

36 roots can be found in the Stochastic Approximation (SA) method by Robbins and
 37 Monro [31], which can be interpreted as a stochastic variant of the gradient descent
 38 method. Convergence in probability of the SA method is ensured if the step-length
 39 sequence $\{\alpha_k\}$ (called gain sequence) is of harmonic type, i.e., it is non-summable
 40 but square-summable. Under suitable assumptions, the method converges in the
 41 mean square sense [31] and almost surely [37]. Its key property is the ability to avoid
 42 zigzagging due to noise when approaching the solution, thanks to the decay of the gain
 43 sequence. However, a significant drawback of the SA method is its slow convergence.

44 Asymptotically ideal step lengths for SA methods include the norm of the inverse
 45 Hessian at the solution [3]. Adjustments to the classical harmonic gain sequence,
 46 including adaptive step lengths based on changes in the sign of the difference between
 47 consecutive iterates, are analyzed in [13, 21] with the aim of speeding up the SA
 48 method. This idea is further developed in [39], where almost sure convergence of the
 49 method is proved. Adaptive step-length schemes are introduced also in [40], with
 50 the objective of reducing the dependence of the behavior of the method on user-
 51 defined parameters. The results in [5] are closely related to the SA method and
 52 concern methods based on search directions that are not necessarily noisy gradients,
 53 but some gradient-related directions. A hybrid approach that combines a line-search
 54 technique with SA is analyzed in [25] for noisy gradient directions and arbitrary
 55 descent directions. General descent directions are also considered in [24]. We also note
 56 that gradient approximations may need to be computed by using finite differences;
 57 an overview of finite-difference methods for stochastic optimization is given in [16].
 58 Variance-reduction SA methods with line search for stochastic variational inequalities
 59 are considered in [19].

60 In the realm of machine learning, many stochastic versions of the gradient method
 61 have been developed. Starting from the basic stochastic and minibatch gradient
 62 methods – see, e.g., [7] and the references therein – variance reduction techniques
 63 for the gradient estimates have been developed, with the aim of improving conver-
 64 gence. Among them we mention SVRG [20], SAGA [12] and its version using Jacobian
 65 sketching [18], which will be considered in section 5. These methods have constant
 66 step lengths and get linear convergence in expectation.

67 Stochastic optimization methods exploiting search directions based on second-
 68 order information have been developed to get better theoretical and practical conver-
 69 gence properties, especially when badly-scaled problems are considered. Stochastic
 70 versions of Newton-type methods are discussed in [2, 6, 8, 9, 33, 34, 35, 36] and a
 71 variant of the adaptive cubic regularization scheme using a dynamic rule for build-
 72 ing inexact Hessian information is proposed in [1]. Stochastic BFGS methods are
 73 analyzed, e.g., in [8, 10, 17, 28, 29, 30]. In particular, in [30] Moritz et al. propose
 74 a stochastic L-BFGS algorithm based on the same inverse Hessian approximation as
 75 in [10], but use SVRG instead of the standard stochastic gradient approximation. This
 76 algorithm, which applies a constant step length, has Q-linear rate of convergence of
 77 the expected value of the error in the objective function. A further modification to this
 78 L-BFGS scheme is proposed by Gower et al. in [17], where a stochastic block BFGS
 79 update is used, in which the vector pairs for updating the inverse Hessian are replaced
 80 by matrix pairs gathering directions and matrix-vector products between subsampled
 81 Hessians and those directions. The resulting algorithm uses constant step length and
 82 has Q-linear convergence rate of the expected value of objective function error, as in
 83 the previous case, but appears more efficient by numerical experiments.

84 **Our contribution.** We propose a Line-search Second-Order Stochastic (LSOS)
 85 algorithmic framework for stochastic optimization problems, where Newton and quasi-
 86 Newton directions in a rather broad meaning are used. Inexactness is allowed in the
 87 sense that the (approximate) Newton direction can be obtained as inexact solution of
 88 the corresponding system of linear equations. We focus on convex problems as they
 89 appear in a wide variety of applications, such as machine learning and least squares.
 90 Furthermore, many stochastic problems need regularization and hence become convex.

91 We prove almost sure convergence of the methods fitting into the LSOS framework
 92 and show by experiments the effectiveness of our approach when using Newton and
 93 inexact Newton directions affected by noise.

94 For finite-sum objective functions such as those arising in machine learning, we
 95 investigate the use of the stochastic L-BFGS Hessian approximations in [10] together
 96 with line searches and the SAGA variance reduction technique for the gradient esti-
 97 mates. The resulting algorithm has almost sure convergence to the solution, while
 98 for the efficient state-of-the-art stochastic L-BFGS methods in [17, 30] it has been
 99 proved only that the function error tends to zero in expectation. We also prove
 100 that the expected function error has linear convergence rate and provide a worst-
 101 case $\mathcal{O}(\log(\varepsilon^{-1}))$ complexity bound. Finally, numerical experiments show that our
 102 algorithm is competitive with the stochastic L-BFGS methods mentioned above.

103 **Notation.** $\mathbb{E}(x)$ denotes the expectation of a random variable x , $\mathbb{E}(x|y)$ the
 104 conditional expectation of x given y , and $\text{var}(x)$ the variance of x . $\|\cdot\|$ indicates
 105 either the Euclidean vector norm or the corresponding induced matrix norm, while
 106 $|\cdot|$ is the cardinality of a set. \mathbb{R}_+ and \mathbb{R}_{++} denote the sets of real non-negative
 107 and positive numbers, respectively. Vectors are written in boldface and subscripts
 108 indicate the elements of a sequence, e.g., $\{\mathbf{x}_k\}$. Throughout the paper M_1, M_2, M_3, \dots
 109 and c_1, c_2, c_3, \dots denote positive constants, without specifying their actual values.
 110 Other constants are defined when they are used. Finally, “a.s.” abbreviates “almost
 111 sure/surely”.

112 **Outline of the paper.** The rest of this article is organized as follows. In [sec-](#)
 113 [tion 2](#), we define the general Stochastic Second-Order (SOS) framework with pre-
 114 defined step-length sequence, which is the basis for the family of algorithms proposed
 115 in this work, and we give preliminary assumptions and results used in the sequel. In
 116 [section 3](#) we provide the convergence theory of the algorithms fitting into the SOS
 117 framework. In [section 4](#) we introduce a SOS version named LSOS, which combines
 118 non-monotone line searches and (if needed) pre-defined step lengths in order to make
 119 the algorithm faster, and provide its convergence analysis. In [section 5](#) we special-
 120 ize LSOS for finite sum objective functions, obtaining a stochastic L-BFGS method
 121 with line search only, and in [section 6](#) we provide its convergence theory, including
 122 convergence rate and complexity results. In [section 7](#), numerical experiments on two
 123 classes of stochastic problems and comparisons with state-of-the art methods show
 124 the effectiveness of our approach. Concluding remarks are given in [section 8](#).

125 **2. Preliminaries.** We assume that for problem (1.1) we can only compute

$$126 \quad (2.1) \quad \begin{aligned} f(\mathbf{x}) &= \phi(\mathbf{x}) + \varepsilon_f(\mathbf{x}), \\ \mathbf{g}(\mathbf{x}) &= \nabla\phi(\mathbf{x}) + \varepsilon_g(\mathbf{x}), \\ B(\mathbf{x}) &= \nabla^2\phi(\mathbf{x}) + \varepsilon_B(\mathbf{x}), \end{aligned}$$

127 with $\varepsilon_f(\mathbf{x})$ being a random number, $\varepsilon_g(\mathbf{x})$ a random vector and $\varepsilon_B(\mathbf{x})$ a symmetric
 128 random matrix. The general algorithmic scheme we analyze in this paper is given in

129 Algorithm 2.1.

Algorithm 2.1 Second-Order Stochastic (SOS) method

1: given $\mathbf{x}^0 \in \mathbb{R}^n$ and $\{\alpha_k\} \subset \mathbb{R}_+$
 2: **for** $k = 0, 1, 2, \dots$ **do**
 3: compute $\mathbf{d}_k \in \mathbb{R}^n$
 4: set $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$
 5: **end for**

130 For now we assume that $\{\alpha_k\}$ is given and it satisfies the conditions stated in
 131 [Assumption 2.2](#) below. We also assume that $f(\mathbf{x})$, $\mathbf{g}(\mathbf{x})$ and $B(\mathbf{x})$ are available for
 132 any $\mathbf{x} \in \mathbb{R}^n$. Although here we do not specify how \mathbf{d}_k is obtained, we call the algorithm
 133 “Second-Order” because in the next sections we will compute \mathbf{d}_k by exploiting noisy
 134 second-order information about $\phi(\mathbf{x})$.

135 We make the following assumptions.

136 **ASSUMPTION 2.1.** *The function ϕ is strongly convex and has Lipschitz-continuous*
 137 *gradient.*

If [Assumption 2.1](#) holds, then there exists a unique $\mathbf{x}_* \in \mathbb{R}^n$ that solves (1.1), with $\nabla\phi(\mathbf{x}_*) = \mathbf{0}$. Furthermore, for some positive constants μ and L and any $\mathbf{x} \in \mathbb{R}^n$ we have

$$\mu I \preceq \nabla^2\phi(\mathbf{x}) \preceq LI,$$

138 where I is the identity matrix, and

139 (2.2)
$$\frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_*\|^2 \leq \phi(\mathbf{x}) - \phi(\mathbf{x}_*) \leq \frac{L}{2} \|\nabla\phi(\mathbf{x})\|^2.$$

ASSUMPTION 2.2. *The gain sequence $\{\alpha_k\}$ satisfies*

$$\alpha_k > 0 \text{ for all } k, \quad \sum_k \alpha_k = \infty, \quad \sum_k \alpha_k^2 < \infty.$$

140 This is a standard assumption for SA methods.

141 Henceforth we denote \mathcal{F}_k the σ -algebra generated by $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_k$.

ASSUMPTION 2.3. *Let $\{\mathbf{x}_k\}$ be a sequence generated by [Algorithm 2.1](#). The gradient noise $\boldsymbol{\varepsilon}_g(\mathbf{x})$ is such that*

$$\mathbb{E}(\boldsymbol{\varepsilon}_g(\mathbf{x}) | \mathcal{F}_k) = \mathbf{0} \text{ and } \mathbb{E}(\|\boldsymbol{\varepsilon}_g(\mathbf{x})\|^2 | \mathcal{F}_k) \leq M_1.$$

142 In other words, we assume that the expected gradient noise is zero and the variance
 143 of gradient errors,

144 (2.3)
$$\text{var}(\|\boldsymbol{\varepsilon}_g(\mathbf{x})\| | \mathcal{F}_k) = \mathbb{E}(\|\boldsymbol{\varepsilon}_g(\mathbf{x})\|^2 | \mathcal{F}_k) - \mathbb{E}^2(\|\boldsymbol{\varepsilon}_g(\mathbf{x})\| | \mathcal{F}_k),$$

is bounded. From (2.3) and [Assumption 2.3](#) it also follows that

$$\mathbb{E}(\|\boldsymbol{\varepsilon}_g(\mathbf{x})\|^2 | \mathcal{F}_k) \geq \mathbb{E}^2(\|\boldsymbol{\varepsilon}_g(\mathbf{x})\| | \mathcal{F}_k)$$

and hence

$$\mathbb{E}(\|\boldsymbol{\varepsilon}_g(\mathbf{x})\| | \mathcal{F}_k) \leq \sqrt{\mathbb{E}(\|\boldsymbol{\varepsilon}_g(\mathbf{x})\|^2 | \mathcal{F}_k)} \leq \sqrt{M_1} := M_2.$$

145 We observe that [Assumptions 2.1](#) and [2.3](#) imply

$$146 \quad (2.4) \quad \|\nabla\phi(\mathbf{x})\|^2 + \mathbb{E}(\|\varepsilon_g(\mathbf{x})\|^2|\mathcal{F}_k) \leq L^2\|\mathbf{x} - \mathbf{x}_*\|^2 + M_1 \leq c_1(1 + \|\mathbf{x} - \mathbf{x}_*\|^2),$$

147 with $c_1 = \max\{L^2, M_1\}$. Moreover, [\(2.4\)](#) and [Assumption 2.3](#) imply

$$148 \quad (2.5) \quad \mathbb{E}(\|\mathbf{g}(\mathbf{x})\|^2|\mathcal{F}_k) \leq c_1(1 + \|\mathbf{x} - \mathbf{x}_*\|^2),$$

149 which can be proved as follows:

$$\begin{aligned} 150 \quad \mathbb{E}(\|\mathbf{g}(\mathbf{x})\|^2|\mathcal{F}_k) &= \mathbb{E}(\|\nabla\phi(\mathbf{x}) + \varepsilon_g(\mathbf{x})\|^2|\mathcal{F}_k) \\ 151 &= \mathbb{E}(\|\nabla\phi(\mathbf{x})\|^2 + 2\nabla\phi(\mathbf{x})^\top \varepsilon_g(\mathbf{x}) + \|\varepsilon_g(\mathbf{x})\|^2|\mathcal{F}_k) \\ 152 &= \|\nabla\phi(\mathbf{x})\|^2 + 2\nabla\phi(\mathbf{x})^\top \mathbb{E}(\varepsilon_g(\mathbf{x})|\mathcal{F}_k) + \mathbb{E}(\|\varepsilon_g(\mathbf{x})\|^2|\mathcal{F}_k) \\ 153 &\leq c_1(1 + \|\mathbf{x} - \mathbf{x}_*\|^2), \end{aligned}$$

154 where the last inequality comes from $\mathbb{E}(\varepsilon_g(\mathbf{x})|\mathcal{F}_k) = 0$.

