# Stochastic Variance-Reduced Prox-Linear Algorithms for Nonconvex Composite Optimization

Junyu Zhang[*]        Lin Xiao[†]

May 12, 2021

### Abstract

We consider the problem of minimizing composite functions of the form $f(g(x)) + h(x)$, where $f$ and $h$ are convex functions (which can be nonsmooth) and $g$ is a smooth vector mapping. In addition, we assume that $g$ is the average of finite number of component mappings or the expectation over a family of random component mappings. We propose a class of stochastic variance-reduced prox-linear algorithms for solving such problems and bound their sample complexities for finding an $\epsilon$-stationary point in terms of the total number of evaluations of the component mappings and their Jacobians. When $g$ is a finite average of $N$ components, we obtain sample complexity $\mathcal{O}(N + N^{4/5}\epsilon^{-1})$ for both mapping and Jacobian evaluations. When $g$ is a general expectation, we obtain sample complexities of $\mathcal{O}(\epsilon^{-5/2})$ and $\mathcal{O}(\epsilon^{-3/2})$ for component mappings and their Jacobians respectively. If in addition $f$ is smooth, then improved sample complexities of $\mathcal{O}(N + N^{1/2}\epsilon^{-1})$ and $\mathcal{O}(\epsilon^{-3/2})$ are derived for $g$ being a finite average and a general expectation respectively, for both component mapping and Jacobian evaluations.

**Keywords:** stochastic composite optimization, nonsmooth optimization, variance reduction, proximal mapping, prox-linear algorithm, sample complexity.

## 1 Introduction

We consider composite optimization problems of the form

$$\operatorname*{minimize}_{x \in \mathbf{R}^n} \quad f(g(x)) + h(x), \tag{1}$$

where $f : \mathbf{R}^m \to \mathbf{R}$ is a convex and possibly nonsmooth function, $g : \mathbf{R}^n \to \mathbf{R}^m$ is a smooth mapping (vector-valued function), and $h : \mathbf{R}^n \to \mathbf{R}$ is a convex and lower-semicontinuous function. Although both $f$ and $h$ are convex, the problem is in general nonconvex due to the composition of $f$ and $g$. In addition, we assume that $g$ is either the average of finite number of component mappings, i.e., $g(x) = \frac{1}{N} \sum_{i=1}^{N} g_i(x)$, or the expectation of a family of random component mappings, i.e., $g(x) = \mathbf{E}_\xi[g_\xi(x)]$ where $\xi$ is a random variable. More explicitly, we consider the problems

$$\operatorname*{minimize}_{x \in \mathbf{R}^n} \quad f\left(\frac{1}{N} \sum_{i=1}^{N} g_i(x)\right) + h(x) \tag{2}$$

[*]Department of Electrical Engineering, Princeton University, Princeton, NJ 08544, USA. Email: `junyuz@princeton.edu`

[†]Facebook AI Research (FAIR), Seattle, WA 98109, USA. Email: `linx@fb.com`

and

$$\underset{x \in \mathbf{R}^n}{\text{minimize}} \quad f\big(\mathbf{E}_\xi[g_\xi(x)]\big) + h(x). \tag{3}$$

Clearly, problem (2) is a special case of (3) where the random variable $\xi$ follows the uniform distribution over the finite set $\{1, 2, \ldots, N\}$. We consider them separately because the sample complexity for solving problem (2) can be much lower than that of the general case (3).

An effective method for solving the composite optimization problem (1) is the (deterministic) *prox-linear* algorithm (e.g., [20, 36], which iteratively minimizes a model of the objective function where $g(x)$ is replaced by a linear approximation. Specifically, let $g' : \mathbf{R}^n \to \mathbf{R}^{m \times n}$ denote the Jacobian of $g$, then each iteration of prox-linear algorithm takes the form

$$x^{k+1} = \underset{x}{\text{argmin}} \left\{ f\big(g(x^k) + g'(x^k)(x - x^k)\big) + h(x) + \frac{M}{2} \|x - x^k\|^2 \right\}, \tag{4}$$

where $M > 0$ is a parameter to penalize the deviation of $x^{k+1}$ from $x^k$ in squared Euclidean distance. Since $f$ and $h$ are convex, the subproblem in (4) is a convex optimization problem. For the algorithm to be efficient in practice, we also need the functions $f$ and $h$ to be relatively simple, meaning that the subproblem in (4) admits a closed-form solution or can be solved efficiently.

For problems (2) and (3), the finite-average and expectation structure of $g$ allow us to use a randomly sampled subset of $g_i$ or $g_\xi$ and their Jacobians to approximate the expectations $g$ and $g'$. Specifically, during each iteration $k$, let $\mathcal{B}^k$ and $\mathcal{S}^k$ be two subsets of $\{1, 2, \ldots, N\}$ sampled uniformly at random or two sets of realizations of $\xi$ sampled from its distribution. A straightforward approach is to construct the mini-batch approximations

$$\tilde{g}^k = \frac{1}{|\mathcal{B}^k|} \sum_{i \in \mathcal{B}^k} g_i(x^k), \qquad \tilde{J}^k = \frac{1}{|\mathcal{S}^k|} \sum_{i \in \mathcal{S}^k} g'_i(x^k), \tag{5}$$

and use them to replace $g(x^k)$ and $g'(x^k)$ in (4), leading to the *stochastic prox-linear* algorithm:

$$x^{k+1} = \underset{x}{\text{argmin}} \left\{ f\big(\tilde{g}^k + \tilde{J}^k(x - x^k)\big) + h(x) + \frac{M}{2} \|x - x^k\|^2 \right\}. \tag{6}$$

While each iteration of (6) uses less samples of $g_\xi$ and $g'_\xi$ than the full-batch method (4), the simple mini-batch construction in (5) may not be able to reduce the overall sample complexity due to increased number of iterations required (see, e.g., [18] and [60, Section 3]).

In this paper, we develop a class of stochastic *variance-reduced* prox-linear algorithms for solving problems (2) and (3). By leveraging the variance reduction techniques of SVRG [31, 57] and SARAH/SPIDER [37, 24], we obtain significantly lower sample complexities than that of the full-batch prox-linear method. Before getting to the details, we first present several applications.

## 1.1 Application examples

Composite optimization problems of the forms (2) and (3) arise from risk-averse optimization (e.g, [48, 51] and a mean-variance tradeoff example in [59]) and stochastic variational inequalities (e.g., [30, 32], through a reformulation in [26]). In machine learning, a well-known example is policy evaluation for reinforcement learning (e.g., [15, 52, 54, 55]). Here we give several additional examples, and explain how the stochastic prox-linear algorithms can be applied.

**Systems of nonlinear equations for ERM**  Solving systems of nonlinear equations is one of the most fundamental problems in computational science and engineering (e.g., [41]). Given a system of nonlinear equations $g(x) = 0$ where $g : \mathbf{R}^n \to \mathbf{R}^m$ is a smooth mapping, a standard approach is to minimize the composite function $f(g(x))$ where $f$ is non-negative merit function and $f(z) = 0$ if only if $z = 0$. A popular choice is the squared Euclidean norm $f(\cdot) = \|\cdot\|^2$. The classical Gauss-Newton method iteratively minimizes a simple model by linearizing $g$ at $x^k$:

$$x^{k+1} = \operatorname*{argmin}_x \; \left\| g(x^k) + g'(x^k)(x - x^k) \right\|^2.$$

Nesterov [36] proposed a modified scheme with sharp merit functions such as $f(\cdot) = \|\cdot\|$ and a quadratic penalty term as in (4). For empirical risk minimization (ERM) problems of the form

$$\operatorname*{minimize}_x \quad F(x) \triangleq \frac{1}{N} \sum_{i=1}^{N} F_i(x),$$

where each $F_i$ is twice differentiable, we can apply Gauss-Newton type of methods by letting $g_i(x) = F_i'(x)$ and $g'(x) = F_i''(x)$ (the gradient and Hessian of $F_i$ respectively) and use either a smooth or a sharp merit function $f$. The resulting optimization problem is of the form (2) and we can exploit the finite-average structure with the sub-sampled prox-linear algorithm (6). This approach can be particularly useful for solving non-convex ERM problems (see, e.g., [50] and [13]). Efficient numerical algorithms for solving the subproblem in each iteration are discussed in [50] for $f(\cdot) = \|\cdot\|^2$ and in [36] for $f(\cdot) = \|\cdot\|$.

**Truncated stochastic gradient method**  Consider the stochastic optimization problem

$$\operatorname*{minimize}_x \quad g(x) \triangleq \mathbf{E}\big[g_\xi(x)\big],$$

where each $g_\xi : \mathbf{R}^n \to \mathbf{R}$ is smooth. Suppose we know the minimum value $g^* = \inf_x g(x)$ or a lower bound of it (in many machine learning problems $g(x) \geq 0$), then the problem is equivalent to

$$\operatorname*{minimize}_x \quad f(g(x)), \qquad \text{where} \quad f(z) = \max\{z, \, g^*\}.$$

In this case, the mini-batch stochastic prox-linear method (6) becomes

$$x^{k+1} = \operatorname*{argmin}_x \left\{ \max \left\{ \tilde{g}^k + \tilde{J}^k(x - x^k), \; g^* \right\} + \frac{M}{2} \|x - x^k\|^2 \right\}, \tag{7}$$

which has a closed-form solution

$$x_{k+1} = x_k - \min \left\{ \frac{1}{M}, \; \frac{\tilde{g}^k - g^*}{\|\tilde{J}^k\|^2} \right\} \cdot \tilde{J}^k.$$

This update has a very similar step-size rule as Polyak's rule for subgradient method [44]. Because the simple model used in (7) truncates the linear model with the known lower bound, it is called the *truncated stochastic gradient method*. Recent studies [1, 2, 16] show that it converges faster and is more stable than the classical stochastic gradient method with a wide range of step sizes. In this paper, we use variance reduction techniques to construct the estimates $\tilde{g}^k$ and $\tilde{J}^k$ and obtain better sample complexity for this method.

**Minimax stochastic optimization** Consider the problem of minimizing the maximum of $m$ expectations:

$$\underset{x \in \mathcal{X}}{\text{minimize}} \ \max_{1 \leq j \leq m} g^{(j)}(x), \qquad \text{where} \quad g^{(j)}(x) = \mathbf{E}_{\xi_j}\big[g^{(j)}_{\xi_j}(x)\big].$$

Here we assume that $\mathcal{X}$ is a closed convex set and the random variables $\xi_i$ follow (slightly) different probability distributions. This is a special case of *distributionally robust optimization* (see [47] and references therein), which has many applications in operations research and statistical machine learning. It can be put into the form of (3) with the definitions $\xi = [\xi_1, \ldots, \xi_m]$ and

$$f(z) = \max_{1 \leq j \leq m} z_j, \qquad g_\xi(x) = \big[g^{(1)}_{\xi_1}(x), \ldots, g^{(m)}_{\xi_m}(x)\big], \qquad h(x) = \delta_{\mathcal{X}}(x),$$

where $\delta_{\mathcal{X}}$ denotes the indicator function of $\mathcal{X}$. In this case, the update in (6) requires solving a convex quadratic programming problem. Similar formulations may apply to other distributionally robust optimization problems.

