

1 **ON THE SAGA ALGORITHM WITH DECREASING STEP\***

2 BERNARD BERCU<sup>†</sup>, LUIS FREDES<sup>†</sup>, AND EMÉRIC GBAGUIDI<sup>†</sup>

3 **Abstract.** Stochastic optimization naturally appear in many application areas, including ma-  
 4 chine learning. Our goal is to go further in the analysis of the Stochastic Average Gradient Acceler-  
 5 ated (SAGA) algorithm. To achieve this, we introduce a new  $\lambda$ -SAGA algorithm which interpolates  
 6 between the Stochastic Gradient Descent ( $\lambda = 0$ ) and the SAGA algorithm ( $\lambda = 1$ ). Firstly, we  
 7 investigate the almost sure convergence of this new algorithm with decreasing step which allows us  
 8 to avoid the restrictive strong convexity and Lipschitz gradient hypotheses associated to the objec-  
 9 tive function. Secondly, we establish a central limit theorem for the  $\lambda$ -SAGA algorithm. Finally, we  
 10 provide the non-asymptotic  $\mathbf{L}^p$  rates of convergence.

11 **Key words.** SAGA algorithm, decreasing step, almost sure convergence, asymptotic normality,  
 12 non-asymptotic rates of convergence

13 **MSC codes.** 68Q25, 68T05, 60F05, 60F15

14 **1. Introduction.** Our goal is to solve the classical optimization problem in  $\mathbb{R}^d$   
 15 which can be written as

16  $(\mathcal{P}) \quad \min_{x \in \mathbb{R}^d} f(x),$

17 where  $f$  is the average of many functions,

18  $(1.1) \quad f(x) = \frac{1}{N} \sum_{k=1}^N f_k(x).$

19 This type of problem is frequently encountered in statistical learning and a standard  
 20 way to solve  $(\mathcal{P})$  is to make use of the Gradient Descent algorithm. However, in a  
 21 large  $N$  context, this approach has a very high computational cost. This limitation  
 22 has led to the development of many stochastic algorithms for optimization [8, 26].

23 These new methods have taken a major role in recent advances of the neural  
 24 networks. Our goal is to go further in the analysis of the Stochastic Gradient Descent  
 25 (SGD) algorithm [33] and the SAGA algorithm [12]. The standard SGD algorithm is  
 26 given for all  $n \geq 1$ , by

27  $(\text{SGD}) \quad X_{n+1} = X_n - \gamma_n \nabla f_{U_{n+1}}(X_n) = X_n - \gamma_n (\nabla f(X_n) + \varepsilon_{n+1}),$

28 where the initial state  $X_1$  is a squared integrable random vector of  $\mathbb{R}^d$  which can be  
 29 arbitrarily chosen,  $\nabla f(X_n)$  is the gradient of the function  $f$  calculated at the value  $X_n$ ,  
 30  $\varepsilon_{n+1} = \nabla f_{U_{n+1}}(X_n) - \nabla f(X_n)$  and  $(U_n)$  is a sequence of independent and identically  
 31 distributed random variables, with uniform distribution on  $\{1, 2, \dots, N\}$ , which is  
 32 also independent from the sequence  $(X_n)$ . Moreover,  $(\gamma_n)$  is a positive deterministic  
 33 sequence decreasing towards zero and satisfying the standard conditions

34  $(1.2) \quad \sum_{n=1}^{\infty} \gamma_n = +\infty \quad \text{and} \quad \sum_{n=1}^{\infty} \gamma_n^2 < +\infty.$

---

\*Submitted to the editors September 11, 2025.

**Funding:** This project has benefited from state support managed by the Agence Nationale de la Recherche (French National Research Agency) under the reference ANR-20-SFRI-0001.

<sup>†</sup>Institut de Mathématiques de Bordeaux, Université de Bordeaux (bernard.bercu@math.u-bordeaux.fr, luis.fredes@math.u-bordeaux.fr, thierry-emerice.gbaguidi@math.u-bordeaux.fr).

35 We clearly have from (1.1) that  $(\varepsilon_n)$  is a martingale difference sequence adapted to  
 36 the filtration  $(\mathcal{F}_n)$  where  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$ .

37 The SAGA algorithm is a stochastic variance reduction algorithm which was proposed  
 38 ten years ago in the pioneering work of [12]. It slightly differs from the SGD algorithm  
 39 as it is given, for all  $n \geq 1$ , by

$$40 \quad (\text{SAGA}) \quad X_{n+1} = X_n - \gamma_n \left( \nabla f_{U_{n+1}}(X_n) - g_{n,U_{n+1}} + \frac{1}{N} \sum_{k=1}^N g_{n,k} \right),$$

41 where the initial states  $X_0$  and  $X_1$  are squared integrable random vectors of  $\mathbb{R}^d$   
 42 which can be arbitrarily chosen, the initial value  $g_{1,k}$  is given, for any  $k = 1, \dots, N$ ,  
 43 by  $g_{1,k} = \nabla f_k(X_0)$ . Moreover, the sequence  $(g_{n,k})$  is updated, for all  $n \geq 1$  and  
 44  $1 \leq k \leq N$ , as

$$45 \quad (1.3) \quad g_{n+1,k} = \begin{cases} \nabla f_k(X_n) & \text{if } U_{n+1} = k, \\ g_{n,k} & \text{otherwise.} \end{cases}$$

46 One can observe that in most of all papers dealing with the SAGA algorithm,  
 47 the step size is a fixed value  $\gamma$  which depends on the strong convexity constant  $\mu$   
 48 and the Lipschitz gradient constant  $L$  associated with  $f$ . This will not be the case  
 49 here at all as the benefit of decreasing step algorithms has been recently investigated  
 50 in [23]. Our work aims to investigate the almost sure convergence as well as the  
 51 asymptotic normality of the SGD and SAGA algorithms with decreasing step sequence  
 52  $(\gamma_n)$  satisfying (1.2).

53 **Our contributions.** The goal of this paper is to answer to several natural questions.

- 54 (a) Is it possible to study the convergence of the SAGA algorithm with decreasing  
 55 step ?
- 56 (b) Can we relax the strong convexity and the Lipschitz gradient assumptions ?
- 57 (c) Can we prove a central limit theorem for our new version of the SAGA algo-  
 58 rithm ?
- 59 (d) Is it possible to provide non-asymptotic  $\mathbf{L}^p$  bounds for the SAGA algorithm ?

60 We shall propose positive answers to all these questions by extending [12] in several  
 61 directions.

62 The paper is structured as follows. Section 2 is devoted to the state of the art  
 63 concerning the SGD and SAGA algorithms. In Section 3, we present our new version  
 64 of the SAGA algorithm which we shall call the  $\lambda$ -SAGA algorithm. Section 4 deals  
 65 with the main results of the paper. We establish the asymptotic properties of our  $\lambda$ -  
 66 SAGA algorithm such as the almost sure convergence and the asymptotic normality.  
 67 Non-asymptotic  $\mathbf{L}^p$  rates of convergence are also provided. In Section 5, we illustrate  
 68 our theoretical results by some numerical experiments on real dataset. All technical  
 69 proofs are postponed to the appendices.

70 **2. Related work.** The stochastic approximations, initiated by [33] and [20],  
 71 have taken a major role in optimization issues. The SGD algorithm, often known  
 72 as a special case of the Robbins-Monro algorithm, is probably the most standard  
 73 stochastic algorithm used in machine learning. The properties of this algorithm were  
 74 investigated in several studies. The almost sure convergence results were established  
 75 in [7, 8, 13, 22, 34, 35, 37]. The convergence rates were proven in [21, 25, 26, 28]. The

76 study of the asymptotic normality of stochastic approximations also appear in several  
77 works such that [13, 14, 29, 36, 40, 41].

78 In a high-dimensional context, many accelerated algorithms were proposed in  
79 literature in order to improve the Robbins-Monro algorithm performances [1, 30, 15,  
80 12, 38, 24]. In this paper, we will focus on the SAGA algorithm first introduced by  
81 [12] for the minimization of the average of many functions and which is a well-known  
82 variance reduction method. This algorithm is a variant of the Stochastic Average  
83 Gradient (SAG) method proposed earlier in [35, 37]. It uses the concept of covariates  
84 to make an unbiased variant of the SAG method that has similar performances but  
85 is easier to implement [17]. The idea behind the SAGA algorithm, is to make use  
86 of the control variates, a well-known technique in Monte-Carlo simulation designed  
87 to reduce the variance of the SGD algorithm in order to accelerate its convergence.  
88 This algorithm incorporates knowledge about gradients on all previous data points  
89 rather than only using the gradient for the sampled data point [12, 27]. This method  
90 requires a storage linear in  $N$  [16]. Several works have studied the convergence of the  
91 SAGA algorithm, which is undoubtedly one of the most celebrated variance reduction  
92 algorithms.

93 Defazio et al. [12] established that the SAGA algorithm converges in  $\mathbf{L}^2$  at ex-  
94ponential rate. This result has been shown by assuming that the function  $f$  is  $\mu$ -  
95strongly convex and with  $L$ -Lipschitz gradient and by considering a fixed constant  
96step  $\gamma$  which tightly depends on the unknown values  $\mu$  and  $L$ . The almost sure con-  
97vergence of the SAGA algorithm was not investigated in [12]. More recently, it was  
98shown by [31] that for a fixed constant step  $\gamma = 1/(3L)$ ,  $f(X_n)$  and  $X_n$  both converge  
99almost surely to  $f(x^*)$  and  $x^*$  respectively, where  $x^*$  is the unique point of  $\mathbb{R}^d$  such  
100that  $\nabla f(x^*) = 0$ . This algorithm has been also investigated in [11, 16, 27, 32] and  
101there are now many variations on the original SAGA algorithm of [12]. For example,  
102[32] proposed a variant of the SAGA algorithm that includes arbitrary importance  
103sampling and minibatching schemes.

104 Despite a decade of research, several issues remain open on the SAGA algorithm.  
105 The choice of the step  $\gamma_n$  is clearly one of them. The vast majority of the theory for  
106 the SAGA algorithm relies on a fixed constant step  $\gamma$  depending on the values  $\mu$  and  
107  $L$  [12, 11, 27, 16, 31, 17]. However, from a practical point of view, the values  $\mu$  and  
108  $L$  are unknown and there is no guarantee on the convergence results established for  
109 this algorithm. We shall propose here to make use of decreasing step sequence  $(\gamma_n)$   
110 which allows us to avoid these constraints and relax some classic assumptions such  
111 that the  $\mu$ -strong convexity. Moreover, to the best of our knowledge, no result about  
112 the asymptotic normality of the SAGA algorithm is available in the literature so far.

113 **3. The  $\lambda$ -SAGA algorithm.** We introduce in this section the  $\lambda$ -SAGA algo-  
114rithm which can be seen as a generalization of the SAGA algorithm. We recall below  
115 the general principle of the Monte Carlo method that gave birth to the  $\lambda$ -SAGA al-  
116gorithm. Suppose that we would like to estimate the expectation  $\mathbb{E}[X]$  of a square  
117integrable real random variable  $X$ . Let us also consider another square integrable  
118real random variable  $Y$  strongly positively correlated to  $X$  and for which we know  
119 how to compute the expectation  $\mathbb{E}[Y]$ . Then, it is possible to find a reduced variance  
120 estimator of  $\mathbb{E}[X]$ , given by  $Z_\lambda = X - \lambda(Y - \mathbb{E}[Y])$  with  $\lambda$  in  $[0, 1]$  (see, e.g., [9, 12]).  
121 One can obviously see that  $\mathbb{E}[Z_\lambda] = \mathbb{E}[X]$ , which means that  $Z_\lambda$  is an unbiased esti-  
122 mator of  $\mathbb{E}[X]$ . Moreover,  $\mathbb{V}[Z_\lambda] = \mathbb{V}[X] + \lambda^2\mathbb{V}[Y] - 2\lambda\text{Cov}(X, Y)$ . Hence, as soon  
123 as  $\text{Cov}(X, Y) > 0$ , we can choose  $\lambda$  in  $[0, 1]$  such that  $\mathbb{V}[Z_\lambda] \leq \mathbb{V}[X]$ . Now, using this

124 principle of variance reduction, the  $\lambda$ -SAGA algorithm is defined, for all  $n \geq 1$ , by

$$125 \quad (\lambda\text{-SAGA}) \quad X_{n+1} = X_n - \gamma_n \left( \nabla f_{U_{n+1}}(X_n) - \lambda \left( g_{n,U_{n+1}} - \frac{1}{N} \sum_{k=1}^N g_{n,k} \right) \right), \quad \blacksquare$$

126 where the initial states  $X_0$  and  $X_1$  are squared integrable random vectors of  $\mathbb{R}^d$   
 127 which can be arbitrarily chosen, the parameter  $\lambda$  belongs to  $[0, 1]$ , and  $(\gamma_n)$  is a  
 128 positive deterministic sequence decreasing towards zero and satisfying (1.2).

