 \mathbb{E}[\|v_{j-1}^{(s)}\|^2].$$

207 Hence, by Lemma 3, we have

$$208 \quad \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] \leq \sum_{j=1}^t \mathbb{E}[\|v_j^{(s)} - v_{j-1}^{(s)}\|^2] \stackrel{(2.6)}{\leq} L^2 \eta^2 \sum_{j=1}^t \mathbb{E}[\|v_{j-1}^{(s)}\|^2].$$

210 Note that $\|\nabla F(w_0^{(s)}) - v_0^{(s)}\|^2 = 0$. By summing over $t = 0, \dots, m$ ($m \geq 1$), we have

$$211 \quad \sum_{t=0}^m \mathbb{E}\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2 \leq L^2\eta^2 \left[m\mathbb{E}\|v_0^{(s)}\|^2 + (m-1)\mathbb{E}\|v_1^{(s)}\|^2 + \dots + \mathbb{E}\|v_{m-1}^{(s)}\|^2 \right]. \\ 212$$

213 By choosing $\eta \leq \frac{2}{L(\sqrt{1+4m}+1)}$, we have

$$214 \quad \begin{aligned} & \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] - (1-L\eta) \sum_{t=0}^m \mathbb{E}[\|v_t^{(s)}\|^2] \\ 215 & \leq L^2\eta^2 \left[m\mathbb{E}\|v_0^{(s)}\|^2 + (m-1)\mathbb{E}\|v_1^{(s)}\|^2 + \dots + \mathbb{E}\|v_{m-1}^{(s)}\|^2 \right] \\ 216 & \quad - (1-L\eta) \left[\mathbb{E}\|v_0^{(s)}\|^2 + \mathbb{E}\|v_1^{(s)}\|^2 + \dots + \mathbb{E}\|v_m^{(s)}\|^2 \right] \\ 217 (2.7) \quad & \leq \left[L^2\eta^2 m - (1-L\eta) \right] \sum_{t=1}^m \mathbb{E}[\|v_{t-1}^{(s)}\|^2] \leq 0, \\ 218 \end{aligned}$$

219 since $\eta = \frac{2}{L(\sqrt{1+4m}+1)}$ is a root of equation $L^2\eta^2 m - (1-L\eta) = 0$. Therefore, with
220 $\eta \leq \frac{2}{L(\sqrt{1+4m}+1)}$, we have

$$221 \quad \mathbb{E}[F(w_{m+1}^{(s)})] \leq \mathbb{E}[F(w_0^{(s)})] - \frac{\eta}{2} \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2]. \quad \square \\ 222$$

223 The above result is for a single outer loop iteration of NC-SARAH, which includes a full
224 gradient step together with the inner loop. Since the outer loop iteration concludes with $\tilde{w}_s =$
225 $w_{m+1}^{(s)}$, and $\tilde{w}_{s-1} = w_0^{(s)}$, we have $\mathbb{E}[F(\tilde{w}_s)] \leq \mathbb{E}[F(\tilde{w}_{s-1})] - \frac{\eta}{2} \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2]$.
226 Summing over $1 \leq s \leq S$ gives

$$227 (2.8) \quad \mathbb{E}[F(\tilde{w}_S)] \leq \mathbb{E}[F(\tilde{w}_0)] - \frac{\eta}{2} \sum_{s=1}^S \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2]. \\ 228$$

229 This proves our main result:

230 THEOREM 1. Suppose that Assumption 1 holds. Consider NC-SARAH (Algorithm 1.1)
231 with $\eta \leq \frac{2}{L(\sqrt{1+4m}+1)}$. Then, for any given \tilde{w}_0 , we have

$$232 \quad \frac{1}{(m+1)S} \sum_{s=1}^S \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \leq \frac{2}{\eta[(m+1)S]} [F(\tilde{w}_0) - F^*], \\ 233$$

234 where F^* is any lower bound of F , and $w_t^{(s)}$ is the t -th iteration in the s -th outer loop.

235 The proof easily follows from (2.8) since F^* is a lower bound of F (that is, $\mathbb{E}[F(\tilde{w}_S)] \geq F^*$).
236 We note that the term $\frac{1}{(m+1)S} \sum_{s=1}^S \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2]$ is simply the average of the
237 expectation of the squared norms of the gradients of all the iteration results generated by
238 NC-SARAH. For nonconvex problems, our goal is to achieve

$$239 (2.9) \quad \frac{1}{(m+1)S} \sum_{s=1}^S \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \leq \epsilon. \\ 240$$

241 We note that, for simplicity, if \bar{w}_s is chosen uniformly at random from all the iterations
 242 generated by NC-SARAH, we are able to have accuracy $\mathbb{E}[\|\nabla F(\bar{w}_s)\|^2] \leq \epsilon$. From Theorem
 243 1 with $\eta = \mathcal{O}(1/\sqrt{m+1})$ we infer that (2.9) can be realized for $S = \mathcal{O}(\frac{1}{\epsilon\sqrt{m+1}} \vee 1)$. The
 244 total complexity of NC-SARAH is equal to $S(n+2m)$ which proves:

245 COROLLARY 1. Suppose that Assumption 1 holds. Let us consider NC-SARAH (Algo-
 246 rithm 1.1) with $\eta = \frac{2}{L(\sqrt{1+4m+1})}$ where m is the inner loop size. Then, in order to achieve
 247 an ϵ -accurate solution, the total complexity is $\mathcal{O}\left(\left[\left(\frac{n+2m}{\sqrt{m+1}}\right)\frac{1}{\epsilon}\right] \vee [n+2m]\right)$.

248 Proof. To achieve $\frac{1}{(m+1)S} \sum_{s=1}^S \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \leq \epsilon$, it is sufficient to prove

$$249 \quad (2.10) \quad \frac{2}{\eta[(m+1)S]} [F(\tilde{w}_0) - F^*] \leq \epsilon.$$

251 Notice that $\eta\sqrt{m+1} = \frac{2}{L} \frac{\sqrt{m+1}}{\sqrt{1+4m+1}} \leq \frac{2}{L}$. Hence, in order to achieve (2.10), we need

$$252 \quad S \geq \frac{2}{\eta[(m+1)\epsilon]} [F(\tilde{w}_0) - F^*] \geq \frac{L[F(\tilde{w}_0) - F^*]}{(\sqrt{m+1})} \frac{1}{\epsilon}. \quad \square$$

253 Together with the requirement $S \geq 1$, we can choose $S = \mathcal{O}\left(\left[\frac{1}{\sqrt{m+1}} \cdot \frac{1}{\epsilon}\right] \vee 1\right)$. Therefore,
 the total complexity to achieve ϵ -accurate solution is

$$(n+2m)S = \mathcal{O}\left(\left[\left(\frac{n+2m}{\sqrt{m+1}}\right)\frac{1}{\epsilon}\right] \vee [n+2m]\right).$$

254 The total complexity can be minimized over the inner loop size m . By choosing $m = n$,
 255 we achieve the minimal total complexity:

256 COROLLARY 2. Suppose that Assumption 1 holds. Let us consider NC-SARAH (Algo-
 257 rithm 1.1) with $\eta = \frac{2}{L(\sqrt{1+4m+1})}$ where m is the inner loop size and chosen equal to $m = n$.
 258 Then, in order to achieve an ϵ -accurate solution, the total complexity is $\mathcal{O}\left(\frac{\sqrt{n}}{\epsilon} \vee n\right)$.

259 REMARK 1. The total complexity in Corollary 1 covers all choices for the inner loop
 260 size m . For example, in the case of $m = 0$, NC-SARAH recovers the Gradient Descent
 261 (GD) algorithm which has total complexity $\mathcal{O}\left(\frac{n}{\epsilon}\right)$. Theorem 1 for $m = 0$ also recovers the
 262 requirement on the learning rate for GD, which is $\eta \leq \frac{1}{L}$.

263 The above results explain the relationship between NC-SARAH and GD and explains
 264 the advantages of the inner loop and outer loop of NC-SARAH. NC-SARAH becomes more
 265 beneficial in ML applications where n is large.

266 **2.2. Mini-batch case.** The above results can be extended to the *mini-batch* case where
 267 instead of (1.2) the update rule (1.3) is used as explained in the introduction.

268 THEOREM 2. Suppose that Assumption 1 holds. Consider NC-SARAH (Algorithm 1.1)
 269 by replacing the update of v_t in the inner loop by (1.3) with

$$270 \quad (2.11) \quad \eta \leq \frac{2}{L \left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} + 1 \right)}.$$

272 Then, for any given \tilde{w}_0 , we have

$$273 \quad \frac{1}{(m+1)S} \sum_{s=1}^S \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \leq \frac{2}{\eta[(m+1)S]} [F(\tilde{w}_0) - F^*],$$

275 where F^* is any lower bound of F , and $w_t^{(s)}$ is the t -th iteration in the s -th outer loop.

276 *Proof.* Consider $v_t^{(s)}$ defined by (1.3) in NC-SARAH (Algorithm 1.1) for any $s \geq 1$.