155 The following theorem (see [\[31\]](#)) will be used in [Section 3](#).

THEOREM 2.4. *Let $U_k, \beta_k, \xi_k, \rho_k \geq 0$ be \mathcal{F}_k -measurable random variables such that*

$$\mathbb{E}(U_{k+1}|\mathcal{F}_k) \leq (1 + \beta_k)U_k + \xi_k - \rho_k, \quad k = 1, 2, \dots$$

156 *If $\sum_k \beta_k < \infty$ and $\sum_k \xi_k < \infty$, then $U_k \rightarrow U$ a.s. and $\sum_k \rho_k < \infty$ a.s..*

3. Convergence theory of Algorithm SOS. The assumptions stated in the previous section generally form a common set of assumptions for SA and related methods. Actually, [Assumption 2.1](#) is different from the commonly used assumption that for some symmetric positive definite matrix B and for all $\eta \in (0, 1)$, we have

$$\inf_{\eta < \|\mathbf{x} - \mathbf{x}_*\| < \frac{1}{\eta}} (\mathbf{x} - \mathbf{x}_*)^\top B \nabla\phi(\mathbf{x}) > 0.$$

157 However, the restriction to strongly convex problems allows us to prove a more general
158 convergence result.

159 While the SA method uses the negative gradient direction, in [\[24\]](#) general descent
160 directions have been considered such that for all k

$$161 \quad (3.1) \quad \mathbf{g}(\mathbf{x}_k)^\top \mathbf{d}_k < 0,$$

$$162 \quad (3.2) \quad (\mathbf{x}_k - \mathbf{x}_*)^\top \mathbb{E}(\mathbf{d}_k|\mathcal{F}_k) \leq -c_3\|\mathbf{x}_k - \mathbf{x}_*\| \quad \text{a.s.},$$

$$163 \quad \|\mathbf{d}_k\| \leq c_4\|\mathbf{g}(\mathbf{x}_k)\| \quad \text{a.s..}$$

164 Here we relax [\(3.1\)](#) and [\(3.2\)](#) so that the direction \mathbf{d}_k need neither be a descent
165 direction nor satisfy [\(3.2\)](#). This relaxation allows us to extend the set of directions
166 covered by the theoretical analysis presenter further on. At each iteration, we allow
167 a deviation from a descent direction proportional to δ_k , where $\{\delta_k\}$ is a predefined
168 sequence of positive numbers that converges to zero with almost arbitrary rate. More
169 precisely, the following condition must hold:

$$170 \quad (3.3) \quad \sum_k \alpha_k \delta_k < \infty.$$

171 Thus, a possible choice could be $\delta_k = \nu^k$, where $\nu \in (0, 1)$, regardless of the choice
172 of the gain sequence. On the other hand, if we choose the standard gain sequence

173 $\alpha_k = 1/k$, then $\delta_k = 1/k^\epsilon$, with arbitrary small $\epsilon > 0$, is a suitable choice. Roughly
 174 speaking, the set of feasible directions is rather wide while we are far away from the
 175 solution, and the descent condition is enforced as we progress towards the solution.
 176 More precisely, we make the following assumptions on the search directions.

ASSUMPTION 3.1. *The direction \mathbf{d}_k satisfies*

$$\nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k | \mathcal{F}_k) \leq \delta_k c_2 - c_3 \|\nabla\phi(\mathbf{x}_k)\|^2.$$

ASSUMPTION 3.2. *The direction \mathbf{d}_k satisfies*

$$\|\mathbf{d}_k\| \leq c_4 \|\mathbf{g}(\mathbf{x}_k)\| \text{ a.s..}$$

We observe that [Assumptions 3.1](#) and [3.2](#) can be seen as a stochastic version of well-known sufficient conditions that guarantee gradient-related directions in the deterministic setting [[4](#), p. 36], i.e.,

$$\nabla\phi(\mathbf{x}_k)^\top \mathbf{d}_k \leq -q_1 \|\nabla\phi(\mathbf{x}_k)\|^{p_1}, \quad \|\mathbf{d}_k\| \leq q_2 \|\nabla\phi(\mathbf{x}_k)\|^{p_2}$$

177 for $q_1, q_2 > 0$ and $p_1, p_2 \geq 0$.

178 In the following theorem we prove almost sure convergence for the general Algo-
 179 rithm SOS.

180 **THEOREM 3.3.** *Let [Assumptions 2.1](#) to [2.3](#) and [Assumptions 3.1](#) and [3.2](#) hold,*
 181 *and let $\{\mathbf{x}_k\}$ be generated by [Algorithm 2.1](#). Assume also that [\(3.3\)](#) holds. Then*
 182 *$\mathbf{x}_k \rightarrow \mathbf{x}_*$ a.s..*

183 *Proof.* Since $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$ we have, by [Assumption 2.1](#) and the descent
 184 lemma [[4](#), Proposition A24],

$$185 \quad \phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*) \leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) + \alpha_k \nabla\phi(\mathbf{x}_k)^\top \mathbf{d}_k + \frac{L}{2} \alpha_k^2 \|\mathbf{d}_k\|^2.$$

186 Therefore, by [Assumption 3.2](#),

$$\begin{aligned} 187 \quad \mathbb{E}(\phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*) | \mathcal{F}_k) &\leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) + \alpha_k \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k | \mathcal{F}_k) \\ 188 \quad &\quad + \frac{L}{2} \alpha_k^2 c_4^2 \mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2 | \mathcal{F}_k) \\ 189 \quad &= \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) + \alpha_k \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k | \mathcal{F}_k) \\ 190 \quad &\quad + \alpha_k^2 c_5 \mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2 | \mathcal{F}_k) \end{aligned}$$

191 where $c_5 = Lc_4^2/2$. From [\(2.5\)](#) (arising from [Assumptions 2.1](#) and [2.3](#)) and [Assump-](#)
 192 [tion 3.1](#) it follows that

$$\begin{aligned} 193 \quad \mathbb{E}(\phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*) | \mathcal{F}_k) &\leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) + \alpha_k^2 c_1 c_5 (1 + \|\mathbf{x}_k - \mathbf{x}_*\|^2) \\ 194 \quad &\quad + \alpha_k (\delta_k c_2 - c_3 \|\nabla\phi(\mathbf{x}_k)\|^2). \end{aligned}$$

195 Since [\(2.2\)](#) holds, we have

$$\begin{aligned} 196 \quad \mathbb{E}(\phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*) | \mathcal{F}_k) &\leq (1 + \alpha_k^2 c_6) (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) + \alpha_k^2 c_1 c_5 \\ 197 \quad &\quad + \alpha_k \delta_k c_2 - \alpha_k c_3 \frac{2}{L} (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)), \end{aligned}$$

with $c_6 = 2c_1 c_5 / \mu$. Taking $\beta_k = \alpha_k^2 c_6$, $U_k = \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)$, $\xi_k = \alpha_k^2 c_1 c_5 + \alpha_k \delta_k c_2$ and $\rho_k = 2\alpha_k c_3 / L (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*))$, we have

$$\sum_k \beta_k < \infty, \quad \sum_k \xi_k < \infty$$

because of [Assumption 2.2](#) and (3.3), and $U_k \geq 0$ as \mathbf{x}_* is the solution of (1.1). Therefore, by [Theorem 2.4](#) we conclude that $\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)$ converges a.s. and $\sum_k \rho_k < \infty$ a.s.. Hence, we have

$$0 = \lim_{k \rightarrow \infty} \rho_k = \lim_{k \rightarrow \infty} \alpha_k c_3 \frac{2}{L} (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \quad a.s..$$

There are two possibilities for the sequence $\{\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)\}$: either there exists an infinite set $\mathcal{K} \subset \mathbb{N}$ such that

$$\lim_{k \in \mathcal{K}, k \rightarrow \infty} \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) = 0 \quad a.s.$$

198 or there exists $\varepsilon > 0$ such that

$$199 \quad (3.4) \quad \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) \geq \varepsilon \quad a.s. \text{ for all } k \text{ sufficiently large.}$$

If \mathcal{K} exists, then we have that the whole sequence $\{\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)\}$ converges to zero a.s., and then $\mathbf{x}_k \rightarrow \mathbf{x}_*$ a.s. because of the continuity of ϕ . On the other hand, if (3.4) holds, then

$$\sum_k \rho_k = \sum_k \alpha_k c_3 \frac{2}{L} (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \geq c_3 \frac{2}{L} \varepsilon \sum_k \alpha_k = \infty \quad a.s.,$$

200 which is a contradiction. Thus we conclude that $\mathbf{x}_k \rightarrow \mathbf{x}_*$ a.s.. \square

201 Now we extend the scope of search directions towards second-order approxima-
202 tions. Since [Assumption 2.1](#) holds, we also assume that the approximate Hessians are
203 positive definite and bounded.

ASSUMPTION 3.4. *For every approximate Hessian $B(\mathbf{x})$,*

$$\mu I \preceq B(\mathbf{x}) \preceq LI.$$

204 This assumption is fulfilled in many significant cases. For example, in binary classi-
205 fication, mini-batch subsampled Hessians are taken as positive definite and bounded
206 matrices, either with a proper choice of the subsample [32], or with regularization [8].
207 The same is true for least squares problems.

[Assumption 3.4](#) implies

$$\frac{1}{L} I \preceq B^{-1}(\mathbf{x}) \preceq \frac{1}{\mu} I,$$

208 and hence $\|B^{-1}(\mathbf{x})\| \leq \mu^{-1}$.

209 We also assume that the noise terms $\varepsilon_f(\mathbf{x})$, $\varepsilon_g(\mathbf{x})$ and $\varepsilon_B(\mathbf{x})$ are mutually inde-
210 pendent, which implies that the same is true for f , \mathbf{g} and B . This independence
211 assumption will be relaxed in [section 5](#) in order to cope with finite-sum problems,
212 where the gradient and Hessian approximations may be taken from the same sample.
213 By defining

$$214 \quad (3.5) \quad \mathbf{d}_k = -D_k \mathbf{g}(\mathbf{x}_k), \quad D_k = B^{-1}(\mathbf{x}_k),$$

we have

$$\|\mathbf{d}_k\| \leq \frac{1}{\mu} \|\mathbf{g}(\mathbf{x}_k)\|,$$

215 thus [Assumption 3.2](#) holds. Furthermore, since D_k is independent of $\mathbf{g}(\mathbf{x}_k)$, we obtain

$$\begin{aligned}
216 \quad \mathbb{E}(\nabla\phi(\mathbf{x}_k)^\top \mathbf{d}_k | \mathcal{F}_k) &= \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(-D_k \mathbf{g}(\mathbf{x}_k) | \mathcal{F}_k) \\
217 &= \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(-D_k | \mathcal{F}_k) \mathbb{E}(\mathbf{g}(\mathbf{x}_k) | \mathcal{F}_k) \\
218 &= \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(-D_k | \mathcal{F}_k) \nabla\phi(\mathbf{x}_k) \\
219 &= \mathbb{E}(-\nabla\phi(\mathbf{x}_k)^\top D_k \nabla\phi(\mathbf{x}_k) | \mathcal{F}_k) \\
220 &\leq \mathbb{E}\left(-\frac{1}{L} \|\nabla\phi(\mathbf{x}_k)\|^2 | \mathcal{F}_k\right) = -\frac{1}{L} \|\nabla\phi(\mathbf{x}_k)\|^2
\end{aligned}$$

221 and hence

$$222 \quad (3.6) \quad \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k | \mathcal{F}_k) \leq -\frac{1}{L} \|\nabla\phi(\mathbf{x}_k)\|^2.$$

223 Then [Assumption 3.1](#) holds with $c_2 = 0$ and $c_3 = \frac{1}{L}$.

224 **COROLLARY 3.5.** *Let [Assumptions 2.1 to 2.3](#) and [Assumption 3.4](#) hold, and let*
225 *(3.3) hold. If $\{\mathbf{x}_k\}$ is a sequence generated by [Algorithm 2.1](#) with \mathbf{d}_k defined in (3.5),*
226 *then $\mathbf{x}_k \rightarrow \mathbf{x}_*$ a.s.*

227 *Proof.* The proof is an immediate consequence of [Theorem 3.3](#) and the previous
228 observations. \square

229 Finally, let us consider the case of inexact Newton methods in the stochastic
230 approximation framework, i.e., when the linear system

$$231 \quad (3.7) \quad B(\mathbf{x}_k) \mathbf{d}_k = -\mathbf{g}(\mathbf{x}_k)$$

232 is solved only approximately, i.e.,

$$233 \quad (3.8) \quad \|B(\mathbf{x}_k) \mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)\| \leq \delta_k \gamma_k,$$

234 where γ_k is a random variable, and δ_k satisfies (3.3).

235 For deterministic inexact Newton methods, global convergence has been proved
236 when $\{\mathbf{x}_k\}$ is bounded and the forcing terms are small enough – see the alternative
237 statement of [Theorem 3.4](#) in [15, page 400]. Thus, we will assume $\{\mathbf{x}_k\}$ bounded in
238 the stochastic case as well. For γ_k we assume bounded variance as follows.

ASSUMPTION 3.6. *The sequence of random variables $\{\gamma_k\}$ is such that*

$$\mathbb{E}(\gamma_k^2 | \mathcal{F}_k) \leq M_3.$$

239 Note that [Assumption 3.6](#) implies

$$240 \quad \mathbb{E}(\gamma_k | \mathcal{F}_k) \leq \sqrt{\mathbb{E}(\gamma_k^2 | \mathcal{F}_k)} \leq \sqrt{M_3} := M_4.$$

241 The main property of the search direction that allows us to prove [Theorem 3.3](#) is
242 stated in [Assumption 3.1](#). Now we prove that [Assumption 3.1](#) holds if the sequence
243 $\{\mathbf{x}_k\}$ is bounded and [Assumption 3.6](#) holds.

244 **LEMMA 3.7.** *Let $\{\mathbf{x}_k\}$ be a sequence generated by [Algorithm 2.1](#) such that (3.8)*
245 *and [Assumption 3.6](#) hold. If $\{\mathbf{x}_k\}$ is bounded, then [Assumption 3.1](#) holds.*

246 *Proof.* If $\{\mathbf{x}_k\}$ is bounded then $\|\nabla\phi(\mathbf{x}_k)\| \leq M_5$ as ϕ is continuously differentiable.
247 Furthermore, [Assumption 3.6](#) implies

$$248 \quad (3.9) \quad \|\nabla\phi(\mathbf{x}_k)\| \mathbb{E}(\gamma_k | \mathcal{F}_k) \leq M_5 M_4 := M_6.$$

Let us denote $\mathbf{r}_k = B(\mathbf{x}_k)\mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)$. Then, by (3.8), $\|\mathbf{r}_k\| \leq \delta_k \gamma_k$. Furthermore,

$$\mathbf{d}_k = B(\mathbf{x}_k)^{-1}\mathbf{r}_k - B(\mathbf{x}_k)^{-1}\mathbf{g}(\mathbf{x}_k).$$

Setting $\mathbf{d}_k^N = -B(\mathbf{x}_k)^{-1}\mathbf{g}(\mathbf{x}_k)$, we have

$$\mathbf{d}_k - \mathbf{d}_k^N = B(\mathbf{x}_k)^{-1}\mathbf{r}_k$$

and

$$\nabla\phi(\mathbf{x}_k)^\top \mathbf{d}_k = \nabla\phi(\mathbf{x}_k)^\top (\mathbf{d}_k - \mathbf{d}_k^N + \mathbf{d}_k^N) = \nabla\phi(\mathbf{x}_k)^\top \mathbf{d}_k^N + \nabla\phi(\mathbf{x}_k)^\top (\mathbf{d}_k - \mathbf{d}_k^N).$$

Taking the conditional expectation, we get

$$\nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k|\mathcal{F}_k) = \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k^N|\mathcal{F}_k) + \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k - \mathbf{d}_k^N|\mathcal{F}_k).$$

It has been shown, see (3.6), that

$$\nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k^N|\mathcal{F}_k) \leq -\frac{1}{L}\|\nabla\phi(\mathbf{x}_k)\|^2,$$

249 thus

$$250 \quad (3.10) \quad \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k|\mathcal{F}_k) \leq -\frac{1}{L}\|\nabla\phi(\mathbf{x}_k)\|^2 + \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(B(\mathbf{x}_k)^{-1}\mathbf{r}_k|\mathcal{F}_k).$$

251 Furthermore,

$$252 \quad \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(B(\mathbf{x}_k)^{-1}\mathbf{r}_k|\mathcal{F}_k) \leq \|\nabla\phi(\mathbf{x}_k)\| \mathbb{E}(\|B(\mathbf{x}_k)^{-1}\| \|\mathbf{r}_k\| |\mathcal{F}_k)$$

$$253 \quad (3.11) \quad \leq \frac{1}{\mu}\|\nabla\phi(\mathbf{x}_k)\| \delta_k \mathbb{E}(\gamma_k|\mathcal{F}_k) \leq \frac{1}{\mu} \delta_k M_6$$

because of (3.8) and (3.9). Putting together (3.10) and (3.11), we get

$$\nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k|\mathcal{F}_k) \leq \delta_k c_2 - c_3 \|\nabla\phi(\mathbf{x}_k)\|^2$$

254 with $c_2 = M_6/\mu$ and $c_3 = 1/L$. Therefore, [Assumption 3.1](#) holds. \square

255 Notice that [Assumption 3.2](#) is not necessarily satisfied by the direction \mathbf{d}_k in (3.8).
256 Therefore, we cannot apply [Theorem 3.3](#). Nevertheless, we can prove the following.