**Exact penalty method for stochastic optimization** Consider the following constrained stochastic optimization problem

$$\begin{aligned}
\underset{x \in \mathcal{X}}{\text{minimize}} \quad & \mathbf{E}_{\xi_0}\big[g^{(0)}_{\xi_0}(x)\big] \\
\text{subject to} \quad & \mathbf{E}_{\xi_j}\big[g^{(j)}_{\xi_j}(x)\big] \geq 0, \quad j = 1, \ldots, m_I, \\
& \mathbf{E}_{\xi_j}\big[g^{(j)}_{\xi_j}(x)\big] = 0, \quad j = m_I + 1, \ldots, m.
\end{aligned}$$

Using the exact penalty approach (see, e.g., [5, 28]), this problem can be reformulated as

$$\begin{aligned}
\underset{x}{\text{minimize}} \quad & \mathbf{E}_{\xi_0}\big[g^{(0)}_{\xi_0}(x)\big] + \sum_{j=1}^{m_I} c_j \max\Big\{0, \mathbf{E}_{\xi_j}\big[g^{(j)}_{\xi_j}(x)\big]\Big\} \\
& + \sum_{j=m_I+1}^{m} c_j \left|\mathbf{E}_{\xi_j}\big[g^{(j)}_{\xi_j}(x)\big]\right| + \delta_{\mathcal{X}}(x),
\end{aligned}$$

where $c_j > 0$ for $j = 1, \ldots, m$ are sufficiently large positive constants (to ensure the penalty terms vanish at optimality). It is straightforward to rewrite the above problem as (3) and we omit the details. The update in (6) also requires solving a convex quadratic programming problem.

## 1.2 Related work

The deterministic composite optimization problem (1) is a classical problem in nonconvex and nonsmooth optimization, and its study can date back to the late 70s in the last century; see, e.g., [4, 25, 43]. Recently, there has been a renewed interest in such problems due to many emerging applications, including the robust phase retrieval problem considered in [23], the low-rank semidefinite programming (SDP) problem considered in [3], and the robust blind deconvolution problem considered in [12], and so on. In fact, many of these applications involve the average or expectation over large amount of component loss functions, similar to those shown in problems (2) and (3).

For solving the nonlinear least-square problems (when $f = \|\cdot\|^2$), the idea of linearizing the inner mapping $g$ is well-known from the classical Gauss-Newton method (e.g, [39, Section 10.3]).

For nonsmooth $f$, the trial of linearizing the inner mapping $g$ was made in [7, 10], where the linearization is used to construct a descent direction for line-search. In [36], Nesterov proposed the Gauss-Newton type of algorithm (4) for nonsmooth $f$, analyzed its general convergence properties and proved local quadratic convergence under a non-degeneracy assumption. More recently, it has received more attention under the name of prox-linear algorithm. The authors of [11, 21, 40] discussed its iteration complexity and the numerical cost of solving the subproblem in each iteration. In [19, 20], the authors studied its fast local convergence property under the quadratic growth or the error-bound conditions. Additional references can be found in [8, 9, 34].

In the stochastic settings, it is worth noting that [16, 17, 22, 27] have considered the problem

$$\underset{x \in \mathcal{X}}{\text{minimize}} \quad \mathbf{E}_\xi \big[ f_\xi(g_\xi(x)) \big],$$

where the expectation is taken outside of the composition (in many cases $f$ does not depend on the random variable $\xi$). This problem is essentially a special case of the classical stochastic programming problem. The problems we consider in (2) and (3) are quite different.

Algorithms for solving stochastic composite optimization problems of the forms (2) and (3) have been studied recently in [6, 29, 35, 46, 54, 55, 58, 59, 61]. Since these are all stochastic or randomized algorithms, a common measure of performance is their sample complexity, i.e., the total number of samples of the component mappings $g_i$ or $g_\xi$ and their Jacobians required to output some point $\bar{x}$ such that $\mathbf{E} \big[ \|\mathcal{G}(\bar{x})\|^2 \big] \leq \epsilon$, where $\epsilon$ is a predefined precision and $\mathcal{G}(\bar{x})$ is the composite gradient mapping at $\bar{x}$ (for a precise definition, see (11) in Section 2). When both $f$ and $g$ are smooth and $g$ is a finite-average, the best sample complexity is $\mathcal{O}(N + N^{1/2}\epsilon^{-1})$ given in [61], which matches the best known complexity for nonconvex finite-sum optimization without composition [24, 38, 42, 56]. When both $f$ and $g$ are smooth and $g$ is a general expectation, the state-of-the-art sample complexity is the $\mathcal{O}(\epsilon^{-3/2})$ obtained in [61]. When $f$ is convex but nonsmooth and $g$ is a finite sum of $N$ smooth mappings, the authors of [46] applied the conjugate function of $f$ and transformed problem (2) to a min-max saddle-point problem. The sample complexity of their method (without counting subproblem cost) is $\mathcal{O}(N\epsilon^{-1})$.

After the initial submission of this paper, we were brought to attention the independent work [53]. The authors also consider problems (2) and (3) and develop stochastic Gauss-Newton methods (same form as prox-linear algorithms) using SARAH [37] for variance reduction. They obtained sample complexity $\mathcal{O}(\epsilon^{-5/2})$ for $g_\xi$ and $\mathcal{O}(\epsilon^{-3/2})$ for $g'_\xi$, but for a slightly different stationarity measure than the one used in this paper. We will comment on the connections to our results at the ends of Sections 3 and 4.

## 1.3 Contributions and outline

In this paper, we develop a class of stochastic *variance-reduced* prox-linear algorithms for solving problems (2) and (3), by constructing the estimates $\tilde{g}^k$ and $\tilde{J}^k$ in (6) with the variance reduction techniques of SVRG [31, 57] and SARAH/Spider [37, 24]. Our main results are summarized below.

- When $f$ is convex and nonsmooth and $g$ is a finite average, we construct an SVRG type estimator augmented with additional first-order correction, and obtain the sample complexity $\mathcal{O}(N + N^{4/5}\epsilon^{-1})$ for both component mapping $(g_i)$ and Jacobian $(g'_i)$ evaluations.

- When $f$ is convex and nonsmooth and $g$ is an expectation of random smooth mappings, we use the SARAH/Spider estimator, and obtain a sample complexity of $\mathcal{O}(\epsilon^{-5/2})$ for the random mappings $(g_\xi)$ and $\mathcal{O}(\epsilon^{-3/2})$ for the Jacobians $(g'_\xi)$.

- When $f$ is smooth, we also adopt the SARAH/SPIDER estimator. For both component mapping and Jacobian evaluations, we obtain the sample complexities $\mathcal{O}(N + \sqrt{N}\epsilon^{-1})$ and $\mathcal{O}(\epsilon^{-3/2})$ for the finite average case and expectation case respectively.

The first result above (with nonsmooth $f$ and finite-sum $g$) appears to be new and our sample complexity improves over the best known in the literature [46]. The second result is among the first in the literature to derive improved sample complexity for nonsmooth $f$ and with $g$ being an expectation (see also [53]). These results can be extended to the cases when $f$ is *weakly convex* (see its definition in, e.g., [16, 21]). We omit details to keep the presentation relatively simple, but will make remarks on the necessary changes where it is applicable.

Note that most work on stochastic composite optimization (SCO) construct the gradient estimators based on chain-rule (see e.g. [54, 61, 60]), and they all fail when $f$ is nonsmooth. The significance of our results (and those in [53]) is to show that using the prox-linear framework, instead of the chain-rule, can take advantage of variance reduction techniques in the nonsmooth composite setting to achieve better sample complexity. Another feature that distinguishes our first two results from the existing smooth SCO literature is the imbalance between the required estimation accuracy of $\tilde{g}$ and $\tilde{J}$. Unlike the chain rule based algrithms for smooth SCO problems where the required accuracy for $\tilde{g}$ and $\tilde{J}$ are of the same order (see e.g. [29, 59, 61]), the nonsmooth SCO problem requires the *order* of estimation accuracy for $\tilde{g}$ to be much higher than $\tilde{J}$. New techniques are required to handle this challenge.

Our results with $f$ being smooth match those in [61], which are obtained by using variance-reduced gradient estimators based on the chain rule, i.e., $(\tilde{J}^k)^T f'(\tilde{g}^k)$, in contrast to using the proximal mapping of $f$ in (6). It is often observed in practice that algorithms based on proximal mappings can be more efficient than those based on gradients, even though in theory they have the same sample complexity (e.g., [1, 2, 16]). Therefore it is very meaningful to establish the convergence and complexity of proximal-mapping based methods even when $f$ is smooth. In addition, we comment on its effectiveness by relating to the classical Gauss-Newton method at the end of this paper.

**Organization** In Section 2, we present a general framework of stochastic variance-reduced prox-linear algorithms using the update formula (6), without specifying how the estimates $\tilde{g}^k$ and $\tilde{J}^k$ are constructed. In Sections 3 and 4, we assume that $f$ can be nonsmooth, and present the constructions of $\tilde{g}^k$ and $\tilde{J}^k$ and the resulting sample complexities for solving problems (2) and (3) respectively. In Sections 5 and 6, we assume that $f$ is smooth and present the estimators and the corresponding sample complexities for solving these two problems respectively. In Section 7, we present preliminary numerical experiments to demonstrate the effectiveness of the proposed algorithms. We conclude the paper in Section 8 with further discussions on different variance reduction techniques for stochastic composite optimization.

## 2　The algorithm framework

In this section, we present a framework of stochastic variance-reduced prox-linear algorithms using the update formula (6). In order to simplify notations, we define

$$\Phi(x) \triangleq f(g(x)) + h(x), \tag{8}$$

where $g$ is either the average of finite number of component mappings as in problem (2), or the expectation of a family of random component mappings as in (3). We make the following assumptions throughout the paper.