277 Then by Lemma 1, for any $t \geq 1$,

(2.12)

$$278 \quad \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] = \sum_{j=1}^t \mathbb{E}[\|v_j^{(s)} - v_{j-1}^{(s)}\|^2] - \sum_{j=1}^t \mathbb{E}[\|\nabla F(w_j^{(s)}) - \nabla F(w_{j-1}^{(s)})\|^2].$$

280 Following the proof of Lemma 1, we would like to determine η such that the expression
281 in (2.5)

$$282 \quad \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] - (1 - L\eta) \sum_{t=0}^m \mathbb{E}[\|v_t^{(s)}\|^2] \leq 0.$$

284 Let

$$285 \quad (2.13) \quad \xi_t = \nabla f_t(w_j^{(s)}) - \nabla f_t(w_{j-1}^{(s)}).$$

287 Let $\mathcal{F}_j = \sigma(w_0^{(s)}, I_1, I_2, \dots, I_{j-1})$ be the σ -algebra generated by $w_0^{(s)}, I_1, I_2, \dots, I_{j-1}$;
288 $\mathcal{F}_0 = \mathcal{F}_1 = \sigma(w_0^{(s)})$. Note that \mathcal{F}_j also contains all the information of $w_0^{(s)}, \dots, w_j^{(s)}$ as well
289 as $v_0^{(s)}, \dots, v_{j-1}^{(s)}$. We have

$$\begin{aligned} 290 \quad & \mathbb{E}[\|v_j^{(s)} - v_{j-1}^{(s)}\|^2 | \mathcal{F}_j] - \|\nabla F(w_j^{(s)}) - \nabla F(w_{j-1}^{(s)})\|^2 \\ 291 \quad & \stackrel{(1.3)}{=} \mathbb{E}\left[\left\|\frac{1}{b} \sum_{i \in I_j} [\nabla f_i(w_j^{(s)}) - \nabla f_i(w_{j-1}^{(s)})]\right\|^2 \middle| \mathcal{F}_j\right] - \left\|\frac{1}{n} \sum_{i=1}^n [\nabla f_i(w_j^{(s)}) - \nabla f_i(w_{j-1}^{(s)})]\right\|^2 \\ 292 \quad & \stackrel{(2.13)}{=} \mathbb{E}\left[\left\|\frac{1}{b} \sum_{i \in I_j} \xi_i\right\|^2 \middle| \mathcal{F}_j\right] - \left\|\frac{1}{n} \sum_{i=1}^n \xi_i\right\|^2 = \frac{1}{b^2} \mathbb{E}\left[\sum_{i \in I_j} \sum_{k \in I_j} \xi_i^T \xi_k \middle| \mathcal{F}_j\right] - \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^n \xi_i^T \xi_k \\ 293 \quad & = \frac{1}{b^2} \mathbb{E}\left[\sum_{i \neq k \in I_j} \xi_i^T \xi_k + \sum_{i \in I_j} \xi_i^T \xi_i \middle| \mathcal{F}_j\right] - \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^n \xi_i^T \xi_k \\ 294 \quad & = \frac{1}{b^2} \left[\frac{b}{n(n-1)} \sum_{i \neq k} \xi_i^T \xi_k + \frac{b}{n} \sum_{i=1}^n \xi_i^T \xi_i \right] - \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^n \xi_i^T \xi_k \\ 295 \quad & = \frac{1}{b^2} \left[\frac{b}{n(n-1)} \sum_{i=1}^n \sum_{k=1}^n \xi_i^T \xi_k + \left(\frac{b}{n} - \frac{b}{n(n-1)} \right) \sum_{i=1}^n \xi_i^T \xi_i \right] - \frac{1}{n^2} \sum_{i=1}^n \sum_{k=1}^n \xi_i^T \xi_k \\ 296 \quad & = \frac{1}{bn} \left[\left(\frac{(b-1)}{(n-1)} - \frac{b}{n} \right) \sum_{i=1}^n \sum_{k=1}^n \xi_i^T \xi_k + \frac{(n-b)}{(n-1)} \sum_{i=1}^n \xi_i^T \xi_i \right] \\ 297 \quad & = \frac{1}{bn} \left(\frac{n-b}{n-1} \right) \left[-\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^n \xi_i^T \xi_k + \sum_{i=1}^n \xi_i^T \xi_i \right] \\ 298 \quad & = \frac{1}{bn} \left(\frac{n-b}{n-1} \right) \left[-n \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \right\|^2 + \sum_{i=1}^n \|\xi_i\|^2 \right] \\ 299 \quad & \leq \frac{1}{b} \left(\frac{n-b}{n-1} \right) \frac{1}{n} \sum_{i=1}^n \|\xi_i\|^2 \end{aligned}$$

300 $\stackrel{(2.13)}{=} \frac{1}{b} \left(\frac{n-b}{n-1} \right) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(w_j^{(s)}) - \nabla f_i(w_{j-1}^{(s)})\|^2 \stackrel{(2.1)}{\leq} \frac{1}{b} \left(\frac{n-b}{n-1} \right) L^2 \eta^2 \|v_{j-1}^{(s)}\|^2.$

301

302 Hence, by taking expectation, we have

303 $\mathbb{E}[\|v_j^{(s)} - v_{j-1}^{(s)}\|^2] - \mathbb{E}[\|\nabla F(w_j^{(s)}) - \nabla F(w_{j-1}^{(s)})\|^2] \leq \frac{1}{b} \left(\frac{n-b}{n-1} \right) L^2 \eta^2 \mathbb{E}[\|v_{j-1}^{(s)}\|^2].$

304

305 By (2.12), for $t \geq 1$,

306 $\mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] = \sum_{j=1}^t \mathbb{E}[\|v_j^{(s)} - v_{j-1}^{(s)}\|^2] - \sum_{j=1}^t \mathbb{E}[\|\nabla F(w_j^{(s)}) - \nabla F(w_{j-1}^{(s)})\|^2]$

307 $\leq \frac{1}{b} \left(\frac{n-b}{n-1} \right) L^2 \eta^2 \sum_{j=1}^t \mathbb{E}[\|v_{j-1}^{(s)}\|^2].$

308

309 Note that $\|\nabla F(w_0^{(s)}) - v_0^{(s)}\|^2 = 0$. By summing over $t = 0, \dots, m$ ($m \geq 1$), we have

310 $\sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2]$

311 $\leq \frac{1}{b} \left(\frac{n-b}{n-1} \right) L^2 \eta^2 \left[m \mathbb{E}\|v_0^{(s)}\|^2 + (m-1) \mathbb{E}\|v_1^{(s)}\|^2 + \dots + \mathbb{E}\|v_{m-1}^{(s)}\|^2 \right].$

312

313 By choosing $\eta \leq \frac{2}{L \left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} + 1 \right)}$, we have

314 $\sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] - (1 - L\eta) \sum_{t=0}^m \mathbb{E}[\|v_t^{(s)}\|^2]$

315 $\leq \frac{1}{b} \left(\frac{n-b}{n-1} \right) L^2 \eta^2 \left[m \mathbb{E}\|v_0^{(s)}\|^2 + (m-1) \mathbb{E}\|v_1^{(s)}\|^2 + \dots + \mathbb{E}\|v_{m-1}^{(s)}\|^2 \right]$

316 $- (1 - L\eta) \left[\mathbb{E}\|v_0^{(s)}\|^2 + \mathbb{E}\|v_1^{(s)}\|^2 + \dots + \mathbb{E}\|v_m^{(s)}\|^2 \right]$

317 (2.14) $\leq \left[\frac{1}{b} \left(\frac{n-b}{n-1} \right) L^2 \eta^2 m - (1 - L\eta) \right] \sum_{t=1}^m \mathbb{E}[\|v_{t-1}^{(s)}\|^2] \leq 0,$

318

319 since $\eta = \frac{2}{L \left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} + 1 \right)}$ is a root of equation $\frac{1}{b} \left(\frac{n-b}{n-1} \right) L^2 \eta^2 m - (1 - L\eta) = 0$.

320 Therefore, with $\eta \leq \frac{2}{L \left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} + 1 \right)}$, we have

321 $\mathbb{E}[F(w_{m+1}^{(s)})] \leq \mathbb{E}[F(w_0^{(s)})] - \frac{\eta}{2} \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2].$

322

323 Following the same derivation of Theorem 1, for any given \tilde{w}_0 , we could achieve the
324 desired result. \square

325 We can again derive similar corollaries as was done for Theorem 1.

COROLLARY 3. For the conditions in Theorem 2 with equality for η in (2.11), in order to achieve an ϵ -accurate solution the total complexity is

$$\mathcal{O} \left(\left[\left(\frac{n+2bm}{m+1} \right) \left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} \right) \frac{1}{\epsilon} \right] \vee [n+2bm] \right).$$

326 *Proof.* By Theorem 2, let $\eta = \frac{2}{L \left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} + 1 \right)}$. Hence, we have

$$327 \quad \frac{1}{(m+1)S} \sum_{s=1}^S \sum_{t=0}^m \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \leq \frac{2}{\eta[(m+1)S]} [F(\tilde{w}_0) - F^*]$$

$$328 \quad = \frac{\left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} + 1 \right)}{(m+1)S} L[F(\tilde{w}_0) - F^*]$$

$$329 \quad 330 \quad \leq \frac{\left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} \right)}{(m+1)S} 2L[F(\tilde{w}_0) - F^*] = \epsilon.$$

331 In order to achieve the ϵ -accurate solution, we need

$$332 \quad S = \frac{\left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} \right)}{(m+1)\epsilon} 2L[F(\tilde{w}_0) - F^*] = \mathcal{O} \left(\frac{\left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} \right)}{(m+1)\epsilon} \vee 1 \right),$$

333

334 since $S \geq 1$. Therefore, the total complexity is

$$335 \quad 336 \quad (n+b \cdot 2m)S = \mathcal{O} \left(\left[\left(\frac{n+2bm}{m+1} \right) \left(\sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} \right) \frac{1}{\epsilon} \right] \vee [n+2bm] \right). \quad \square$$

337 COROLLARY 4. For the conditions in Corollary 3 with $b = n^{1/2-\gamma}$ and $m = n^{1/2+\gamma}$,
 338 where $0 \leq \gamma \leq 1/2$, in order to achieve an ϵ -accurate solution the total complexity is
 339 $\mathcal{O} \left(\frac{\sqrt{n}}{\epsilon} \vee n \right)$.