257 **THEOREM 3.8.** *Let [Assumptions 2.1 to 2.3](#) and [Assumptions 3.4](#) and [3.6](#) hold. Let*
258 *$\{\mathbf{x}_k\}$ be a sequence generated by [Algorithm 2.1](#) with search direction \mathbf{d}_k satisfying (3.8)*
259 *with δ_k such that (3.3) holds. If $\{\mathbf{x}_k\}$ is bounded, then $\mathbf{x}_k \rightarrow \mathbf{x}_*$ a.s..*

Proof. The direction \mathbf{d}_k satisfies

$$\|\mathbf{d}_k\| = \|B(\mathbf{x}_k)^{-1}(\mathbf{r}_k - \mathbf{g}(\mathbf{x}_k))\| \leq \frac{1}{\mu}(\|\mathbf{r}_k\| + \|\mathbf{g}(\mathbf{x}_k)\|) \leq \frac{1}{\mu}(\delta_k \gamma_k + \|\mathbf{g}(\mathbf{x}_k)\|),$$

thanks to (3.8). Therefore,

$$\|\mathbf{d}_k\|^2 \leq \frac{2}{\mu^2} (\delta_k^2 \gamma_k^2 + \|\mathbf{g}(\mathbf{x}_k)\|^2)$$

260 and

$$261 \quad \mathbb{E}(\|\mathbf{d}_k\|^2|\mathcal{F}_k) \leq \frac{2}{\mu^2} (\delta_k^2 \mathbb{E}(\gamma_k^2|\mathcal{F}_k) + \mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2|\mathcal{F}_k))$$

$$262 \quad \leq \frac{2}{\mu^2} (\delta_k^2 M_3 + c_1(1 + \|\mathbf{x}_k - \mathbf{x}_*\|^2)),$$

263 because of (2.5) and Assumption 3.6. Therefore,

$$264 \quad (3.12) \quad \mathbb{E}(\|\mathbf{d}_k\|^2|\mathcal{F}_k) \leq c_7 + c_8\|\mathbf{x}_k - \mathbf{x}_*\|^2,$$

for $c_7 = (2/\mu^2)(\delta_{\max}^2 M_3 + c_1)$ and $c_8 = 2c_1/\mu^2$, where $\delta_{\max} = \max_k \delta_k$. Using the descent lemma and Assumption 2.1 as in Theorem 3.3, we get

$$\phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*) \leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) + \alpha_k \nabla \phi(\mathbf{x}_k)^\top \mathbf{d}_k + \frac{L}{2} \alpha_k^2 \|\mathbf{d}_k\|^2$$

265 and, by (3.12) and Lemma 3.7,

$$\begin{aligned} 266 \quad \mathbb{E}(\phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*)|\mathcal{F}_k) &\leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) + \alpha_k \nabla \phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k|\mathcal{F}_k) \\ 267 \quad &+ \frac{L}{2} \alpha_k^2 \mathbb{E}(\|\mathbf{d}_k\|^2|\mathcal{F}_k) \\ 268 \quad &\leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) + \alpha_k \left(\delta_k \frac{M_6}{\mu} - \frac{1}{L} \|\nabla \phi(\mathbf{x}_k)\|^2 \right) \\ 269 \quad &+ \alpha_k^2 \frac{L}{2} (c_7 + c_8 \|\mathbf{x}_k - \mathbf{x}_*\|^2). \end{aligned}$$

270 Using (2.2), we get

$$\begin{aligned} 271 \quad \mathbb{E}(\phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*)|\mathcal{F}_k) &\leq (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \left(1 + \alpha_k^2 \frac{L}{\mu} c_8 \right) + \alpha_k^2 \frac{L}{2} c_7 \\ 272 \quad &+ \alpha_k \delta_k \frac{M_6}{\mu} - \alpha_k \frac{2}{L^2} (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)). \end{aligned}$$

Now we define

$$\beta_k = \alpha_k^2 \frac{L}{\mu} c_8, \quad \xi_k = \alpha_k^2 \frac{L}{2} c_7 + \alpha_k \delta_k \frac{M_6}{\mu}, \quad \rho_k = \alpha_k \frac{2}{L^2} (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*))$$

and

$$U_k = \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*).$$

273 As the hypotheses of Theorem 2.4 are fulfilled due to (3.3) and Assumption 2.2, we
274 have that $\mathbf{x}_k \rightarrow \mathbf{x}_*$ a.s. \square

So far we have required only that γ_k is a random variable with bounded variance. Following inexact Newton methods in the deterministic case [15], the norm of the left-hand side in (3.7) can be used to define the inexactness in the solution of the linear system. By setting

$$\gamma_k = c_9 \|\mathbf{g}(\mathbf{x}_k)\|,$$

275 we have that (3.8) implies

$$\begin{aligned} 276 \quad \|\mathbf{d}_k\| &= \|B(\mathbf{x}_k)^{-1}(\mathbf{r}_k - \mathbf{g}(\mathbf{x}_k))\| \leq \frac{1}{\mu} (\|\mathbf{r}_k\| + \|\mathbf{g}(\mathbf{x}_k)\|) \\ 277 \quad &\leq \frac{1}{\mu} (\delta_k c_9 \|\mathbf{g}(\mathbf{x}_k)\| + \|\mathbf{g}(\mathbf{x}_k)\|) = \frac{1}{\mu} (\delta_{\max} c_9 + 1) \|\mathbf{g}(\mathbf{x}_k)\|. \end{aligned}$$

Therefore, for this choice of γ_k we have that Assumption 3.2 holds as well. Furthermore, by (2.5) we get

$$\mathbb{E}(\gamma_k^2|\mathcal{F}_k) = c_9^2 \mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2|\mathcal{F}_k) \leq c_9^2 c_1 (1 + \|\mathbf{x}_k - \mathbf{x}_*\|^2).$$

Assuming that $\{\mathbf{x}_k\}$ is bounded, we get

$$\mathbb{E}(\gamma_k^2 | \mathcal{F}_k) \leq c_9^2 c_1 (1 + \|\mathbf{x}_k - \mathbf{x}_*\|^2) \leq c_9^2 c_1 (1 + M_7) := M_8$$

and hence [Assumption 3.6](#) holds as well. The previous observations imply the following convergence statement, whose proof is straightforward.

COROLLARY 3.9. *Let [Assumptions 2.1 to 2.3](#) and [Assumption 3.4](#) hold. Let $\{\mathbf{x}_k\}$ be a sequence generated by [Algorithm 2.1](#) where \mathbf{d}_k satisfies*

$$\|B(\mathbf{x}_k)\mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)\| \leq \delta_k \|\mathbf{g}(\mathbf{x}_k)\|,$$

and δ_k satisfies [\(3.3\)](#) holds. If $\{\mathbf{x}_k\}$ is bounded then $\mathbf{x}_k \rightarrow \mathbf{x}_*$ a.s..

The next theorem considers the most general case, extending the tolerance for inexact solutions of the Newton linear system [\(3.7\)](#) even further. Let us define

$$(3.13) \quad \gamma_k = \omega_1 \eta_k + \omega_2 \|\mathbf{g}(\mathbf{x}_k)\|$$

for some $\omega_1, \omega_2 \geq 0$ and a random variable η_k such that

$$(3.14) \quad \mathbb{E}(\eta_k^2 | \mathcal{F}_k) \leq M_9,$$

i.e., with bounded variance.

THEOREM 3.10. *Let [Assumptions 2.1 to 2.3](#) and [Assumption 3.4](#) hold. Let $\{\mathbf{x}_k\}$ be a sequence generated by [Algorithm 2.1](#) where \mathbf{d}_k satisfies*

$$\|B(\mathbf{x}_k)\mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)\| \leq \delta_k \gamma_k,$$

with γ_k defined by [\(3.13\)](#) and δ_k such that [\(3.3\)](#) holds. If $\{\mathbf{x}_k\}$ is bounded, then $\mathbf{x}_k \rightarrow \mathbf{x}_*$ a.s..

Proof. First, note that the search direction \mathbf{d}_k satisfies

$$\begin{aligned} \|\mathbf{d}_k\| &= \|B(\mathbf{x}_k)^{-1}(\mathbf{r}_k - \mathbf{g}(\mathbf{x}_k))\| \leq \frac{1}{\mu} (\delta_k \gamma_k + \|\mathbf{g}(\mathbf{x}_k)\|) \\ &= \frac{1}{\mu} (\omega_1 \delta_k \eta_k + (1 + \omega_2 \delta_k) \|\mathbf{g}(\mathbf{x}_k)\|), \end{aligned}$$

and then, by [\(3.14\)](#), [\(2.5\)](#), and [Assumption 2.1](#),

$$\begin{aligned} \mathbb{E}(\|\mathbf{d}_k\|^2 | \mathcal{F}_k) &\leq \frac{2}{\mu^2} (\omega_1^2 \delta_k^2 \mathbb{E}(\eta_k^2 | \mathcal{F}_k) + (1 + \omega_2 \delta_k)^2 \mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2 | \mathcal{F}_k)) \\ &\leq \frac{2}{\mu^2} \omega_1^2 \delta_k^2 M_9 + \frac{2}{\mu^2} (1 + \omega_2 \delta_k)^2 c_1 (1 + \|\mathbf{x}_k - \mathbf{x}_*\|^2) \\ &= \frac{2}{\mu^2} (\omega_1^2 \delta_k^2 M_9 + c_1 (1 + \omega_2 \delta_k)^2) + \frac{2}{\mu^2} c_1 (1 + \omega_2 \delta_k)^2 \frac{2}{\mu} (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \\ &\leq M_{10} + M_{11} (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)), \end{aligned}$$

where $M_{10} = 2/\mu^2 (\omega_1^2 \delta_{\max}^2 M_9 + c_1 (1 + \omega_2 \delta_{\max})^2)$ and $M_{11} = 4c_1/\mu^3 (1 + \omega_2 \delta_{\max})^2$.

Since $\mathbf{d}_k = B(\mathbf{x}_k)^{-1}(\mathbf{r}_k - \mathbf{g}(\mathbf{x}_k))$, using the same arguments as for [\(3.10\)](#) we obtain

$$(3.15) \quad \nabla \phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k | \mathcal{F}_k) \leq -\frac{1}{L} \|\nabla \phi(\mathbf{x}_k)\|^2 + \nabla \phi(\mathbf{x}_k)^\top \mathbb{E}(B(\mathbf{x}_k)^{-1} \mathbf{r}_k | \mathcal{F}_k).$$

301 Furthermore,

$$\begin{aligned}
302 \quad \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(B(\mathbf{x}_k)^{-1}\mathbf{r}_k|\mathcal{F}_k) &\leq \|\nabla\phi(\mathbf{x}_k)\| \mathbb{E}(\|B(\mathbf{x}_k)^{-1}\|\|\mathbf{r}_k\||\mathcal{F}_k) \\
303 \quad &\leq \frac{1}{\mu}\|\nabla\phi(\mathbf{x}_k)\| \delta_k \mathbb{E}(\gamma_k|\mathcal{F}_k)
\end{aligned}$$

304 and

$$\begin{aligned}
305 \quad \mathbb{E}(\gamma_k|\mathcal{F}_k) &= \mathbb{E}(\omega_1\eta_k + \omega_2\|\mathbf{g}(\mathbf{x}_k)\||\mathcal{F}_k) \\
306 \quad &= \omega_1\mathbb{E}(\eta_k|\mathcal{F}_k) + \omega_2\mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\||\mathcal{F}_k) \\
307 \quad &\leq \omega_1\sqrt{\mathbb{E}(\eta_k^2|\mathcal{F}_k)} + \omega_2\sqrt{\mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2|\mathcal{F}_k)} \\
308 \quad &\leq \omega_1\sqrt{M_9} + \omega_2\sqrt{c_1(1 + \|\mathbf{x}_k - \mathbf{x}_*\|^2)} \\
309 \quad &\leq \omega_1\sqrt{M_9} + \omega_2\sqrt{c_1(1 + M_7)} := M_{12}.
\end{aligned}$$

310 Putting together the above estimates and using the descent lemma as in the previous
311 proofs, we get

$$\begin{aligned}
312 \quad \mathbb{E}(\phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*)|\mathcal{F}_k) &\leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) + \alpha_k \nabla\phi(\mathbf{x}_k)^\top \mathbb{E}(\mathbf{d}_k|\mathcal{F}_k) + \alpha_k^2 \frac{L}{2} \mathbb{E}(\|\mathbf{d}_k\|^2|\mathcal{F}_k) \\
313 \quad &\leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) + \alpha_k^2 \frac{L}{2} (M_{10} + M_{11}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*))) \\
314 \quad &\quad + \alpha_k \left(-\frac{1}{L} \|\nabla\phi(\mathbf{x}_k)\|^2 + \|\nabla\phi(\mathbf{x}_k)\| \frac{1}{\mu} \delta_k M_{12} \right) \\
315 \quad &\leq (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \left(1 + \alpha_k^2 \frac{L}{2} M_{11} \right) + \alpha_k^2 \frac{L}{2} M_{10} \\
316 \quad &\quad + \alpha_k \delta_k \frac{1}{\mu} M_{13} M_{12} - \alpha_k \frac{2}{L^2} (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)),
\end{aligned}$$

where $\|\nabla\phi(\mathbf{x}_k)\| \leq M_{13}$ because of the boundedness of $\{\mathbf{x}_k\}$ and the continuity of $\nabla\phi$. By defining

$$\beta_k = \alpha_k^2 \frac{L}{2} M_{11}, \quad \xi_k = \alpha_k^2 \frac{L}{2} M_{10} + \alpha_k \delta_k \frac{1}{\mu} M_{13} M_{12}, \quad \rho_k = \alpha_k \frac{2}{L^2} (\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)),$$

317 and observing that (3.3) and Assumption 2.2 hold, we can apply Theorem 2.4 to get
318 the thesis. \square

319 **4. SOS method with line search.** It is well known that in practice a gain
320 sequence that satisfies Assumption 2.2 is usually too conservative and makes the
321 algorithm slow because the step length becomes too small soon. In order to avoid this
322 drawback, we propose a practical version of Algorithm SOS that uses a line search in
323 the initial phase and then reduces to SOS if the step length obtained with the line
324 search becomes too small, e.g., smaller than some predetermined threshold $t_{\min} > 0$.