**Assumption 1.** *The function $f : \mathbf{R}^m \to \mathbf{R} \cup \{+\infty\}$ is convex and $\ell_f$-Lipschitz continuous, i.e.,*

$$|f(u) - f(v)| \le \ell_f \|u - v\|, \qquad \forall\, u, v \in \mathbf{R}^m.$$

*The function $h : \mathbf{R}^n \to \mathbf{R} \cup \{+\infty\}$ is convex and lower semi-continuous.*

**Assumption 2.** *The vector mapping $g : \mathbf{R}^n \to \mathbf{R}^m$ is $\ell_g$-Lipschitz continuous and its Jacobian $g' : \mathbf{R}^n \to \mathbf{R}^{m \times n}$ is $L_g$-Lipschitz continuous, i.e.,*

$$\begin{aligned}
\|g(x) - g(y)\| &\le \ell_g \|x - y\|, \\
\|g'(x) - g'(y)\| &\le L_g \|x - y\|,
\end{aligned}$$

*for all $x, y \in \operatorname{dom} h$, where $\| \cdot \|$ for matrices denotes the spectral norm.*

A direct consequence of the Lipschitz condition on $g'$ in Assumption 2 is

$$\left\| g(x) - g(y) - g(y)'(x - y) \right\| \le \frac{L_g}{2} \|x - y\|^2. \tag{9}$$

(See, e.g., [41, Theorem 3.2.12].) Throughout the paper, we also assume the objective function $\Phi$ is lower bounded, as stated in the following assumption.

**Assumption 3.** *There exists $\Phi_*$ such that $\Phi_* = \inf_x \Phi(x) > -\infty$.*

Under these assumptions, we have the following result.

**Lemma 1.** *Suppose Assumption 1 and 2 hold, then for any $x, y \in \operatorname{dom} h$,*

$$f(g(x)) \le f(g(y) + g'(y)(x - y)) + \frac{\ell_f L_g}{2} \|x - y\|^2. \tag{10}$$

*Proof.* By the Lipschitz continuity of $f$ and $g'$, we have

$$\begin{aligned}
f(g(x)) &= f\big(g(y) + g'(y)(x - y)\big) + f(g(x)) - f\big(g(y) + g'(y)(x - y)\big) \\
&\le f\big(g(y) + g'(y)(x - y)\big) + \big| f(g(x)) - f\big(g(y) + g'(y)(x - y)\big) \big| \\
&\le f\big(g(y) + g'(y)(x - y)\big) + \ell_f \big\| g(x) - g(y) - g'(y)(x - y) \big\| \\
&\le f\big(g(y) + g'(y)(x - y)\big) + \frac{\ell_f L_g}{2} \|x - y\|^2,
\end{aligned}$$

where the last inequality is due to (9). $\qquad \square$

As a result of Lemma 1, $f(g(y) + g'(y)(x - y)) + h(x) + \frac{M}{2} \|x - y\|^2$ is an upper bound of the objective function $f(g(x)) + h(x)$ as long as $M \ge \ell_f L_g$. This is exactly the principle of majorization used in the update (4). In order to exploit the finite-average structure of problem (2), we can approximate the full average $g(x^k)$ and $g'(x^k)$ with randomly sampled mini-batch estimators $\tilde{g}^k$ and $\tilde{J}^k$ as in (5). For problem (3), sampling based methods are the only choices because the full expectations $\mathbf{E}_\xi[\cdot]$ are impossible to evaluate in most cases. As shown in several previous work

---
**Algorithm 1:** Stochastic variance-reduced prox-linear algorithm
---
**1 input:** initial point $x_0^1$, $M > 0$, number of outer and inner iterations $K$ and $\tau$.

**2 for** $k = 1, \ldots, K$ **do**

**3**      **for** $i = 0, \ldots, \tau - 1$ **do**

**4**          **if** $i == 0$ **then**

**5**             compute $\tilde{g}_0^k$ and $\tilde{J}_0^k$ using *large* batches $\mathcal{B}_0^k$ and $\mathcal{S}_0^k$ respectively.

**6**          **else**

**7**             compute $\tilde{g}_i^k$ and $\tilde{J}_i^k$ using *small* batches $\mathcal{B}_i^k$ and $\mathcal{S}_i^k$ respectively.

**8**          **end**

**9**          $x_{i+1}^k = \underset{x}{\operatorname{argmin}} \ \left\{ f\big(\tilde{g}_i^k + \tilde{J}_i^k(x - x_i^k)\big) + h(x) + \frac{M}{2}\|x - x_i^k\|^2 \right\}.$

**10**      **end**

**11**      Set $x_0^{k+1} = x_\tau^k$.

**12 end**

**13 output:** choose $x_{i*}^{k*}$ from $\{x_i^k\}_{i=0,\ldots,\tau-1}^{k=1,\ldots,K}$ uniformly at random.
---

(see, e.g., [18] and [60, Section 3]), the simple mini-batching scheme (5) usually does not reduce the overall sample complexity for problems with similar structure, compared with using the full-batch in the finite-average case and using a single sample in the expectation case.

In this paper, we propose a class of stochastic variance-reduced prox-linear algorithms, outlined in Algorithm 1, and shown that they achieve better sample complexities than simple mini-batching. Following the celebrated SVRG method [31, 57], our framework employs an outer loop of $K$ stages and an inner loop of $\tau$ iterations. During the first iteration of each inner loop, the mapping and Jacobian approximations $\tilde{g}_0^k$ and $\tilde{J}_0^k$ are computed using relatively large sample batches. In the rest of inner iterations, they are computed with relatively small sample batches. It turns out that different variance-reduced estimators are needed to obtain the best sample complexity under different assumptions on $f$ and the structure of $g$. We will present the details of constructing different estimators and their convergence analysis in the remaining sections of this paper.

In order to characterize the sample complexity of different algorithms, we first define what is an $\epsilon$-stationary point. For any $x \in \operatorname{dom} h$, we define the proximal point

$$x_+ \triangleq \underset{y}{\operatorname{argmin}} \left\{ f\big(g(x) + g'(x)(y - x)\big) + h(y) + \frac{M}{2}\|y - x\|^2 \right\}$$

and the composite gradient mapping at $x$,

$$\mathcal{G}_M(x) \triangleq M(x - x_+). \tag{11}$$

Given any $\epsilon > 0$, we call $\bar{x}$ an $\epsilon$-stationary point of $\Phi$ defined in (8) if $\|\mathcal{G}_M(\bar{x})\|^2 \leq \epsilon$. Note that when $h = 0$ and $f$ is the identity mapping, we have $\mathcal{G}_M(x) = \nabla\Phi(x)$ for any $M > 0$ and the definition of $\epsilon$-stationary point reduces to its classical form $\|\nabla\Phi(x)\|^2 \leq \epsilon$ for smooth optimization. For the validity of $\|\mathcal{G}_M(\cdot)\|^2$ as an optimality measure under nontrivial $h$ and nonsmooth $f$, the readers are referred to [20]. To simplify notation, we will omit the subscript $M$ (which is a constant throughout this paper) and denote the composite gradient mapping as $\mathcal{G}(x)$.

The sample complexity of a randomized algorithm, such as Algorithm 1, is the total number of evaluations of the component mappings $g_i$ or $g_\xi$ and their Jacobians required in order to output

some $\bar{x}$ satisfying

$$\mathbf{E}\big[\|\mathcal{G}(\bar{x})\|^2\big] \le \epsilon, \tag{12}$$

where the expectation is taken over all the random samplings during the iterations of the algorithm.

Notice that the proximal point $x_+$ used in the definition of $\mathcal{G}(x)$ is computed with $g(x)$ and $g'(x)$, which can be very costly if not impossible to evaluate. In Algorithm 1, the proximal point $x_{i+1}^k$ is computed using the estimates $\tilde{g}_i^k$ and $\tilde{J}_i^k$, i.e.,

$$x_{i+1}^k = \underset{x}{\operatorname{argmin}} \left\{ f\big(\tilde{g}_i^k + \tilde{J}_i^k(x - x_i^k)\big) + h(x) + \frac{M}{2}\|x - x_i^k\|^2 \right\}. \tag{13}$$

This leads to a convenient approximation,

$$\widetilde{\mathcal{G}}(x_i^k) \triangleq M(x_i^k - x_{i+1}^k), \tag{14}$$

of the true gradient mapping $\mathcal{G}(x_i^k) = M(x_i^k - \hat{x}_{i+1}^k)$, where

$$\hat{x}_{i+1}^k = \underset{x}{\operatorname{argmin}} \left\{ f\big(g(x_i^k) + g'(x_i^k)(x - x_i^k)\big) + h(x) + \frac{M}{2}\|x - x_i^k\|^2 \right\}. \tag{15}$$

Since the definitions of $\epsilon$-stationary point and sample complexity are based on the true gradient mapping $\mathcal{G}$ but computationally we only have access to the approximation $\widetilde{\mathcal{G}}$, we need to derive a bound between them for the purpose of complexity analysis. Not surprisingly, such a bound depends on the approximation quality of the estimators $\tilde{g}_i^k$ and $\tilde{J}_i^k$, as shown in the following lemma.

**Lemma 2.** *Under Assumptions 1 and 2, the iterates generated by Algorithm 1 satisfy*

$$\begin{aligned}
\frac{M - \ell_f L_g}{M^2}\big\|\mathcal{G}(x_i^k)\big\|^2 \quad \le \quad & \frac{2M + \ell_f L_g}{M^2}\big\|\widetilde{\mathcal{G}}(x_i^k)\big\|^2 \\
& + 4\ell_f\big\|\tilde{g}_i^k - g(x_i^k)\big\| + \frac{2\ell_f}{L_g}\big\|\tilde{J}_i^k - g'(x_i^k)\big\|^2.
\end{aligned}$$

*Proof.* For the ease of notation, we denote

$$\begin{aligned}
F(x; x_i^k) &= f\big(g(x_i^k) + g'(x_i^k)(x - x_i^k)\big), \tag{16} \\
\widetilde{F}(x; x_i^k) &= f\big(\tilde{g}_i^k + \tilde{J}_i^k(x - x_i^k)\big). \tag{17}
\end{aligned}$$

Since both $f$ and $h$ are convex (Assumption 1), the following two functions are $M$-strongly convex:

$$F(x; x_i^k) + h(x) + \frac{M}{2}\|x - x_i^k\|^2, \tag{18}$$

$$\widetilde{F}(x; x_i^k) + h(x) + \frac{M}{2}\|x - x_i^k\|^2. \tag{19}$$

According to (15) and (13), $\hat{x}_{i+1}^k$ and $x_{i+1}^k$ are the minimizers of these two functions respectively. Therefore

$$\begin{aligned}
F(\hat{x}_{i+1}^k; x_i^k) + h(\hat{x}_{i+1}^k) + \frac{M}{2}\|\hat{x}_{i+1}^k - x_i^k\|^2 \quad \le \quad & F(x_{i+1}^k; x_i^k) + h(x_{i+1}^k) + \frac{M}{2}\|x_{i+1}^k - x_i^k\|^2 \\
& - \frac{M}{2}\|\hat{x}_{i+1}^k - x_{i+1}^k\|^2,
\end{aligned}$$

and

$$\widetilde{F}(x_{i+1}^k; x_i^k) + h(x_{i+1}^k) + \frac{M}{2}\|x_{i+1}^k - x_i^k\|^2 \;\; \leq \;\; \widetilde{F}(\hat{x}_{i+1}^k; x_i^k) + h(\hat{x}_{i+1}^k) + \frac{M}{2}\|\hat{x}_{i+1}^k - x_i^k\|^2$$
$$- \frac{M}{2}\|\hat{x}_{i+1}^k - x_{i+1}^k\|^2.$$