340 *Proof.* Let $b = n^\alpha$ and $m = n^\beta$ where $0 \leq \alpha, \beta \leq 1$, we have

$$341 \quad \left(\frac{n+2bm}{m+1} \right) \sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1} \right)} = \left(\frac{n+2n^{\alpha+\beta}}{n^\beta+1} \right) \sqrt{1 + 4n^{\beta-\alpha} \left(\frac{n-n^\alpha}{n-1} \right)} \\ 342 \quad 343 \quad \leq \frac{n+2n^{\alpha+\beta}}{n^\beta} 2\sqrt{1+n^{\beta-\alpha}}.$$

344 If $\beta \geq \alpha$, we have

$$345 \quad \frac{n+2n^{\alpha+\beta}}{n^\beta} 2\sqrt{1+n^{\beta-\alpha}} \leq 2\sqrt{2} \left(\frac{n+2n^{\alpha+\beta}}{n^\beta} \right) n^{(\beta-\alpha)/2} \\ 346 \quad = 2\sqrt{2}(n^{1-\alpha/2-\beta/2} + 2n^{\alpha/2+\beta/2}).$$

348 In order to minimize the order of n , we need to choose $1 - \alpha/2 - \beta/2 = \alpha/2 + \beta/2$, which
 349 is equivalent to $\alpha + \beta = 1$ with $\beta \geq \alpha$. The best option is to choose $\alpha + \beta = 1$ with $\beta \geq 1/2$
 350 and $0 \leq \alpha \leq 1/2$ in order to achieve $\mathcal{O}(n^{1/2})$.

351 If $\beta \leq \alpha$, we have $\frac{n+2n^{\alpha+\beta}}{n^\beta} 2\sqrt{1+n^{\beta-\alpha}} \leq 2\sqrt{2}(n^{1-\beta} + 2n^\alpha)$. In order to minimize
 352 the order of n , we need to choose $1-\beta = \alpha$, which is equivalent to $\alpha + \beta = 1$ with $\beta \leq \alpha$.
 353 The best option is to choose $\beta = 1/2$ and $\alpha = 1/2$ in order to achieve $\mathcal{O}(n^{1/2})$.

354 Therefore, with $b = n^\alpha$ and $m = n^\beta$ where $\alpha + \beta = 1$ with $\beta \geq 1/2$ and $0 \leq \alpha \leq 1/2$,
 355 we have $\left(\frac{n+2bm}{m+1}\right) \sqrt{1 + \frac{4m}{b} \left(\frac{n-b}{n-1}\right)} = \mathcal{O}(n^{1/2})$.

356 By Corollary 3 with $bm = n^{\alpha+\beta} = n$, it implies the total complexity

$$357 \quad 358 \quad (n + b \cdot 2m)S = \mathcal{O}\left(\frac{\sqrt{n}}{\epsilon} \vee n\right).$$

359 Hence, by setting $\alpha = 1/2 - \gamma$ and $\beta = 1/2 + \gamma$ with $0 \leq \gamma \leq 1/2$, we obtain the corollary. \square

360 **REMARK 2.** *The choice of η in Theorem 2 is more general than in Theorem 1. For $b = n$
 361 and $m = m_0$, for some non-negative integer m_0 , it recovers the convergence rate of Gradient
 362 Descent with learning rate $\eta \leq \frac{1}{L}$ and total complexity $\mathcal{O}\left(\frac{n}{\epsilon}\right)$ (see Corollary 3).*

363 **3. Comparison of NC-SARAH, SPIDER, and SpiderBoost.** As shown in the previous
 364 section, like SPIDER [5] and SpiderBoost [18], also NC-SARAH enjoys the same asymptotic
 365 total complexity of $\mathcal{O}(n + \sqrt{n}/\epsilon)$.

366 In this section, we show practical advantages of NC-SARAH over SPIDER and Spider-
 367 Boost. By using our notation of b and m , the three algorithms have the following properties:

- **SPIDER:** For $0 \leq \gamma \leq 1/2$,

$$369 \quad 370 \quad (3.1) \quad b = n^{1/2-\gamma}, \quad m = n^{1/2+\gamma}, \quad \eta_t^{(s)} = \min\left\{\frac{\epsilon}{Ln^\gamma \|v_t^{(s)}\|}, \frac{1}{2Ln^\gamma}\right\},$$

371 where $v_t^{(s)}$ is the SARAH update (1.3), and $\eta_t^{(s)}$ denotes the learning rate of the t -th
 372 iteration in the inner loop of the s -th outer loop.

- **SpiderBoost:**

$$374 \quad 375 \quad (3.2) \quad b = n^{1/2}, \quad m = n^{1/2}, \quad \eta = \frac{1}{2L}.$$

- **NC-SARAH:** For $0 \leq \gamma \leq 1/2$,

$$377 \quad (3.3) \quad b = n^{1/2-\gamma}, \quad m = n^{1/2+\gamma}, \quad \eta = \frac{2}{L \left(\sqrt{1 + 4n^{2\gamma} \left(\frac{n-n^{1/2-\gamma}}{n-1} \right)} + 1 \right)}.$$

378 The following subsections analyze NC-SARAH in comparison to SPIDER and SpiderBoost,
 379 respectively.

380 **3.1. NC-SARAH vs SPIDER.** NC-SARAH and SPIDER have the same flexibility of
 381 choosing mini-batch size $b \in [1, \sqrt{n}]$. However, the learning rate of SPIDER can be quite
 382 small compared to the learning rate of NC-SARAH because SPIDER's learning rate scales
 383 linearly with ϵ for learning rates $\leq \frac{1}{2Ln^\gamma}$ and ϵ will be small especially when we want a small
 384 ϵ -accurate solution.

385 **COROLLARY 5.** *For the same mini-batch size b for the inner loop and the same number
 386 of inner loop iterations m , the learning rate choice of NC-SARAH is strictly larger than that
 387 of SPIDER when $n > 1$.*

389 *Proof.* We recall the SPIDER's setting:

390 $b = n^{1/2-\gamma}$, $m = n^{1/2+\gamma}$, $\eta_t^{(s)} = \min \left\{ \frac{\epsilon}{Ln^\gamma \|v_t^{(s)}\|}, \frac{1}{2Ln^\gamma} \right\}$,

391

392 where $v_t^{(s)}$ is the SARAH update (1.3), and $\eta_t^{(s)}$ is the learning rate of the s -outer loop and
393 t -th iteration in the inner loop; and the NC-SARAH's setting:

394 $b = n^{1/2-\gamma}$, $m = n^{1/2+\gamma}$, $\eta = \frac{2}{L \left(\sqrt{1 + 4n^{2\gamma} \left(\frac{n-n^{1/2-\gamma}}{n-1} \right)} + 1 \right)}$,

395

396 where $0 \leq \gamma \leq 1/2$.

397 Since $\eta_t^{(s)} = \min \left\{ \frac{\epsilon}{Ln^\gamma \|v_t^{(s)}\|}, \frac{1}{2Ln^\gamma} \right\} \leq \frac{1}{2Ln^\gamma}$. In order to achieve the desired result, it
398 is sufficient to show that, for $0 \leq \gamma \leq 1/2$,

399 (3.4)
$$\frac{2}{L \left(\sqrt{1 + 4n^{2\gamma} \left(\frac{n-n^{1/2-\gamma}}{n-1} \right)} + 1 \right)} > \frac{1}{2Ln^\gamma}.$$

400

401 This is equivalent to showing

402
$$4n^\gamma > \sqrt{1 + 4n^{2\gamma} \left(\frac{n-n^{1/2-\gamma}}{n-1} \right)} + 1$$

403
$$16n^{2\gamma} - 8n^\gamma + 1 > 1 + 4n^{2\gamma} \left(\frac{n-n^{1/2-\gamma}}{n-1} \right)$$

404
$$4 - \frac{2}{n^\gamma} > \frac{n-n^{1/2-\gamma}}{n-1} = 1 - \frac{n^{1/2-\gamma} - 1}{n-1}.$$

405

406 The last inequality clearly holds since $4 - \frac{2}{n^\gamma} \geq 2$. Therefore, we obtain (3.4). \square

407 **3.2. NC-SARAH vs SpiderBoost.** It is clear that, compared to SpiderBoost with only
408 $b = \sqrt{n}$, NC-SARAH has more flexibility of choosing mini-batch size $b = n^{1/2-\gamma}$ for some
409 $0 \leq \gamma \leq 1/2$. In order to compare the learning rate of NC-SARAH and SpiderBoost for the
410 same mini-batch size, we let $\gamma = 0$ in NC-SARAH's parameter setting; we obtain

411
$$b = \sqrt{n}$$
, $m = \sqrt{n}$, $\eta = \frac{2}{L \left(\sqrt{1 + 4 \left(\frac{n-n^{1/2}}{n-1} \right)} + 1 \right)}.$

412

413 We have the following corollary.

414 **COROLLARY 6.** For the same mini-batch size b for the inner loop and the same number
415 of inner loop iterations m , the learning rate choice of NC-SARAH is at least a factor $4/(\sqrt{5} + 1) = 1.236$ larger than that of SpiderBoost when $n > 1$.

417 *Proof.* We need to choose $\gamma = 0$ for NC-SARAH in order to have the same option as
418 SpiderBoost, i.e., $b = \sqrt{n}$ and $m = \sqrt{m}$. Hence, we have the learning rate of NC-SARAH

419
$$\frac{2}{L \left(\sqrt{1 + 4 \left(\frac{n-n^{1/2}}{n-1} \right)} + 1 \right)}.$$

420

421 We notice that $\frac{n-n^{1/2}}{n-1} < 1$ for $n > 1$. Hence, we have

$$422 \quad \frac{2}{L \left(\sqrt{1 + 4 \left(\frac{n-n^{1/2}}{n-1} \right)} + 1 \right)} > \frac{2}{L(\sqrt{5}+1)} = \frac{4}{(\sqrt{5}+1)} \cdot \frac{1}{2L}.$$

423 This completes the proof. \square

424 In Theorem 2 we can choose a smaller learning rate than the RHS of (2.11) in NC-SARAH. In this sense, the above corollary shows that SpiderBoost is a special case of NC-SARAH. The following subsection confirms numerically that the choice of $b = \sqrt{n}$ and $m = \sqrt{n}$ in SpiderBoost is not the best choice.