325 Since the search directions considered in the previous sections do not have to be
326 descent directions (not even for the current objective function approximation), and
327 the line search can be performed considering only the approximate objective function,
328 we choose a nonmonotone line-search strategy.

329 We state the new algorithmic framework in Algorithm 4.1. Note that this algo-
330 rithm remains well defined even with the monotone (classical Armijo) line search – if
331 the search direction is not a descent one, we shift to the predefined gain sequence.

332 We prove the a.s. convergence of Algorithm LSOS under a mild additional as-
333 sumption.

Algorithm 4.1 Line-search Second-Order Stochastic (LSOS) method

- 1: given $\mathbf{x}^0 \in \mathbb{R}^n$, $\eta \in (0, 1)$, $t_{\min} > 0$ and $\{\alpha_k\}, \{\delta_k\}, \{\zeta_k\} \subset \mathbb{R}_+$
 - 2: set LSphase = *active*
 - 3: **for** $k = 0, 1, 2, \dots$ **do**
 - 4: compute a search direction \mathbf{d}_k such that

$$(4.1) \quad \|B(\mathbf{x}_k)\mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)\| \leq \delta_k \|\mathbf{g}(\mathbf{x}_k)\|.$$
 - 5: find a step length t_k as follows:
 - 6: **if** LSphase = *active* **then** find t_k that satisfies

$$(4.2) \quad f(\mathbf{x}_k + t_k \mathbf{d}_k) \leq f(\mathbf{x}_k) + \eta t_k \mathbf{g}(\mathbf{x}_k)^\top \mathbf{d}_k + \zeta_k$$
 - 7: **if** $t_k < t_{\min}$ **then** set LSphase = *inactive*
 - 8: **if** LSphase = *inactive* **then** set $t_k = \alpha_k$
 - 9: set $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$
 - 10: **end for**
-

ASSUMPTION 4.1. *The objective function estimator f is unbiased, i.e.,*

$$\mathbb{E}(\varepsilon_f(\mathbf{x}) | \mathcal{F}_k) = 0.$$

334 THEOREM 4.2. *Let Assumptions 2.1 to 2.3, Assumption 3.4 and Assumption 4.1*
 335 *hold. Assume also that the sequence $\{\zeta_k\}$ is summable and the forcing term sequence*
 336 *$\{\delta_k\}$ satisfies (3.3). If the sequence $\{\mathbf{x}_k\}$ generated by Algorithm 4.1 is bounded, then*
 337 *$\mathbf{x}_k \rightarrow \mathbf{x}_*$ a.s..*

Proof. If there exists an iteration k such that $t_k < t_{\min}$, then Algorithm LSOS reduces to SOS and the thesis follows from Corollary 3.9. Let us consider the case $t_k \geq t_{\min}$ for all k . Using $\mathbf{r}_k = B(\mathbf{x}_k)\mathbf{d}_k + \mathbf{g}(\mathbf{x}_k)$ we obtain

$$\mathbf{d}_k = B(\mathbf{x}_k)^{-1}\mathbf{r}_k - B(\mathbf{x}_k)^{-1}\mathbf{g}(\mathbf{x}_k).$$

338 Furthermore, Assumption 3.4 together with (4.1) implies

$$\begin{aligned} 339 \quad \mathbf{g}(\mathbf{x}_k)^\top \mathbf{d}_k &= \mathbf{g}(\mathbf{x}_k)^\top B(\mathbf{x}_k)^{-1}\mathbf{r}_k - \mathbf{g}(\mathbf{x}_k)^\top B(\mathbf{x}_k)^{-1}\mathbf{g}(\mathbf{x}_k) \\ 340 &\leq \|\mathbf{g}(\mathbf{x}_k)\| \|B(\mathbf{x}_k)^{-1}\| \|\mathbf{r}_k\| - \frac{1}{L} \|\mathbf{g}(\mathbf{x}_k)\|^2 \\ 341 &\leq \frac{1}{\mu} \delta_k \|\mathbf{g}(\mathbf{x}_k)\|^2 - \frac{1}{L} \|\mathbf{g}(\mathbf{x}_k)\|^2 \\ 342 &= \left(\frac{\delta_k}{\mu} - \frac{1}{L} \right) \|\mathbf{g}(\mathbf{x}_k)\|^2. \end{aligned}$$

Assumption 2.2 together with (3.3) implies $\delta_k \rightarrow 0$. Therefore, there exists \bar{k} such that

$$\delta_k \leq \frac{\mu}{2L} \quad \text{for all } k \geq \bar{k}$$

343 and hence

$$344 \quad (4.3) \quad \mathbf{g}(\mathbf{x}_k)^\top \mathbf{d}_k \leq -\frac{1}{2L} \|\mathbf{g}(\mathbf{x}_k)\|^2.$$

345 Furthermore, LSphase is active at each iteration and for $k \geq \bar{k}$ we get

$$\begin{aligned}
346 \quad f(\mathbf{x}_k + t_k \mathbf{d}_k) &\leq f(\mathbf{x}_k) + \eta t_k \mathbf{g}(\mathbf{x}_k)^\top \mathbf{d}_k + \zeta_k \\
347 \quad &\leq f(\mathbf{x}_k) - \eta t_k \frac{1}{2L} \|\mathbf{g}(\mathbf{x}_k)\|^2 + \zeta_k \\
348 \quad &\leq f(\mathbf{x}_k) - \eta t_{\min} \frac{1}{2L} \|\mathbf{g}(\mathbf{x}_k)\|^2 + \zeta_k.
\end{aligned}$$

349 Setting $c_{10} = \eta t_{\min}/(2L)$, taking the conditional expectation and using [Assump-](#)
350 [tion 4.1](#), we get

$$351 \quad (4.4) \quad \phi(\mathbf{x}_{k+1}) \leq \phi(\mathbf{x}_k) - c_{10} \mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2 | \mathcal{F}_k) + \zeta_k.$$

352 [Assumption 2.3](#) implies

$$353 \quad (4.5) \quad \mathbb{E}(\mathbf{g}(\mathbf{x}_k) | \mathcal{F}_k) = \nabla \phi(\mathbf{x}_k),$$

354 and thus we get

$$355 \quad (4.6) \quad \|\nabla \phi(\mathbf{x}_k)\|^2 = \|\mathbb{E}(\mathbf{g}(\mathbf{x}_k) | \mathcal{F}_k)\|^2 \leq \mathbb{E}^2(\|\mathbf{g}(\mathbf{x}_k)\|^2 | \mathcal{F}_k) \leq \mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2 | \mathcal{F}_k),$$

356 which, together with [Assumption 2.1](#), implies

$$357 \quad (4.7) \quad \frac{\mu}{L} \|\mathbf{x}_k - \mathbf{x}_*\|^2 \leq \|\nabla \phi(\mathbf{x}_k)\|^2 \leq \mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2 | \mathcal{F}_k).$$

358 Combining (4.7) with (4.4) we have

$$359 \quad (4.8) \quad \phi(\mathbf{x}_{k+1}) \leq \phi(\mathbf{x}_k) - c_{11} \|\mathbf{x}_k - \mathbf{x}_*\|^2 + \zeta_k \quad \text{for all } k \geq \bar{k}$$

for a suitable \bar{k} , where $c_{11} = c_{10}\mu/L$. The boundedness of the iterates and the continuity of ϕ imply the existence of a constant Q such that $\phi(\mathbf{x}_k) \geq Q$ for all k . Furthermore, (4.8) implies that, for all $p \in \mathbb{N}$,

$$Q \leq \phi(\mathbf{x}_{\bar{k}+p}) \leq \phi(\mathbf{x}_{\bar{k}}) - c_{11} \sum_{j=0}^{p-1} \|\mathbf{x}_{\bar{k}+j} - \mathbf{x}_*\|^2 + \sum_{j=0}^{p-1} \zeta_{\bar{k}+j}.$$

360 Taking the expectation, letting p tend to infinity and using the summability of ζ_k , we
361 conclude that

$$362 \quad \sum_{k=0}^{\infty} \mathbb{E}(\|\mathbf{x}_k - \mathbf{x}_*\|^2) < \infty.$$

Finally, using Markov's inequality we have that for any $\epsilon > 0$

$$P(\|\mathbf{x}_k - \mathbf{x}_*\| \geq \epsilon) \leq \frac{\mathbb{E}(\|\mathbf{x}_k - \mathbf{x}_*\|^2)}{\epsilon^2}$$

and therefore

$$\sum_{k=0}^{\infty} P(\|\mathbf{x}_k - \mathbf{x}_*\| \geq \epsilon) < \infty.$$

363 The almost sure convergence follows from Borel-Cantelli Lemma [[22](#), Theorem 2.7],
364 which completes the proof. \square

365 **5. Specializing LSOS for finite sums.** Now we consider finite-sum problems,
 366 where the objective function is, e.g., the sample mean of a finite family of convex
 367 functions. This is the case, for example, of machine learning problems in which the
 368 logistic loss, the quadratic loss or other loss functions are used, usually coupled with ℓ_2 -
 369 regularization terms. Recently, much attention has been devoted to the development
 370 of methods for the solution of problems of this type. Therefore, we analyze extensions
 371 to this setting of the LSOS algorithmic framework.

372 Specifically, we focus on objective functions of the form

$$373 \quad (5.1) \quad \phi(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \phi_i(\mathbf{x}),$$

374 where each $\phi_i(\mathbf{x})$ is twice continuously differentiable and $\bar{\mu}$ -strongly convex, and has
 375 Lipschitz-continuous gradient with Lipschitz constant \bar{L} . It is straightforward to show
 376 that these assumptions imply that ϕ satisfies [Assumption 2.1](#).

377 We assume that at each iteration k a sample \mathcal{N}_k of size $N_k \ll N$ is chosen
 378 randomly and uniformly from $\mathcal{N} = \{1, \dots, N\}$. Then, we consider

$$379 \quad f_{\mathcal{N}_k}(\mathbf{x}) = \frac{1}{N_k} \sum_{i \in \mathcal{N}_k} \phi_i(\mathbf{x}),$$

380 which is an unbiased estimator of $\phi(\mathbf{x})$, i.e., [Assumption 4.1](#) holds.

381 By considering the first and second derivatives of $f_{\mathcal{N}_k}$, we obtain the following
 382 subsampled gradient and Hessian of ϕ :

$$383 \quad (5.2) \quad \mathbf{g}_{\mathcal{N}_k}(\mathbf{x}) = \frac{1}{N_k} \sum_{i \in \mathcal{N}_k} \nabla \phi_i(\mathbf{x}), \quad B_{\mathcal{N}_k}(\mathbf{x}) = \frac{1}{N_k} \sum_{i \in \mathcal{N}_k} \nabla^2 \phi_i(\mathbf{x}),$$

384 which are unbiased estimators of the gradient and the Hessian of ϕ as well. More
 385 precisely, the first equality in [Assumption 2.3](#) holds (i.e., $\mathbb{E}(\varepsilon_g(\mathbf{x})|\mathcal{F}_k) = 0$) together
 386 with [Assumption 3.4](#).

387 The derivative estimates in (5.2) can be replaced by more sophisticated ones,
 388 with the aim of improving the performance of second-order stochastic optimization
 389 methods. The Hessian approximation $B_{\mathcal{N}_k}(\mathbf{x})$ only needs to satisfy [Assumption 3.4](#) in
 390 order to prove the results contained in this section. Therefore, the theory we develop
 391 still holds if one replaces the subsampled Hessian approximation with a quasi-Newton
 392 approximation. For example, in [10] Byrd et al. propose to use subsampled gradients
 393 and an approximation of the inverse of the Hessian $\nabla^2 \phi(\mathbf{x})$, say H_k , built by means of
 394 a stochastic variant of limited memory BFGS (L-BFGS). Given a memory parameter
 395 m , H_k is defined by applying m BFGS updates to an initial matrix, using the m
 396 most recent correction pairs $(\mathbf{s}_j, \mathbf{y}_j) \in \mathbb{R}^n \times \mathbb{R}^n$ like in the deterministic version of the
 397 L-BFGS method. The pairs are obtained by averaging iterates, i.e., every l steps the
 398 following vectors are computed

$$399 \quad (5.3) \quad \mathbf{w}_j = \frac{1}{l} \sum_{i=k-l+1}^k \mathbf{x}_i, \quad \mathbf{w}_{j-1} = \frac{1}{l} \sum_{i=k-2l+1}^{k-l} \mathbf{x}_i,$$

400 where $j = \frac{k}{l}$, and they are used to build \mathbf{s}_j and \mathbf{y}_j as specified next:

$$401 \quad (5.4) \quad \mathbf{s}_j = \mathbf{w}_j - \mathbf{w}_{j-1}, \quad \mathbf{y}_j = B_{\mathcal{T}_j}(\mathbf{w}_j) \mathbf{s}_j,$$

where $\mathcal{T}_j \subset \{1, \dots, N\}$. By defining the set of the m most recent correction pairs as

$$\{(\mathbf{s}_j, \mathbf{y}_j), j = 1, \dots, m\},$$

402 the inverse Hessian approximation is computed as

$$403 \quad (5.5) \quad H_k = H_k^{(m)},$$

404 where for $j = 1, \dots, m$

$$405 \quad (5.6) \quad H_k^{(j)} = \left(I - \frac{\mathbf{s}_j \mathbf{y}_j^\top}{\mathbf{s}_j^\top \mathbf{y}_j} \right)^\top H_k^{(j-1)} \left(I - \frac{\mathbf{y}_j \mathbf{s}_j^\top}{\mathbf{s}_j^\top \mathbf{y}_j} \right) + \frac{\mathbf{s}_j \mathbf{s}_j^\top}{\mathbf{s}_j^\top \mathbf{y}_j},$$

and $H_k^{(0)} = (\mathbf{s}_m^\top \mathbf{y}_m / \|\mathbf{y}_m\|^2) I$. It can be proved (see [10, Lemma 3.1] and [30, Lemma 4]) that for approximate inverse Hessians of the form (5.5) there exist constants $\lambda_2 \geq \lambda_1 > 0$ such that

$$\lambda_1 I \preceq H_k \preceq \lambda_2 I,$$

i.e., [Assumption 3.4](#) holds with $\mu = \min\{\bar{\mu}, 1/\lambda_2\}$ and $L = \max\{\bar{L}, 1/\lambda_1\}$. The authors of [10] propose a version of [Algorithm 2.1](#) in which the direction is computed as

$$\mathbf{d}_k = -H_k \mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k),$$

406 and prove R-linear decrease of the expected value of the error in the function value.