Summing the two inequalities above and rearranging the terms, we obtain

$$M\|\hat{x}_{i+1}^k - x_{i+1}^k\|^2 \;\; \leq \;\; F(x_{i+1}^k; x_i^k) - \widetilde{F}(x_{i+1}^k; x_i^k) + \widetilde{F}(\hat{x}_{i+1}^k; x_i^k) - F(\hat{x}_{i+1}^k; x_i^k). \qquad (20)$$

Using the Lipschitz property of $f$, we have

$$
\begin{aligned}
\left| F(x_{i+1}^k; x_i^k) - \widetilde{F}(x_{i+1}^k; x_i^k) \right| 
&= \left| f\big(g(x_i^k) + g'(x_i^k)(x_{i+1}^k - x_i^k)\big) - f\big(\tilde{g}_i^k + \tilde{J}_i^k(x_{i+1}^k - x_i^k)\big) \right| \\
&\leq \ell_f \left\| \big(g(x_i^k) - \tilde{g}_i^k\big) + \big(g'(x_i^k) - \tilde{J}_i^k\big)(x_{i+1}^k - x_i^k) \right\| \\
&\leq \ell_f \left( \|\tilde{g}_i^k - g(x_i^k)\| + \|\tilde{J}_i^k - g'(x_i^k)\| \|x_{i+1}^k - x_i^k\| \right) \\
&\leq \ell_f \left( \|\tilde{g}_i^k - g(x_i^k)\| + \frac{1}{2L_g}\|\tilde{J}_i^k - g'(x_i^k)\|^2 + \frac{L_g}{2}\|x_{i+1}^k - x_i^k\|^2 \right).
\end{aligned}
$$

Replacing $x_{i+1}^k$ in the above inequality with $\hat{x}_{i+1}^k$, we get

$$\left| F(\hat{x}_{i+1}^k; x_i^k) - \widetilde{F}(\hat{x}_{i+1}^k; x_i^k) \right| \;\; \leq \;\; \ell_f \left( \|\tilde{g}_i^k - g(x_i^k)\| + \frac{1}{2L_g}\|\tilde{J}_i^k - g'(x_i^k)\|^2 + \frac{L_g}{2}\|\hat{x}_{i+1}^k - x_i^k\|^2 \right).$$

Combining the two bounds above with (20) gives

$$
\begin{aligned}
M\|\hat{x}_{i+1}^k - x_{i+1}^k\|^2 &\leq 2\ell_f \|\tilde{g}_i^k - g(x_i^k)\| + \frac{\ell_f}{L_g}\|\tilde{J}_i^k - g'(x_i^k)\|^2 \\
&\quad + \frac{\ell_f L_g}{2}\|x_{i+1}^k - x_i^k\|^2 + \frac{\ell_f L_g}{2}\|\hat{x}_{i+1}^k - x_i^k\|^2.
\end{aligned}
$$

Next, using the fact that $\|a+b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$ and the above inequality, we have

$$
\begin{aligned}
M\|\hat{x}_{i+1}^k - x_i^k\|^2] &\leq 2M\|x_{i+1}^k - x_i^k\|^2 + 2M\|\hat{x}_{i+1}^k - x_{i+1}^k\|^2 \\
&\leq 2M\|x_{i+1}^k - x_i^k\|^2 + 4\ell_f\|\tilde{g}_i^k - g(x_i^k)\| + \frac{2\ell_f}{L_g}\|\tilde{J}_i^k - g'(x_i^k)\|^2 \\
&\quad + \ell_f L_g\|x_{i+1}^k - x_i^k\|^2 + \ell_f L_g\|\hat{x}_{i+1}^k - x_i^k\|^2.
\end{aligned}
$$

Rearranging the terms yields

$$
\begin{aligned}
(M - \ell_f L_g)\|\hat{x}_{i+1}^k - x_i^k\|^2 &\leq (2M + \ell_f L_g)\|x_{i+1}^k - x_i^k\|^2 \\
&\quad + 4\ell_f\|\tilde{g}_i^k - g(x_i^k)\| + \frac{2\ell_f}{L_g}\|\tilde{J}_i^k - g'(x_i^k)\|^2.
\end{aligned}
$$

Finally, using the definitions $\mathcal{G}(x_i^k) = M(x_i^k - \hat{x}_{i+1}^k)$ and $\widetilde{\mathcal{G}}(x_i^k) = M(x_i^k - x_{i+1}^k)$, we obtain the desired result. $\qquad\square$

**Extension to the weakly convex case.** The function $f$ is $\rho$-weakly convex if $f(x) + \frac{\rho}{2}\|x\|^2$ is convex. In order to extends results in this paper for weakly convex $f$, we need to increase $M$ to ensure that the functions in (18) and (19) are strongly convex (in fact, strong convexity in expectation is sufficient).

## 3 The nonsmooth and finite-average case

In this section, we consider the composite finite-average problem (2) with nonsmooth $f$ and smooth $g_i$'s. In particular, we replace Assumption 2 with the following more structured one, which implies Assumption 2.

**Assumption 4.** *For each $i = 1, \ldots, N$, the mapping $g_i : \mathbf{R}^n \to \mathbf{R}^m$, is $\ell_{g,i}$-Lipschitz continuous and its Jacobian matrix $g_i' : \mathbf{R}^n \to \mathbf{R}^{m\times n}$ is $L_{g,i}$-Lipschitz continuous. Namely,*

$$\|g_i(x) - g_i(x)\| \le \ell_{g,i}\|x - y\|,$$
$$\|g_i'(x) - g_i'(x)\| \le L_{g,i}\|x - y\|,$$

*for all $x, y \in \operatorname{dom} h$ and $i = 1, \ldots, N$.*

A direct consequence of this assumption is that $g$ is $\left(\frac{1}{N}\sum_i^N \ell_{g,i}\right)$-Lipschitz continuous and $g'$ is $\left(\frac{1}{N}\sum_i^N L_{g,i}\right)$-Lipschitz continuous. Due to the root-mean square inequality $\frac{z_1+\ldots+z_N}{N} \le \sqrt{\frac{z_1^2+\ldots+z_N^2}{N}}$, We define

$$\ell_g = \sqrt{\frac{1}{N}\sum_i^N \ell_{g,i}^2}, \qquad L_g = \sqrt{\frac{1}{N}\sum_i^N L_{g,i}^2}, \tag{21}$$

which can serve as the Lipschitz constants of $g$ and $g'$ respectively as in Assumption 2.

In this case, we construct the estimates $\tilde{g}_0^k$ and $\tilde{J}_0^k$ using the full batch. In other words, we let $\mathcal{B}_0^k = \mathcal{S}_0^k = \{1, 2, \ldots, N\}$ and replace Line 5 in Algorithm 1 with

$$\tilde{g}_0^k = g(x_0^k) = \frac{1}{N}\sum_{i=1}^N g_i(x_0^k), \tag{22}$$

$$\tilde{J}_0^k = g'(x_0^k) = \frac{1}{N}\sum_{i=1}^N g_i'(x_0^k). \tag{23}$$

For $i > 0$, we sample with replacement from $\{1, 2, \ldots, N\}$ to obtain smaller sets $\mathcal{B}_i^k$ and $\mathcal{S}_i^k$ (whose cardinalities will be determined later), and apply the following construction:

$$\tilde{g}_i^k = \frac{1}{|\mathcal{B}_i^k|}\sum_{j\in\mathcal{B}_i^k}\left(g_j(x_i^k) - g_j(x_0^k) - g_j'(x_0^k)(x_i^k - x_0^k)\right) + g(x_0^k) + g'(x_0^k)(x_i^k - x_0^k), \tag{24}$$

$$\tilde{J}_i^k = \frac{1}{|\mathcal{S}_i^k|}\sum_{j\in\mathcal{S}_i^k}\left(g_j'(x_i^k) - g_j'(x_0^k)\right) + g'(x_0^k). \tag{25}$$

It is worth noting that here we use the standard SVRG estimator [31] to construct $\tilde{J}_i^k$, but the estimator for $\tilde{g}_i^k$ is augmented with a first-order correction (a similar estimator was proposed in [62]).

We remark that for nonsmooth $f$, the first-order correction scheme in (24) is essential for achieving a sample complexity that is sublinear in $N$, whereas purely applying the SVRG estimator will only result in a sample complexity linear in $N$. This is very different from the case with smooth $f$ (see e.g. [29, 35]). The main reason for such distinction is that nonsmooth SCO problem requires the estimation accuracy for $\tilde{g}_i^k$ to be much higher than $\tilde{J}_i^k$. In addition, the SARAH/SPIDER estimators seem to be not compatible with the first-order correction technique and we are not able to combine them together in order to obtain a sample complexity that is sublinear in $N$.

The following lemma bounds the approximation errors of these estimators.