129 One can establish a link between the SGD, SAGA and  $\lambda$ -SAGA algorithms. Indeed,  
 130 the  $\lambda$ -SAGA algorithm with  $\lambda = 0$  corresponds to the absence of variance reduction  
 131 and reduces to the SGD algorithm. Furthermore, one can easily see that we find  
 132 again the SAGA algorithm by choosing  $\lambda = 1$ . The motivation to introduce and  
 133 study the  $\lambda$ -SAGA algorithm comes from our desire to propose a unified convergence  
 134 analysis for the SGD and SAGA algorithms and to investigate what happens in the  
 135 intermediate cases  $0 < \lambda < 1$ . We shall now state the general assumptions which we  
 136 will use in all the sequel.

137 **ASSUMPTION 3.1.** *Assume that function  $f$  is continuously differentiable with a*  
 138 *unique equilibrium point  $x^*$  in  $\mathbb{R}^d$  such that  $\nabla f(x^*) = 0$ .*

139 **ASSUMPTION 3.2.** *Suppose that for all  $x \in \mathbb{R}^d$  with  $x \neq x^*$ ,*

$$140 \quad \langle x - x^*, \nabla f(x) \rangle > 0.$$

141 **ASSUMPTION 3.3.** *Assume there exists a positive constant  $L$  such that, for all*  
 142  *$x \in \mathbb{R}^d$ ,*

$$143 \quad \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x) - \nabla f_k(x^*)\|^2 \leq L \|x - x^*\|^2.$$

144 These assumptions are not really restrictive and they are fulfilled in many appli-  
 145 cations. One can observe that Assumption 3.2 is obviously weaker than the standard  
 146 hypothesis that each function  $f_k$  for  $1 \leq k \leq N$  is  $\mu$ -strongly convex with  $\mu > 0$ . Note  
 147 also that Assumption 3.3 ensures that at  $x^*$ , the gradient of all functions  $f_k$  for any  
 148  $1 \leq k \leq N$ , does not change arbitrarily with respect to the vector  $x \in \mathbb{R}^d$ . Such an  
 149 assumption is essential for convergence of most gradient-based algorithms; without  
 150 it, the gradient would not provide a good indicator of how far to move to decrease  $f$ .  
 151 One can also observe that if each function  $f_k$  has Lipschitz continuous gradient with  
 152 constant  $\sqrt{L_k}$ , then Assumption 3.3 is satisfied by taking  $L$  as the average value of  
 153 all  $L_k$ . The most interesting improvement here is that both conditions are local in  $x^*$   
 154 and sufficient for all of our analysis.

155 **4. Main results.** In this section, we present the main results of the paper. First  
 156 of all, we provide an almost sure convergence analysis for the  $\lambda$ -SAGA algorithm with  
 157 decreasing step. After that, we establish its asymptotic normality. Lastly, we conclude  
 158 this section by focusing on non-asymptotic  $L^p$  rates of convergence of this stochastic  
 159 algorithm.

160 **4.1. Almost sure convergence.** Our first result deals with the almost sure  
 161 convergence of the  $\lambda$ -SAGA algorithm.

162 THEOREM 4.1. Consider a fixed  $\lambda \in [0, 1]$ . Assume that  $(X_n)$  is the sequence  
 163 generated by the  $\lambda$ -SAGA algorithm with decreasing step sequence  $(\gamma_n)$  satisfying (1.2).  
 164 In addition, suppose that Assumptions 3.1, 3.2 and 3.3 are satisfied. Then, we have

$$165 \quad (4.1) \quad \lim_{n \rightarrow +\infty} X_n = x^* \quad a.s.$$

166 and

$$167 \quad (4.2) \quad \lim_{n \rightarrow +\infty} f(X_n) = f(x^*) \quad a.s.$$

168 *Proof.* Recall that for all  $n \geq 1$ ,

$$169 \quad X_{n+1} = X_n - \gamma_n \left( \nabla f_{U_{n+1}}(X_n) - \lambda \left( g_{n, U_{n+1}} - \frac{1}{N} \sum_{k=1}^N g_{n,k} \right) \right).$$

170 Hence, the  $\lambda$ -SAGA algorithm can be rewritten as

$$171 \quad (4.3) \quad X_{n+1} = X_n - \gamma_n (Y_{n+1} - \lambda Z_{n+1}),$$

172 where

$$173 \quad \begin{cases} Y_{n+1} &= \nabla f_{U_{n+1}}(X_n), \\ Z_{n+1} &= \nabla f_{U_{n+1}}(\phi_{n, U_{n+1}}) - \frac{1}{N} \sum_{k=1}^N \nabla f_k(\phi_{n,k}), \end{cases}$$

174 and  $\phi_{n,k}$  is the point such that  $g_{n,k} = \nabla f_k(\phi_{n,k})$ . As  $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$  and the  
 175 sequence  $(U_n)$  is independent of the sequence  $(X_n)$ , we clearly have from (1.1) that

$$176 \quad (4.4) \quad \mathbb{E}[Y_{n+1} | \mathcal{F}_n] = \nabla f(X_n) \quad \text{and} \quad \mathbb{E}[Z_{n+1} | \mathcal{F}_n] = 0 \quad a.s.$$

177 As a consequence,  $(Z_n)$  is a martingale difference sequence adapted to the filtration  
 178  $(\mathcal{F}_n)$ .

179 Hereafter, define for all  $n \geq 1$ ,

$$180 \quad V_n = \|X_n - x^*\|^2.$$

181 We obtain from (4.3) that for all  $n \geq 1$ ,

$$\begin{aligned} 182 \quad V_{n+1} &= \|X_{n+1} - x^*\|^2, \\ 183 \quad &= \|X_n - x^* - \gamma_n(Y_{n+1} - \lambda Z_{n+1})\|^2, \\ 184 \quad &= V_n - 2\gamma_n \langle X_n - x^*, Y_{n+1} - \lambda Z_{n+1} \rangle + \gamma_n^2 \|Y_{n+1} - \lambda Z_{n+1}\|^2. \end{aligned}$$

185 Moreover, we have from Jensen's inequality and the fact that  $\lambda$  belongs to  $[0, 1]$  that

$$\begin{aligned} 186 \quad (4.5) \quad &\mathbb{E}[\|Y_{n+1} - \lambda Z_{n+1}\|^2 | \mathcal{F}_n] \\ &= \mathbb{E}[\|(Y_{n+1} - \nabla f_{U_{n+1}}(x^*)) - \lambda(Z_{n+1} - \nabla f_{U_{n+1}}(x^*)) + (1 - \lambda)\nabla f_{U_{n+1}}(x^*)\|^2 | \mathcal{F}_n] \\ &\leq 2\mathbb{E}[\|Y_{n+1} - \nabla f_{U_{n+1}}(x^*)\|^2 | \mathcal{F}_n] + 2\mathbb{E}[\|Z_{n+1} - \nabla f_{U_{n+1}}(x^*)\|^2 | \mathcal{F}_n] \\ &\quad + 2\mathbb{E}[\|\nabla f_{U_{n+1}}(x^*)\|^2 | \mathcal{F}_n]. \end{aligned}$$

187 First of all, we clearly have

$$188 \quad (4.6) \quad \mathbb{E}[\|\nabla f_{U_{n+1}}(x^*)\|^2 | \mathcal{F}_n] = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x^*)\|^2 \quad a.s.$$

189 In addition, denote

$$190 \quad A_n = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(\phi_{n,k}) - \nabla f_k(x^*)\|^2 \quad \text{and} \quad \Sigma_n = \frac{1}{N} \sum_{k=1}^N \nabla f_k(\phi_{n,k}).$$

191 Since  $\nabla f(x^*) = 0$ , we obtain by expanding the norm that

$$192 \quad (4.7) \quad \mathbb{E}[\|Z_{n+1} - \nabla f_{U_{n+1}}(x^*)\|^2 | \mathcal{F}_n] = A_n - \|\Sigma_n\|^2 \quad a.s.$$

193 Furthermore, define for all  $x \in \mathbb{R}^d$ ,

$$194 \quad \tau^2(x) = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x) - \nabla f_k(x^*)\|^2.$$

195 One can observe that

$$196 \quad (4.8) \quad \mathbb{E}[\|Y_{n+1} - \nabla f_{U_{n+1}}(x^*)\|^2 | \mathcal{F}_n] = \tau^2(X_n) \quad a.s.$$

197 Putting together the three contributions (4.6), (4.7) and (4.8), we deduce from (4.5)  
198 that

$$199 \quad (4.9) \quad \mathbb{E}[\|Y_{n+1} - \lambda Z_{n+1}\|^2 | \mathcal{F}_n] \leq 2(\tau^2(X_n) + A_n + \theta^*) \quad a.s.$$

200 where

$$201 \quad \theta^* = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x^*)\|^2.$$

202 Consequently, it follows from (4.4) and (4.9) that for all  $n \geq 1$ ,

$$203 \quad (4.10) \quad \mathbb{E}[V_{n+1} | \mathcal{F}_n] \leq V_n - 2\gamma_n \langle X_n - x^*, \nabla f(X_n) \rangle + 2\gamma_n^2 (\tau^2(X_n) + A_n + \theta^*) \quad a.s.$$

204 Furthermore, let  $(T_n)$  be the sequence of Lyapunov functions defined, for all  $n \geq 2$ ,  
205 by

$$206 \quad (4.11) \quad T_n = V_n + 2N\gamma_{n-1}^2 A_n.$$

207 It follows from the very definition of the sequence  $(\phi_{n,k})$  associated with (1.3) that

$$\begin{aligned} 208 \quad \mathbb{E}[A_{n+1} | \mathcal{F}_n] &= \frac{1}{N} \sum_{k=1}^N \mathbb{E}[\|\nabla f_k(\phi_{n+1,k}) - \nabla f_k(x^*)\|^2 | \mathcal{F}_n], \\ 209 \quad &= \frac{1}{N} \sum_{k=1}^N \left( \frac{1}{N} \|\nabla f_k(X_n) - \nabla f_k(x^*)\|^2 + \left(1 - \frac{1}{N}\right) \|\nabla f_k(\phi_{n,k}) - \nabla f_k(x^*)\|^2 \right), \\ (4.12) \quad & \\ 210 \quad &= \frac{1}{N} \tau^2(X_n) + \left(1 - \frac{1}{N}\right) A_n, \end{aligned}$$

211 almost surely. Hence, we obtain from (4.10) and (4.12) that

$$\begin{aligned}
212 \quad \mathbb{E}[T_{n+1}|\mathcal{F}_n] &= \mathbb{E}[V_{n+1}|\mathcal{F}_n] + 2N\gamma_n^2\mathbb{E}[A_{n+1}|\mathcal{F}_n], \\
213 \quad &\leq V_n - 2\gamma_n\langle X_n - x^*, \nabla f(X_n) \rangle + 2\gamma_n^2(\tau^2(X_n) + A_n + \theta^*) + 2N\gamma_n^2\mathbb{E}[A_{n+1}|\mathcal{F}_n] \\
214 \quad &= V_n + 2N\gamma_n^2A_n - 2\gamma_n\langle X_n - x^*, \nabla f(X_n) \rangle + 2\gamma_n^2(2\tau^2(X_n) + \theta^*), \\
215 \quad &\leq V_n + 2N\gamma_{n-1}^2A_n - 2\gamma_n\langle X_n - x^*, \nabla f(X_n) \rangle + 2\gamma_n^2(2\tau^2(X_n) + \theta^*), \\
(4.13) \quad &\leq T_n - 2\gamma_n\langle X_n - x^*, \nabla f(X_n) \rangle + 2\gamma_n^2(2\tau^2(X_n) + \theta^*). \quad \blacksquare
\end{aligned}$$

217 Additionally, we clearly have  $V_n \leq T_n$  almost surely and it follows from Assumption  
218 3.3 that

$$219 \quad \tau^2(X_n) \leq LV_n \leq LT_n.$$

220 Finally, we deduce from (4.13) that

$$221 \quad (4.14) \quad \mathbb{E}[T_{n+1}|\mathcal{F}_n] \leq (1 + 4L\gamma_n^2)T_n - 2\gamma_n\langle X_n - x^*, \nabla f(X_n) \rangle + 2\gamma_n^2\theta^* \quad a.s.,$$

222 which can be rewritten as

$$223 \quad \mathbb{E}[T_{n+1}|\mathcal{F}_n] \leq (1 + a_n)T_n + \mathcal{A}_n - \mathcal{B}_n \quad a.s.$$

224 where  $a_n = 4L\gamma_n^2$ ,  $\mathcal{A}_n = 2\gamma_n^2\theta^*$  and  $\mathcal{B}_n = 2\gamma_n\langle X_n - x^*, \nabla f(X_n) \rangle$ . The four sequences  
225  $(T_n)$ ,  $(a_n)$ ,  $(\mathcal{A}_n)$  and  $(\mathcal{B}_n)$  are positive sequences of random variables adapted to  $(\mathcal{F}_n)$ .