407 As regards the gradient estimate, we observe that the second part of [Assump-](#)
 408 [tion 2.3](#) is not required by the method presented in this section. Notice that we can
 409 replace the subsampled gradient estimate in (5.2) with alternative estimates coming,
 410 e.g., from variance reduction techniques, which have gained much attention in the lit-
 411 erature. This is the case of the stochastic L-BFGS algorithm by Moritz et al. [30] and
 412 the stochastic block L-BFGS by Gower et al. [17], where SVRG gradient approxima-
 413 tions are used. The former method computes the same inverse Hessian approximation
 414 as in [10], while the latter uses an adaptive sketching technique exploiting the action
 415 of a sub-sampled Hessian on a set of random vectors rather than just on a single
 416 vector. Both stochastic BFGS algorithms use constant step lengths and have Q-linear
 417 rate of convergence of the expected value of the error in the objective function, but
 418 the block L-BFGS one appears more efficient than the other in most of the numerical
 419 experiments reported in [17].

420 Instead of choosing the SVRG approximation, we apply a mini-batch variant of
 421 the SAGA algorithm [12], used in [18]. Starting from the matrix $J^0 \in \mathbb{R}^{n \times N}$ whose
 422 columns are defined as $J_0^{(i)} = \nabla \phi_i(\mathbf{x}^0)$, at each iteration we compute the gradient
 423 approximation as

$$424 \quad (5.7) \quad \mathbf{g}_{\mathcal{N}_k}^{\text{SAGA}}(\mathbf{x}_k) = \frac{1}{N_k} \sum_{i \in \mathcal{N}_k} \left(\nabla \phi_i(\mathbf{x}_k) - J_k^{(i)} \right) + \frac{1}{N} \sum_{l=1}^N J_k^{(l)},$$

425 and, after updating the iterate, we set

$$426 \quad (5.8) \quad J_{k+1}^{(i)} = \begin{cases} J_k^{(i)} & \text{if } i \notin \mathcal{N}_k, \\ \nabla \phi_i(\mathbf{x}_{k+1}) & \text{if } i \in \mathcal{N}_k. \end{cases}$$

427 As in SVRG, the set $\{1, \dots, N\}$ is partitioned into a fixed-number n_b of random mini-
 428 batches which are used in order. One advantage of SAGA over SVRG is that it only
 429 requires a full gradient computation at the beginning of the algorithm, while SVRG
 430 requires a full gradient evaluation each n_b iterations.

431 *Remark 5.1.* By assuming that all the ϕ_i 's have Lipschitz-continuous gradients
 432 with Lipschitz constant \bar{L} , we have that the gradient estimates $\mathbf{g}_{\mathcal{N}_k}(\mathbf{x})$ and $\mathbf{g}_{\mathcal{N}_k}^{\text{SAGA}}(\mathbf{x})$
 433 are Lipschitz continuous with the same Lipschitz constant.

Our algorithmic framework for objective functions of the form (5.1) is called LSOS-FS (where FS stands for Finite Sums) and is outlined in Algorithm 5.1. For the sake of generality, we refer to generic gradient and Hessian approximations, denoted $\mathbf{g}(\mathbf{x}_k)$ and $B(\mathbf{x}_k)$, respectively. We consider the possibility of introducing inexactness in the computation of the direction

$$\mathbf{d}_k = -B(\mathbf{x}_k)^{-1}\mathbf{g}(\mathbf{x}_k),$$

434 even if for the L-BFGS strategy mentioned above, where $H_k = B(\mathbf{x}_k)^{-1}$, the direction
 435 can be computed exactly by a matrix-vector product with H_k .

Algorithm 5.1 LSOS for Finite Sums (LSOS-FS)

- 1: given $\mathbf{x}^0 \in \mathbb{R}^n$, $\eta, \beta \in (0, 1)$, $\{\delta_k\} \subset \mathbb{R}_+$ and $\{\zeta_k\} \subset \mathbb{R}_{++}$
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: compute $f_{\mathcal{N}_k}(\mathbf{x}_k)$, $\mathbf{g}(\mathbf{x}_k)$ and $B(\mathbf{x}_k)$
- 4: find a search direction \mathbf{d}_k such that

$$(5.9) \quad \|B(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k)\| \leq \delta_k \|\mathbf{g}(\mathbf{x}_k)\|$$

- 5: find the smallest integer $j \geq 0$ such that the step length $t_k = \beta^j$ satisfies

$$(5.10) \quad f_{\mathcal{N}_k}(\mathbf{x}_k + t_k \mathbf{d}_k) \leq f_{\mathcal{N}_k}(\mathbf{x}_k) + \eta t_k \mathbf{g}(\mathbf{x}_k)^\top \mathbf{d}_k + \zeta_k$$

- 6: set $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$
 - 7: **end for**
-

436 **6. Convergence theory of Algorithm LSOS-FS.** We assume

$$437 \quad (6.1) \quad \sum_k \zeta_k < \infty.$$

438 In the initial phase of the computation, nondescent directions are likely to occur;
 439 however, by requiring $\zeta_k > 0$ we ensure that the line search remains well defined.
 440 Furthermore, by (6.1) it is $\zeta_k \rightarrow 0$, which, by reasoning as in the proof of Theorem 4.2,
 441 implies that Algorithm 5.1 will eventually determine a descent direction for the current
 442 approximation of the objective function.

443 Algorithm LSOS-FS computes the step length t_k by applying a backtracking
 444 line-search to the approximate function $f_{\mathcal{N}_k}(\mathbf{x})$. In the next lemma we prove that
 445 the sequence $\{t_k\}$ is bounded away from zero for all k large enough, if the gradient
 446 approximation is the subsampled gradient $\mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)$. Throughout this section we use
 447 δ_{\max} defined at the beginning of page 10.

448 **LEMMA 6.1.** *Let Algorithm 5.1 be applied to problem (5.1) with $\mathbf{g}(\mathbf{x}_k) = \mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)$,*
 449 *and let $\delta_k \rightarrow 0$. Then the step-length sequence $\{t_k\}$ is such that*

$$450 \quad (6.2) \quad t_k \geq \frac{\beta(1-\eta)\mu^2}{L^2(1+\delta_{\max})^2} := t_{\min} \in (0, 1),$$

451 *for all k large enough.*

452 *Proof.* If $t_k = 1$, then (6.2) holds. If $t_k < 1$, then there exists $t'_k = t_k/\beta$ such that

$$453 \quad (6.3) \quad f_{\mathcal{N}_k}(\mathbf{x}_k + t'_k \mathbf{d}_k) > f_{\mathcal{N}_k}(\mathbf{x}_k) + \eta t'_k \mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)^\top \mathbf{d}_k.$$

454 Furthermore, by the descent lemma applied to $f_{\mathcal{N}_k}$ and the Lipschitz continuity of
455 $\mathbf{g}_{\mathcal{N}_k}$ we have

$$456 \quad (6.4) \quad f_{\mathcal{N}_k}(\mathbf{x}_k + t'_k \mathbf{d}_k) \leq f_{\mathcal{N}_k}(\mathbf{x}_k) + t'_k \mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)^\top \mathbf{d}_k + \frac{L}{2} (t'_k)^2 \|\mathbf{d}_k\|^2.$$

457 Combining (6.3) and (6.4) we obtain

$$458 \quad (6.5) \quad t_k = \beta t'_k > \frac{-2\beta(1-\eta) \mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)^\top \mathbf{d}_k}{L \|\mathbf{d}_k\|^2}.$$

459 Following the proof of Theorem 4.2, we can show that (4.3) holds for all $k \geq \bar{k}$ with
460 $\mathbf{g} = \mathbf{g}_{\mathcal{N}_k}$ and thus

$$461 \quad (6.6) \quad -\mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)^\top \mathbf{d}_k \geq \frac{\|\mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)\|^2}{2L}.$$

462 On the other hand,

$$463 \quad \|\mathbf{d}_k\| = \|(B_{\mathcal{N}_k}(\mathbf{x}_k))^{-1}(\mathbf{r}_k - \mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k))\| \leq \frac{1}{\mu} (\|\mathbf{r}_k\| + \|\mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)\|)$$

$$464 \quad \leq \frac{\delta_k + 1}{\mu} \|\mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)\|,$$

465 where the last inequality comes from (5.9). Therefore, since $\delta_k \leq \delta_{\max}$, we obtain

$$466 \quad \|\mathbf{d}_k\|^2 \leq \frac{(\delta_{\max} + 1)^2}{\mu^2} \|\mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)\|^2.$$

467 This, together with (6.5) and (6.6), gives the thesis. \square

468 In the following theorem we state the convergence of the LSOS-FS method. The
469 proof is omitted since it follows the steps of the proof of Theorem 4.2. The Lemma
470 above exploits $\mathbf{g}(\mathbf{x}_k) = \mathbf{g}_{\mathcal{N}_k}(\mathbf{x}_k)$; for general $\mathbf{g}(\mathbf{x}_k)$ we have to assume that the step
471 lengths are bounded away from zero. Notice that we do not need the assumption of
472 bounded iterates since the line search is performed at each iteration and the function
473 is strongly convex and thus bounded from below.

474 **THEOREM 6.2.** *Let $\{\mathbf{x}_k\}$ be the sequence generated by Algorithm 5.1 applied to*
475 *problem (5.1). Assume that (6.1) and Assumption 3.4 hold, and $\mathbf{g}(\mathbf{x})$ is a Lipschitz-*
476 *continuous unbiased gradient estimate. Moreover, assume $\{t_k\}$ is bounded away from*
477 *zero. Then $\{\mathbf{x}_k\}$ converges a.s. to the unique minimizer of ϕ .*

478 Finally, we provide the convergence rate analysis of LSOS-FS. We prove that
479 the expected function error converges R -linearly provided that ζ_k vanishes R -linearly.
480 We also prove that a Q -linear rate of convergence can be achieved if the monotone
481 (Armijo) line search is employed and the descent direction is ensured. The latter
482 condition can be provided by putting an upper bound on the forcing term, which is in
483 line with the classical (deterministic) analysis. The results are stated in the following
484 three theorems, whose proofs rely on the steps of the proof of Theorem 4.2. Since L is
485 an upper bound of the spectrum of the Hessian estimates, without loss of generality
486 we can assume $L \geq 1$.

487 **THEOREM 6.3.** *Let $\{\mathbf{x}_k\}$ be a sequence generated by [Algorithm 5.1](#) applied to prob-*
 488 *lem (5.1). Let $\delta_k \rightarrow 0$ and let $\zeta_k \rightarrow 0$ R-linearly. Let [Assumption 3.4](#) hold, $\mathbf{g}(\mathbf{x})$ be*
 489 *a Lipschitz-continuous unbiased gradient estimate and the sequence $\{t_k\}$ be bounded*
 490 *away from zero. Then there exist constants $\rho_1 \in (0, 1)$ and $C > 0$ such that*

$$491 \quad (6.7) \quad \mathbb{E}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \leq \rho_1^k C.$$

492 *Proof.* Let t_{\min} be a lower bound for the sequence $\{t_k\}$. Following the steps of
 493 the proof of [Theorem 4.2](#) we obtain (4.4) with $c_{10} = \eta t_{\min}/(2L)$, or equivalently,

$$494 \quad \phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*) \leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) - c_{10} \mathbb{E}(\|\mathbf{g}(\mathbf{x}_k)\|^2 | \mathcal{F}_k) + \zeta_k.$$

495 Moreover, using (4.5), (4.6) and the right-hand inequality in (2.2), we have

$$496 \quad \phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*) \leq \phi(\mathbf{x}_k) - \phi(\mathbf{x}_*) - \frac{2c_{10}}{L}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) + \zeta_k.$$

497 Applying the expectation we get

$$498 \quad (6.8) \quad \mathbb{E}(\phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*)) \leq \rho \mathbb{E}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) + \zeta_k,$$

499 where $\rho = 1 - 2c_{10}/L = 1 - \eta t_{\min}/L^2 \in (0, 1)$. Applying the induction argument we
 500 obtain

$$501 \quad \mathbb{E}(\phi(\mathbf{x}_j) - \phi(\mathbf{x}_*)) \leq \rho^j \mathbb{E}(\phi(\mathbf{x}_0) - \phi(\mathbf{x}_*)) + v_j,$$

502 where $v_j = \sum_{i=1}^{j-1} \rho^{i-1} \zeta_{j-i}$. The thesis follows by recalling that Lemma 4.2 from [23]
 503 implies $v_j \rightarrow 0$ R-linearly, with a factor $\rho_v = \frac{1}{2}(1 + \max\{\rho, \rho_\zeta\}) \in (0, 1)$, where
 504 $\rho_\zeta \in (0, 1)$ is an R-linear convergence factor of the sequence ζ_k . Finally, the statement
 505 holds with $\rho_1 = \max\{\rho, \rho_v\}$. \square

506 Notice that the condition $\delta_k \rightarrow 0$ can be relaxed with $0 < \delta_k \rightarrow \delta_{\min}$ where $\delta_{\min} <$
 507 $\mu/(2L)$. The reason is that, eventually, the inexact second-order direction becomes a
 508 descent direction if (6.6) holds for all k large enough. Under the same argument we
 509 can prove [Lemma 6.1](#) and the proof is essentially the same as for [Theorem 6.3](#). Thus,
 510 the R-linear convergence is attainable under the persistent inexactness in solving the
 511 Newton equation.