425 **3.3. Numerical Experiments.** In this subsection, we numerically verify the advantages
426 of NC-SARAH over SPIDER and SpiderBoost. We consider the binary classification problem
427 with non-convex loss function in [18] as follows

$$428 \quad (3.5) \quad \min_{w \in \mathbb{R}^d} \left\{ F(w) = \frac{1}{n} \sum_{i=1}^n \left[f_i(w) := \log(1 + \exp(-y_i x_i^T w)) + \frac{\lambda}{2} \sum_{j=1}^d \frac{w_j^2}{1+w_j^2} \right] \right\},$$

where $\{x_i, y_i\}_{i=1}^n$ is the training data with $x_i \in \mathbb{R}^d$ and $y_i \in \{-1, +1\}$, $i \in [n]$.

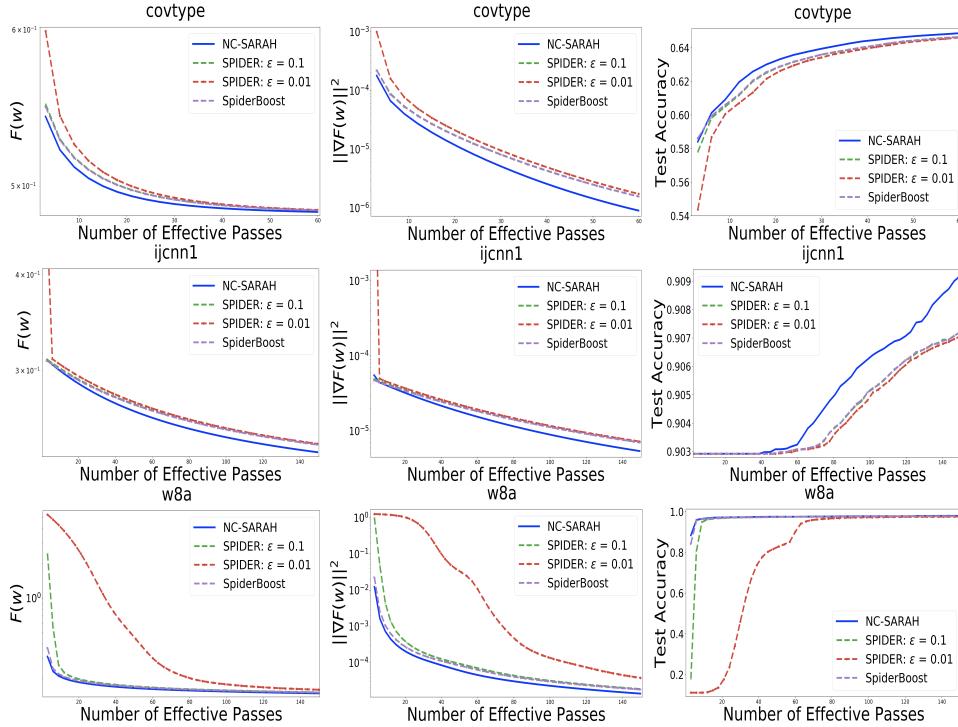


Fig. 1: Comparisons of $F(w)$, $\|\nabla F(w)\|^2$, and Test Accuracy among NC-SARAH, SPIDER with $\epsilon = 0.1$, SPIDER with $\epsilon = 0.01$, and SpiderBoost on covtype, ijcn1 and w8a datasets

434 We conducted experiments to demonstrate the advantage in performance of NC-SARAH
435 over SPIDER and SpiderBoost on the popular classification data sets covtype ($n = 406$, 708
436

437 training data; estimated $L \simeq 1.90$, *ijcnn1* ($n = 91,701$ training data; estimated $L \simeq 1.77$,
 438 and *w8a* ($n = 49,749$ training data, estimated $L \simeq 7.05$) from LIBSVM [3]. Since we
 439 only care about the non-convexity of each f_i , we can simply choose $\lambda = 0.01$. Additional
 440 experiments on more data sets are reported in the supplemental material.

441 Figure 1 shows comparisons of the values of $F(w)$, $\|\nabla F(w)\|^2$, and Test Accuracy
 442 among NC-SARAH, SPIDER with $\epsilon = 0.1$, SPIDER with $\epsilon = 0.01$, and SpiderBoost. In
 443 order to fit Spiderboost’s mini-batch size of $b = \sqrt{n}$ we choose $\gamma = 0$ in NC-SARAH
 444 and SPIDER. In this scenario, SPIDER with $\epsilon = 0.1$ performs similarly to SpiderBoost.
 445 We observe that NC-SARAH has better performance than both SPIDER and SpiderBoost
 446 since NC-SARAH is able to adopt a larger learning rate than those used in SPIDER and
 447 SpiderBoost (Corollaries 5 and 6) as shown in Figure 1. We experimented with 10 runs and
 448 reported the average results with the same initial point w_0 for all the algorithms.

449 SpiderBoost only allows a mini-batch size of $b = \sqrt{n}$ while NC-SARAH allows $b =$
 450 $n^{1/2-\gamma}$ with $m = n^{1/2+\gamma}$ for $\gamma \in [0, 0.5]$. Figure 2 shows the sensitivity of γ for NC-
 451 SARAH. We observe that the choice of $b = \sqrt{n}$ and $m = \sqrt{n}$ (or equivalently $\gamma = 0$) is
 452 not a good choice; for *covtype* $b \in [n^{0.1}, n^{0.2}]$ (or equivalently $\gamma = 0.3, 0.4$) leads to the best
 453 performance. This demonstrates that allowing a flexible range of mini-batch sizes beyond
 454 $b = \sqrt{n}$ is beneficial.

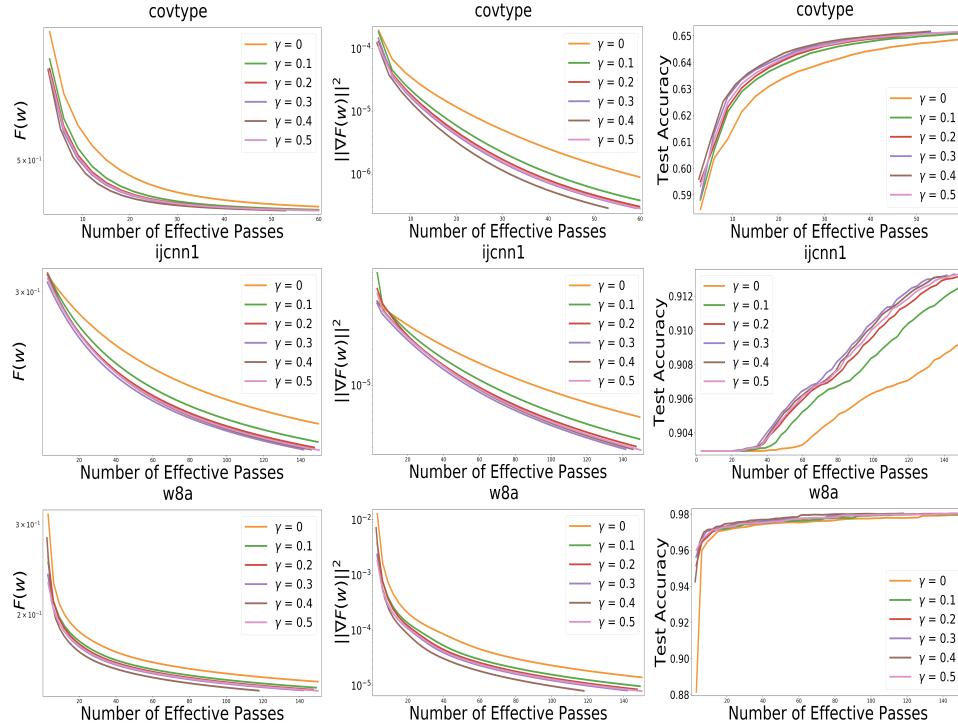


Fig. 2: Comparisons of $F(w)$, $\|\nabla F(w)\|^2$, and Test Accuracy for NC-SARAH with different values of $\gamma = \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$ on *covtype*, *ijcnn1* and *w8a* datasets

455 **4. Convex Case: SARAH++.** In this section, we propose a new variant of SARAH+
 456 (Algorithm 4.1) [11], called SARAH++ (Algorithm 4.2), for convex problems of form (1.1).
 Different from SARAH, SARAH++ provides a stopping criteria for the inner loop; as soon

as

$$\|v_{t-1}^{(s)}\|^2 \leq \gamma \|v_0^{(s)}\|^2,$$

457 the inner loop finishes. This idea originates from the property of SARAH that, for each outer
458 loop iteration s , $\mathbb{E}[\|v_t^{(s)}\|^2] \rightarrow 0$ as $t \rightarrow \infty$ in the strongly convex case (Theorems 1a and
459 1b in [11]). Therefore, it does not make any sense to update with tiny steps when $\|v_t^{(s)}\|^2$ is
460 small. (We note that SVRG [6] does not have this property.) SARAH+ suggests to empirically
choose parameter $\gamma = 1/8$ [11] without theoretical guarantee.

Algorithm 4.1 SARAH+ [11]

Parameters: the learning rate $\eta > 0$, $0 < \gamma \leq 1$, the maximum inner loop size m , and the outer loop size S

Initialize: \tilde{w}_0

Iterate:

for $s = 1, 2, \dots, S$ **do**

$w_0^{(s)} = \tilde{w}_{s-1}$

$v_0^{(s)} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_0^{(s)})$

$w_1^{(s)} = w_0^{(s)} - \eta v_0^{(s)}$

$t = 1$

while $\|v_{t-1}^{(s)}\|^2 > \gamma \|v_0^{(s)}\|^2$ **and** $t \leq m$ **do**

Sample i_t uniformly at random from $[n]$

$v_t^{(s)} = \nabla f_{i_t}(w_t^{(s)}) - \nabla f_{i_t}(w_{t-1}^{(s)}) + v_{t-1}^{(s)}$

$w_{t+1}^{(s)} = w_t^{(s)} - \eta v_t^{(s)}$

$t \leftarrow t + 1$

end while

Set $\tilde{w}_s = w_t^{(s)}$

end for

461

Here, we modify SARAH+ (Algorithm 4.1) into SARAH++ (Algorithm 4.2) by choosing the stopping criteria for the inner loop as

$$\|v_{t-1}^{(s)}\|^2 < \gamma \|v_0^{(s)}\|^2 \text{ where } \gamma \geq L\eta$$

462 and by introducing a stopping criteria for the outer loop.