226 We clearly have from (1.2) that

$$227 \quad \sum_{n=1}^{\infty} a_n < +\infty \quad \text{and} \quad \sum_{n=1}^{\infty} \mathcal{A}_n < +\infty.$$

228 Then, it follows from the Robbins-Siegmund Theorem [34] given by Theorem A.1 that  
229  $(T_n)$  converges a.s. towards a finite random variable  $T$  and the series

$$230 \quad (4.15) \quad \sum_{n=1}^{\infty} \mathcal{B}_n < +\infty \quad a.s.$$

231 Consequently,  $(V_n)$  also converges a.s. to a finite random variable  $V$ . It only remains  
232 to show that  $V = 0$  almost surely. Assume by contradiction that  $V > 0$ . For some  
233 positive constants  $a < b$ , denote by  $\Omega$  the annulus of  $\mathbb{R}^d$ ,

$$234 \quad \Omega = \{x \in \mathbb{R}^d, \quad 0 < a < \|x - x^*\|^2 < b\}.$$

235 Let  $F$  be the function defined, for all  $x \in \mathbb{R}^d$ , by

$$236 \quad F(x) = \langle x - x^*, \nabla f(x) \rangle.$$

237 We have from Assumption 3.1 that  $F$  is a continuous function in  $\Omega$  compact. It implies  
238 that there exists a positive constant  $c$  such that  $F(x) > c$  for all  $x \in \Omega$ . However, for  $n$   
239 large enough,  $X_n \in \Omega$ , which ensures that  $\gamma_n\langle X_n - x^*, \nabla f(X_n) \rangle > c\gamma_n$ . Consequently,  
240 it follows from (4.15) that

$$241 \quad \sum_{n=1}^{\infty} \gamma_n < +\infty.$$

242 This is of course in contradiction with assumption (1.2). Finally, we obtain that  $V = 0$   
 243 almost surely, leading to

$$244 \quad \lim_{n \rightarrow +\infty} X_n = x^* \quad a.s.$$

245 By continuity of the function  $f$ , we also have (4.2), which completes the proof of  
 246 Theorem 4.1.  $\square$

247 **4.2. Asymptotic normality.** We now focus our attention on the asymptotic  
 248 normality of the  $\lambda$ -SAGA algorithm with decreasing step. In this subsection, we  
 249 assume that  $f$  is twice differentiable and we denote by  $H = \nabla^2 f(x^*)$  the Hessian  
 250 matrix of  $f$  at the point  $x^*$ .

251 **ASSUMPTION 4.2.** *Suppose that  $f$  is twice differentiable with a unique equilibrium*  
 252 *point  $x^*$  in  $\mathbb{R}^d$  such that  $\nabla f(x^*) = 0$ . Denote by  $\rho = \lambda_{\min}(H)$  the minimum eigen-*  
 253 *value of  $H$ . We assume that  $\rho > 1/2$ .*

254 The central limit theorem for the  $\lambda$ -SAGA algorithm is as follows.

255 **THEOREM 4.3.** *Consider a fixed  $\lambda \in [0, 1]$ . Let  $(X_n)$  be the sequence generated by*  
 256 *the  $\lambda$ -SAGA algorithm with decreasing step  $\gamma_n = 1/n$ . Suppose that Assumption 4.2*  
 257 *is satisfied. Assume also that*

$$258 \quad (4.16) \quad \lim_{n \rightarrow +\infty} X_n = x^* \quad a.s.$$

259 *Then, we have the asymptotic normality*

$$260 \quad (4.17) \quad \sqrt{n}(X_n - x^*) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}_d(0, \Sigma)$$

261 *where the asymptotic covariance matrix is given by*

$$262 \quad \Sigma = (1 - \lambda)^2 \int_0^\infty (e^{-(H - \mathbb{I}_d/2)u})^T \Gamma e^{-(H - \mathbb{I}_d/2)u} du,$$

263 *with*

$$264 \quad \Gamma = \frac{1}{N} \sum_{k=1}^N \nabla f_k(x^*) (\nabla f_k(x^*))^T.$$

265 *Proof.* The proof of Theorem 4.3 can be found in Appendix B.  $\square$

266 **Remark 4.4.** It was proven in Theorem 4.1 that the almost sure convergence (4.16)  
 267 holds under Assumptions 3.1, 3.2 and 3.3. It is obvious to see that Assumption 4.2  
 268 implies Assumption 3.1. Consequently, Theorem 4.3 is also true when replacing (4.16)  
 269 by Assumptions 3.2 and 3.3.

270 Many conclusions can be drawn from Theorem 4.3. First of all, if we assume  
 271 that  $\Gamma$  is a positive definite matrix, we obtain that our  $\lambda$ -SAGA algorithm with  
 272  $0 \leq \lambda < 1$ , converges towards a centered normal distribution with positive definite  
 273 variance. However, as soon as  $\lambda = 1$ , the limit distribution becomes a centered  
 274 normal with variance  $\Sigma = 0$ , in other words a Dirac measure. Thus, the asymptotic  
 275 distribution of the SAGA algorithm has zero variance and one can therefore try to  
 276 understand it. In fact, the conditional variances of the two terms of the martingale  
 277 difference  $(\varepsilon_n)$  extracted from this algorithm, converge almost surely to exactly the  
 278 same matrix. Therefore, the conditional variance of  $(\varepsilon_n)$  vanishes which explains  
 279 the final result for the SAGA algorithm. Moreover, Theorem 4.3 clearly shows the

280 asymptotic variance reduction effect. Indeed, when  $\lambda$  grows to 1, we observe that  
 281 the variance  $\Sigma$  decreases and converges towards 0. Hence, for statistical inference  
 282 purposes such that hypothesis test and confidence interval, we can take  $\lambda$  just a little  
 283 smaller than 1 to reduce the variance with respect to SGD, but without canceling it.

284 **4.3. Non-asymptotic convergence rates.** In the same vein as [4] for the  
 285 Robbins-Monro algorithm, we shall now establish non-asymptotic  $\mathbf{L}^p$  convergence  
 286 rates. Hence, our goal is to investigate, for all integer  $p \geq 1$ , the convergence rate of  
 287  $\mathbb{E}[\|X_n - x^*\|^{2p}]$  for the  $\lambda$ -SAGA algorithm where the decreasing step is defined, for  
 288 all  $n \geq 1$  by,

$$289 \quad (4.18) \quad \gamma_n = \frac{c}{n^\alpha},$$

290 for some positive constant  $c$  and  $1/2 < \alpha \leq 1$ . First of all, we focus our attention  
 291 on the standard case  $p = 1$  by analyzing our algorithm with a little more stringent  
 292 condition than Assumption 3.2.

293 **ASSUMPTION 4.5.** *Assume there exists a positive constant  $\mu$  such that for all  $x \in$   
 294  $\mathbb{R}^d$  with  $x \neq x^*$ ,*

$$295 \quad \langle x - x^*, \nabla f(x) \rangle \geq \mu \|x - x^*\|^2.$$

296 Although this is a strengthened version of Assumption 3.2, it is still weaker than the  
 297 usual  $\mu$ -strong convexity assumption on the function  $f$ . This condition is sometimes  
 298 called in the literature the Restricted Secant Inequality [18, 19, 39].

299 **THEOREM 4.6.** *Consider a fixed  $\lambda \in [0, 1]$ . Let  $(X_n)$  be the sequence generated by  
 300 the  $\lambda$ -SAGA algorithm with decreasing step sequence  $(\gamma_n)$  defined by (4.18). Suppose  
 301 that Assumptions 3.1, 3.3 and 4.5 are satisfied with  $2c\mu \leq 2^\alpha$  and  $2c\mu > 1$  if  $\alpha = 1$ .  
 302 Then, there exists a positive constant  $K$  such that for all  $n \geq 1$ ,*

$$303 \quad (4.19) \quad \mathbb{E}[\|X_n - x^*\|^2] \leq \frac{K}{n^\alpha}.$$

304 *Proof.* The proof of Theorem 4.6 can be found in Appendix C. □

305 Next, we carry out our analysis in the general case  $p \geq 1$ . It requires a strengthened  
 306 version of Assumption 3.3 given as follows.

307 **ASSUMPTION 4.7.** *Assume that for some integer  $p \geq 1$ , there exists a positive  
 308 constant  $L_p$  such that for all  $x \in \mathbb{R}^d$ ,*

$$309 \quad \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x) - \nabla f_k(x^*)\|^{2p} \leq L_p \|x - x^*\|^{2p}.$$

310 **THEOREM 4.8.** *Consider a fixed  $\lambda \in [0, 1]$ . Let  $(X_n)$  be the sequence generated by  
 311 the  $\lambda$ -SAGA algorithm with decreasing step sequence  $(\gamma_n)$  defined by (4.18) and such  
 312 that the initial state  $X_1$  belongs to  $\mathbf{L}^{2p}$ . Suppose that Assumptions 3.1, 4.5 and 4.7  
 313 are satisfied with  $pc\mu \leq 2^\alpha$  and  $c\mu > 1$  if  $\alpha = 1$ . Then, there exists a positive constant  
 314  $K_p$  such that for all  $n \geq 1$ ,*

$$315 \quad (4.20) \quad \mathbb{E}[\|X_n - x^*\|^{2p}] \leq \frac{K_p}{n^{p\alpha}}.$$

316 *Proof.* The proof of Theorem 4.8 can be found in Appendix D. □

317 *Remark 4.9.* It is easy to see that Assumption 4.7 also implies that for all  $x \in \mathbb{R}^d$ ,  
 318

$$319 \quad (4.21) \quad (f(x) - f(x^*))^p \leq \frac{\sqrt{L_p}}{2^p} \|x - x^*\|^{2p}.$$

320 Then, it follows from Theorem 4.8 together with inequality (4.21) that there exists a  
 321 positive constant  $M_p = 2^{-p} K_p \sqrt{L_p}$  such that for all  $n \geq 1$ ,

$$322 \quad (4.22) \quad \mathbb{E}[(f(X_n) - f(x^*))^p] \leq \frac{M_p}{n^{p\alpha}}.$$

323 **5. Numerical experiments.** Consider the logistic regression model (see, e.g.,  
 324 [3, 6]) associated with the classical minimization problem ( $\mathcal{P}$ ) of the convex function  
 325  $f$  given, for all  $x \in \mathbb{R}^d$ , by

$$326 \quad f(x) = \frac{1}{N} \sum_{k=1}^N f_k(x) = \frac{1}{N} \sum_{k=1}^N (\log(1 + \exp(\langle x, w_k \rangle)) - y_k \langle x, w_k \rangle),$$

327 where  $x \in \mathbb{R}^d$  is a vector of unknown parameters,  $w_k \in \mathbb{R}^d$  is a vector of features  
 328 and the binary output  $y_k \in \{0, 1\}$ . As stated, this problem is equivalent to the log-  
 329 likelihood maximization problem, where the aim is to find the parameter  $x^*$  that  
 330 maximizes the probability of a given sample  $((w_1, y_1), \dots, (w_N, y_N))$ , which follows a  
 331 model depending only on the unknown parameter  $x$ . To be more precise, our model  
 332 has a Bernoulli probability with parameter  $p_k(x)$  following a logistic function for each  
 333  $k$ ,

$$334 \quad p_k(x) = \frac{\exp(\langle x, w_k \rangle)}{1 + \exp(\langle x, w_k \rangle)}.$$

335 It is easy to see that  $f$  is twice differentiable and its Hessian matrix is given by

$$336 \quad \nabla^2 f(x) = \frac{1}{N} \sum_{k=1}^N p_k(x)(1 - p_k(x))w_k w_k^T.$$

337 Consequently,  $f$  has an unique equilibrium point  $x^*$  and if we assume that the min-  
 338 imum eigenvalue of  $H = \nabla^2 f(x^*)$  is greater than  $1/2$ , Assumption 4.2 will be au-  
 339 tomatically satisfied, and therefore Assumption 3.1 too. Moreover, one can observe  
 340 that Assumptions 3.3 and 4.7 hold with

$$341 \quad L_p = \frac{1}{4^p N} \sum_{k=1}^N \|w_k\|^{4p}.$$

342 We conducted experiments on the MNIST dataset in order to present a visualisation  
 343 of the almost sure convergence in Theorem 4.1, the asymptotic normality in Theorem  
 344 4.3 and the  $\mathbf{L}^2$  bound in Theorem 4.6. For the almost sure convergence, the training  
 345 database considered here includes  $N = 60000$  images in gray-scale format and size  
 346  $28 \times 28$ . Each image  $w_k$  is therefore a vector of dimension  $d = 28 \times 28 = 784$ . Each  
 347 of these images is identified with a number from 0 to 9 and we divide it into a binary  
 348 classification so that  $y_k = 0$  if  $\{0, 1, 2, 3, 4\}$  and  $y_k = 1$  if  $\{5, 6, 7, 8, 9\}$ . The results

349 concerning the convergence of the estimator  $X_n$  of  $x^*$  are illustrated in Figure 1. The  
 350 convergence is ordered from slowest to fastest in an increasing order with respect to  
 351  $\lambda \in \{0, 0.5, 0.9, 1\}$ .