**Lemma 3.** *Suppose Assumption 4 holds and $\tilde{g}_i^k$ and $\tilde{J}_i^k$ are constructed according to (24) and (25) respectively, then*

$$\mathbf{E}\Big[\big\|\tilde{g}_i^k - g(x_i^k)\big\| \,\big|\, x_i^k\Big] \leq \frac{L_g}{2\sqrt{|\mathcal{B}_i^k|}}\big\|x_i^k - x_0^k\big\|^2,$$

$$\mathbf{E}\Big[\big\|\tilde{J}_i^k - g'(x_i^k)\big\|^2 \,\big|\, x_i^k\Big] \leq \frac{L_g^2}{|\mathcal{S}_i^k|}\big\|x_i^k - x_0^k\big\|^2,$$

*where $\mathbf{E}[\cdot|x_i^k]$ denotes conditional expectation given $x_i^k$, i.e., expectation with respect to the random indices in $\mathcal{B}_i^k$ and $\mathcal{S}_i^k$.*

*Proof.* To prove the first inequality, we start with (24) and write

$$\tilde{g}_i^k - g(x_i^k) = \frac{1}{|\mathcal{B}_i^k|}\sum_{j\in\mathcal{B}_i^k} Z_j, \tag{26}$$

where

$$Z_j = g_j(x_i^k) - g_j(x_0^k) - g_j'(x_0^k)(x_i^k - x_0^k) + g(x_0^k) + g'(x_0^k)(x_i^k - x_0^k) - g(x_i^k).$$

Since $j$ is randomly sampled from $\{1, 2, \ldots, N\}$, we have $\mathbf{E}[g_j(x_i^k)] = g(x_i^k)$ and $\mathbf{E}[g_j'(x_i^k)] = g'(x_i^k)$, which implies $\mathbf{E}[Z_j|x_i^k] = 0$. That is, $\tilde{g}_i^k$ is an unbiased estimate of $g(x_i^k)$. In addition, we have

$$\mathbf{E}\big[g_j(x_i^k) - g_j(x_0^k) - g_j'(x_0^k)(x_i^k - x_0^k) \,\big|\, x_i^k\big] = g(x_i^k) - g(x_0^k) - g'(x_0^k)(x_i^k - x_0^k).$$

This allows us to bound the variance of $Z_j$ as follows:

$$\begin{aligned}
\mathbf{E}\big[\|Z_j\|^2 \,\big|\, x_i^k\big] &= \mathbf{E}\Big[\big\|g_j(x_i^k) - g_j(x_0^k) - g_j'(x_0^k)(x_i^k - x_0^k)\big\|^2 \,\big|\, x_i^k\Big] \\
&\quad - \big\|g(x_i^k) - g(x_0^k) - g'(x_0^k)(x_i^k - x_0^k)\big\|^2 \\
&\leq \mathbf{E}\Big[\big\|g_j(x_i^k) - g_j(x_0^k) - g_j'(x_0^k)(x_i^k - x_0^k)\big\|^2 \,\big|\, x_i^k\Big] \\
&\leq \frac{1}{N}\sum_{j=1}^N \left(\frac{L_{g,j}}{2}\|x_i^k - x_0^k\|^2\right)^2 \\
&= \frac{L_g^2}{4}\|x_i^k - x_0^k\|^4,
\end{aligned}$$

where the last inequality is due to (9) and Assumption 4 respectively. In the last equality, we used the definition of $L_g$ in (21).

12

Combining the above inequality with (26) yields

$$\mathbf{E}\Big[\big\|\tilde{g}_i^k - g(x_i^k)\big\|^2 \mid x_i^k\Big] \leq \frac{L_g^2}{4|\mathcal{B}_i^k|}\big\|x_i^k - x_0^k\big\|^4.$$

Next, using the concavity of $\sqrt{\cdot}$ and Jensen's inequality, we obtain the desired result:

$$\mathbf{E}\Big[\big\|\tilde{g}_i^k - g(x_i^k)\big\| \mid x_i^k\Big] \leq \sqrt{\mathbf{E}\Big[\big\|\tilde{g}_i^k - g(x_i^k)\big\|^2 \mid x_i^k\Big]} \leq \frac{L_g}{2\sqrt{|\mathcal{B}_i^k|}}\big\|x_i^k - x_0^k\big\|^2.$$

To prove the second inequality, we define $Z_j = g_j'(x_i^k) - g_j'(x_0^k) + g'(x_0^k) - g'(x_i^k)$ and follow a similar line of arguments. $\qquad\square$

Next, we prove a descent property of the algorithm, which is a crucial step for the convergence analysis.

**Lemma 4.** *Suppose Assumptions 1 and 4 hold and the estimates $\tilde{g}_0^k$, $\tilde{J}_0^k$, $\tilde{g}_i^k$ and $\tilde{J}_i^k$ in Algorithm 1 are constructed as in (22)-(25) respectively. Then for $k = 1, \ldots, K$ and $i = 0, \ldots, \tau - 1$,*

$$\begin{aligned}
\Phi(x_{i+1}^k) &\leq \Phi(x_i^k) - \frac{M - 2\ell_f L_g}{2M^2}\big\|\widetilde{\mathcal{G}}(x_i^k)\big\|^2 \\
&\quad + 2\ell_f\big\|\tilde{g}_i^k - g(x_i^k)\big\| + \frac{\ell_f}{2L_g}\big\|\tilde{J}_i^k - g'(x_i^k)\big\|^2.
\end{aligned} \tag{27}$$

*Proof.* By the definition of $\Phi$ in (8) and Lemma 1, we have

$$\begin{aligned}
\Phi(x_{i+1}^k) &\leq f\big(g(x_i^k) + g'(x_i^k)(x_{i+1}^k - x_i^k)\big) + \frac{\ell_f L_g}{2}\big\|x_{i+1}^k - x_i^k\big\|^2 + h(x_{i+1}^k) \\
&= f\big(\tilde{g}_i^k + \tilde{J}_i^k(x_{i+1}^k - x_i^k)\big) + h(x_{i+1}^k) + \frac{M}{2}\big\|x_{i+1}^k - x_i^k\big\|^2 - \frac{M - \ell_f L_g}{2}\big\|x_{i+1}^k - x_i^k\big\|^2 \\
&\quad + \underbrace{f\big(g(x_i^k) + g'(x_i^k)(x_{i+1}^k - x_i^k)\big) - f\big(\tilde{g}_i^k + g'(x_i^k)(x_{i+1}^k - x_i^k)\big)}_{T_1} \\
&\quad + \underbrace{f\big(\tilde{g}_i^k + g'(x_i^k)(x_{i+1}^k - x_i^k)\big) - f\big(\tilde{g}_i^k + \tilde{J}_i^k(x_{i+1}^k - x_i^k)\big)}_{T_2}.
\end{aligned} \tag{28}$$

According to (13), we have

$$f\big(\tilde{g}_i^k + \tilde{J}_i^k(x_{i+1}^k - x_i^k)\big) + h(x_{i+1}^k) + \frac{M}{2}\big\|x_{i+1}^k - x_i^k\big\|^2 \leq f(\tilde{g}_i^k) + h(x_i^k).$$

Therefore,

$$\begin{aligned}
\Phi(x_{i+1}^k) &\leq f(\tilde{g}_i^k) + h(x_i^k) - \frac{M - \ell_f L_g}{2}\big\|x_{i+1}^k - x_i^k\big\|^2 + T_1 + T_2 \\
&\leq f(g(x_i^k)) + h(x_i^k) - \frac{M - \ell_f L_g}{2}\big\|x_{i+1}^k - x_i^k\big\|^2 + T_1 + T_2 + \underbrace{f(\tilde{g}_i^k) - f(g(x_i^k))}_{T_3} \\
&= \Phi(x_i^k) - \frac{M - \ell_f L_g}{2}\big\|x_{i+1}^k - x_i^k\big\|^2 + T_1 + T_2 + T_3.
\end{aligned} \tag{29}$$

13

By the Lipschitz property of $f$, we have

$$T_1 \le \ell_f \|\tilde{g}_i^k - g(x_i^k)\|, \qquad T_3 \le \ell_f \|\tilde{g}_i^k - g(x_i^k)\|,$$

and

$$
\begin{aligned}
T_2 &\le \ell_f \|(\tilde{J}_i^k - g'(x_i^k))(x_{i+1}^k - x_i^k)\| \le \ell_f \|\tilde{J}_i^k - g'(x_i^k)\| \cdot \|x_{i+1}^k - x_i^k\| \\
&\le \frac{\ell_f}{2L_g} \|\tilde{J}_i^k - g'(x_i^k)\|^2 + \frac{\ell_f L_g}{2} \|x_{i+1}^k - x_i^k\|^2.
\end{aligned}
\tag{30}
$$

Combining the bounds on $T_1$, $T_2$ and $T_3$ and the inequality (29) yields

$$
\begin{aligned}
\Phi(x_{i+1}^k) &\le \Phi(x_i^k) - \left(\frac{M}{2} - \ell_f L_g\right) \|x_{i+1}^k - x_i^k\|^2 \\
&\quad + 2\ell_f \|\tilde{g}_i^k - g(x_i^k)\| + \frac{\ell_f}{2L_g} \|\tilde{J}_i^k - g'(x_i^k)\|^2,
\end{aligned}
$$

which, upon noticing $\widetilde{\mathcal{G}}(x_i^k) = -M(x_{i+1}^k - x_i^k)$, is equivalent to the desired result. $\qquad\square$

Recall the definition that $\widetilde{\mathcal{G}}(x_i^k) := M(x_i^k - x_{i+1}^k)$. In order to complete the convergence analysis, we define a stochastic Lyapunov function

$$
R_i^k = \mathbf{E}\left[\Phi(x_i^k) + c_i \left\|\sum_{t=0}^{i-1} \widetilde{\mathcal{G}}(x_t^k)\right\|^2\right], \quad k = 1, \ldots, K, \quad i = 0, \ldots, \tau,
\tag{31}
$$

where the coefficients $c_i$ for $i = 0, 1, \ldots, \tau$ are obtained through the recursion:

$$
\begin{aligned}
c_\tau &= 0, \\
c_i &= c_{i+1}\left(1 + \frac{1}{\tau}\right) + \frac{1}{3M\sqrt{|\mathcal{B}_i^k|}} + \frac{1}{5M|\mathcal{S}_i^k|}, \quad i = \tau - 1, \ldots, 0.
\end{aligned}
\tag{32}
$$

(Our choices or $|\mathcal{B}_i^k|$ and $|\mathcal{S}_i^k|$ will not depend on $k$.) In addition, we define the following constant

$$
\gamma \triangleq \min_{0 \le i \le \tau-1} \frac{1}{3}\left(\frac{1}{4M} - c_{i+1}(1 + \tau)\right).
$$

We can ensure $\gamma > 0$ by choosing $\tau$, $|\mathcal{B}_i^k|$ and $|\mathcal{S}_i^k|$ appropriately. We will discuss how to set these values after the following lemma, where we simply assume $\gamma > 0$.