4.1. SARAH++ and Its Convergence Analysis. Before analyzing and explaining

463 SARAH++ in detail, we introduce the following assumptions used in this section.

465 ASSUMPTION 2 (L -smooth). *Each $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, $i \in [n]$, is L -smooth, i.e., there exists
466 a constant $L > 0$ such that, $\forall w, w' \in \mathbb{R}^d$,*

468 (4.1)
$$\|\nabla f_i(w) - \nabla f_i(w')\| \leq L\|w - w'\|.$$

469 ASSUMPTION 3 (μ -strongly convex). *The objective function $F : \mathbb{R}^d \rightarrow \mathbb{R}$, is μ -strongly
470 convex, i.e., there exists a constant $\mu > 0$ such that $\forall w, w' \in \mathbb{R}^d$,*

472
$$F(w) \geq F(w') + \nabla F(w')^T(w - w') + \frac{\mu}{2}\|w - w'\|^2.$$

473 Under Assumption 3, let us define the (unique) optimal solution of (1.1) as w_* . Then
474 strong convexity of F implies that

475 (4.2)
$$2\mu[F(w) - F(w_*)] \leq \|\nabla F(w)\|^2, \forall w \in \mathbb{R}^d.$$

476 We note here, for future use, that for strongly convex functions of the form (1.1), arising in
 477 machine learning applications, the condition number is defined as $\kappa \stackrel{\text{def}}{=} L/\mu$. Assumption 3
 478 covers a wide range of problems, e.g. l_2 -regularized empirical risk minimization problems
 479 with convex losses.

480 We separately assume the special case of strong convexity of all f_i 's with $\mu = 0$, called
 481 the general convexity assumption, which we will use for convergence analysis.

482 ASSUMPTION 4. *Each function $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, $i \in [n]$, is convex, i.e.,*

$$483 \quad f_i(w) \geq f_i(w') + \nabla f_i(w')^T(w - w').$$

485 SARAH++ is motivated by the following lemma.

486 LEMMA 2. *Suppose that Assumptions 2 and 4 hold. Consider a single outer loop iteration*
 487 *in SARAH (Algorithm 1) with $\eta \leq \frac{1}{L}$. Then, for $t \geq 0$ and any $s \geq 1$, we have*

$$488 \quad \mathbb{E}[F(w_{t+1}^{(s)}) - F(w_*)] \leq \mathbb{E}[F(w_t^{(s)}) - F(w_*)] - \frac{\eta}{2}\mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \\ 489 \quad (4.3) \quad + \frac{\eta}{2} \left(L\eta\mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_t^{(s)}\|^2] \right), \\ 490$$

491 where w_* is any optimal solution of F .

492 Proof. By using (2.4) and adding $-F(w_*)$ for both sides, where $w_* = \arg \min_w F(w)$,
 493 we have

$$494 \quad \mathbb{E}[F(w_{t+1}^{(s)}) - F(w_*)] \\ 495 \quad \leq \mathbb{E}[F(w_t^{(s)}) - F(w_*)] - \frac{\eta}{2}\mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] + \frac{\eta}{2}\mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] \\ 496 \quad - \left(\frac{\eta}{2} - \frac{L\eta^2}{2} \right) \mathbb{E}[\|v_t^{(s)}\|^2] \\ 497 \quad \stackrel{(2)}{\leq} \mathbb{E}[F(w_t^{(s)}) - F(w_*)] - \frac{\eta}{2}\mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \\ 498 \quad + \frac{\eta}{2} \frac{\eta L}{(2 - \eta L)} \left(\mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_t^{(s)}\|^2] \right) - \left(\frac{\eta}{2} - \frac{L\eta^2}{2} \right) \mathbb{E}[\|v_t^{(s)}\|^2] \\ 499 \quad = \mathbb{E}[F(w_t^{(s)}) - F(w_*)] - \frac{\eta}{2}\mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \\ 500 \quad + \frac{\eta}{2} \left(\frac{\eta L}{(2 - \eta L)} \left(\mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_t^{(s)}\|^2] \right) - (1 - L\eta)\mathbb{E}[\|v_t^{(s)}\|^2] \right) \\ 501 \quad \stackrel{\eta \leq \frac{1}{L}}{\leq} \mathbb{E}[F(w_t^{(s)}) - F(w_*)] - \frac{\eta}{2}\mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \\ 502 \quad + \frac{\eta}{2} \left(\eta L \left(\mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_t^{(s)}\|^2] \right) - (1 - L\eta)\mathbb{E}[\|v_t^{(s)}\|^2] \right) \\ 503 \quad = \mathbb{E}[F(w_t^{(s)}) - F(w_*)] - \frac{\eta}{2}\mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] + \frac{\eta}{2} \left(L\eta\mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_t^{(s)}\|^2] \right). \quad \square \\ 504$$

Clearly, if

$$L\eta\mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_t^{(s)}\|^2] \leq \gamma\mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_t^{(s)}\|^2] \leq 0,$$

505 where $\eta \leq \frac{\gamma}{L}$, inequality (4.3) implies

$$506 \quad 507 \quad \mathbb{E}[F(w_{t+1}^{(s)}) - F(w_*)] \leq \mathbb{E}[F(w_t^{(s)}) - F(w_*)] - \frac{\eta}{2}\mathbb{E}[\|\nabla F(w_t^{(s)})\|^2].$$

Algorithm 4.2 SARAH++

Parameters: The controlled factor $0 < \gamma \leq 1$, the learning rate $0 < \eta \leq \frac{\gamma}{L}$, the total iteration $T > 0$, and the maximum inner loop size $m \leq T$.

Initialize: \tilde{w}_0

$G = 0$

Iterate:

$s = 0$

while $G < T$ **do**

- $s \leftarrow s + 1$
- $w_0^{(s)} = \tilde{w}_{s-1}$
- $v_0^{(s)} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_0^{(s)})$
- $t = 0$
- while** $\|v_t^{(s)}\|^2 \geq \gamma \|v_0^{(s)}\|^2$ **and** $t \leq m$ **do**

 - $w_{t+1}^{(s)} = w_t^{(s)} - \eta v_t^{(s)}$
 - $t \leftarrow t + 1$
 - if** $m \neq 0$ **then**

 - Sample i_t uniformly at random from $[n]$
 - $v_t^{(s)} = \nabla f_{i_t}(w_t^{(s)}) - \nabla f_{i_t}(w_{t-1}^{(s)}) + v_{t-1}^{(s)}$

 - end if**

- end while**
- $T_s = t$
- $\tilde{w}_s = w_{T_s}^{(s)}$
- $G \leftarrow G + T_s$
- end while**
- $S = s$
- Set $\hat{w} = \tilde{w}_S$

508 For this reason, we choose the stopping criteria for the inner loop in SARAH++ as $\|v_t^{(s)}\|^2 <$
 509 $\gamma \|v_0^{(s)}\|^2$ with $\gamma \geq L\eta$. Unlike SARAH+, for analyzing the convergence rate γ can be as
 510 small as $L\eta$.

511 The above discussion leads to SARAH++ (Algorithm 4.2). In order to analyze its con-
 512 vergence for convex problems, we define random variable T_s as the stopping time of the inner
 513 loop in the s -th outer iteration:

514
$$T_s = \min \left\{ \min_{t \geq 0} \left\{ t : \|v_t^{(s)}\|^2 < \gamma \|v_0^{(s)}\|^2 \right\}, m + 1 \right\}, \quad s = 1, 2, \dots$$

 515

516 Note that T_s is at least 1 since at $t = 0$, the condition $\|v_0^{(s)}\|^2 \geq \gamma \|v_0^{(s)}\|^2$ always holds (and
 517 $m \geq 0$).

518 Let random variable S be the stopping time of the outer iterations as a function of an
 519 algorithm parameter $T > 0$:

520
$$S = \min_{\hat{S}} \left\{ \hat{S} : \sum_{s=1}^{\hat{S}} T_s \geq T \right\}.$$

 521

522 Notice that SARAH++ maintains a running sum $G = \sum_{j=1}^s T_i$ against which parameter T is
 523 compared in the stopping criteria of the outer loop.