512 **THEOREM 6.4.** *Let $\{\mathbf{x}_k\}$ be a sequence generated by [Algorithm 5.1](#) applied to prob-*
 513 *lem (5.1). Assume that $\zeta_k \rightarrow 0$ R-linearly and $\delta_k \rightarrow \delta_{\min}$, where $\delta_{\min} < \mu/(2L)$.*
 514 *Moreover, let [Assumption 3.4](#) be satisfied, $\mathbf{g}(\mathbf{x})$ be a Lipschitz-continuous unbiased*
 515 *gradient estimate and $\{t_k\}$ be bounded away from zero. Then there exist $\rho_1 \in (0, 1)$*
 516 *and $C > 0$ such that (6.7) holds.*

517 An immediate consequence of the previous theorem is the following worst-case
 518 complexity result.

COROLLARY 6.5. *Let $\{\mathbf{x}_k\}$ be a sequence generated by [Algorithm 5.1](#) applied to*
problem (5.1). Assume that $\zeta_k \rightarrow 0$ R-linearly and $\delta_k \rightarrow \delta_{\min}$, where $\delta_{\min} < \mu/(2L)$.
Moreover, Let [Assumption 3.4](#) be satisfied, $\mathbf{g}(\mathbf{x})$ be a Lipschitz-continuous unbiased
*gradient estimate and $\{t_k\}$ be bounded away from zero. Then, to achieve $\mathbb{E}(\phi(\mathbf{x}_k) -$
 $\phi(\mathbf{x}_*)) \leq \varepsilon$ for some $\varepsilon \in (0, e^{-1})$, [Algorithm 5.1](#) takes at most*

$$k_{\max} = \left\lceil \frac{|\log(C)| + 1}{|\log(\rho_1)|} \log(\varepsilon^{-1}) \right\rceil,$$

519 where $\rho_1 \in (0, 1)$ and $C > 0$ satisfy (6.7).

Proof. Theorem 6.4 implies (6.7). Thus, $\mathbb{E}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \leq \varepsilon$ for all

$$k \geq \frac{\log(C) - \log(\varepsilon)}{-\log(\rho_1)}.$$

Now, using the fact that $\log(\varepsilon) < -1$ and $\log(\rho_1) < 0$ we can provide an upper bound to the right-hand side of the previous inequality as follows

$$\frac{\log(C) - \log(\varepsilon)}{-\log(\rho_1)} \leq \frac{|\log(C)|\log(\varepsilon^{-1}) + \log(\varepsilon^{-1})}{|\log(\rho_1)|} = \frac{|\log(C)| + 1}{|\log(\rho_1)|} \log(\varepsilon^{-1})$$

520 and the thesis holds. \square

521 In order to achieve a Q-linear rate of convergence, the standard Armijo line search
522 has to be used, i.e., $\zeta_k = 0$ has to be set in (5.10). Again, the forcing terms δ_k need
523 not vanish in order to achieve the desired rate (i.e., Newton's equation can be solved
524 inexactly), but it must be bounded above away from one. More in detail, it must be
525 $\delta_{\max} \leq \mu/(2L)$, as stated in the following theorem. A sequence $\{\delta_k\}$ satisfying the
526 requirement of the theorem can be defined as $\delta_k = \mu/(2L)$ for all k .

527 **THEOREM 6.6.** *Let $\{\mathbf{x}_k\}$ be a sequence generated by Algorithm 5.1 applied to prob-*
528 *lem (5.1). Assume that $\delta_{\max} \leq \mu/(2L)$ and $\zeta_k = 0$ for all k . Moreover, suppose that*
529 *Assumption 3.4 is satisfied, $\mathbf{g}(\mathbf{x})$ is a Lipschitz-continuous unbiased gradient estimate*
530 *and the sequence $\{t_k\}$ is bounded away from zero. Then there exists $\rho_2 \in (0, 1)$ such*
531 *that for all k*

532 (6.9)
$$\mathbb{E}(\phi(\mathbf{x}_{k+1}) - \phi(\mathbf{x}_*)) \leq \rho_2 \mathbb{E}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)).$$

Proof. Notice that the Lipschitz continuity of the gradient estimate implies that (6.6) holds for every k since $\delta_k \leq \mu/(2L)$. Let t_{\min} be a lower bound for the sequence $\{t_k\}$. By following the steps of the proof of Theorem 6.3, we have that (6.8) holds with $\zeta_k = 0$. Therefore, by setting

$$\rho_2 = \rho = 1 - \frac{\eta t_{\min}}{L^2} \leq 1 - \frac{\eta(1-\eta)\beta\mu^2}{L^2(2L+\mu)}$$

533 the thesis holds. \square

534 Since Theorem 6.6 implies $\mathbb{E}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_*)) \leq \rho_2^k(\phi(\mathbf{x}_0) - \phi(\mathbf{x}_*))$, following the
535 same reasoning as in Corollary 6.5, we obtain the following complexity result.

COROLLARY 6.7. *Let $\{\mathbf{x}_k\}$ be a sequence generated by Algorithm 5.1 applied to*
problem (5.1). Assume that $\delta_{\max} \leq \mu/(2L)$ and $\zeta_k = 0$ for all k . Moreover, sup-
pose that Assumption 3.4 is satisfied, $\mathbf{g}(\mathbf{x})$ is a Lipschitz-continuous unbiased gra-
dent estimate and the sequence $\{t_k\}$ is bounded away from zero. Then, in order
to achieve $\mathbb{E}(\phi(\mathbf{x}_k) - \phi(\mathbf{x}_)) \leq \varepsilon$ for some $\varepsilon \in (0, e^{-1})$, LSOS-FS takes at most*
 $k_{\max} = \mathcal{O}(\log(\varepsilon^{-1}))$ iterations. More precisely,

$$k_{\max} = \left\lceil \frac{|\log(\phi(\mathbf{x}_0) - \phi(\mathbf{x}_*))| + 1}{|\log(\rho_2)|} \log(\varepsilon^{-1}) \right\rceil,$$

536 where ρ_2 satisfies (6.9).

537 **7. Numerical experiments.** We developed MATLAB implementations of the
 538 algorithms discussed in the previous sections and tested them on two sets of stochastic
 539 optimization problems. The first set consists of general convex problems with the
 540 addition of random noise in the evaluation of the objective function and its derivatives.
 541 On these problems we tested Algorithms SOS and LSOS discussed in sections 2 to 4.
 542 The second set consists of finite-sum problems arising in training linear classifiers
 543 with regularized logistic regression models. On these problems we tested a specialized
 544 version of LSOS-FS. All the tests were run with MATLAB R2019b on a server available
 545 at the University of Campania “L. Vanvitelli”, equipped with 8 Intel Xeon Platinum
 546 8168 CPUs, 1536 GB of RAM and Linux CentOS 7.5 operating system.

547 **7.1. Convex random problems.** The first set of test problems was defined by
 548 setting

$$549 \quad (7.1) \quad \phi(\mathbf{x}) = \sum_{i=1}^n \lambda_i (e^{x_i} - x_i) + (\mathbf{x} - \mathbf{e})^\top A(\mathbf{x} - \mathbf{e}),$$

550 where, given a scalar $\kappa \gg 1$, the coefficients λ_i are logarithmically spaced between 1
 551 and κ , $A \in \mathbb{R}^{n \times n}$ is symmetric positive definite with eigenvalues λ_i , and $\mathbf{e} \in \mathbb{R}^n$ has all
 552 entries equal to 1. Changing the values of n and κ allows us to have strongly convex
 553 problems with variable size and conditioning. In order to obtain unbiased estimates
 554 of ϕ and its gradient, we considered $\varepsilon_f(\mathbf{x}) \sim \mathcal{N}(0, \sigma)$ and $(\varepsilon_g(\mathbf{x}))_i \sim \mathcal{N}(0, \sigma)$ for all i ,
 555 where $\mathcal{N}(0, \sigma)$ is the normal distribution with mean 0 and standard deviation σ . We
 556 considered $\sigma \in (0, 1]$. Since the Hessian estimate can be biased, we set it equal to the
 557 diagonal matrix $\varepsilon_B(\mathbf{x}) = \text{diag}(\mu_1, \dots, \mu_n)$, where $\mu_j \sim \mathcal{N}(0, \sigma)$ for all j .

558 In applying Algorithm 4.1 to this set of problems, we introduced a small modi-
 559 fication in the switching criterion at line 7 of the algorithm, by deactivating the line
 560 search whenever $t_k \|\mathbf{d}_k\| < t_{\min}$ instead of deactivating it when $t_k < t_{\min}$.

We first ran Algorithm LSOS with exact solution of the noisy Newton systems,
 i.e., $\delta_k = 0$ in (4.1). The parameters were set as $n = 10^3, \kappa = 10^2, 10^3, 10^4$,
 $\sigma = 0.1\% \kappa, 0.5\% \kappa, 1\% \kappa$, and A was generated by using the MATLAB `sprandsym`
 function with density 0.5 and eigenvalues $\lambda_1, \dots, \lambda_n$. It was verified experimentally
 that the condition number of the Hessian of ϕ is close to κ at the solution. This
 solution was computed with high accuracy by using the (deterministic) L-BFGS im-
 plementation by Mark Schmidt, available from <https://www.cs.ubc.ca/~schmidtm/Software/minFunc.html>. The starting point was set as a random vector with distri-
 bution of entries $\mathcal{N}(0, 5)$. The noisy Newton systems were solved by the MATLAB
`backslash` operator. The parameter used to switch between the line search and the
 pre-defined gain sequence was set as $t_{\min} = 10^{-3}$. The gain sequence $\{\alpha_k\}$ used after
 the deactivation of the line search was defined as

$$\alpha_k = \alpha_{k_\tau} \frac{T}{T + k - k_\tau} \quad \text{for all } k > k_\tau,$$

561 where k_τ is the first iteration such that $t_{k_\tau} \|\mathbf{d}_{k_\tau}\| < t_{\min}$, $\alpha_{k_\tau} = t_{\min} / \|\mathbf{d}_{k_\tau}\|$ and
 562 $T = 10^6$. In the nonmonotone line search we set $\eta = 10^{-4}$ and $\zeta_k = \vartheta^k$ for all k ,
 563 where $\vartheta = 0.9$.

564 LSOS was compared with the following algorithms:

- 565 • SOS (Algorithm 2.1) with exact solution of the noisy Newton systems and
 566 gain sequence defined as

$$567 \quad (7.2) \quad \alpha_k = \frac{1}{\|\mathbf{d}^0\|} \frac{T}{T + k}$$

568 • Stochastic Gradient Descent with step length (7.2), referred to as SGD.
 569 For both SOS and SGD the choice of the starting point was the same as for LSOS.
 570 The comparison was performed in terms of the absolute error of the objective
 571 function value (with respect to the optimal value computed by the deterministic L-
 572 BFGS algorithm) versus the execution time. We ran each algorithm 20 times on each
 573 problem and computed the average error and the average execution time spent until
 574 each iteration k . The results are shown in Figure 1, where each error line is plotted
 575 with its 95% confidence interval (which does not appear in the pictures because its
 576 size is negligible). The time interval on the x axis is the average time required by
 577 LSOS to perform 50 iterations.

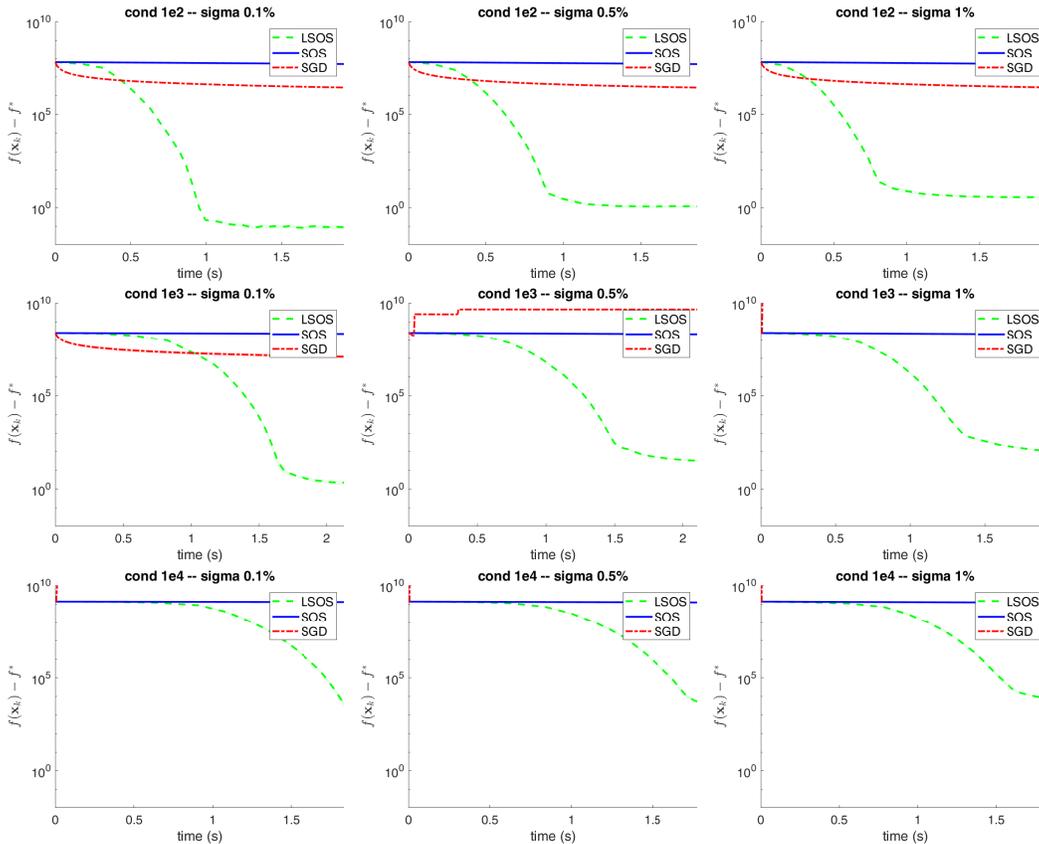


FIG. 1. Test set 1: comparison of LSOS, SOS and SGD. The condition number increases from top to bottom, the noise increases from left to right.

578 The figure shows that the introduction of the line search yields much better
 579 exploitation of the second-order directions, thus enabling the method to approach the
 580 solution faster. The line search also allows us to overcome typical problems associated
 581 with the choice of a pre-defined gain sequence, which may strongly affect the speed
 582 of the algorithm and possibly lead to divergence in practice.

583 We also investigated the effect of the inexactness in the solution of the noisy
 584 Newton systems. To this aim, we considered problems of the form (7.1) with size
 585 $n = 2 \cdot 10^4$, where, following [14], the symmetric positive definite matrix A was

586 defined as

$$587 \quad A = V D V^T.$$

Here D is a diagonal matrix with diagonal entries $\lambda_1, \dots, \lambda_n$ and

$$V = (I - 2 \mathbf{v}_3 \mathbf{v}_3^T)(I - 2 \mathbf{v}_2 \mathbf{v}_2^T)(I - 2 \mathbf{v}_1 \mathbf{v}_1^T),$$

588 with \mathbf{v}_j random vectors of unit norm. Since for these problems the Hessian is available
 589 in factorized form, we solved the noisy Newton systems with the Conjugate Gradient
 590 (CG) method implemented in the MATLAB `pcg` function, exploiting the factorization
 591 to compute matrix-vector products. In this case, we compared three versions of
 592 Algorithm LSOS:

- 593 • LSOS with $\delta_k = 0$ in (4.1);
- 594 • LSOS with $\delta_k = \varrho^k$ and $\varrho = 0.95$, referred to as LSOS-I (where I denotes the
 595 inexact solution of the Newton systems according to (4.1));
- 596 • a line-search version of the SGD algorithm (corresponding to LSOS with
 597 $\mathbf{d}_k = -\mathbf{g}(\mathbf{x}_k)$), referred to as SGD-LS.