**Lemma 5.** *Suppose Assumptions 1 and 4 hold and the estimates $\tilde{g}_0^k$, $\tilde{J}_0^k$, $\tilde{g}_i^k$ and $\tilde{J}_i^k$ in Algorithm 1 are constructed as in (22)-(25) respectively. In addition, we assume $M \ge 4\ell_f L_g$ and $\gamma > 0$. Then for each $k = 1, \ldots, K$,*

$$
\sum_{i=0}^{\tau-1} \mathbf{E}\left[\|\mathcal{G}(x_i^k)\|^2\right] \le \frac{R_0^k - R_\tau^k}{\gamma} = \frac{\mathbf{E}[\Phi(x_0^k)] - \mathbf{E}[\Phi(x_0^{k+1})]}{\gamma}.
\tag{33}
$$

*Proof.* For the ease of notation, we write the stochastic Lyapunov function as

$$R_i^k = \mathbf{E}\big[\Phi(x_i^k) + c_i\|G_i^k\|^2\big],$$

where

$$G_i^k = \sum_{t=0}^{i-1} \widetilde{\mathcal{G}}(x_t^k) = -M(x_i^k - x_0^k). \tag{34}$$

In particular, we have $G_0^k = -M(x_0^k - x_0^k) = 0$. Moreover, we have

$$
\begin{aligned}
\mathbf{E}\big[\|G_{i+1}^k\|^2\big] &= \mathbf{E}\big[\|G_i^k + \widetilde{\mathcal{G}}(x_i^k)\|^2\big] \\
&\leq \left(1 + \frac{1}{\tau}\right)\mathbf{E}\big[\|G_i^k\|^2\big] + (1+\tau)\mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big].
\end{aligned}
\tag{35}
$$

Combining Lemmas 2 and 4 yields

$$
\begin{aligned}
\mathbf{E}\big[\Phi(x_{i+1}^k)\big] &\leq \mathbf{E}\big[\Phi(x_i^k)\big] - \frac{M - 2\ell_f L_g}{2M^2}\mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big] + \left(\frac{\ell_f L_g}{\sqrt{|\mathcal{B}_i^k|}} + \frac{\ell_f L_g}{2|\mathcal{S}_i^k|}\right)\mathbf{E}\big[\|x_i^k - x_0^k\|^2\big] \\
&= \mathbf{E}\big[\Phi(x_i^k)\big] - \frac{M - 2\ell_f L_g}{2M^2}\mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big] + \frac{1}{M^2}\left(\frac{\ell_f L_g}{\sqrt{|\mathcal{B}_i^k|}} + \frac{\ell_f L_g}{2|\mathcal{S}_i^k|}\right)\mathbf{E}\big[\|G_i^k\|^2\big],
\end{aligned}
$$

where in the last equality we used (34). Adding both sides of (35) to that of the above inequality and using the assumption $M \geq 4\ell_f L_g$, we obtain

$$
\begin{aligned}
\mathbf{E}\big[\Phi(x_{i+1}^k) + c_{i+1}\|G_{i+1}^k\|^2\big] &\leq \mathbf{E}\big[\Phi(x_i^k)\big] - \left(\frac{M - 2\ell_f L_g}{2M^2} - c_{i+1}(1+\tau)\right)\mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big] \\
&\quad + \left(\frac{1}{M^2}\left(\frac{\ell_f L_g}{\sqrt{|\mathcal{B}_i^k|}} + \frac{\ell_f L_g}{2|\mathcal{S}_i^k|}\right) + c_{i+1}\left(1 + \frac{1}{\tau}\right)\right)\mathbf{E}\big[\|G_i^k\|^2\big] \\
&\leq \mathbf{E}\big[\Phi(x_i^k)\big] - \left(\frac{1}{4M} - c_{i+1}(1+\tau)\right)\mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big] \tag{36} \\
&\quad + \left(\frac{1}{4M}\left(\frac{1}{\sqrt{|\mathcal{B}_i^k|}} + \frac{1}{2|\mathcal{S}_i^k|}\right) + c_{i+1}\left(1 + \frac{1}{\tau}\right)\right)\mathbf{E}\big[\|G_i^k\|^2\big].
\end{aligned}
$$

Next, combining Lemma 2 with Lemmas 3 yields

$$\frac{M - \ell_f L_g}{M^2}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \leq \frac{2M + \ell_f L_g}{M^2}\mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big] + \left(\frac{2\ell_f L_g}{\sqrt{|\mathcal{B}_i^k|}} + \frac{2\ell_f L_g}{|\mathcal{S}_i^k|}\right)\mathbf{E}\big[\|x_i^k - x_0^k\|^2\big].$$

Using the equality $G_i^k = -M(x_i^k - x_0^k)$ and the assumption $M \geq 4\ell_f L_g$, the above inequality implies

$$\frac{3}{4M}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \leq \frac{9}{4M}\mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big] + \frac{1}{2M}\left(\frac{1}{\sqrt{|\mathcal{B}_i^k|}} + \frac{1}{|\mathcal{S}_i^k|}\right)\mathbf{E}\big[\|G_i^k\|^2\big]. \tag{37}$$

Multiplying both sides of (37) by $\left( \frac{1}{4M} - c_{i+1}(1 + \tau) \right) \big/ \frac{9}{4M}$, which is positive by the assumption $\gamma > 0$, and adding the resulting inequality to (36), we get

$$\mathbf{E}\big[\Phi(x_{i+1}^k) + c_{i+1}\|G_{i+1}^k\|^2\big] \ \le \ \mathbf{E}\left[\Phi(x_i^k) + \left( c_{i+1}\left(1 + \frac{1}{\tau}\right) + \frac{1}{3M\sqrt{|\mathcal{B}_i^k|}} + \frac{1}{5M|\mathcal{S}_i^k|} \right)\|G_i^k\|^2\right]$$
$$- \frac{1}{3}\left( \frac{1}{4M} - c_{i+1}(1 + \tau) \right) \mathbf{E}[\|\mathcal{G}(x_i^k)\|^2].$$

Now, using the definitions in (31) and (32), the above inequality is the same as

$$\frac{1}{3}\left( \frac{1}{4M} - c_{i+1}(1 + \tau) \right) \mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \ \le \ R_i^k - R_{i+1}^k.$$

Recalling the definition of $\gamma$ and summing up the above inequality over $i$ from 0 to $\tau - 1$, we get

$$\gamma \sum_{i=0}^{\tau-1} \mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \ \le \ R_0^k - R_\tau^k \ = \ \mathbf{E}\big[\Phi(x_0^k)\big] - \mathbf{E}\big[\Phi(x_\tau^k)\big],$$

where the last equality is due to the observations that $c_\tau = 0$ and $G_0^k = 0$. Finally, dividing both sides by $\gamma$ and using $x_0^{k+1} = x_\tau^k$ give the desired result. $\qquad\square$

The next lemma shows how to choose the inner loop length $\tau$ and the two mini-batch sizes $|\mathcal{B}_i^k|$ and $|\mathcal{S}_i^k|$ to ensure $\gamma > 0$. We use $\lceil\cdot\rceil$ to denote the nearest integer from above.

**Lemma 6.** *If we choose $\tau = \lceil \frac{1}{2}N^{1/5} - 1 \rceil$, $|\mathcal{B}_i^k| = \lceil 4N^{4/5} \rceil$ and $|\mathcal{S}_i^k| = \lceil N^{2/5} \rceil$ for $i = 1, \ldots, \tau - 1$, then $\gamma \ge \frac{1}{15M}$.*

*Proof.* To simplify notation, let $B = |\mathcal{B}_i^k|$ and $S = |\mathcal{S}_i^k|$ for $i = 1, \ldots, \tau - 1$. From (32), we deduce

$$(c_i + C) = (c_{i+1} + C)\left(1 + \frac{1}{\tau}\right), \qquad \text{where} \quad C = \frac{\tau}{3M\sqrt{B}} + \frac{\tau}{5MS}.$$

Consequently, with $c_\tau = 0$, we have for all $i = 1, \ldots, \tau$,

$$c_i \ = \ (c_\tau + C)\left(1 + \frac{1}{\tau}\right)^{\tau - i} - C \ \le \ C\left(1 + \frac{1}{\tau}\right)^{\tau} - C \ \le \ Ce - C \ = \ C(e - 1),$$

where the last inequality is due to the fact that $(1 + 1/\tau)^\tau \le e$ with $e$ is Euler's number (the basis of natural logarithm). Therefore,

$$\begin{aligned}
\gamma \ &= \ \min_{0 \le i \le \tau - 1} \frac{1}{3}\left( \frac{1}{4M} - c_{i+1}(1 + \tau) \right) \\
&\ge \ \frac{1}{3}\left( \frac{1}{4M} - C(e-1)(1 + \tau) \right) \\
&= \ \frac{1}{3M}\left( \frac{1}{4} - \left( \frac{1}{3\sqrt{B}} + \frac{1}{5S} \right)(e-1)\tau(1 + \tau) \right) \\
&\ge \ \frac{1}{3M}\left( \frac{1}{4} - \left( \frac{1}{3\sqrt{B}} + \frac{1}{5S} \right)2(1 + \tau)^2 \right).
\end{aligned}$$

Finally, setting $\tau = \frac{1}{2}N^{1/5} - 1$, $B = 4N^{4/5}$ and $S = N^{2/5}$ yields $\gamma \ge \frac{1}{15M}$. $\qquad\square$

16

Combining Lemma 5 and Lemma 6, we arrive at the main result of this section.

**Theorem 1.** *Suppose Assumptions 1, 3 and 4 hold for problem (2). Let the estimates $\tilde{g}_0^k$, $\tilde{J}_0^k$, $\tilde{g}_i^k$ and $\tilde{J}_i^k$ in Algorithm 1 be given in (22)-(25) respectively. If we choose $M \geq 4\ell_f L_g$ and $\tau = \lceil \frac{1}{2} N^{1/5} - 1 \rceil$, and*

$$|\mathcal{B}_i^k| = \lceil 4N^{4/5} \rceil, \qquad |\mathcal{S}_i^k| = \lceil N^{2/5} \rceil, \qquad i = 1, \dots, \tau - 1, \qquad k = 1, \dots, K,$$

*then the output of Algorithm 1 satisfies*

$$\mathbf{E}\big[\|\mathcal{G}(x_{i*}^{k*})\|^2\big] \leq \frac{15M\big(\Phi(x_0^1) - \Phi_*\big)}{K\tau}. \tag{38}$$

*To get an $\epsilon$-stationary point in expectation, the total sample complexity for the component mappings $g_j$ and their Jacobians are both $\mathcal{O}(N + N^{4/5}\epsilon^{-1})$.*

*Proof.* Summing up the inequality (33) over $k$ from 1 to $K$ and using the fact $\Phi(x_\tau^K) > \Phi_*$, we get

$$\sum_{k=1}^K \sum_{i=0}^{\tau-1} \mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \leq \frac{\Phi(x_0^1) - \Phi(x^*)}{\gamma}.$$

By the random choice of the output $x_{i*}^{k*}$, we can get the inequality (38).