524 For the general convex case which supposes Assumption 4 in addition to smoothness we
 525 have the next theorem.

526 THEOREM 3 (Smooth general convex). *Suppose that Assumptions 2 and 4 hold. Consider SARAH++ (Algorithm 4.2) with $\eta \leq \frac{\gamma}{L}$, $0 < \gamma \leq 1$. Then,*

$$528 \quad \mathbb{E} \left[\frac{1}{T_1 + \dots + T_S} \sum_{s=1}^S \sum_{t=0}^{T_s-1} \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2 | T_1, \dots, T_S] \right] \leq \frac{2}{T\eta} [F(\tilde{w}_0) - F(w_*)].$$

530 *Proof.* We recall the following definitions. T_s is the stopping time (a random variable)
531 of the s -th outer iteration such that

$$532 \quad T_s = \min \left\{ \min_{t \geq 0} \left\{ t : \|v_t^{(s)}\|^2 < \gamma \|v_0^{(s)}\|^2 \right\}, m+1 \right\}, \quad s = 1, 2, \dots$$

534 and S is the stopping time of the outer iterations (a random variable) and such that for some
535 $T > 0$

$$536 \quad S = \min_{\hat{S}} \left\{ \hat{S} : \sum_{s=1}^{\hat{S}} T_s \geq T \right\}.$$

538 Note that $T_s \geq 1$ is the first time such that $\|v_{T_s}^{(s)}\|^2 < \gamma \|v_0^{(s)}\|^2$. Hence, for a given T_s ,
539 we have $\|v_t^{(s)}\|^2 \geq \gamma \|v_0^{(s)}\|^2$, for $0 \leq t \leq T_s - 1$, and

$$\begin{aligned} 540 \quad \mathbb{E}[F(w_{T_s}^{(s)}) - F(w_*)] &\leq \mathbb{E}[F(w_{T_s-1}^{(s)}) - F(w_*)] - \frac{\eta}{2} \mathbb{E}[\|\nabla F(w_{T_s-1}^{(s)})\|^2] \\ 541 \quad &\quad + \frac{\eta}{2} \left(L\eta \mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_{T_s-1}^{(s)}\|^2] \right) \\ 542 \quad &\leq \mathbb{E}[F(w_{T_s-1}^{(s)}) - F(w_*)] - \frac{\eta}{2} \mathbb{E}[\|\nabla F(w_{T_s-1}^{(s)})\|^2] \\ 543 \quad &\quad + \frac{\eta}{2} \left(\gamma \mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_{T_s-1}^{(s)}\|^2] \right) \\ 544 \quad &\leq \mathbb{E}[F(w_{T_s-1}^{(s)}) - F(w_*)] - \frac{\eta}{2} \mathbb{E}[\|\nabla F(w_{T_s-1}^{(s)})\|^2] \\ 545 \quad &\leq \mathbb{E}[F(w_0^{(s)}) - F(w_*)] - \frac{\eta}{2} \sum_{t=0}^{T_s-1} \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2]. \end{aligned}$$

547 Since $\tilde{w}_s = w_{T_s}^{(s)}$ and $\tilde{w}_{s-1} = w_0^{(s)}$, for given T_1, \dots, T_S , we have

$$\begin{aligned} 548 \quad \mathbb{E}[F(\tilde{w}_S) - F(w_*)] &\leq \mathbb{E}[F(\tilde{w}_{S-1}) - F(w_*)] - \frac{\eta}{2} \sum_{t=0}^{T_S-1} \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2] \\ 549 \quad &\leq \mathbb{E}[F(\tilde{w}_0) - F(w_*)] - \frac{\eta}{2} \sum_{s=1}^S \sum_{t=0}^{T_s-1} \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2]. \end{aligned}$$

551 Since $F(\tilde{w}_S) \geq F(w_*)$, bringing the second term of the RHS to the LHS. For any given
552 \tilde{w}_0 , we have

$$553 \quad \frac{\eta}{2} \sum_{s=1}^S \sum_{t=0}^{T_s-1} \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2 | T_1, \dots, T_S] \leq [F(\tilde{w}_0) - F(w_*)],$$

555 which is equivalent to

$$556 \quad \frac{1}{T_1 + \dots + T_S} \sum_{s=1}^S \sum_{t=0}^{T_s-1} \mathbb{E}[\|\nabla F(w_t^{(s)})\|^2 | T_1, \dots, T_S] \leq \frac{1}{T_1 + \dots + T_S} \frac{2}{\eta} [F(\tilde{w}_0) - F(w_*)]$$

$$\leq \frac{2}{\eta T} [F(\tilde{w}_0) - F(w_*)],$$

559 where the last inequality follows since $\sum_{s=1}^S T_s \geq T$. Hence, by taking the expectation to
560 both sides, we could achieve the desired result. \square

561 The theorem leads to the next corollary about iteration complexity, i.e., we bound T
562 which is the total number of iterations performed by the inner loop across all outer loop
563 iterations. This is different from the total complexity since T does not separately count the n
564 gradient evaluations when the full gradient is computed in the outer loop.

565 **COROLLARY 7** (Smooth general convex). *For the conditions in Theorem 3 with $\eta =$
566 $\mathcal{O}(\frac{1}{L})$, we achieve an ϵ -accurate solution after $\mathcal{O}(\frac{1}{\epsilon})$ inner loop iterations.*

567 *Proof.* The proof is trivial since we want $\frac{2}{\eta T} [F(\tilde{w}_0) - F(w_*)] = \epsilon$, which requires
568 $T = \frac{2[F(\tilde{w}_0) - F(w_*)]}{\eta} \cdot \frac{1}{\epsilon} = \mathcal{O}(\frac{1}{\epsilon})$ iterations, where we could choose $\eta = \mathcal{O}(\frac{1}{L})$. \square

569 By supposing Assumption 3 in addition to the smoothness and general convexity as-
570 sumptions, we can prove a linear convergence rate. For strongly convex objective functions
571 we have the following result.

572 **THEOREM 4** (Smooth strongly convex). *Suppose that Assumptions 2, 3 and 4 hold.
573 Consider SARAH++ (Algorithm 4.2) with $\eta \leq \frac{\gamma}{L}$, $0 < \gamma \leq 1$. Then, for the final output \hat{w} of
574 SARAH++, we have*

$$575 \quad (4.4) \quad \mathbb{E}[F(\hat{w}) - F(w_*)] \leq (1 - \mu\eta)^T [F(\tilde{w}_0) - F(w_*)].$$

577 *Proof.* Following the beginning part of the proof of Theorem 3, we have, for a given T_s ,

$$\begin{aligned} 578 \quad \mathbb{E}[F(w_{T_s}^{(s)}) - F(w_*)] &\leq \mathbb{E}[F(w_{T_{s-1}}^{(s)}) - F(w_*)] - \frac{\eta}{2} \mathbb{E}[\|\nabla F(w_{T_{s-1}}^{(s)})\|^2] \\ 579 &\stackrel{(4.2)}{\leq} (1 - \mu\eta) \mathbb{E}[F(w_{T_{s-1}}^{(s)}) - F(w_*)] \\ 580 &\leq (1 - \mu\eta)^{T_s} \mathbb{E}[F(w_0^{(s)}) - F(w_*)] \end{aligned}$$

582 Since $\tilde{w}_s = w_{T_s}^{(s)}$ and $\tilde{w}_{s-1} = w_0^{(s)}$, for given T_1, \dots, T_S , we have

$$\begin{aligned} 583 \quad \mathbb{E}[F(\hat{w}) - F(w_*)|T_1, \dots, T_S] &= \mathbb{E}[F(\tilde{w}_S) - F(w_*)|T_1, \dots, T_S] \\ 584 &\leq (1 - \mu\eta)^{T_1 + \dots + T_S} [F(\tilde{w}_0) - F(w_*)] \\ 585 &\leq (1 - \mu\eta)^T [F(\tilde{w}_0) - F(w_*)], \end{aligned}$$

587 where the last inequality follows since $\sum_{s=1}^S T_s \geq T$. Hence, by taking the expectation to
588 both sides, we could have $\mathbb{E}[F(\hat{w}) - F(w_*)] \leq (1 - \mu\eta)^T [F(\tilde{w}_0) - F(w_*)]$. \square

589 This leads to the following iteration complexity.

590 **COROLLARY 8** (Smooth strongly convex). *For the conditions in Theorem 4 with $\eta =$
591 $\mathcal{O}(\frac{1}{L})$, we achieve $\mathbb{E}[F(\hat{w}) - F(w_*)] \leq \epsilon$ after $\mathcal{O}(\kappa \log(\frac{1}{\epsilon}))$ total iterations, where $\kappa = L/\mu$
592 is the condition number.*

593 *Proof.* We want $(1 - \mu\eta)^T [F(\tilde{w}_0) - F(w_*)] = \epsilon$. Hence,

$$594 \quad T = -\frac{1}{\log(1 - \mu\eta)} \log \left(\frac{[F(\tilde{w}_0) - F(w_*)]}{\epsilon} \right).$$

596 Note that: $-\frac{1}{x} - 1 \leq -\frac{1}{\log(1+x)} \leq -\frac{1}{x}$, $-1 < x < 0$. We can have

$$597 \quad 598 \quad \left(\frac{1}{\mu\eta} - 1 \right) \log \left(\frac{[F(\tilde{w}_0) - F(w_*)]}{\epsilon} \right) \leq T \leq \frac{1}{\mu\eta} \log \left(\frac{[F(\tilde{w}_0) - F(w_*)]}{\epsilon} \right).$$

599 By choosing $\eta = \mathcal{O}(\frac{1}{L})$, we have $T = \mathcal{O}(\kappa \log(\frac{1}{\epsilon}))$. \square

600 REMARK 3. *The proofs of the above results hold for any $m \leq T$. If we choose $m = 0$,*
 601 *then SARAH++ reduces to the Gradient Descent algorithm since the inner “while” loop stops*
 602 *right after updating $w_1^{(s)} = w_0^{(s)} - \eta v_0^{(s)}$. In this case, Corollaries 7 and 8 recover the rate*
 603 *of convergence and complexity of GD.*

604 In this section, we showed that SARAH++ has a guarantee of theoretical convergence
 605 (see Theorems 3 and 4) while SARAH+ does not have such a guarantee.

606 An interesting open question we would like to discuss here is the total complexity of
 607 SARAH++. Although we have shown the convergence results of SARAH++ in terms of the
 608 iteration complexity, the total complexity which is computed as the total number of evalua-
 609 tions of the component gradient functions still remains an open question. It is clear that the
 610 total complexity must depend on the learning rate η (or γ) – the factor that decides when to
 611 stop the inner iterations.

612 We note that T can be “closely” understood as the total number of updates $w_{t+1}^{(s)}$ of the
 613 algorithm. The total complexity is equal to $\sum_{i=1}^S (n + 2(T_i - 1))$. For the special case
 614 $T_i = 1$, $i = 1, \dots, S$, the algorithm recovers the GD algorithm with $T = \sum_{i=1}^S T_s = S$.
 615 Since each full gradient takes n gradient evaluations, the total complexity for this case is
 616 equal to $nS = \mathcal{O}(\frac{n}{\epsilon})$ (in the general convex case) and $nS = \mathcal{O}(n\kappa \log(\frac{1}{\epsilon}))$ (in the strongly
 617 convex case).