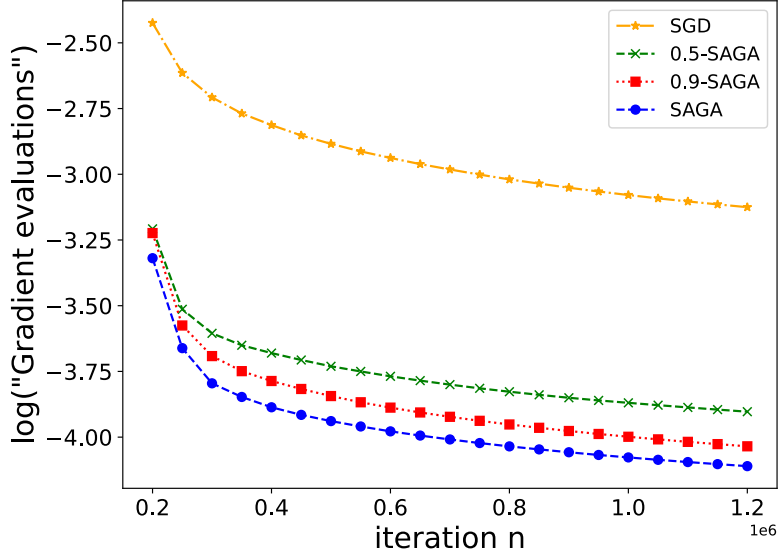


Fig. 1: Convergence with  $\gamma_n = 1/n$  for 1.2M of iterations. Here we put “Gradient evaluations” since instead of using  $\|\nabla f(X_n)\|$ , we use the norm of the mean associated to the lines in the matrix  $g_n$ ,  $\|\sum_{k=1}^N g_{n,k}\|/N$ . This quantity keeps track of the convergence since it also converges to 0 and its lines converge to the gradients of the functions  $f_k$ , that is for each  $1 \leq k \leq N$ ,  $\lim g_{n,k} = \nabla f_k(x^*)$  as  $n$  goes to infinity.

352 Moreover, to illustrate the asymptotic normality result, we use  $N = 100$ , the first  
 353 100 images in the MNIST dataset, and the distributional convergence

$$354 \quad \lim_{n \rightarrow \infty} \mathbb{E}(h(\sqrt{n}(X_n - x^*))) = \mathbb{E}(h(\mathcal{N}_d(0, \Sigma))),$$

355 where  $h$  is defined, for all  $x \in \mathbb{R}^d$ , by  $h(x) = \sum_{i=1}^d x_i$ . It follows from Theorem 4.3 that

$$356 \quad h(\sqrt{n}(X_n - x^*)) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, \sigma^2(\lambda)).$$

357 As the equilibrium point  $x^*$  and the asymptotic variance  $\sigma^2(\lambda)$  are unknown, we  
 358 use estimators from the standard Monte Carlo procedure. We denote for each fixed  
 359 lambda  $\hat{\sigma}_n^2(\lambda)$  the estimator of  $\sigma^2(\lambda)$ . Given the form of our function  $h$ , we deduce  
 360 that the limiting variances should be related as  $(1 - \lambda)^2 \sigma^2(0) = \sigma^2(\lambda)$ . The results  
 361 are shown in Figure 2.

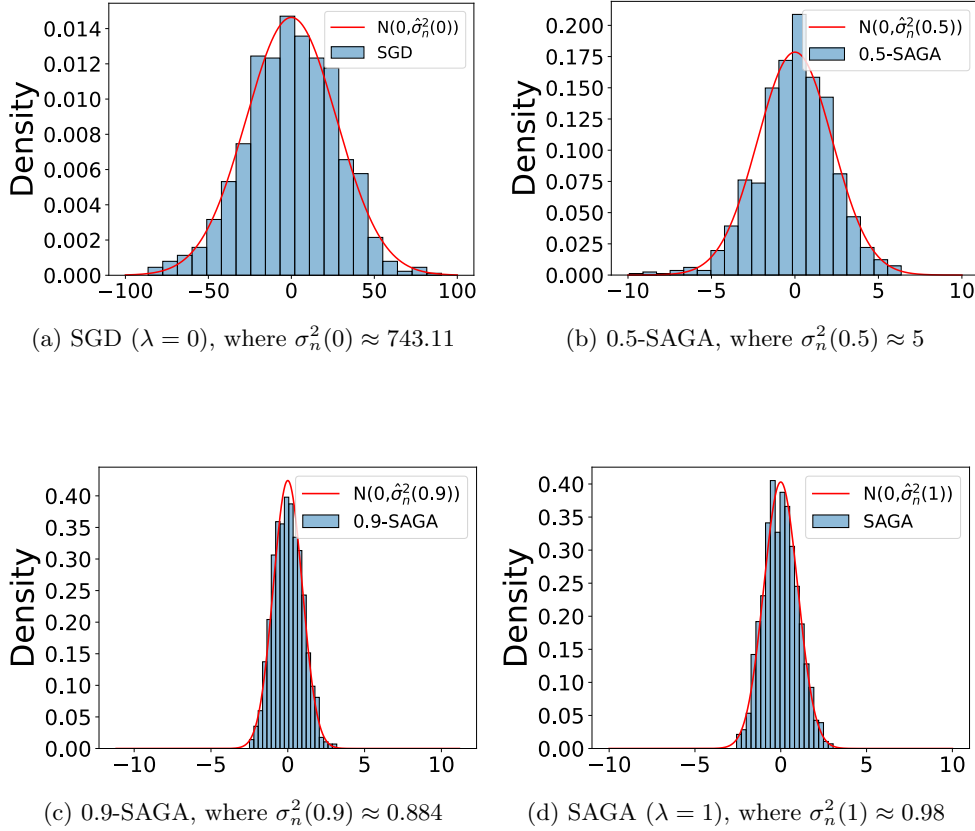


Fig. 2: We used 1000 samples, where each one was obtained by running the associated algorithm for  $n = 500000$  iterations.

362 The main purpose of this plot is to represent the decreasing behavior of the  
 363 variance with respect to the parameter  $\lambda$ . Even though for the SAGA ( $\lambda = 1$ )  
 364 we know that its variance converges to zero, for finite  $n$  we can only see that it is  
 365 shrinking with respect to  $n$  to obtain at the limit a Dirac mass at 0. Here, the sample  
 366 variances satisfy  $\hat{\sigma}_n^2(0.9) < \hat{\sigma}_n^2(1)$ . Nevertheless, they are still very close in the scale  
 367 of the sample variance  $\hat{\sigma}_n^2(0)$  of the SGD. We explain this as a consequence of the  
 368 approximations and the fact that the models  $\lambda = 0.9$  and  $\lambda = 1$  are intimately related.

369 Finally, we present approximate results of the mean squared error  $\mathbb{E}[\|X_n - x^*\|^2]$ .  
 370 For that purpose, we suppose that Assumption 4.5 is satisfied so that Theorem 4.6  
 371 holds. We run each algorithm for 100 epochs, where each epoch consists of 1000  
 372 iterations. In order to approximate the expectation, we apply the standard Monte  
 373 Carlo procedure with 1000 samples. Here, the approximation of  $x^*$  is the result of  
 374 running the SAGA algorithm for  $40M$  iterations. The results are illustrated in Figure  
 375 3. The plot just gives an intuition on the behavior of the mean squared error, since  
 376 the constant  $K$  in Theorem 4.6 is unknown.

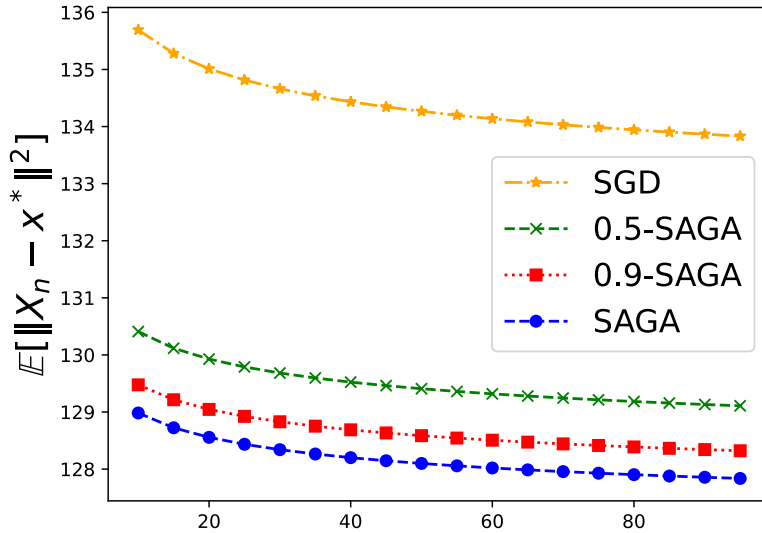


Fig. 3: Mean squared error with respect to epochs. We confirm the decreasing order of the mean squared error of  $X_n - x^*$  with respect to  $\lambda$  and  $n$ .

377 **6. Conclusion.** Stochastic optimization is one of the main challenges of ma-  
 378 chine learning touching almost every aspect of the discipline. Thus, in order to meet  
 379 expectations, the SGD algorithm has been studied at length. However, the advent of  
 380 Big Data for model learning led to the development of more sophisticated stochastic  
 381 methods. In our study, we therefore highlight the properties of the new  $\lambda$ -SAGA al-  
 382 gorithm which is a generalization of the SAGA algorithm. We were able to establish  
 383 the almost sure convergence and the asymptotic normality of this novel algorithm  
 384 by using a decreasing step and without the strong convexity and Lipschitz gradient  
 385 assumptions. The other major contribution of our paper concerns the convergence  
 386 rates in  $\mathbf{L}^p$  of the  $\lambda$ -SAGA algorithm. Finally, stochastic algorithms offer multiple  
 387 guarantees in terms of convergence and certainly promise to continue to have profound  
 388 impacts on the fast development of the machine learning field.

389 **Appendix A. Some useful existing results.** We first recall the well-known  
 390 Robbins-Siegmund Theorem [34].

391 **THEOREM A.1** (Robbins-Siegmund theorem). *Let  $(V_n), (a_n), (\mathcal{A}_n), (\mathcal{B}_n)$  be four*  
 392 *positive sequences of random variables adapted to a filtration  $(\mathcal{F}_n)$  such that*

$$393 \quad \mathbb{E}[V_{n+1} | \mathcal{F}_n] \leq (1 + a_n)V_n + \mathcal{A}_n - \mathcal{B}_n,$$

394 *where*

$$395 \quad \sum_{n=1}^{\infty} a_n < +\infty \quad \text{and} \quad \sum_{n=1}^{\infty} \mathcal{A}_n < +\infty \quad \text{a.s.}$$

396 *Then,  $(V_n)$  converges almost surely towards a finite random variable  $V$  and*

$$397 \quad \sum_{n=1}^{\infty} \mathcal{B}_n < +\infty \quad \text{a.s.}$$

398 The next two lemma provide very useful inequality for non-asymptotic convergence  
 399 rates. The first lemma is given by Lemma A.3 in supplementary material of [5], see  
 400 also Theorem 1 in [4].

401 LEMMA A.2. [5]. Let  $(Z_n)$  be a sequence of positive real numbers satisfying, for  
 402 all  $n \geq 1$ , the recursive inequality

$$403 \quad (A.1) \quad Z_{n+1} \leq \left(1 - \frac{a}{(n+1)^\alpha}\right) Z_n + \frac{b}{(n+1)^\beta},$$

404 where  $a, b, \alpha$  and  $\beta$  are positive constants satisfying  $a \leq 2^\alpha$ ,  $\alpha \leq 1$ ,  $1 < \beta < 2$  and  
 405  $\beta \leq 2\alpha$  with  $\beta < a + 1$  in the special case where  $\alpha = 1$ . Then, there exists a positive  
 406 constant  $C$  such that, for any  $n \geq 1$ ,

$$407 \quad (A.2) \quad Z_n \leq \frac{C}{n^{\beta-\alpha}}.$$

408 The second lemma is given without proof in [10] in the special case  $p \in (0, 2]$ , see  
 409 Lemma B.3 as well as the seminal paper [2]. We extend it to the case  $p$  even and we  
 410 propose a short proof for the sake of completeness.