598 The CG method in LSOS and LSOS-I was run until the residual norm of the Newton
 599 system had been reduced by $\max(\delta_k, 10^{-6})$ with respect to the initial residual norm.

600 In Figure 2 we report the results obtained with the three algorithms, in terms of
 601 average error on the objective function versus average execution time over 20 runs,
 602 with 95% confidence intervals (not visible, as in the previous tests). In this case the
 603 time interval on the x axis is the average time required by LSOS-I to perform 250
 604 iterations. The plots clearly show that LSOS-I outperforms the other methods.

605 **7.2. Binary classification problems.** The second set of test problems models
 606 the training a linear classifier by minimization of the ℓ_2 -regularized logistic regres-
 607 sion. Given N pairs (\mathbf{a}_i, b_i) , where $\mathbf{a}_i \in \mathbb{R}^n$ is a training point and $b_i \in \{-1, 1\}$ the
 608 corresponding class label, an unbiased hyperplane approximately separating the two
 609 classes can be found by minimizing the function

$$610 \quad (7.3) \quad \phi(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \phi_i(\mathbf{x}),$$

611 where

$$612 \quad \phi_i(\mathbf{x}) = \log \left(1 + e^{-b_i \mathbf{a}_i^\top \mathbf{x}} \right) + \frac{\mu}{2} \|\mathbf{x}\|^2$$

613 and $\mu > 0$. By setting $z_i(\mathbf{x}) = 1 + e^{-b_i \mathbf{a}_i^\top \mathbf{x}}$, the gradient and the Hessian of ϕ_i are

$$614 \quad \nabla \phi_i(\mathbf{x}) = \frac{1 - z_i(\mathbf{x})}{z_i(\mathbf{x})} b_i \mathbf{a}_i + \mu \mathbf{x} \quad \text{and} \quad \nabla^2 \phi_i(\mathbf{x}) = \frac{z_i(\mathbf{x}) - 1}{z_i^2(\mathbf{x})} \mathbf{a}_i \mathbf{a}_i^\top + \mu I.$$

615 From $\frac{z_i(\mathbf{x}) - 1}{z_i^2(\mathbf{x})} \in (0, 1)$ it follows that ϕ_i is μ -strongly convex and

$$616 \quad \mu I \preceq \nabla^2 \phi_i(\mathbf{x}) \preceq LI, \quad L = \mu + \max_{i=1, \dots, N} \|\mathbf{a}_i\|^2.$$

617 We applied the L-BFGS version of Algorithm LSOS-FS described in section 5,
 618 which is sketched in Algorithm 7.1.

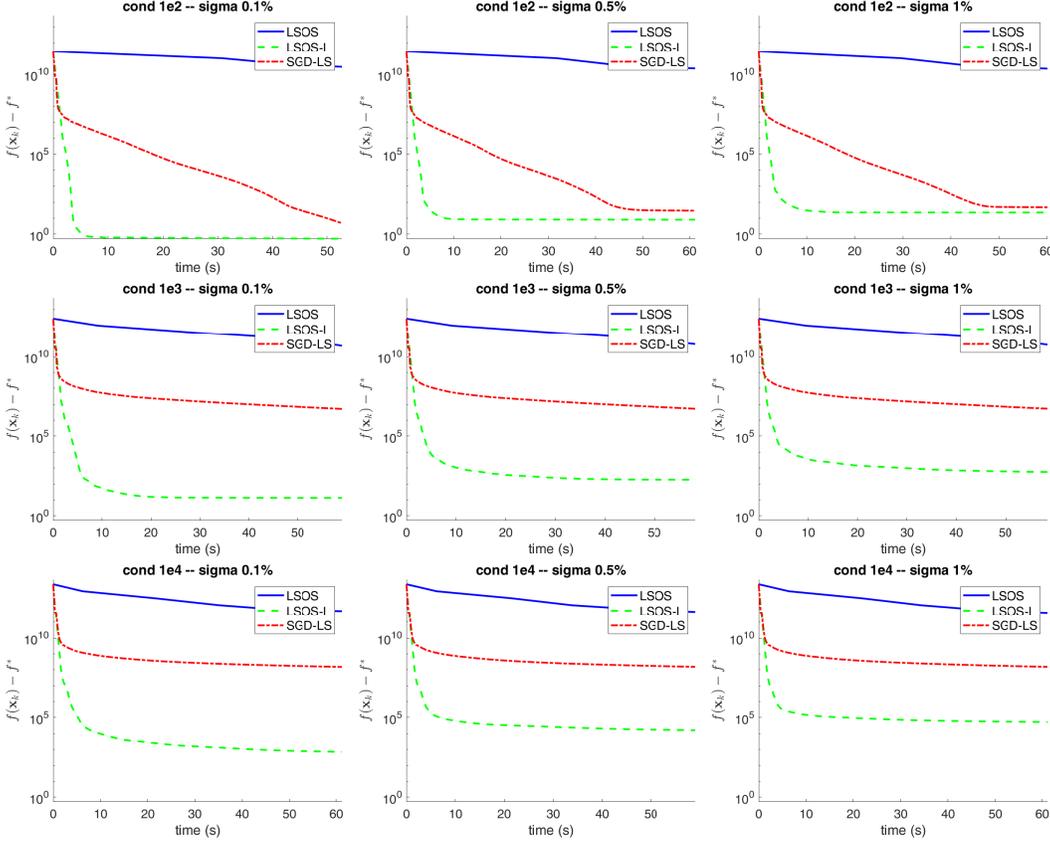


FIG. 2. Test set 2: comparison of LSOS, LSOS-I and SGD-LS. The condition number increases from top to bottom, the noise increases from left to right.

Algorithm 7.1 LSOS-BFGS

- 1: given $\mathbf{x}^0 \in \mathbb{R}^n$, $m, l \in \mathbb{N}$, $\eta, \vartheta \in (0, 1)$
 - 2: **for** $k = 0, 1, 2, \dots$ **do**
 - 3: compute a partition $\{\mathcal{K}_0, \mathcal{K}_1, \dots, \mathcal{K}_{n_b-1}\}$ of $\{1, \dots, N\}$
 - 4: **for** $r = 0, \dots, n_b - 1$ **do**
 - 5: choose $\mathcal{N}_k = \mathcal{K}_r$ and compute $\mathbf{g}(\mathbf{x}_k) = \mathbf{g}_{\mathcal{N}_k}^{\text{SAGA}}(\mathbf{x}_k)$ as in (5.7)-(5.8)
 - 6: compute $\mathbf{d}_k = -H_k \mathbf{g}(\mathbf{x}_k)$ with H_k defined in (5.5)-(5.6)
 - 7: find a step length t_k satisfying

$$f_{\mathcal{N}_k}(\mathbf{x}_k + t_k \mathbf{d}_k) \leq f_{\mathcal{N}_k}(\mathbf{x}_k) + \eta t_k \mathbf{g}(\mathbf{x}_k)^\top \mathbf{d}_k + \vartheta^k$$
 - 8: set $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$;
 - 9: **if** $\text{mod}(k, l) = 0$ and $k \geq 2l$ **then**
 - 10: update the L-BFGS correction pairs by using (5.3)-(5.4)
 - 11: **end if**
 - 12: **end for**
 - 13: **end for**
-

619 To test the effectiveness of LSOS-BFGS we considered six binary classification
 620 datasets from the LIBSVM collection available from [https://www.csie.ntu.edu.tw/
 621 ~cjlin/libsvmtools/datasets/](https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/), which we list in Table 1.

TABLE 1

Datasets from LIBSVM. For each dataset the number of training points and the number of features (space dimension) are reported; the datasets are sorted by the increasing number of features. Whenever a training set was not specified in LIBSVM, we selected it by using the MATLAB `crossvalind` function so that it contained 70% of the available data.

name	N	n
covtype	406709	54
w8a	49749	300
epsilon	400000	2000
gisette	6000	5000
real-sim	50617	20958
rcv1	20242	47236

We compared [Algorithm 7.1](#) with the stochastic L-BFGS algorithms proposed in [17] and [30] (referred to as GGR and MNJ, respectively), both using a constant step length selected by means of a grid search over the set $\{1, 5 \cdot 10^{-1}, 10^{-1}, 5 \cdot 10^{-2}, 10^{-2}, \dots, 5 \cdot 10^{-5}, 10^{-5}\}$, and with a mini-batch variant of the SAGA algorithm equipped with the same line search used in LSOS-BFGS. The implementations of GGR and MNJ were taken from the MATLAB StochBFGS code available from https://perso.telecom-paristech.fr/rgower/software/StochBFGS_dist-0.0.zip. In [Algorithm 7.1](#) we set $\vartheta = 0.999$ and started the line searches from a value t_{ini} selected by means of a grid search over $\{1, 5 \cdot 10^{-1}, 10^{-1}, 5 \cdot 10^{-2}, 10^{-2}, \dots, 5 \cdot 10^{-5}, 10^{-5}\}$. In particular, we set $t_{\text{ini}} = 5 \cdot 10^{-3}$ for `epsilon`, $t_{\text{ini}} = 5 \cdot 10^{-2}$ for `covtype` and `w8a`, and $t_{\text{ini}} = 1 \cdot 10^{-2}$ for `gisette`, `rcv1` and `real-sim`. We adopted the same strategy as the line-search version of SAGA used for the comparison, setting $t_{\text{ini}} = 5 \cdot 10^{-1}$ for `epsilon` and $t_{\text{ini}} = 1$ for the other datasets. Furthermore, we set $m = 10$ and $l = 5$. Since the first L-BFGS update pair is available after the first $2l = 10$ iterations, following [10] we take $\mathbf{d}_k = -\mathbf{g}(\mathbf{x}_k)$ for the first 10 iterations. The same values of m and l were used in the MNJ algorithm proposed in [30]. For GGR, following the indications coming from the results in [17], we set $m = 5$ and used the sketching based on the previous directions (indicated as `prev` in [17]), with sketch size $l = \lceil \sqrt[3]{n} \rceil$. We chose the sample size equal to $\lceil \sqrt{N} \rceil$ and the regularization parameter $\mu = 1/N$, as in the experiments reported in [17]. We decided to stop the algorithms when a maximum execution time was reached, i.e., 60 seconds for `covtype`, `w8a` and `gisette`, and 300 seconds for `epsilon`, `real-sim` and `rcv1`.

Figure 3 shows a comparison among the four algorithms in terms of the average absolute error of the objective function (with respect to the optimal value computed with the L-BFGS code by Mark Schmidt) versus the average execution time. As in the previous experiments, the error and the times were averaged over 20 runs and the plots show their 95% confidence interval (shaded lines, when visible). For all the algorithms, the grid search for defining or initializing the step lengths was performed on the first of the 20 runs and then fixed for the remaining 19 runs.

The results show that LSOS-BFGS algorithm outperforms the other stochastic L-BFGS algorithms on `w8a` and `gisette`, and outperforms GGR on `real-sim` and `rcv1`. It is worth noting that for `covtype` and `rcv1` the error for GGR tends to increase after a certain iteration, while the other algorithms seem to keep a much less “swinging” decrease. Furthermore, LSOS-BFGS seems to have a less oscillatory behavior with respect to GGR and MNJ. We conjecture that this behavior is due to the use of the line-search strategy. Since, in general, stopping criteria on this type of problems rely on the number of iterations, the number of epochs or the computational time, we

659 believe that a smoother behaviour could be associated with more consistent results if
 660 one decides to stop the execution in advance (see, e.g., the behavior of MNJ on `epsilon`).
 661 Finally, we observe that LSOS-BFGS is more efficient than the line-search-based mini-
 662 batch SAGA on all the problems, showing that the introduction of stochastic second-
 663 order information is crucial for the performance of the algorithm.

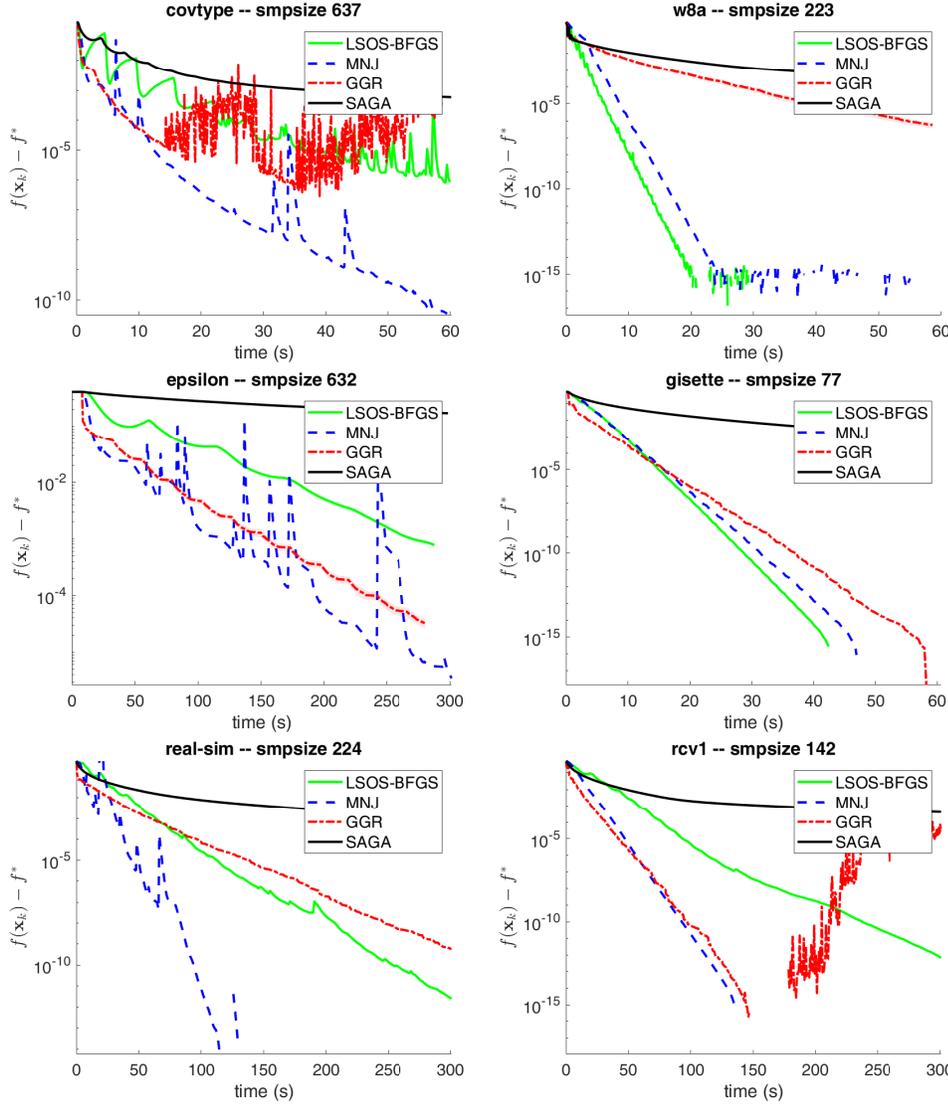


FIG. 3. Binary classification problems: comparison of LSOS-BFGS, MNJ, GGR and SAGA.