To get an $\epsilon$-stationary point in expectation, we need to set $K\tau = \mathcal{O}(\epsilon^{-1})$, which implies

$$K = \mathcal{O}(\tau^{-1}\epsilon^{-1}) = \mathcal{O}(N^{-1/5}\epsilon^{-1}).$$

Consequently, the sample complexity of the component mappings (the $g_i$'s) is

$$KN + K\tau B \;=\; \mathcal{O}(N^{-1/5}\epsilon^{-1}) \cdot N + \mathcal{O}(\epsilon^{-1}) \cdot 4N^{4/5} \;=\; \mathcal{O}(N + N^{4/5}\epsilon^{-1}).$$

and the sample complexity for the component Jacobians is

$$KN + K\tau S \;=\; \mathcal{O}(N^{-1/5}\epsilon^{-1}) \cdot N + \mathcal{O}(\epsilon^{-1}) \cdot N^{2/5} \;=\; \mathcal{O}(N + N^{4/5}\epsilon^{-1}).$$

This completes the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Remark 1.** *Up to this point, we notice that all the analysis leading to Theorem 1 only requires the sample batches between iterations to be independent. Whereas within each iteration we do not require the independence between $\mathcal{B}_i^k$ and $\mathcal{S}_i^k$. Therefore, in practice one can simply use the same mini-batch $\mathcal{B}_i^k = \mathcal{S}_i^k$ to estimate both $\tilde{g}_i^k$ and $\tilde{J}_i^k$, with batch size equal to $\max\{S, B\} = \lceil 4N^{4/5} \rceil$. Or, we can use a random subset of $\mathcal{B}_i^k$ of size $\lceil N^{2/5} \rceil$ to compute $\tilde{J}_i^k$ in order to save computation.*

**Remark 2.** *The nonsmooth and finite-sum case is also considered in [53]. But their results are limited to using the simple mini-batch scheme for both component mapping and Jacobian estimation. As a consequence, their sample complexities for the component mappings and their Jacobians are $\mathcal{O}(\epsilon^{-3})$ and $\mathcal{O}(\epsilon^{-2})$ respectively, without explicit dependence on $N$. They are similar to our results in Section 4.1 on using mini-batches when $g$ is a general expectation.*

---

**Algorithm 2:** Simple mini-batch prox-linear algorithm

---

**1 input:** initial point $x_0$, parameter $M > 0$, and number of iterations $T$.

**2 for** $i = 0, \ldots, T - 1$ **do**

**3** $\quad$ sample mini-batches $\mathcal{B}_i$ and $\mathcal{S}_i$ from distribution of $\xi$, and compute $\tilde{g}_i$ and $\tilde{J}_i$ as in (39).

**4** $\quad$ $x_{i+1} = \underset{x}{\mathrm{argmin}} \left\{ f(\tilde{g}_i + \tilde{J}_i(x - x_i)) + h(x) + \frac{M}{2}\|x - x_i\|^2 \right\}.$

**5 end**

**6 output:** choose $x_{i^*}$ from $\{x_0, x_1, \ldots, x_{T-1}\}$ uniformaly at random.

---

# 4 The nonsmooth and expectation case

In this section, we consider the composite stochastic optimization problem (3), which we repeat here for convenience:

$$\underset{x}{\text{minimize}} \quad \Phi(x) \triangleq f(g(x)) + h(x), \quad \text{where} \quad g(x) = \mathbf{E}_\xi\big[g_\xi(x)\big].$$

We assume that $f$ and $h$ satisfy Assumption 1 and the $g_\xi$'s satisfy the following assumption.

**Assumption 5.** *The random mappings $g_\xi : \mathbf{R}^n \to \mathbf{R}^m$ and their Jacobians are mean-squares Lipschitz continuous, i.e., there exist constants $\ell_g$ and $L_g$ such that for all $x, y \in \mathrm{dom}\, h$,*

$$
\begin{aligned}
\mathbf{E}\big[\|g_\xi(x) - g_\xi(y)\|^2\big] &\leq \ell_g^2 \|x - y\|^2, \\
\mathbf{E}\big[\|g'_\xi(x) - g'_\xi(y)\|^2\big] &\leq L_g^2 \|x - y\|^2.
\end{aligned}
$$

*Furthermore, there exist constants $\sigma_g^2$ and $\sigma_{g'}^2$ such that for all $x \in \mathrm{dom}\, h$,*

$$
\begin{aligned}
\mathbf{E}\big[\|g_\xi(x) - g(x)\|^2\big] &\leq \sigma_g^2, \\
\mathbf{E}\big[\|g'_\xi(x) - g'(x)\|^2\big] &\leq \sigma_{g'}^2.
\end{aligned}
$$

Assumption 5 implies Assumption 2, but is weaker than assuming that $g_\xi$ and $g'_\xi$ are almost surely $\ell_g$- and $L_g$-Lipschitz respectively.

In this case, the first-order correction used in (24) is no longer useful in reducing the estimation errors because we cannot evaluate $g(x_0^k)$ or $g'(x_0^k)$ accurately. Instead, we turn to the SARAH/SPIDER estimator developed in [37, 24]. But before doing that, we first examine the simple mini-batch scheme outlined in (5) and (6).

## 4.1 The simple mini-batch method

The simple mini-batch method is to run Algorithm 1 with only one epoch ($K = 1$) and $\tau = T$ iterations, where during each iteration we set

$$\tilde{g}_i = \frac{1}{|\mathcal{B}_i|} \sum_{\xi \in \mathcal{B}_i} g_\xi(x_i), \qquad \text{and} \qquad \tilde{J}_i = \frac{1}{|\mathcal{S}_i|} \sum_{\xi \in \mathcal{S}_i} g'_\xi(x_i). \tag{39}$$

Since there is only one epoch, we omit the superscript $k$ on $x_i^k$, $\tilde{g}_i^k$ and $\tilde{J}_i^k$ to write $x_i$, $\tilde{g}_i$ and $\tilde{J}_i$. Similar to Remark 1, we do not require the independence between $\mathcal{B}_i$ and $\mathcal{S}_i$. For clarity, we present the resulting method as Algorithm 2. The following complexity result holds.

**Theorem 2.** *Suppose Assumptions 1, 3 and 5 hold for problem* (3). *If we choose $M \geq 4\ell_f L_g$ and the batch sizes $|\mathcal{B}_i| = B \geq \frac{36\ell_f^2 \sigma_g^2}{\epsilon^2}$ and $|\mathcal{S}_i| = S \geq \frac{2\ell_f \sigma_{g'}^2}{L_g \epsilon}$, then the output $x_{i^*}$ of Algorithm 2 satisfies*

$$\mathbf{E}\big[\|\mathcal{G}(x_{i^*})\|^2\big] \leq 12M \left( \frac{\Phi(x_0) - \Phi_*}{T} + \epsilon \right). \tag{40}$$

*Consequently by setting $T = \mathcal{O}(\epsilon^{-1})$, the sample complexities for the component mappings $g_\xi$ and their Jacobians for getting an $\epsilon$-solution are $\mathcal{O}(\epsilon^{-3})$ and $\mathcal{O}(\epsilon^{-2})$ respectively.*

*Proof.* From the construction of $\tilde{g}_i$ and $\tilde{J}_i$ in (39), we have $\mathbf{E}[\tilde{g}_i] = g(x_i)$ and $\mathbf{E}[\tilde{J}_i] = g'(x_i)$. Moreover, by Assumption 5, we have

$$\mathbf{E}\big[\|\tilde{g}_i - g(x_i)\|^2\big] \leq \frac{\sigma_g^2}{B}, \qquad \mathbf{E}\big[\|\tilde{J}_i - g'(x_i)\|^2\big] \leq \frac{\sigma_{g'}^2}{S}.$$

Using Jensen's inequality, the variance bound on $\tilde{g}_i$ further implies that $\mathbf{E}\big[\|\tilde{g}_i - g(x_i)\|\big] \leq \frac{\sigma_g}{\sqrt{B}}$. Together with Lemma 4, we have

$$\frac{M - 2\ell_f L_g}{2M^2} \mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big] \ \leq \ \mathbf{E}[\Phi(x_i)] - \mathbf{E}[\Phi(x_{i+1})] + \frac{2\ell_f \sigma_g}{\sqrt{B}} + \frac{\ell_f \sigma_{g'}^2}{2L_g S}. \tag{41}$$

On the other hand, applying Lemma 2 yields

$$\frac{M - \ell_f L_g}{M^2} \mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \ \leq \ \frac{2M + \ell_f L_g}{M^2} \mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big] + \frac{4\ell_f \sigma_g}{\sqrt{B}} + \frac{2\ell_f \sigma_{g'}^2}{L_g S}. \tag{42}$$

Next, we multiply both sides of (42) by $\frac{M - 2\ell_f L_g}{2(2M + \ell_f L_g)}$ and add them to (41) to cancel the terms containing $\mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big]$. Then with $M \geq 4\ell_f L_g$, we have $\frac{M - 2\ell_f L_g}{2(2M + \ell_f L_g)} \in \left[\frac{1}{9}, \frac{1}{4}\right]$ and obtain

$$\frac{1}{12M} \mathbf{E}[\|\mathcal{G}(x_i)\|^2] \leq \mathbf{E}[\Phi(x_i)] - \mathbf{E}[\Phi(x_{i+1})] + \frac{3\ell_f \sigma_g}{\sqrt{B}} + \frac{\ell_f \sigma_{g'}^2}{L_g S}.$$

Summing up the above inequality over $i$ from 0 to $T - 1$ and dividing by $T$, we obtain

$$\frac{1}{T} \sum_{i=0}^{T-1} \mathbf{E}\big[\|\mathcal{G}(x_i)\|^2\big] \ \leq \ 12M \left( \frac{\Phi(x_0) - \Phi_*}{T} + \frac{3\ell_f \sigma_g}{\sqrt{B}} + \frac{\ell_f \sigma_{g'}^2}{L_g S} \right).$$

Finally, using $|\mathcal{B}_i| = B \geq \frac{36\ell_f^2 \sigma_g^2}{\epsilon^2}$ and $|\mathcal{S}_i| = S \geq \frac{2\ell_f \sigma_{g'}^2}{L_g \epsilon}$ yields (40). The sample complexities for $g_\xi$ and $g'_\xi$ can be obtained as $TB = \mathcal{O}(\epsilon^{-3})$ and $TS = \mathcal{O}(\epsilon^{-2})$ respectively. □

## 4.2 Using the SARAH/SPIDER estimator

In this section, we show that by using the SARAH/SPIDER estimator [37, 24], the sample complexities for the component mappings and Jacobians can be improved to $\mathcal{O}(\epsilon^{-5/2})$ and $\mathcal{O}(\epsilon^{-3/2})$, respectively. We note that for solving problem (3) when $f$ is nonsmooth and convex (more generally weakly convex), even the $\mathcal{O}(\epsilon^{-3})$ and $\mathcal{O}(\epsilon^{-2})$ sample complexities established in Theorem 2 seem to be new in the literature.