411 LEMMA A.3. Let  $p$  be a positive even integer. There exist two positive constant  
 412  $C_p$  and  $D_p$  such that for any  $a, b \in \mathbb{R}^d$ ,

$$413 \quad (A.3) \quad \|a + b\|^{2+p} \leq \|a\|^{2+p} + (2+p)\langle a, b \rangle \|a\|^p + C_p \|a\|^p \|b\|^2 + D_p \|b\|^{2+p}.$$

414 *Proof.* We prove Lemma A.3 by induction. For the base case  $p = 2$ , we have

$$415 \quad \|a + b\|^4 = (\|a\|^2 + 2\langle a, b \rangle + \|b\|^2)^2 \\ 416 \quad = \|a\|^4 + 4(\langle a, b \rangle)^2 + \|b\|^4 + 4\langle a, b \rangle \|a\|^2 + 4\langle a, b \rangle \|b\|^2 + 2\|a\|^2 \|b\|^2.$$

417 It follows from Cauchy–Schwarz inequality that  $(\langle a, b \rangle)^2 \leq \|a\|^2 \|b\|^2$ . Moreover, we  
 418 also have  $2\langle a, b \rangle \leq \|a\|^2 + \|b\|^2$  which implies that  $4\langle a, b \rangle \|b\|^2 \leq 2\|a\|^2 \|b\|^2 + 2\|b\|^4$ .  
 419 Hence, we obtain from these two inequalities that

$$420 \quad \|a + b\|^4 \leq \|a\|^4 + 4\langle a, b \rangle \|a\|^2 + 8\|a\|^2 \|b\|^2 + 3\|b\|^4,$$

421 which leads to  $C_2 = 8$  and  $D_2 = 3$ . Hereafter, assume that inequality (A.3) holds up  
 422 to  $q \geq 2$  and let  $p = 2 + q$ . We have by induction

$$423 \quad \|a + b\|^{2+p} = \|a + b\|^2 \|a + b\|^p = \|a + b\|^2 \|a + b\|^{2+q} \\ 424 \quad \leq \left(\|a\|^2 + 2\langle a, b \rangle + \|b\|^2\right) \left(\|a\|^{2+q} + (2+q)\langle a, b \rangle \|a\|^q + C_q \|a\|^q \|b\|^2 + D_q \|b\|^{2+q}\right) \\ 425 \quad \leq \|a\|^{2+p} + (2+p)\langle a, b \rangle \|a\|^p + (C_q + 2p + 1)\|a\|^p \|b\|^2 + (2C_q + p)\|a\|^{q+1} \|b\|^3 \\ 426 \quad + D_q \|a\|^2 \|b\|^p + C_q \|a\|^q \|b\|^4 + 2D_q \|a\| \|b\|^{p+1} + D_q \|b\|^{2+p}$$

427 where the last inequality is the result of applying Cauchy–Schwarz inequality  $\langle a, b \rangle \leq$   
 428  $\|a\| \|b\|$  for all the terms multiplied by  $\langle a, b \rangle$  but the ones isolated in the second term.

429 Furthermore, it follows from Young's inequality for products that

$$\begin{aligned}
430 \quad \|a\|^{q+1}\|b\|^3 &= \|a\|^{q+1}\|b\|^{2(q+1)/p} \times \|b\|^{3-2(q+1)/p} \leq \frac{\|a\|^p\|b\|^2}{p/(p-1)} + \frac{\|b\|^{p+2}}{p}, \\
431 \quad \|a\|^2\|b\|^p &= \|a\|^2\|b\|^{4/p} \times \|b\|^{p-4/p} \leq \frac{\|a\|^p\|b\|^2}{p/2} + \frac{\|b\|^{p+2}}{p/q}, \\
432 \quad \|a\|^q\|b\|^4 &= \|a\|^q\|b\|^{2q/p} \times \|b\|^{4-2q/p} \leq \frac{\|a\|^p\|b\|^2}{p/q} + \frac{\|b\|^{p+2}}{p/2}, \\
433 \quad \|a\|\|b\|^{p+1} &= \|a\|\|b\|^{2/p} \times \|b\|^{p+1-2/p} \leq \frac{\|a\|^p\|b\|^2}{p} + \frac{\|b\|^{p+2}}{p/(p-1)}.
\end{aligned}$$

434 Finally, we obtain (A.3) with  $C_p$  and  $D_p$  satisfying the system defined, for  $p \geq 4$ , by

$$\begin{cases}
435 \quad C_p &= 3p + \frac{4}{p}((p-1)C_{p-2} + D_{p-2}) \\
D_p &= 1 + \frac{4}{p}(C_{p-2} + (p-1)D_{p-2})
\end{cases}$$

436 with initial values  $C_2 = 8$  and  $D_2 = 3$ , which achieves the proof of Lemma A.3.  $\square$

437 *Remark A.4.* One can easily compute  $C_4 = 39$  and  $D_4 = 18$ . Moreover, one can  
438 observe that we always have  $D_p \leq C_p$ . Consequently, we can make use of (A.3) with  
439  $C_p$  instead of  $D_p$ .

#### 440 **Appendix B. Proof of Theorem 4.3.**

441 *Proof.* The  $\lambda$ -SAGA algorithm can be rewritten as

$$442 \quad X_{n+1} = X_n - \gamma_n(\nabla f(X_n) + \varepsilon_{n+1}),$$

443 where  $\varepsilon_{n+1} = \mathcal{Y}_{n+1} - \lambda Z_{n+1}$  with

$$\begin{cases}
444 \quad \mathcal{Y}_{n+1} &= \nabla f_{U_{n+1}}(X_n) - \nabla f(X_n), \\
Z_{n+1} &= \nabla f_{U_{n+1}}(\phi_{n,U_{n+1}}) - \frac{1}{N} \sum_{k=1}^N \nabla f_k(\phi_{n,k}).
\end{cases}$$

445 We already saw that  $(\varepsilon_n)$  is a martingale difference adapted to the filtration  $(\mathcal{F}_n)$ .

446 Moreover,

$$\begin{aligned}
447 \quad \mathbb{E}[\varepsilon_{n+1}\varepsilon_{n+1}^T | \mathcal{F}_n] &= \mathbb{E}[\mathcal{Y}_{n+1}\mathcal{Y}_{n+1}^T | \mathcal{F}_n] - \lambda \mathbb{E}[\mathcal{Y}_{n+1}Z_{n+1}^T | \mathcal{F}_n] - \lambda \mathbb{E}[Z_{n+1}\mathcal{Y}_{n+1}^T | \mathcal{F}_n] \\
448 \quad &+ \lambda^2 \mathbb{E}[Z_{n+1}Z_{n+1}^T | \mathcal{F}_n].
\end{aligned}$$

449 In addition, we clearly have that almost surely

$$\begin{aligned}
450 \quad \mathbb{E}[\mathcal{Y}_{n+1}\mathcal{Y}_{n+1}^T | \mathcal{F}_n] &= \frac{1}{N} \sum_{k=1}^N \nabla f_k(X_n) (\nabla f_k(X_n))^T - \nabla f(X_n) (\nabla f(X_n))^T, \\
451 \quad \mathbb{E}[Z_{n+1}Z_{n+1}^T | \mathcal{F}_n] &= \frac{1}{N} \sum_{k=1}^N \nabla f_k(\phi_{n,k}) (\nabla f_k(\phi_{n,k}))^T - \frac{1}{N^2} \left( \sum_{k=1}^N \nabla f_k(\phi_{n,k}) \right) \left( \sum_{k=1}^N (\nabla f_k(\phi_{n,k}))^T \right), \\
452 \quad \mathbb{E}[Z_{n+1}\mathcal{Y}_{n+1}^T | \mathcal{F}_n] &= \frac{1}{N} \sum_{k=1}^N \nabla f_k(\phi_{n,k}) (\nabla f_k(X_n))^T - \frac{1}{N} \sum_{k=1}^N \nabla f_k(\phi_{n,k}) (\nabla f(X_n))^T.
\end{aligned}$$

453 We now claim that for all  $1 \leq k \leq N$

$$454 \quad (\text{B.1}) \quad \lim_{n \rightarrow +\infty} \phi_{n,k} = x^* \quad a.s.$$

455 As a matter of fact, for a fixed value  $1 \leq k \leq N$ , the probability that  $U_n = k$  occurs  
 456 for infinitely many  $n$ . Consequently,  $(\phi_{n,k})$  is a sub-sequence of  $(X_n)$ , since  $\phi_{n,k}$  is  
 457 updated to  $X_n$  each time  $\{U_n = k\}$ . Hence, the almost sure convergence (B.1) follows  
 458 from (4.16). Combining the almost sure convergence of  $(X_n)$  and  $(\phi_{n,k})$  towards  $x^*$   
 459 with the continuity of  $\nabla f$  given by Assumption 4.2, it follows that almost surely

$$460 \quad \begin{cases} \lim_{n \rightarrow +\infty} \mathbb{E}[\mathcal{Y}_{n+1} \mathcal{Y}_{n+1}^T | \mathcal{F}_n] = \Gamma, \\ \lim_{n \rightarrow +\infty} \mathbb{E}[\mathcal{Y}_{n+1} Z_{n+1}^T | \mathcal{F}_n] = \Gamma, \\ \lim_{n \rightarrow +\infty} \mathbb{E}[Z_{n+1} \mathcal{Y}_{n+1}^T | \mathcal{F}_n] = \Gamma, \\ \lim_{n \rightarrow +\infty} \mathbb{E}[Z_{n+1} Z_{n+1}^T | \mathcal{F}_n] = \Gamma, \end{cases}$$

461 where

$$462 \quad \Gamma = \frac{1}{N} \sum_{k=1}^N \nabla f_k(x^*) (\nabla f_k(x^*))^T,$$

463 which leads to

$$464 \quad \lim_{n \rightarrow +\infty} \mathbb{E}[\varepsilon_{n+1} \varepsilon_{n+1}^T | \mathcal{F}_n] = (1 - \lambda)^2 \Gamma \quad a.s.$$

465 Therefore, we obtain from Toeplitz's lemma that

$$466 \quad \lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{k=1}^n \mathbb{E}[\varepsilon_k \varepsilon_k^T | \mathcal{F}_{k-1}] = (1 - \lambda)^2 \Gamma \quad a.s.$$

467 In addition, we have for all  $n \geq 1$

$$468 \quad (\text{B.2}) \quad \|\varepsilon_{n+1}\| \leq 2 \left( \max_{k=1, \dots, N} \|\nabla f_k(X_n)\| + \lambda \max_{k=1, \dots, N} \|\nabla f_k(\phi_{n,k})\| \right).$$

469 Hence, it follows from (B.2) that

$$470 \quad (\text{B.3}) \quad \|\varepsilon_{n+1}\|^4 \leq 128 \left( \max_{k=1, \dots, N} \|\nabla f_k(X_n)\|^4 + \lambda^4 \max_{k=1, \dots, N} \|\nabla f_k(\phi_{n,k})\|^4 \right).$$

471 However, for all  $k = 1, \dots, N$ , we have

$$472 \quad \|\nabla f_k(X_n)\|^2 \leq 2 \|\nabla f_k(X_n) - \nabla f_k(x^*)\|^2 + 2 \|\nabla f_k(x^*)\|^2,$$

473 which implies that

$$474 \quad \max_{k=1, \dots, N} \|\nabla f_k(X_n)\|^2 \leq 2N \left( \tau^2(X_n) + \theta^* \right).$$

475 Consequently,

$$476 \quad (\text{B.4}) \quad \max_{k=1, \dots, N} \|\nabla f_k(X_n)\|^4 \leq 4N^2 \left( \tau^2(X_n) + \theta^* \right)^2.$$

477 By the same token,

$$478 \quad (\text{B.5}) \quad \max_{k=1, \dots, N} \|\nabla f_k(\phi_{n,k})\|^4 \leq 4N^2 (A_n + \theta^*)^2.$$

479 Hence, we obtain from (B.3), (B.4) and (B.5) that

$$480 \quad \|\varepsilon_{n+1}\|^4 \leq 512N^2 \left( (\tau^2(X_n) + \theta^*)^2 + \lambda^4 (A_n + \theta^*)^2 \right),$$

481 which immediately implies that

$$482 \quad (\text{B.6}) \quad \mathbb{E}[\|\varepsilon_{n+1}\|^4 | \mathcal{F}_n] \leq 512N^2 \left( (\tau^2(X_n) + \theta^*)^2 + \lambda^4 (A_n + \theta^*)^2 \right) \quad a.s.$$

483 Moreover, since  $X_n$  converges towards  $x^*$ , it follows that  $\tau^2(X_n)$  converges to 0 almost  
484 surely. Combining this result with the almost sure convergence of  $A_n$  towards 0 and  
485 (B.6), we find that

$$486 \quad \sup_{n \geq 1} \mathbb{E}[\|\varepsilon_{n+1}\|^4 | \mathcal{F}_n] < +\infty,$$

487 which implies that for all  $\epsilon > 0$ ,

$$488 \quad \lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{k=1}^n \mathbb{E}[\|\varepsilon_k\|^2 \mathbf{1}_{\{\|\varepsilon_k\| \geq \epsilon \sqrt{n}\}} | \mathcal{F}_{k-1}] = 0 \quad a.s.$$

489 Finally, it follows from the central limit theorem for stochastic algorithms given by  
490 Theorem 2.3 in [41] that

$$491 \quad \sqrt{n}(X_n - x^*) \xrightarrow[n \rightarrow +\infty]{\mathcal{L}} \mathcal{N}_d(0, \Sigma),$$

492 where

$$493 \quad \Sigma = (1 - \lambda)^2 \int_0^\infty (e^{-(H - \mathbb{I}_d/2)u})^T \Gamma e^{-(H - \mathbb{I}_d/2)u} du,$$

494 which completes the proof of Theorem 4.3.  $\square$

### 495 **Appendix C. Proof of Theorem 4.6.**

496 *Proof.* We already saw in (4.10) that for all  $n \geq 1$ ,

$$497 \quad \mathbb{E}[V_{n+1} | \mathcal{F}_n] \leq V_n - 2\gamma_n \langle X_n - x^*, \nabla f(X_n) \rangle + 2\gamma_n^2 (\tau^2(X_n) + A_n + \theta^*) \quad a.s.$$