618 However, it is non-trivial to derive the total complexity of SARAH++ since it should
 619 depend on the learning rate η . We leave this question as an open direction for future research.

620 **4.2. Numerical Experiments.** Paper [11] provides experiments showing good overall
 621 performance of SARAH over other algorithms such as SGD [14], SAG [15], SVRG [6], etc.
 622 For this reason, we provide experiments comparing SARAH++ directly with SARAH. We
 623 notice that SARAH (with multiple outer loops) like SARAH++ has theoretical guarantees
 624 with sublinear convergence for general convex and linear convergence for strongly convex
 625 problems as proved in [11]. Because of these theoretical guarantees (which SARAH+ does
 626 not have), SARAH itself may already perform well for convex problems and the question is
 627 whether SARAH++ offers an advantage.

628 We consider ℓ_2 -regularized logistic regression problems with

$$630 \quad (4.5) \quad f_i(w) = \log(1 + \exp(-y_i \langle x_i, w \rangle)) + \frac{\lambda}{2} \|w\|^2,$$

631 where $\{x_i, y_i\}_{i=1}^n$ is the training data and the regularization parameter λ is set to $1/n$, a
 632 widely-used value in literature [15, 11]. The condition number is equal to $\kappa = L/\mu = n$.
 633 We conducted experiments to demonstrate the advantage in performance of SARAH++ over
 634 SARAH for convex problems on popular data sets including *covtype*, *ijcnn1*, *w8a* (introduced
 635 in Section 3.3), and *phishing* ($n = 7,738$ training data, estimated $L \simeq 7.49$) from LIBSVM.

636 Figure 3 shows comparisons between SARAH++ and SARAH for different values of
 637 learning rate η . We depicted the value of $\log[F(w) - F(w_*)]$ (i.e. $F(w) - F(w_*)$ in log
 638 scale) for the y -axis and “number of effective passes” (or number of epochs, where an epoch
 639 is the equivalent of n component gradient evaluations or one full gradient computation) for
 640 the x -axis. For SARAH, we choose the outer loop size $S = 10$ and tune the inner loop size
 641 $m = \{0.5n, n, 2n, 3n, 4n\}$ to achieve the best performance. The optimal solution w_* of the

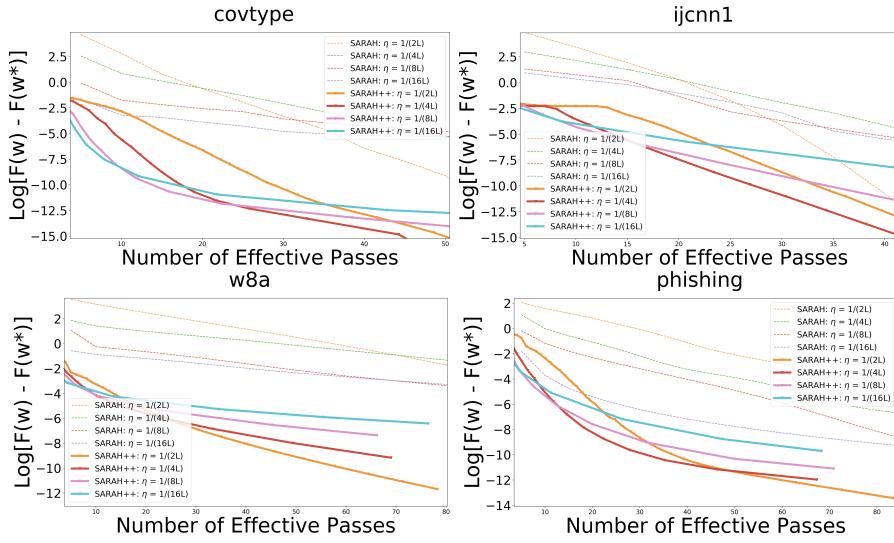


Fig. 3: Comparisons of $\log[F(w) - F(w^*)]$ between SARAH++ and SARAH with different learning rates on *covtype*, *ijcn1*, *w8a*, and *phishing* datasets

strongly convex problem in (4.5) is found by using Gradient Descent with stopping criterion $\|\nabla F(w)\|^2 \leq 10^{-15}$. We observe that, SARAH++ achieves improved overall performance compared to regular SARAH as shown in Figure 3. From the experiments we see that the stopping criteria $\|v_t^{(s)}\|^2 < \gamma \|v_0^{(s)}\|^2$ ($\gamma = L\eta$) of SARAH++ is indeed important. The stopping criteria helps the inner loop to prevent updating tiny redundant steps.

4.3. SARAH Adaptive: A New Practical Variant. We now propose a practical adaptive method which aims to improve performance. Although we do not have any theoretical result for this adaptive method, numerical experiments are very promising and they heuristically show the improved performance on different data sets.

The motivation of this algorithm comes from the intuition of Lemma 2 for convex optimization. For a single outer loop with $\eta \leq \frac{1}{L}$, (4.3) holds for SARAH (Algorithm 1.1). Hence, for any s , we intentionally choose $\eta = \eta_t^{(s)} = \frac{\|v_t^{(s)}\|^2}{L\|v_0^{(s)}\|^2}$ such that $L\eta \mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_t^{(s)}\|^2] = 0$. Since $\|v_t^{(s)}\|^2 \leq \|v_0^{(s)}\|^2$, $t \geq 0$, in [11] for convex problems, we have $\eta_t^{(s)} \leq \frac{1}{L}$, $t \geq 0$. We also stop the inner loop when $\|v_t^{(s)}\|^2 < \gamma \|v_0^{(s)}\|^2$ for some $0 < \gamma \leq 1$. SARAH Adaptive is given in detail in Algorithm 4.3 without convergence analysis.

We have conducted numerical experiments on the same datasets and problems as introduced in the previous subsection. Figures 4 and 5 show the comparison between SARAH Adaptive and SARAH and SARAH++ for different values of η . We observe that SARAH Adaptive has an improved performance over SARAH and SARAH++ (without tuning learning rate). In Figure 6 we present the numerical performance of SARAH Adaptive for different values of $\gamma = \{\frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \frac{1}{6}, \frac{1}{8}, \frac{1}{10}, \frac{1}{12}, \frac{1}{16}\}$.

5. Conclusion and Future Research. In this paper, we address almost important open problems for the original SARAH algorithm, i.e., SARAH for convex and non-convex optimization problems. For non-convex optimization, we propose NC-SARAH, which achieves the state-of-the-art asymptotic total complexity for finding a first-order stationary point based on only the average smooth assumption as SPIDER and SpiderBoost. The total complexity matches the lower-bound worst case complexity up to a constant factor when $n \leq \mathcal{O}(\epsilon^{-2})$. In-

Algorithm 4.3 SARAH Adaptive

Parameters: The maximum inner loop size m , the outer loop size S , the factor $0 < \gamma \leq 1$.

```

Initialize:  $\tilde{w}_0$ 
Iterate:
for  $s = 1, 2, \dots, S$  do
     $w_0^{(s)} = \tilde{w}_{s-1}$ 
     $v_0^{(s)} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_0^{(s)})$ 
     $t = 0$ 
    while  $\|v_t^{(s)}\|^2 \geq \gamma \|v_0^{(s)}\|^2$  and  $t \leq m$  do
         $\eta_t^{(s)} = \frac{1}{L} \cdot \frac{\|v_t^{(s)}\|^2}{\|v_0^{(s)}\|^2}$  (adaptive)
         $w_{t+1}^{(s)} = w_t^{(s)} - \eta_t^{(s)} v_t^{(s)}$ 
         $t \leftarrow t + 1$ 
    if  $m \neq 0$  then
        Sample  $i_t$  uniformly at random from  $[n]$ 
         $v_t^{(s)} = \nabla f_{i_t}(w_t^{(s)}) - \nabla f_{i_t}(w_{t-1}^{(s)}) + v_{t-1}^{(s)}$ 
    end if
    end while
    Set  $\tilde{w}_s = w_t^{(s)}$ 
end for

```

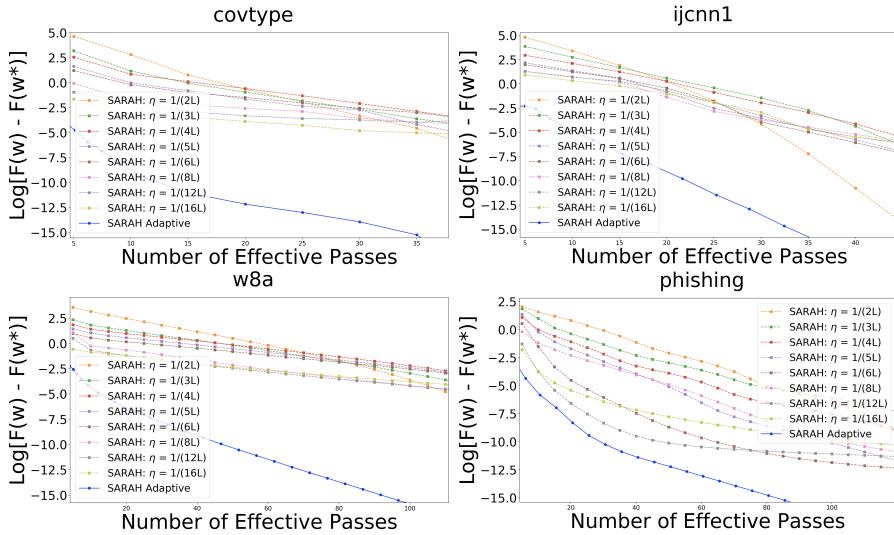


Fig. 4: Comparisons of $\log[F(w) - F(w^*)]$ between SARAH Adaptive and SARAH with different learning rates on *covtype*, *ijcn1*, *w8a*, and *phishing* datasets

669 indeed, NC-SARAH has advantages over SPIDER and SpiderBoost in both theory and practice,
670 i.e., NC-SARAH allows larger learning rates as well as a range of mini-batch sizes; numerical
671 experiments show how NC-SARAH outperforms SPIDER and SpiderBoost. In theory, our
672 proof is significantly simpler and more intuitive. Moreover, we showed the promising numer-
673 ical results for SARAH++ and SARAH Adaptive in the convex case. The total complexity of
674 SARAH++ and the convergence analysis of SARAH Adaptive could be potential for future