664 **8. Conclusions.** The proposed LSOS framework includes a variety of second-
 665 order stochastic optimization algorithms, using Newton, inexact Newton and, for
 666 finite-sum problems, limited-memory quasi-Newton directions. Almost sure conver-
 667 gence of the sequences generated by all the LSOS variants has been proved. For
 668 finite-sum problems, R-linear and Q-linear convergence rates of the expected objective

669 function error have been proved for stochastic L-BFGS Hessian approximations and
 670 any Lipschitz-continuous unbiased gradient estimates. In this case, an $\mathcal{O}(\log(\varepsilon^{-1}))$
 671 complexity bound has been also provided.

672 Numerical experiments have confirmed that line-search techniques in second-order
 673 stochastic methods yield a significant improvement over predefined step-length se-
 674 quences. Furthermore, in the case of finite-sum problems, the experiments have
 675 shown that combining stochastic L-BFGS Hessian approximations with the SAGA
 676 variance reduction technique and with line searches produces methods that are highly
 677 competitive with state-of-the art second-order stochastic optimization methods.

678 A challenging future research agenda includes the extension of (some) of these
 679 results to problems that do not satisfy the strong convexity assumption, as well as
 680 extensions to constrained stochastic problems.

681

REFERENCES

- 682 [1] S. BELLAVIA, G. GURIOLI, AND B. MORINI, *Adaptive cubic regularization methods with dynamic*
 683 *inexact Hessian information and applications to finite-sum minimization*, IMA Journal of
 684 Numerical Analysis, (2020), <https://doi.org/10.1093/imanum/drz076>.
- 685 [2] S. BELLAVIA, N. KREJIĆ, AND N. KRKLEC JERINKIĆ, *Subsampled inexact Newton methods for*
 686 *minimizing large sums of convex functions*, IMA Journal of Numerical Analysis, (2019),
 687 <https://doi.org/10.1093/imanum/drz027>.
- 688 [3] A. BENVENISTE, M. MÉTIVIER, AND P. PRIOURET, *Adaptive algorithms and stochastic approx-*
 689 *imations*, vol. 22 of Applications of Mathematics (New York), Springer-Verlag, Berlin,
 690 1990, <https://doi.org/10.1007/978-3-642-75894-2>. Translated from the French by Stephen
 691 S. Wilson.
- 692 [4] D. P. BERTSEKAS, *Nonlinear programming*, Athena Scientific Optimization and Computation
 693 Series, Athena Scientific, Belmont, MA, second ed., 1999.
- 694 [5] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Gradient convergence in gradient methods with errors*,
 695 SIAM J. Optim., 10 (2000), pp. 627–642, <https://doi.org/10.1137/S1052623497331063>.
- 696 [6] R. BOLLAPRAGADA, R. H. BYRD, AND J. NOCEDAL, *Exact and inexact subsampled Newton*
 697 *methods for optimization*, IMA J. Numer. Anal., 39 (2019), pp. 545–578, [https://doi.org/](https://doi.org/10.1093/imanum/dry009)
 698 [10.1093/imanum/dry009](https://doi.org/10.1093/imanum/dry009).
- 699 [7] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, *Optimization methods for large-scale machine*
 700 *learning*, SIAM Rev., 60 (2018), pp. 223–311, <https://doi.org/10.1137/16M1080173>.
- 701 [8] R. H. BYRD, G. M. CHIN, W. NEVEITT, AND J. NOCEDAL, *On the use of stochastic Hessian*
 702 *information in optimization methods for machine learning*, SIAM J. Optim., 21 (2011),
 703 pp. 977–995, <https://doi.org/10.1137/10079923X>.
- 704 [9] R. H. BYRD, G. M. CHIN, J. NOCEDAL, AND Y. WU, *Sample size selection in optimization*
 705 *methods for machine learning*, Math. Program., 134 (2012), pp. 127–155, [https://doi.org/](https://doi.org/10.1007/s10107-012-0572-5)
 706 [10.1007/s10107-012-0572-5](https://doi.org/10.1007/s10107-012-0572-5).
- 707 [10] R. H. BYRD, S. L. HANSEN, J. NOCEDAL, AND Y. SINGER, *A stochastic Quasi-Newton method*
 708 *for large-scale optimization*, SIAM J. Optim., 26 (2016), pp. 1008–1031, [https://doi.org/](https://doi.org/10.1137/140954362)
 709 [10.1137/140954362](https://doi.org/10.1137/140954362).
- 710 [11] P. J. CARRINGTON, J. SCOTT, AND S. WASSERMAN, eds., *Models and Methods in Social Network*
 711 *Analysis*, Structural Analysis in the Social Sciences, Cambridge University Press, 2005,
 712 <https://doi.org/10.1017/CBO9780511811395>.
- 713 [12] A. DEFAZIO, F. BACH, AND S. LACOSTE-JULIEN, *SAGA: A fast incremental gradient method*
 714 *with support for non-strongly convex composite objectives*, in Advances in Neural Informa-
 715 tion Processing Systems 27, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and
 716 K. Q. Weinberger, eds., Curran Associates, Inc., 2014, pp. 1646–1654, [https://dl.acm.org/](https://dl.acm.org/doi/10.5555/2968826.2969010)
 717 [doi/10.5555/2968826.2969010](https://dl.acm.org/doi/10.5555/2968826.2969010).
- 718 [13] B. DELYON AND A. JUDITSKY, *Accelerated stochastic approximation*, SIAM J. Optim., 3 (1993),
 719 pp. 868–881, <https://doi.org/10.1137/0803045>.
- 720 [14] D. DI SERAFINO, G. TORALDO, M. VIOLA, AND J. BARLOW, *A two-phase gradient method for*
 721 *quadratic programming problems with a single linear constraint and bounds on the vari-*
 722 *ables*, SIAM J. Optim., 28 (2018), pp. 2809–2838, <https://doi.org/10.1137/17M1128538>.
- 723 [15] S. C. EISENSTAT AND H. F. WALKER, *Globally convergent inexact Newton methods*, SIAM J.
 724 Optim., 4 (1994), pp. 393–422, <https://doi.org/10.1137/0804022>.

- 725 [16] M. C. FU, *Chapter 19 Gradient Estimation*, in Simulation, S. G. Henderson and B. L. Nelson,
726 eds., vol. 13 of Handbooks in Operations Research and Management Science, Elsevier,
727 2006, pp. 575–616, [https://doi.org/10.1016/S0927-0507\(06\)13019-4](https://doi.org/10.1016/S0927-0507(06)13019-4).
- 728 [17] R. M. GOWER, D. GOLDFARB, AND P. RICHTÁRIK, *Stochastic block BFGS: Squeezing more*
729 *curvature out of data*, in Proceedings of the 33rd International Conference on International
730 Conference on Machine Learning - Volume 48, JMLR.org, 2016, pp. 1869–1878, <https://dl.acm.org/doi/10.5555/3045390.3045588>.
- 731 [18] R. M. GOWER, P. RICHTÁRIK, AND F. BACH, *Stochastic quasi-gradient methods: variance*
732 *reduction via Jacobian sketching*, Mathematical Programming, (2020), [https://doi.org/10.](https://doi.org/10.1007/s10107-020-01506-0)
733 [1007/s10107-020-01506-0](https://doi.org/10.1007/s10107-020-01506-0).
- 734 [19] A. N. IUSEM, A. JOFRÉ, R. I. OLIVEIRA, AND P. THOMPSON, *Variance-based extragradient*
735 *methods with line search for stochastic variational inequalities*, SIAM J. Optim., 29 (2019),
736 pp. 175–206, <https://doi.org/10.1137/17M1144799>, <https://doi.org/10.1137/17M1144799>.
- 737 [20] R. JOHNSON AND T. ZHANG, *Accelerating stochastic gradient descent using predictive variance*
738 *reduction*, in Advances in Neural Information Processing Systems 26, C. J. C. Burges,
739 L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, eds., Curran Associates,
740 Inc., 2013, pp. 315–323, <https://dl.acm.org/doi/10.5555/2999611.2999647>.
- 741 [21] H. KESTEN, *Accelerated stochastic approximation*, Ann. Math. Statist., 29 (1958), pp. 41–59,
742 <https://doi.org/10.1214/aoms/1177706705>.
- 743 [22] A. KLENKE, *Probability theory*, Universitext, Springer, London, second ed., 2014, [https://doi.](https://doi.org/10.1007/978-1-4471-5361-0)
744 [org/10.1007/978-1-4471-5361-0](https://doi.org/10.1007/978-1-4471-5361-0). A comprehensive course.
- 745 [23] N. KREJIĆ AND N. KRKLEC JERINKIĆ, *Nonmonotone line search methods with variable*
746 *sample size*, Numer. Algorithms, 68 (2015), pp. 711–739, [https://doi.org/10.1007/](https://doi.org/10.1007/s11075-014-9869-1)
747 [s11075-014-9869-1](https://doi.org/10.1007/s11075-014-9869-1).
- 748 [24] N. KREJIĆ, Z. LUŽANIN, Z. OVCIN, AND I. STOJKOVSKA, *Descent direction method with line*
749 *search for unconstrained optimization in noisy environment*, Optim. Methods Softw., 30
750 (2015), pp. 1164–1184, <https://doi.org/10.1080/10556788.2015.1025403>.
- 751 [25] N. KREJIĆ, Z. LUŽANIN, AND I. STOJKOVSKA, *A gradient method for unconstrained optimization*
752 *in noisy environment*, Appl. Numer. Math., 70 (2013), pp. 1–21, [https://doi.org/10.1016/](https://doi.org/10.1016/j.apnum.2013.02.006)
753 [j.apnum.2013.02.006](https://doi.org/10.1016/j.apnum.2013.02.006).
- 754 [26] K. MARTI, *Stochastic optimization methods*, Springer, Heidelberg, third ed., 2015, [https://doi.](https://doi.org/10.1007/978-3-662-46214-0)
755 [org/10.1007/978-3-662-46214-0](https://doi.org/10.1007/978-3-662-46214-0). Applications in engineering and operations research.
- 756 [27] L. MARTÍNEZ, R. ANDRADE, E. G. BIRGIN, AND J. M. MARTÍNEZ, *Packmol: A package for build-*
757 *ing initial configurations for molecular dynamics simulations*, Journal of Computational
758 Chemistry, 30 (2009), pp. 2157–2164, <https://doi.org/10.1002/jcc.21224>.
- 759 [28] A. MOKHTARI AND A. RIBEIRO, *RES: regularized stochastic BFGS algorithm*, IEEE Trans.
760 Signal Process., 62 (2014), pp. 6089–6104, <https://doi.org/10.1109/TSP.2014.2357775>.
- 761 [29] A. MOKHTARI AND A. RIBEIRO, *Global convergence of online limited memory BFGS*, J. Mach.
762 Learn. Res., 16 (2015), pp. 3151–3181.
- 763 [30] P. MORITZ, R. NISHIHARA, AND M. JORDAN, *A linearly-convergent stochastic L-BFGS algo-*
764 *rithm*, in Proceedings of the 19th International Conference on Artificial Intelligence and
765 Statistics, A. Gretton and C. C. Robert, eds., vol. 51 of Proceedings of Machine Learning
766 Research, Cadiz, Spain, 09–11 May 2016, PMLR, pp. 249–258, [http://proceedings.mlr.](http://proceedings.mlr.press/v51/moritz16.html)
767 [press/v51/moritz16.html](http://proceedings.mlr.press/v51/moritz16.html).
- 768 [31] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, Ann. Math. Statistics, 22
769 (1951), pp. 400–407, <https://doi.org/10.1214/aoms/1177729586>.
- 770 [32] F. ROOSTA-KHORASANI AND M. W. MAHONEY, *Sub-sampled Newton methods*, Math. Program.,
771 174 (2019), pp. 293–326, <https://doi.org/10.1007/s10107-018-1346-5>.
- 772 [33] D. RUPPERT, *A Newton-Raphson version of the multivariate Robbins-Monro procedure*, Ann.
773 Statist., 13 (1985), pp. 236–245, <https://doi.org/10.1214/aos/1176346589>.
- 774 [34] J. C. SPALL, *A second order stochastic approximation algorithm using only function measure-*
775 *ments*, in Proceedings of 1994 33rd IEEE Conference on Decision and Control, vol. 3, 1994,
776 pp. 2472–2477.
- 777 [35] J. C. SPALL, *Stochastic version of second-order (Newton-Raphson) optimization using only*
778 *function measurements*, in Proceedings of the 27th Conference on Winter Simulation, WSC
779 '95, USA, 1995, IEEE Computer Society, p. 347–352, [https://doi.org/10.1145/224401.](https://doi.org/10.1145/224401.224633)
780 [224633](https://doi.org/10.1145/224401.224633).
- 781 [36] J. C. SPALL, *Accelerated second-order stochastic optimization using only function measure-*
782 *ments*, in Proceedings of the 36th IEEE Conference on Decision and Control, vol. 2, 1997,
783 pp. 1417–1424.
- 784 [37] J. C. SPALL, *Introduction to stochastic search and optimization*, Wiley-Interscience Series in
785 Discrete Mathematics and Optimization, Wiley-Interscience [John Wiley & Sons], Hobo-
786

- 787 ken, NJ, 2003, <https://doi.org/10.1002/0471722138>. Estimation, simulation, and control.
- 788 [38] D. VICARI, A. OKADA, G. RAGOZINI, AND C. WEIHS, eds., *Analysis and Modeling of Complex*
789 *Data in Behavioral and Social Sciences*, Springer, Cham, 2014, [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-319-06692-9)
790 [978-3-319-06692-9](https://doi.org/10.1007/978-3-319-06692-9).
- 791 [39] Z. XU AND Y.-H. DAI, *New stochastic approximation algorithms with adaptive step sizes*, *Optim. Lett.*, 6 (2012), pp. 1831–1846, <https://doi.org/10.1007/s11590-011-0380-5>.
- 792
- 793 [40] F. YOUSEFIAN, A. NEDIĆ, AND U. V. SHANBHAG, *On stochastic gradient and subgradient*
794 *methods with adaptive steplength sequences*, *Automatica J. IFAC*, 48 (2012), pp. 56–67,
795 <https://doi.org/10.1016/j.automatica.2011.09.043>.