498 Hence, it follows from Assumption 4.5 that  $\langle X_n - x^*, \nabla f(X_n) \rangle \geq \mu V_n$ , which leads to

$$499 \quad \mathbb{E}[V_{n+1} | \mathcal{F}_n] \leq (1 - 2\mu\gamma_n)V_n + 2\gamma_n^2 (\tau^2(X_n) + A_n + \theta^*) \quad a.s.$$

500 By taking the expectation on both side of this inequality, we obtain that for all  $n \geq 1$ ,  
501

$$502 \quad (\text{C.1}) \quad \mathbb{E}[V_{n+1}] \leq (1 - 2\mu\gamma_n)\mathbb{E}[V_n] + 2\gamma_n^2 (\mathbb{E}[\tau^2(X_n)] + \mathbb{E}[A_n] + \theta^*).$$

503 Furthermore, we deduce from Corollary E.2 in Appendix E below that there exist  
504 positive constants  $b_1$  and  $b_2$  such that, for all  $n \geq 1$ ,  $\mathbb{E}[\tau^2(X_n)] \leq b_1$  and  $\mathbb{E}[A_n] \leq b_2$ .  
505 Consequently, (C.1) immediately leads, for all  $n \geq 1$ , to

$$506 \quad \mathbb{E}[V_{n+1}] \leq \left( 1 - \frac{a}{(n+1)^\alpha} \right) \mathbb{E}[V_n] + \frac{b}{(n+1)^{2\alpha}}$$

507 where  $a = 2\mu c$  and  $b = c^2 2^{2\alpha+1}(b_1 + b_2 + \theta^*)$ . Finally, we can conclude from Lemma  
508 A.2 that there exists a positive constant  $K$  such that for any  $n \geq 1$ ,

$$509 \quad \mathbb{E}[\|X_n - x^*\|^2] \leq \frac{K}{n^\alpha},$$

510 which completes the proof of Theorem 4.6.  $\square$

#### 511 **Appendix D. Proof of Theorem 4.8.**

512 *Proof.* First of all, Theorem 4.8 follows from Theorem 4.6 in the special case  
513  $p = 1$ . Hence, we are going to prove Theorem 4.8 by induction on  $n \geq 1$  for some  
514 integer  $p \geq 2$  satisfying Assumption 4.7. As the initial state  $X_1$  belongs to  $\mathbf{L}^{2p}$ , the  
515 base case is immediately true. Next, assume by induction that for some integer  $m$   
516 which will be fixed soon, there exists a positive constant  $K_p$  such that for all  $n \leq m$ ,

$$517 \quad (\text{D.1}) \quad \mathbb{E}[\|X_n - x^*\|^{2p}] \leq \frac{K_p}{n^{p\alpha}}.$$

518 We have from (4.3) together with Lemma A.3 that it exists a positive constant  $C_p$   
519 such that for all  $n \geq 1$ ,

$$\begin{aligned} 520 \quad V_{n+1}^p &= \|X_{n+1} - x^*\|^{2p}, \\ 521 \quad &= \|X_n - x^* - \gamma_n(Y_{n+1} - \lambda Z_{n+1})\|^{2p}, \\ 522 \quad &\leq V_n^p - 2p\gamma_n V_n^{p-1} \langle X_n - x^*, Y_{n+1} - \lambda Z_{n+1} \rangle + C_p \gamma_n^2 V_n^{p-1} \|Y_{n+1} - \lambda Z_{n+1}\|^2 \\ 523 \quad &+ C_p \gamma_n^{2p} \|Y_{n+1} - \lambda Z_{n+1}\|^{2p}. \end{aligned}$$

524 Hence, it follows from (4.4) that

$$\begin{aligned} 525 \quad (\text{D.2}) \quad \mathbb{E}[V_{n+1}^p | \mathcal{F}_n] &\leq V_n^p - 2p\gamma_n V_n^{p-1} \langle X_n - x^*, \nabla f(X_n) \rangle \\ 526 \quad &+ C_p \gamma_n^2 V_n^{p-1} \mathbb{E}[\|Y_{n+1} - \lambda Z_{n+1}\|^2 | \mathcal{F}_n] + C_p \gamma_n^{2p} \mathbb{E}[\|Y_{n+1} - \lambda Z_{n+1}\|^{2p} | \mathcal{F}_n]. \end{aligned}$$

527 We already saw from (4.9) that

$$528 \quad \mathbb{E}[\|Y_{n+1} - \lambda Z_{n+1}\|^2 | \mathcal{F}_n] \leq 2(\tau^2(X_n) + A_n + \theta^*) \quad a.s.$$

529 which leads, via Assumption 4.7, to

$$530 \quad \mathbb{E}[\|Y_{n+1} - \lambda Z_{n+1}\|^2 | \mathcal{F}_n] \leq 2(L_p^{1/p} V_n + A_n + \theta^*) \quad a.s.$$

531 Moreover, as in the proof of (4.5), we deduce from Jensen's inequality that

$$\begin{aligned} 532 \quad (\text{D.3}) \quad \mathbb{E}[\|Y_{n+1} - \lambda Z_{n+1}\|^{2p} | \mathcal{F}_n] &\leq 3^{2p-1} \mathbb{E}[\|Y_{n+1} - \nabla f_{U_{n+1}}(x^*)\|^{2p} | \mathcal{F}_n] \\ 533 \quad &+ 3^{2p-1} \mathbb{E}[\|Z_{n+1} - \nabla f_{U_{n+1}}(x^*)\|^{2p} | \mathcal{F}_n] + 3^{2p-1} \mathbb{E}[\|\nabla f_{U_{n+1}}(x^*)\|^{2p} | \mathcal{F}_n]. \end{aligned}$$

534 Hereafter, we clearly have

$$535 \quad (\text{D.4}) \quad \mathbb{E}[\|\nabla f_{U_{n+1}}(x^*)\|^{2p} | \mathcal{F}_n] = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x^*)\|^{2p} \quad a.s.$$

536 Moreover, denote

$$537 \quad A_{p,n} = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(\phi_{n,k}) - \nabla f_k(x^*)\|^{2p} \quad \text{and} \quad \Sigma_n = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(\phi_{n,k})\|^{2p}.$$

538 It follows once again from Jensen's inequality that

$$539 \quad (D.5) \quad \mathbb{E}[\|Z_{n+1} - \nabla f_{U_{n+1}}(x^*)\|^{2p} | \mathcal{F}_n] \leq 2^{2p-1} (A_{p,n} + \|\Sigma_n\|^{2p}) \quad a.s.$$

540 However, we obtain from Holder's inequality that  $\|\Sigma_n\|^{2p} \leq A_{p,n}$ . Consequently,  
541 inequality (D.5) immediately leads to

$$542 \quad (D.6) \quad \mathbb{E}[\|Z_{n+1} - \nabla f_{U_{n+1}}(x^*)\|^{2p} | \mathcal{F}_n] \leq 4^p A_{p,n} \quad a.s.$$

543 Furthermore, define for all  $x \in \mathbb{R}^d$ ,

$$544 \quad \tau^{2p}(x) = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x) - \nabla f_k(x^*)\|^{2p}.$$

545 As in the proof of Theorem 4.1, one can observe that

$$546 \quad (D.7) \quad \mathbb{E}[\|Y_{n+1} - \nabla f_{U_{n+1}}(x^*)\|^{2p} | \mathcal{F}_n] = \tau^{2p}(X_n) \quad a.s.$$

547 Putting together the three contributions (D.4), (D.6) and (D.7), we obtain from (D.3)  
548 that

$$549 \quad (D.8) \quad \mathbb{E}[\|Y_{n+1} - \lambda Z_{n+1}\|^{2p} | \mathcal{F}_n] \leq 3^{2p-1} (\tau^{2p}(X_n) + 4^p A_{p,n} + \theta_p^*) \quad a.s.$$

550 where

$$551 \quad \theta_p^* = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x^*)\|^{2p}.$$

552 Hence, Assumption 4.7 implies that

$$553 \quad (D.9) \quad \mathbb{E}[\|Y_{n+1} - \lambda Z_{n+1}\|^{2p} | \mathcal{F}_n] \leq 3^{2p-1} (L_p V_n^p + 4^p A_{p,n} + \theta_p^*) \quad a.s.$$

554 Therefore, we deduce from (D.2) and (D.9) that for all  $n \geq 1$ ,

$$555 \quad \mathbb{E}[V_{n+1}^p | \mathcal{F}_n] \leq \left(1 + 2L_p^{1/p} C_p \gamma_n^2 + 3^{2p-1} L_p C_p \gamma_n^{2p}\right) V_n^p \\ 556 \quad \quad \quad - 2p\gamma_n V_n^{p-1} \langle X_n - x^*, \nabla f(X_n) \rangle + 2C_p \gamma_n^2 V_n^{p-1} (A_n + \theta^*) \\ 557 \quad \quad \quad + 3^{2p-1} C_p \gamma_n^{2p} (4^p A_{p,n} + \theta_p^*) \quad a.s.$$

558 Furthermore, it follows from Assumption 4.5 that  $\langle X_n - x^*, \nabla f(X_n) \rangle \geq \mu V_n$ , which  
559 leads to

$$560 \quad \mathbb{E}[V_{n+1}^p | \mathcal{F}_n] \leq \left(1 + 2L_p^{1/p} C_p \gamma_n^2 + 3^{2p-1} L_p C_p \gamma_n^{2p} - 2\mu p \gamma_n\right) V_n^p \\ 561 \quad (D.10) \quad \quad \quad + 2C_p \gamma_n^2 V_n^{p-1} (A_n + \theta^*) + 3^{2p-1} C_p \gamma_n^{2p} (4^p A_{p,n} + \theta_p^*) \quad a.s.$$

562 By taking the expectation on both side of this inequality, we obtain that for all  $n \geq 1$ ,

$$563 \quad \mathbb{E}[V_{n+1}^p] \leq \left(1 + 2L_p^{1/p} C_p \gamma_n^2 + 3^{2p-1} L_p C_p \gamma_n^{2p} - 2\mu p \gamma_n\right) \mathbb{E}[V_n^p] + 2C_p \gamma_n^2 \mathbb{E}[A_n V_n^{p-1}] \\ 564 \quad (D.11) \quad \quad \quad + 2C_p \gamma_n^2 \theta^* \mathbb{E}[V_n^{p-1}] + 3^{2p-1} C_p \gamma_n^{2p} (4^p \mathbb{E}[A_{p,n}] + \theta_p^*)$$

565 We deduce from Corollary F.2 in Appendix F below that there exists a positive constant  
566  $d_p$  such that for all  $n \geq 1$ ,  $\mathbb{E}[A_{p,n}] \leq d_p$ . The main difficulty arising here is

567 to find a sharp upper bound for the crossing term  $\mathbb{E}[A_n V_n^{p-1}]$ . By using once again  
568 Holder's inequality, we have for all  $n \geq 1$ ,

$$569 \quad (\text{D.12}) \quad \mathbb{E}[A_n V_n^{p-1}] \leq \left(\mathbb{E}[A_n^p]\right)^{\frac{1}{p}} \left(\mathbb{E}[V_n^p]\right)^{\frac{p-1}{p}} \quad \text{and} \quad \mathbb{E}[V_n^{p-1}] \leq \left(\mathbb{E}[V_n^p]\right)^{\frac{p-1}{p}}.$$

570 Hence, it is necessary to compute an upper bound for  $\mathbb{E}[A_n^p]$ . Nevertheless, one can  
571 observe from Jensen's inequality that we always have  $A_n^p \leq A_{p,n}$ , which leads that  
572  $\mathbb{E}[A_n^p] \leq d_p$ . Therefore, it follows from the induction hypothesis (D.1) together with  
573 (D.11) and (D.12) that

$$574 \quad (\text{D.13}) \quad \mathbb{E}[V_{n+1}^p] \leq \left(1 + \xi_n - 2\mu p \gamma_n\right) \mathbb{E}[V_n^p] + \frac{b}{(n+1)^{\alpha(p+1)}}$$

where  $\xi_n = 2L_p^{1/p} C_p \gamma_n^2 + 3^{2p-1} L_p C_p \gamma_n^{2p}$  and

$$b = 2^{\alpha(p+1)} c^2 C_p \left(2K_p^{(p-1)/p} (d_p^{1/p} + \theta^*) + c^{p-1} 3^{2p-1} (4^p d_p + \theta_p^*)\right).$$

Hereafter, denote by  $m$  the integer part of the real number

$$\left(\frac{cC_p}{\mu p}\right)^{1/\alpha} \left(2L_p^{1/p} + 3^{2p-1} L_p\right)^{1/\alpha}.$$

575 One can easily check that as soon as  $n \geq m$ ,  $\xi_n \leq \mu p \gamma_n$ . Consequently, we find from  
576 (D.13) that as soon as  $n \geq m$ ,

$$577 \quad (\text{D.14}) \quad \mathbb{E}[V_{n+1}^p] \leq \left(1 - \frac{a}{(n+1)^\alpha}\right) \mathbb{E}[V_n^p] + \frac{b}{(n+1)^{\alpha(p+1)}}$$