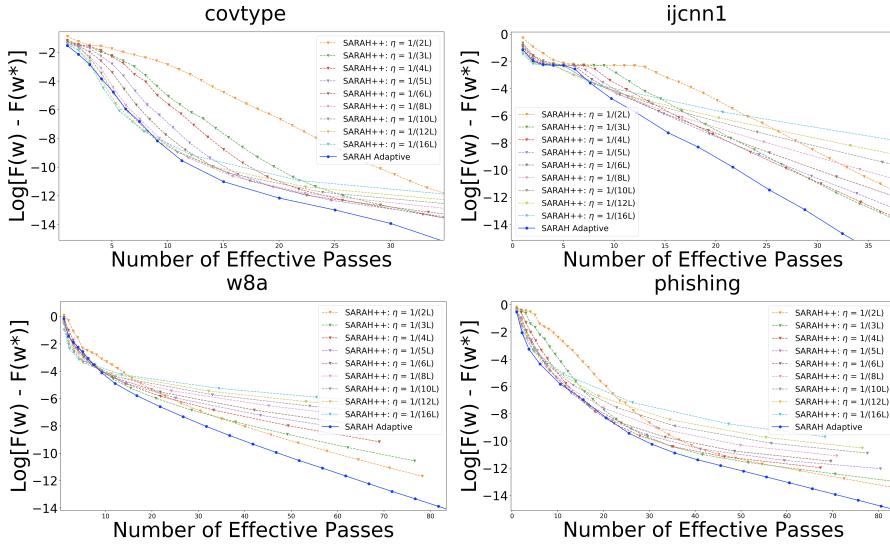


Fig. 5: Comparisons of $\log[F(w) - F(w^*)]$ between SARAH Adaptive and SARAH++ with different learning rates on *covtype*, *ijcn1*, *w8a*, and *phishing* datasets

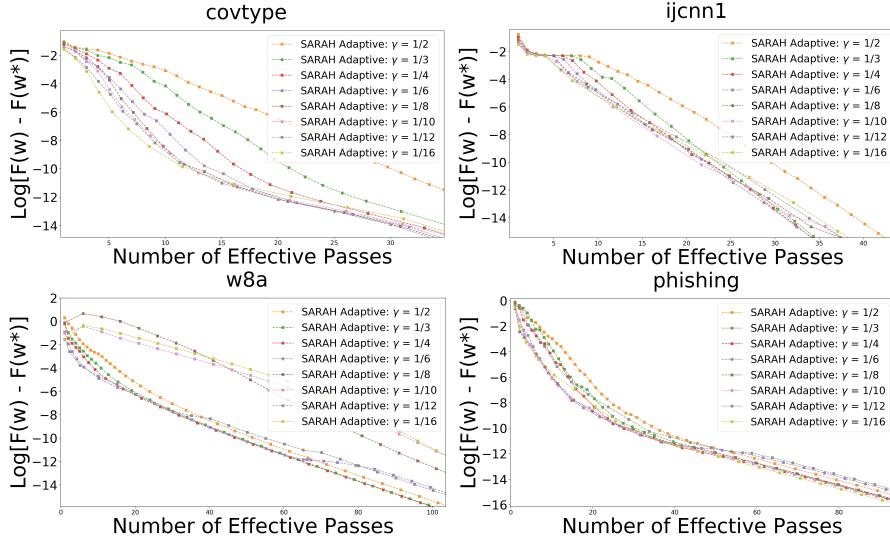


Fig. 6: Comparisons of $\log[F(w) - F(w^*)]$ with different value of γ for SARAH Adaptive on *covtype*, *ijcn1*, *w8a*, and *phishing* datasets

675 research. In addition, we show that SARAH (NC-SARAH and SARAH++) can be reduced to
 676 Gradient Descent - an open problem since 2012 - which may open new research directions.

677 Appendix.

678 **Some Useful Existing Results.** The following existing results are used in the proofs of
 679 our main results.

680 LEMMA 3 (Lemma 2 in [11] (or in [12])). Consider $v_t^{(s)}$ defined by (1.2) (or (1.3)) in

681 SARAH (Algorithm 1.1) for any $s \geq 1$. Then for any $t \geq 1$,

(.1)

$$682 \quad \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] = \sum_{j=1}^t \mathbb{E}[\|v_j^{(s)} - v_{j-1}^{(s)}\|^2] - \sum_{j=1}^t \mathbb{E}[\|\nabla F(w_j^{(s)}) - \nabla F(w_{j-1}^{(s)})\|^2]. \\ 683$$

684 LEMMA 4 (Lemma 3 in [11]). Suppose that Assumptions 2 and 4 hold. Consider $v_t^{(s)}$
 685 defined as (1.2) in SARAH (Algorithm 1.1) with $\eta < 2/L$ for any $s \geq 1$. Then we have that
 686 for any $t \geq 0$,

$$687 \quad (2) \quad \mathbb{E}[\|\nabla F(w_t^{(s)}) - v_t^{(s)}\|^2] \leq \frac{\eta L}{2 - \eta L} \left[\mathbb{E}[\|v_0^{(s)}\|^2] - \mathbb{E}[\|v_t^{(s)}\|^2] \right]. \\ 688$$

689

REFERENCES

- 690 [1] Z. ALLEN-ZHU, *Natasha: Faster non-convex stochastic optimization via strongly non-convex parameter*, in
 691 Proceedings of the 34th International Conference on Machine Learning-Volume 70, JMLR. org, 2017,
 692 pp. 89–97.
- 693 [2] Z. ALLEN-ZHU, *Natasha 2: Faster non-convex optimization than sgd*, in Advances in Neural Information
 694 Processing Systems, 2018, pp. 2675–2686.
- 695 [3] C.-C. CHANG AND C.-J. LIN, *LIBSVM: A library for support vector machines*, ACM Transactions on Intel-
 696 ligent Systems and Technology, 2 (2011), pp. 27:1–27:27.
- 697 [4] A. DEFazio, F. BACH, AND S. LACOSTE-JULIEN, *Saga: A fast incremental gradient method with support*
 698 *for non-strongly convex composite objectives*, in Advances in Neural Information Processing Systems,
 699 2014, pp. 1646–1654.
- 700 [5] C. FANG, C. J. LI, Z. LIN, AND T. ZHANG, *Spider: Near-optimal non-convex optimization via stochastic*
 701 *path-integrated differential estimator*, in Advances in Neural Information Processing Systems, 2018,
 702 pp. 689–699.
- 703 [6] R. JOHNSON AND T. ZHANG, *Accelerating stochastic gradient descent using predictive variance reduction*,
 704 in Advances in Neural Information Processing Systems, 2013, pp. 315–323.
- 705 [7] J. KONEČNÝ AND P. RICHTÁRIK, *Semi-stochastic gradient descent methods*, Frontiers in Applied Mathe-
 706 matics and Statistics, 3 (2017), p. 9.
- 707 [8] L. LEI, C. JU, J. CHEN, AND M. I. JORDAN, *Non-convex finite-sum optimization via SCSG methods*, in
 708 Advances in Neural Information Processing Systems 30, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach,
 709 R. Fergus, S. Vishwanathan, and R. Garnett, eds., Curran Associates, Inc., 2017, pp. 2348–2358.
- 710 [9] J. MAIRAL, *Optimization with first-order surrogate functions*, in International Conference on Machine Learn-
 711 ing, 2013, pp. 783–791.
- 712 [10] Y. NESTEROV, *Introductory lectures on convex optimization : a basic course*, Applied optimization, Kluwer
 713 Academic Publ., Boston, Dordrecht, London, 2004.
- 714 [11] L. M. NGUYEN, J. LIU, K. SCHEINBERG, AND M. TAKÁČ, *Sarah: A novel method for machine learning*
 715 *problems using stochastic recursive gradient*, in Proceedings of the 34th International Conference on
 716 Machine Learning-Volume 70, JMLR. org, 2017, pp. 2613–2621.
- 717 [12] L. M. NGUYEN, J. LIU, K. SCHEINBERG, AND M. TAKÁČ, *Stochastic recursive gradient algorithm for*
 718 *nonconvex optimization*, CoRR, abs/1705.07261 (2017).
- 719 [13] S. J. REDDI, A. HEFNY, S. SRA, B. POCZOS, AND A. SMOLA, *Stochastic variance reduction for nonconvex*
 720 *optimization*, in International conference on machine learning, 2016, pp. 314–323.
- 721 [14] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, The Annals of Mathematical Statistics,
 722 22 (1951), pp. 400–407.
- 723 [15] N. L. ROUX, M. SCHMIDT, AND F. R. BACH, *A stochastic gradient method with an exponential convergence*
 724 *rate for finite training sets*, in Advances in Neural Information Processing Systems, 2012, pp. 2663–
 725 2671.
- 726 [16] M. SCHMIDT, N. LE ROUX, AND F. BACH, *Minimizing finite sums with the stochastic average gradient*,
 727 Mathematical Programming, (2016), pp. 1–30.
- 728 [17] S. SHALEV-SHWARTZ AND T. ZHANG, *Stochastic dual coordinate ascent methods for regularized loss*, Jour-
 729 *nal of Machine Learning Research*, 14 (2013), pp. 567–599.
- 730 [18] Z. WANG, K. JI, Y. ZHOU, Y. LIANG, AND V. TAROKH, *Spiderboost: A class of faster variance-reduced*
 731 *algorithms for nonconvex optimization*, Advances in Neural Information Processing Systems, (2019).
- 732 [19] D. ZHOU, P. XU, AND Q. GU, *Stochastic nested variance reduced gradient descent for nonconvex optimiza-
 733 tion*, in Advances in Neural Information Processing Systems, 2018, pp. 3921–3932.