578 where  $a = p\mu c$ . Finally, we deduce from Lemma A.2 that there exists a positive  
579 constant  $K_p$  such that for all  $n \geq 1$ ,

$$580 \quad \mathbb{E}[\|X_{n+1} - x^*\|^{2p}] \leq \frac{K_p}{(n+1)^{\alpha p}},$$

581 which achieves the induction on  $n$  and completes the proof of Theorem 4.8.  $\square$

### 582 Appendix E. Additional asymptotic result on the convergence in $\mathbf{L}^2$ .

583 The goal of this appendix is to provide additional asymptotic properties of the  $\lambda$ -  
584 SAGA algorithm that will be useful in the proofs of our main results. First of all, we  
585 recall that  $V_n = \|X_n - x^*\|^2$ ,

$$586 \quad A_n = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(\phi_{n,k}) - \nabla f_k(x^*)\|^2 \quad \text{and} \quad \tau^2(x) = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x) - \nabla f_k(x^*)\|^2.$$

587 **THEOREM E.1.** *Consider a fixed  $\lambda \in [0, 1]$ . Let  $(X_n)$  be the sequence generated by  
588 the  $\lambda$ -SAGA algorithm with decreasing step sequence  $(\gamma_n)$  satisfying (1.2). Suppose  
589 that Assumptions 3.1, 3.3 and 4.5 are satisfied. Then, we have almost surely*

$$590 \quad (\text{E.1}) \quad \sum_{n=1}^{\infty} \gamma_n V_n < +\infty, \quad \sum_{n=1}^{\infty} \gamma_n A_n < +\infty, \quad \sum_{n=1}^{\infty} \gamma_n \tau^2(X_n) < +\infty.$$

591 *In addition, we also have*

$$592 \quad (\text{E.2}) \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[V_n] < +\infty, \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[A_n] < +\infty, \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[\tau^2(X_n)] < +\infty.$$

593 *Proof.* Let us consider the same Lyapounov function used in the proof of Theorem  
594 4.1 and given by (4.11). We recall from inequality (4.14) that for all  $n \geq 1$ ,

$$595 \quad (\text{E.3}) \quad \mathbb{E}[T_{n+1}|\mathcal{F}_n] \leq (1 + 4L\gamma_n^2)T_n - 2\gamma_n \langle X_n - x^*, \nabla f(X_n) \rangle + 2\theta^* \gamma_n^2 \quad a.s.$$

596 Let  $(\mathcal{T}_n)$  be the sequence of Lyapunov functions defined, for all  $n \geq 2$ , by  $\mathcal{T}_n = b_n T_n$   
597 where

$$598 \quad b_n = \prod_{k=1}^{n-1} (1 + 4L\gamma_k^2)^{-1}.$$

599 Since  $b_n = b_{n+1}(1 + 4L\gamma_n^2)$ , we obtain from (E.3) that for all  $n \geq 1$ ,

$$600 \quad \mathbb{E}[\mathcal{T}_{n+1}|\mathcal{F}_n] \leq \mathcal{T}_n - 2\gamma_n b_{n+1} \langle X_n - x^*, \nabla f(X_n) \rangle + 2\theta^* b_{n+1} \gamma_n^2 \quad a.s.$$

601 Hence, it follows from Assumption 4.5 that for all  $n \geq 1$ ,

$$602 \quad (\text{E.4}) \quad \mathbb{E}[\mathcal{T}_{n+1}|\mathcal{F}_n] \leq \mathcal{T}_n - 2\mu\gamma_n b_{n+1} V_n + 2\theta^* b_{n+1} \gamma_n^2 \quad a.s.$$

603 Moreover, we clearly have from the right-hand side of (1.2) that  $(b_n)$  converges to a  
604 positive real number  $b$ , which implies that

$$605 \quad (\text{E.5}) \quad \sum_{n=1}^{\infty} b_{n+1} \gamma_n^2 < +\infty.$$

606 Therefore, we deduce from the Robbins-Siegmund Theorem [34] given by Theorem  
607 A.1 that  $(\mathcal{T}_n)$  converges almost surely towards a finite random variable  $\mathcal{T}$  and the  
608 series

$$609 \quad \sum_{n=2}^{\infty} \gamma_n b_{n+1} V_n < +\infty \quad a.s.$$

610 which leads to

$$611 \quad (\text{E.6}) \quad \sum_{n=1}^{\infty} \gamma_n V_n < +\infty \quad a.s.$$

612 We also obtain from relation (E.6) and Assumption 3.3 that

$$613 \quad (\text{E.7}) \quad \sum_{n=1}^{\infty} \gamma_n \tau^2(X_n) < +\infty \quad a.s.$$

614 In addition, by taking the expectation of both sides of (E.4) and using a standard  
615 telescoping argument, we find that

$$616 \quad (\text{E.8}) \quad 2\mu \sum_{n=2}^{\infty} \gamma_n b_{n+1} \mathbb{E}[V_n] \leq \mathbb{E}[\mathcal{T}_2] + 2\theta^* \sum_{n=2}^{\infty} b_{n+1} \gamma_n^2.$$

617 Then, it follows from (E.5) and (E.8) that

$$618 \quad \sum_{n=2}^{\infty} \gamma_n b_{n+1} \mathbb{E}[V_n] < +\infty,$$

619 which implies that

$$620 \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[V_n] < +\infty.$$

621 Consequently, we get from Assumption 3.3 that

$$622 \quad (\text{E.9}) \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[\tau^2(X_n)] < +\infty.$$

623 Furthermore, we already saw from (4.12) that for all  $n \geq 1$ ,

$$624 \quad (\text{E.10}) \quad \mathbb{E}[A_{n+1} | \mathcal{F}_n] = \frac{1}{N} \tau^2(X_n) + \left(1 - \frac{1}{N}\right) A_n \quad a.s.$$

625 For all  $n \geq 1$ , denote  $\mathcal{A}_n = \gamma_n A_n$ . Since  $\gamma_{n+1} \leq \gamma_n$ , we obtain from (E.10) that for  
626 all  $n \geq 1$ ,

$$627 \quad (\text{E.11}) \quad \mathbb{E}[\mathcal{A}_{n+1} | \mathcal{F}_n] \leq \mathcal{A}_n + \frac{1}{N} \gamma_n \tau^2(X_n) - \frac{1}{N} \gamma_n \mathcal{A}_n \quad a.s.$$

628 By considering the almost sure convergence (E.7), it follows once again from the  
629 Robbins-Siegmund Theorem [34] given by Theorem A.1 that  $(\mathcal{A}_n)$  converges almost  
630 surely towards a finite random variable  $\mathcal{A}$  and the series

$$631 \quad \sum_{n=1}^{\infty} \gamma_n \mathcal{A}_n < +\infty \quad a.s.$$

632 Moreover, by taking the expectation of both sides of (E.11) and using a standard  
633 telescoping argument, we obtain that

$$634 \quad (\text{E.12}) \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[A_n] \leq N \mathbb{E}[A_1] + \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[\tau^2(X_n)].$$

635 Finally, we deduce from (E.9) and (E.12) that

$$636 \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[A_n] < +\infty,$$

637 which completes the proof of Theorem E.1. □

638 A straightforward application of Theorem E.1, using the left-hand side of (1.2), is as  
639 follows.

640 **COROLLARY E.2.** *Assume that the conditions of Theorem E.1 hold. Then, we*  
641 *have*

$$642 \quad \lim_{n \rightarrow +\infty} V_n = \lim_{n \rightarrow +\infty} A_n = \lim_{n \rightarrow +\infty} \tau^2(X_n) = 0 \quad a.s.$$

643 *Moreover, we also have*

$$644 \quad \lim_{n \rightarrow +\infty} \mathbb{E}[V_n] = \lim_{n \rightarrow +\infty} \mathbb{E}[A_n] = \lim_{n \rightarrow +\infty} \mathbb{E}[\tau^2(X_n)] = 0.$$

645 **Appendix F. Additional asymptotic result on the convergence in  $L^p$ .**

646 As in the previous Appendix, our purpose is to establish additional asymptotic  
 647 properties of the  $\lambda$ -SAGA algorithm that will be useful in the proofs of our main  
 648 results. First of all, we recall that  $V_n^p = \|X_n - x^*\|^{2p}$ ,

$$649 A_{p,n} = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(\phi_{n,k}) - \nabla f_k(x^*)\|^{2p} \quad \text{and} \quad \tau^{2p}(x) = \frac{1}{N} \sum_{k=1}^N \|\nabla f_k(x) - \nabla f_k(x^*)\|^{2p}.$$

650

651 **THEOREM F.1.** *Consider a fixed  $\lambda \in [0, 1]$ . Let  $(X_n)$  be the sequence generated  
 652 by the  $\lambda$ -SAGA algorithm with decreasing step  $\gamma_n$  satisfying (1.2). Suppose that As-  
 653 sumptions 3.1, 4.5 and 4.7 are satisfied. Then, we have almost surely*

$$654 \text{(F.1)} \quad \sum_{n=1}^{\infty} \gamma_n V_n^p < +\infty, \quad \sum_{n=1}^{\infty} \gamma_n A_{p,n} < +\infty, \quad \sum_{n=1}^{\infty} \gamma_n \tau^{2p}(X_n) < +\infty.$$

655 *In addition, we also have*

$$656 \text{(F.2)} \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[V_n^p] < +\infty, \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[A_{p,n}] < +\infty, \quad \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[\tau^{2p}(X_n)] < +\infty.$$

657 *Proof.* We are going to prove Theorem F.1 by induction on  $p \geq 1$ . First of all,  
 658 Theorem F.1 follows from Theorem E.1 in the special case  $p = 1$ . Hence, the base  
 659 case is immediately true. Next, assume by induction that Theorem F.1 holds for some  
 660 integer  $p - 1$  with  $p \geq 2$ . We recall from inequality (D.10) in the proof of Theorem  
 661 4.8 that for all  $n \geq 1$ ,

$$662 \mathbb{E}[V_{n+1}^p | \mathcal{F}_n] \leq \left(1 + 2L_p^{1/p} C_p \gamma_n^2 + 3^{2p-1} L_p C_p \gamma_n^{2p} - 2\mu p \gamma_n\right) V_n^p \\
 663 \text{(F.3)} \quad + 2C_p \gamma_n^2 V_n^{p-1} (A_n + \theta^*) + 3^{2p-1} C_p \gamma_n^{2p} (4^p A_{p,n} + \theta_p^*) \quad a.s.$$

664 However, it follows from Young's inequality for products that almost surely

$$665 \text{(F.4)} \quad V_n^{p-1} A_n \leq \frac{A_n^p}{p} + \frac{(p-1)V_n^p}{p}.$$

666 Moreover, one can observe from Jensen's inequality that  $A_n^p \leq A_{p,n}$  almost surely.  
 667 Combining the previous inequality with (F.4), it implies that

$$668 \text{(F.5)} \quad V_n^{p-1} A_n \leq \frac{A_{p,n}}{p} + \frac{(p-1)V_n^p}{p}.$$

669 Furthermore, by putting together the inequalities (F.3) and (F.5), we obtain that

$$670 \mathbb{E}[V_{n+1}^p | \mathcal{F}_n] \leq \left(1 + 2(L_p^{1/p} + p^{-1}(p-1))C_p \gamma_n^2 + 3^{2p-1} L_p C_p \gamma_n^{2p}\right) V_n^p - 2\mu p \gamma_n V_n^p \\
 671 \text{(F.6)} \quad + 2C_p \theta^* \gamma_n^2 V_n^{p-1} + \left(3^{2p-1} 4^p + 2p^{-1}\right) C_p \gamma_n^2 A_{p,n} + 3^{2p-1} C_p \gamma_n^{2p} \theta_p^* \quad a.s.$$

672 Let  $(T_{p,n})$  be the sequence of Lyapunov functions defined, for all  $n \geq 2$ , by

$$673 \text{(F.7)} \quad T_{p,n} = V_n^p + N e_p \gamma_{n-1}^2 A_{p,n},$$

674 where  $e_p = C_p(3^{2p-1}4^p + 2p^{-1})$ . By the definition (F.7), we have

$$675 \quad (\text{F.8}) \quad \mathbb{E}[T_{p,n+1}|\mathcal{F}_n] = \mathbb{E}[V_{n+1}^p|\mathcal{F}_n] + Ne_p\gamma_n^2\mathbb{E}[A_{p,n+1}|\mathcal{F}_n] \quad a.s.$$

676 However, we deduce by the same arguments as in (4.12) that

$$677 \quad (\text{F.9}) \quad \mathbb{E}[A_{p,n+1}|\mathcal{F}_n] = \frac{1}{N}\tau^{2p}(X_n) + \left(1 - \frac{1}{N}\right)A_{p,n} \quad a.s.$$

678 Hence, it follows from (F.6), (F.8) and (F.9) that

$$679 \quad \mathbb{E}[T_{p,n+1}|\mathcal{F}_n] \leq T_{p,n} + C_p \left(2(L_p^{1/p} + p^{-1}(p-1))\gamma_n^2 + 3^{2p-1}L_p\gamma_n^{2p}\right)V_n^p - 2\mu p\gamma_n V_n^p \\ 680 \quad (\text{F.10}) \quad + e_p\gamma_n^2\tau^{2p}(X_n) + 2C_p\theta^*\gamma_n^2V_n^{p-1} + 3^{2p-1}C_p\gamma_n^{2p}\theta_p^* \quad a.s.$$

681 Additionally, we clearly have  $V_n^p \leq T_{p,n}$  and Assumption 4.7 leads to

$$682 \quad \tau^{2p}(X_n) \leq L_p V_n^p \leq L_p T_{p,n}.$$

683 Finally, we obtain from (F.10) that

$$684 \quad (\text{F.11}) \quad \mathbb{E}[T_{p,n+1}|\mathcal{F}_n] \leq (1 + a_n)T_{p,n} + \Delta_n - 2\mu p\gamma_n V_n^p \quad a.s.,$$

685 where  $\Delta_n = 2C_p\theta^*\gamma_n^2V_n^{p-1} + 3^{2p-1}C_p\gamma_n^{2p}\theta_p^*$  and

$$686 \quad a_n = e_p L_p \gamma_n^2 + 2C_p(L_p^{1/p} + p^{-1}(p-1))\gamma_n^2 + 3^{2p-1}L_p C_p \gamma_n^{2p}.$$

687 Since the sequence  $(\gamma_n)$  satisfies (1.2), it is easy to see that

$$688 \quad \sum_{n=1}^{\infty} a_n < +\infty.$$

689 Moreover, by the induction hypothesis, we have that

$$690 \quad (\text{F.12}) \quad \left\{ \begin{array}{l} \sum_{n=1}^{\infty} \gamma_n V_n^{p-1} < +\infty \quad a.s., \\ \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[V_n^{p-1}] < +\infty. \end{array} \right.$$

691 From (F.12), one can immediately deduce that

$$692 \quad (\text{F.13}) \quad \left\{ \begin{array}{l} \sum_{n=1}^{\infty} \Delta_n < +\infty \quad a.s., \\ \sum_{n=1}^{\infty} \mathbb{E}[\Delta_n] < +\infty. \end{array} \right.$$

693 Therefore, one uses exactly the same lines as in the proof of Theorem E.1 and the  
694 Robbins-Siegmund Theorem [34] given by Theorem A.1 to show that

$$695 \quad (\text{F.14}) \quad \left\{ \begin{array}{l} \sum_{n=1}^{\infty} \gamma_n V_n^p < +\infty \quad a.s., \\ \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[V_n^p] < +\infty. \end{array} \right.$$

696 Hence, combining (F.14) with Assumption 4.7, one immediately deduces that

$$697 \quad (F.15) \quad \left\{ \begin{array}{l} \sum_{n=1}^{\infty} \gamma_n \tau^{2p}(X_n) < +\infty \quad a.s., \\ \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[\tau^{2p}(X_n)] < +\infty. \end{array} \right.$$

698 Finally, using once again the same arguments as in the proof of Theorem E.1, we  
699 obtain that

$$700 \quad (F.16) \quad \left\{ \begin{array}{l} \sum_{n=1}^{\infty} \gamma_n A_{p,n} < +\infty \quad a.s., \\ \sum_{n=1}^{\infty} \gamma_n \mathbb{E}[A_{p,n}] < +\infty. \end{array} \right. ,$$

701 which achieves the proof of Theorem F.1.  $\square$

702 A useful consequence of Theorem F.1, using the left-hand side of (1.2), is as follows.

703 **COROLLARY F.2.** *Assume that the conditions of Theorem F.1 hold. Then, for all*  
704  *$p \geq 1$ , we have*

$$705 \quad \lim_{n \rightarrow +\infty} V_n^p = \lim_{n \rightarrow +\infty} A_{p,n} = \lim_{n \rightarrow +\infty} \tau^{2p}(X_n) = 0 \quad a.s.,$$

706 *and*

$$707 \quad \lim_{n \rightarrow +\infty} \mathbb{E}[V_n^p] = \lim_{n \rightarrow +\infty} \mathbb{E}[A_{p,n}] = \lim_{n \rightarrow +\infty} \mathbb{E}[\tau^{2p}(X_n)] = 0.$$

708

#### REFERENCES

- 709 [1] Z. ALLEN-ZHU, Katyusha: The first direct acceleration of stochastic gradient methods, Journal  
710 of Machine Learning Research, 18 (2018), pp. 1–51.
- 711 [2] P. ASSOUD, Espaces  $p$ -lisses et  $q$ -convexes. Inégalités de Burkholder, Séminaire Maurey-  
712 Schwartz, (1975), pp. 1–7. talk:15.
- 713 [3] F. BACH, Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic  
714 regression, The Journal of Machine Learning Research, 15 (2014), pp. 595–627.
- 715 [4] F. BACH AND E. MOULINES, Non-asymptotic analysis of stochastic approximation algorithms  
716 for machine learning, Advances in neural information processing systems, 24 (2011).
- 717 [5] B. BERCU AND J. BIGOT, Asymptotic distribution and convergence rates of stochastic  
718 algorithms for entropic optimal transportation between probability measures, The Annals  
719 of Statistics, 49 (2021), pp. 968–987.
- 720 [6] B. BERCU, A. GODICHON, AND B. PORTIER, An efficient stochastic Newton algorithm for  
721 parameter estimation in logistic regressions, SIAM J. Control Optim., 58 (2020), pp. 348–  
722 367.
- 723 [7] D. P. BERTSEKAS AND J. N. TSITSIKLIS, Gradient convergence in gradient methods with errors,  
724 SIAM Journal on Optimization, 10 (2000), pp. 627–642.
- 725 [8] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, Optimization methods for large-scale machine  
726 learning, SIAM review, 60 (2018), pp. 223–311.
- 727 [9] N. CHATTERJI, N. FLAMMARION, Y. MA, P. BARTLETT, AND M. JORDAN, On the theory of  
728 variance reduction for stochastic gradient monte carlo, in International Conference on Ma-  
729 chine Learning, PMLR, 2018, pp. 764–773.
- 730 [10] X. CHEN, Z. LAI, H. LI, AND Y. ZHANG, Online statistical inference for stochastic optimization  
731 via kiefer-wolfowitz methods, Journal of the American Statistical Association, (2024),  
732 pp. 1–24.
- 733 [11] A. DEFAZIO, A simple practical accelerated method for finite sums, Advances in neural infor-  
734 mation processing systems, 29 (2016).

- 735 [12] A. DEFAZIO, F. BACH, AND S. LACOSTE-JULIEN, Saga: A fast incremental gradient method  
736 with support for non-strongly convex composite objectives, Advances in neural information  
737 processing systems, 27 (2014).
- 738 [13] M. DUFLO, Algorithmes stochastiques, Mathématiques et Applications, Springer Berlin Heidel-  
739 berg, 1996.
- 740 [14] V. FABIAN, On asymptotic normality in stochastic approximation, The Annals of Mathematical  
741 Statistics, (1968), pp. 1327–1332.
- 742 [15] O. FERCOQ AND P. RICHTÁRIK, Optimization in high dimensions via accelerated, parallel, and  
743 proximal coordinate descent, Siam review, 58 (2016), pp. 739–771.
- 744 [16] R. GOWER, N. LE ROUX, AND F. BACH, Tracking the gradients using the hessian: A new look at  
745 variance reducing stochastic methods, in International Conference on Artificial Intelligence  
746 and Statistics, PMLR, 2018, pp. 707–715.
- 747 [17] R. M. GOWER, M. SCHMIDT, F. BACH, AND P. RICHTÁRIK, Variance-reduced methods for  
748 machine learning, Proceedings of the IEEE, 108 (2020), pp. 1968–1983.
- 749 [18] C. GUILLE-ESCURET, A. IBRAHIM, B. GOUJAUD, AND I. MITLIAGKAS, Gradient descent is  
750 optimal under lower restricted secant inequality and upper error bound, Advances in Neural  
751 Information Processing Systems, 35 (2022), pp. 24893–24904.
- 752 [19] H. KARIMI, J. NUTINI, AND M. SCHMIDT, Linear convergence of gradient and proximal-gradient  
753 methods under the polyak-lojasiewicz condition, in Joint European conference on machine  
754 learning and knowledge discovery in databases, Springer, 2016, pp. 795–811.
- 755 [20] J. KIEFER AND J. WOLFOWITZ, Stochastic estimation of the maximum of a regression function,  
756 The Annals of Mathematical Statistics, (1952), pp. 462–466.
- 757 [21] H. J. KUSHNER AND H. HUANG, Rates of convergence for stochastic approximation type  
758 algorithms, SIAM Journal on Control and Optimization, 17 (1979), pp. 607–617.
- 759 [22] H. J. KUSHNER AND G. YIN, Stochastic Approximation and Recursive Algorithms and  
760 Applications, Springer New York, NY, 2003.
- 761 [23] C. K. LAUAND, I. KONTOYIANNIS, AND S. MEYN, The case for and against fixed step-size:  
762 Stochastic approximation algorithms in optimization and machine learning, 2025, <https://arxiv.org/abs/2309.02944>, <https://arxiv.org/abs/2309.02944>.
- 763 [24] R. LELUC AND F. PORTIER, Sgd with coordinate sampling: Theory and practice, Journal of  
764 Machine Learning Research, 23 (2022), pp. 1–47.
- 765 [25] J. LIU AND Y. YUAN, On almost sure convergence rates of stochastic gradient methods, in  
766 Conference on Learning Theory, PMLR, 2022, pp. 2963–2983.
- 767 [26] L. NGUYEN, P. H. NGUYEN, M. DIJK, P. RICHTÁRIK, K. SCHEINBERG, AND M. TAKÁC, Sgd  
768 and hogwild! convergence without the bounded gradients assumption, in International  
769 Conference on Machine Learning, PMLR, 2018, pp. 3750–3758.
- 770 [27] B. PALANIAPPAN AND F. BACH, Stochastic variance reduction methods for saddle-point  
771 problems, Advances in Neural Information Processing Systems, 29 (2016).
- 772 [28] M. PELLETTIER, On the almost sure asymptotic behaviour of stochastic algorithms, Stochastic  
773 processes and their applications, 78 (1998), pp. 217–244.
- 774 [29] M. PELLETTIER, Weak convergence rates for stochastic approximation with application to  
775 multiple targets and simulated annealing, Annals of Applied Probability, (1998), pp. 10–44.
- 776 [30] B. T. POLYAK AND A. B. JUDITSKY, Acceleration of stochastic approximation by averaging,  
777 SIAM journal on control and optimization, 30 (1992), pp. 838–855.
- 778 [31] C. POON, J. LIANG, AND C. SCHOENLIEB, Local convergence properties of saga/prox-svrg and  
779 acceleration, in International Conference on Machine Learning, PMLR, 2018, pp. 4124–  
780 4132.
- 781 [32] X. QIAN, Z. QU, AND P. RICHTÁRIK, Saga with arbitrary sampling, in International Conference  
782 on Machine Learning, PMLR, 2019, pp. 5190–5199.
- 783 [33] H. ROBBINS AND S. MONRO, A stochastic approximation method, The annals of mathematical  
784 statistics, (1951), pp. 400–407.
- 785 [34] H. ROBBINS AND D. SIEGMUND, A convergence theorem for non negative almost  
786 supermartingales and some applications, in Optimizing methods in statistics, Elsevier,  
787 1971, pp. 233–257.
- 788 [35] N. ROUX, M. SCHMIDT, AND F. BACH, A stochastic gradient method with an exponential  
789 convergence rate for finite training sets, Advances in neural information processing sys-  
790 tems, 25 (2012).
- 791 [36] J. SACKS, Asymptotic distribution of stochastic approximation procedures, The Annals of  
792 Mathematical Statistics, 29 (1958), pp. 373–405.
- 793 [37] M. SCHMIDT, N. LE ROUX, AND F. BACH, Minimizing finite sums with the stochastic average  
794 gradient, Mathematical Programming, 162 (2017), pp. 83–112.
- 795 [38] L. XIAO AND T. ZHANG, A proximal stochastic gradient method with progressive variance  
796

- 797        reduction, SIAM Journal on Optimization, 24 (2014), pp. 2057–2075.
- 798 [39] X. YI, S. ZHANG, T. YANG, T. CHAI, AND K. H. JOHANSSON, Exponential convergence  
799 for distributed optimization under the restricted secant inequality condition, IFAC-  
800 PapersOnLine, 53 (2020), pp. 2672–2677.
- 801 [40] L. ZHANG AND L. GUO, Asymptotically efficient adaptive identification under saturated output  
802 observations, SIAM Journal on Control and Optimization, 63 (2025), pp. 2338–2368.
- 803 [41] L.-X. ZHANG, Central limit theorems of a recursive stochastic algorithm with applications to  
804 adaptive designs, The Annals of Applied Probability, 26 (2016), pp. 3630–3658.