The SARAH/SPIDER estimators for Algorithm 1 are constructed as follows. For $i = 0$, we set

$$\tilde{g}_0^k = \frac{1}{|\mathcal{B}_0^k|} \sum_{\xi \in \mathcal{B}_0^k} g_\xi(x_0^k), \qquad \text{and} \qquad \tilde{J}_0^k = \frac{1}{|\mathcal{S}_0^k|} \sum_{j \in \mathcal{S}_0^k} g_\xi'(x_0^k). \tag{43}$$

For the rest iterations with $i = 1, \ldots, \tau - 1$,

$$\tilde{g}_i^k = \tilde{g}_{i-1}^k + \frac{1}{|\mathcal{B}_i^k|} \sum_{\xi \in \mathcal{B}_i^k} \left( g_\xi(x_i^k) - g_\xi(x_{i-1}^k) \right), \tag{44}$$

$$\tilde{J}_i^k = \tilde{J}_{i-1}^k + \frac{1}{|\mathcal{S}_i^k|} \sum_{\xi \in \mathcal{S}_i^k} \left( g_\xi'(x_i^k) - g_\xi'(x_{i-1}^k) \right), \tag{45}$$

Here $\mathcal{B}_i^k$ and $\mathcal{S}_i^k$ for $i = 0, 1, \ldots, \tau - 1$ are mini-batches sampled from the underlying distribution of the random variable $\xi$. We require the batches $\mathcal{B}_i^k$ (and $\mathcal{S}_i^k$) to be independently sampled for different iterations, whereas in each iteration $\mathcal{B}_i^k$ and $\mathcal{S}_i^k$ can be dependent or even identical. The mean-squared estimation errors of the above estimators are bounded via the following lemma, which is adapted from [37, Lemma 2] or [24, Lemma 1]. A complete proof can be found in [61, Lemma 1].

**Lemma 7.** *Suppose Assumption 5 holds and $\tilde{g}_i^k$ and $\tilde{J}_i^k$ are constructed through (43)-(45). Then we have for $k = 1, \ldots, K$ and $\tau = 0, 1, \ldots, \tau - 1$,*

$$\mathbf{E}\left[\|\tilde{g}_i^k - g(x_i^k)\|^2\right] \leq \mathbf{E}\left[\|\tilde{g}_0^k - g(x_0^k)\|^2\right] + \sum_{r=1}^{i} \frac{\ell_g^2}{|\mathcal{B}_r^k|} \mathbf{E}\left[\|x_r^k - x_{r-1}^k\|^2\right], \tag{46}$$

$$\mathbf{E}\left[\|\tilde{J}_i^k - g'(x_i^k)\|^2\right] \leq \mathbf{E}\left[\|\tilde{J}_0^k - g'(x_0^k)\|^2\right] + \sum_{r=1}^{i} \frac{L_g^2}{|\mathcal{S}_r^k|} \mathbf{E}\left[\|x_r^k - x_{r-1}^k\|^2\right]. \tag{47}$$

The following theorem establishes the convergence of Algorithm 1 by specifying the batch sizes used in the SARAH/SPIDER estimators, and gives the sample complexities for $g_\xi$ and $g_\xi'$.

**Theorem 3.** *Suppose Assumptions 1, 3 and 5 hold for problem (3). Let the estimates $\tilde{g}_0^k$, $\tilde{J}_0^k$, $\tilde{g}_i^k$ and $\tilde{J}_i^k$ in Algorithm 1 be given in (43)-(45). If we choose $M \geq 4\ell_f L_g$ and $\tau = \lceil \epsilon^{-1/2} \rceil$, and the batch sizes as*

$$|\mathcal{B}_0^k| = \left\lceil \frac{25\ell_f^2 \sigma_g^2}{4\epsilon^2} \right\rceil, \quad |\mathcal{S}_0^k| = \left\lceil \frac{3\ell_f \sigma_{g'}^2}{4L_g\epsilon} \right\rceil, \quad |\mathcal{B}_i^k| = \left\lceil \frac{25\ell_f^2 \ell_g^2}{M\epsilon^{3/2}} \right\rceil, \quad |\mathcal{S}_i^k| = \left\lceil \frac{12\ell_f L_g}{M\epsilon^{1/2}} \right\rceil,$$

*for $i = 1, \ldots, \tau - 1$, then the output $x_{i*}^{k*}$ satisfies*

$$\mathbf{E}\left[\|\mathcal{G}(x_{i*}^{k*})\|^2\right] \leq 24M \left( \frac{\Phi(x_0^1) - \Phi_*}{K\tau} + 3\epsilon \right). \tag{48}$$

*Consequently by setting $K = \mathcal{O}(\epsilon^{-\frac{1}{2}})$, then we get an output $\mathbf{E}[\|\mathcal{G}(x_{i*}^{k*})\|^2] \leq \mathcal{O}(\epsilon)$ with a function evaluation complexity of $\mathcal{O}(\epsilon^{-5/2})$ and a Jacobian evaluation complexity of $\mathcal{O}(\epsilon^{-3/2})$.*

*Proof.* We will choose batch sizes that do not depend on $k$. For the ease of notation, we set $|\mathcal{B}_0^k| = B$, $|\mathcal{S}_0^k| = S$, and $|\mathcal{B}_i^k| = b$ and $|\mathcal{S}_i^k| = s$ for $i = 1, \ldots, \tau - 1$. First, by Assumption 5 and (43), we have

$$\mathbf{E}\big[\|\tilde{g}_0^k - g(x_0^k)\|^2\big] = \frac{\sigma_g^2}{B}, \qquad \mathbf{E}\big[\|\tilde{J}_0^k - g'(x_0^k)\|^2\big] = \frac{\sigma_{g'}^2}{S},$$

which can be substituted into Lemma 7. Then by Lemma 7, we know that

$$\mathbf{E}\big[\|\tilde{g}_i^k - g(x_i^k)\|\big] \;\leq\; \sqrt{\mathbf{E}\big[\|\tilde{g}_i^k - g(x_i^k)\|^2\big]} \;\leq\; \frac{\sigma_g}{\sqrt{B}} + \sqrt{\frac{\ell_g^2}{b}\sum_{r=1}^i \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big]}.$$

Moreover, for any $\delta > 0$, we have

$$\sqrt{\frac{\ell_g^2}{b}\sum_{r=1}^i \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big]} \;\leq\; \frac{\delta}{2} + \frac{\ell_g^2}{2b\delta}\sum_{r=1}^i \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big].$$

Now we invoke Lemma 4. Taking expectation on both sides of (27) and applying (47) and the above bounds, we obtain

$$
\begin{aligned}
\mathbf{E}\big[\Phi(x_{i+1}^k)\big] \;\leq\;\; & \mathbf{E}[\Phi(x_i^k)] - \frac{M - 2\ell_f L_g}{2}\mathbf{E}\big[\|x_{i+1}^k - x_i^k\|^2\big] + \frac{\ell_f \sigma_{g'}^2}{2L_g S} + \frac{2\ell_f \sigma_g}{\sqrt{B}} \\
& + \left(\frac{\ell_f L_g}{2s} + \frac{\ell_f \ell_g^2}{b\delta}\right)\sum_{r=1}^i \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big] + \ell_f \delta,
\end{aligned}
\tag{49}
$$

Similarly, with Lemma 2, we have

$$
\begin{aligned}
\frac{M - \ell_f L_g}{M^2}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \;\leq\;\; & (2M + \ell_f L_g)\mathbf{E}\big[\|x_{i+1}^k - x_i^k\|^2\big] + \frac{2\ell_f \sigma_{g'}^2}{L_g S} + \frac{4\ell_f \sigma_g}{\sqrt{B}} \\
& + \left(\frac{2\ell_f L_g}{s} + \frac{2\ell_f \ell_g^2}{b\delta}\right)\sum_{r=1}^i \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big] + 2\ell_f \delta.
\end{aligned}
\tag{50}
$$

Because $M \geq 4\ell_f L_g$, we have $\frac{1}{2} \cdot \frac{M - 2\ell_f L_g}{2(2M + \ell_f L_g)} \in \left[\frac{1}{18}, \frac{1}{8}\right]$. Therefore, multiplying (50) by $\frac{1}{2} \cdot \frac{M - 2\ell_f L_g}{2(2M + \ell_f L_g)}$ and adding to (49) gives

$$
\begin{aligned}
\frac{1}{24M}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \;\leq\;\; & \mathbf{E}\big[\Phi(x_i^k)\big] - \mathbf{E}\big[\Phi(x_{i+1}^k)\big] - \frac{M - 2\ell_f L_g}{4}\mathbf{E}\big[\|x_{i+1}^k - x_i^k\|^2\big] \\
& + \left(\frac{3\ell_f L_g}{4s} + \frac{5\ell_f \ell_g^2}{4b\delta}\right)\sum_{r=1}^i \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big] + \frac{3\ell_f \sigma_{g'}^2}{4L_g S} + \frac{5\ell_f \sigma_g}{2\sqrt{B}} + \frac{5}{4}\ell_f \delta.
\end{aligned}
$$

Next, we replace $\sum_{r=1}^i \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big]$ in the above inequality by $\sum_{r=1}^\tau \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big]$. Then summing up the above inequality for $i = 0, \ldots, \tau - 1$ gives

$$
\begin{aligned}
\frac{1}{24M}\sum_{i=0}^{\tau-1}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \;\leq\;\; & \mathbf{E}[\Phi(x_0^k)] - \mathbf{E}[\Phi(x_\tau^k)] \\
& - \left(\frac{M}{8} - \frac{3\tau\ell_f L_g}{4s} - \frac{5\tau\ell_f \ell_g^2}{4b\delta}\right)\sum_{r=1}^\tau \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big] \\
& + \left(\frac{3\ell_f \sigma_{g'}^2}{4L_g S} + \frac{5\ell_f \sigma_g}{2\sqrt{B}} + \frac{5}{4}\ell_f \delta\right)\tau.
\end{aligned}
$$

If we set $\delta = \frac{4\epsilon}{5\ell_f}$, $B = \frac{25\ell_f^2 \sigma_g^2}{4\epsilon^2}$, $S = \frac{3\ell_f \sigma_{g'}^2}{4L_g \epsilon}$, $s = \frac{12\ell_f L_g}{M}\tau$ and $b = \frac{20\ell_f \ell_g^2}{M\delta}\tau = \frac{25\ell_f^2 \ell_g^2}{M\epsilon}\tau$, then

$$\frac{M}{8} - \frac{3\tau \ell_f L_g}{4s} - \frac{5\tau \ell_g^2 \ell_f}{4b\delta} \geq 0 \quad \text{and} \quad \frac{3\ell_f \sigma_{g'}^2}{4L_g S} + \frac{5\ell_f \sigma_g}{2\sqrt{B}} + \frac{4}{5}\ell_f \delta \leq 3\epsilon.$$

Therefore,

$$\frac{1}{24M}\sum_{i=0}^{\tau-1}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \leq \mathbf{E}[\Phi(x_0^k)] - \mathbf{E}[\Phi(x_\tau^k)] + 3\tau\epsilon.$$

Summing up the above inequality for $k = 1, \ldots, K$, and dividing by $K\tau$, we obtain

$$\frac{1}{K\tau}\sum_{k=1}^{K}\sum_{i=0}^{\tau-1}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \leq 24M\left(\frac{\Phi(x_0^1) - \Phi_*}{K\tau} + 3\epsilon\right).$$

Since $x_{i*}^{k*}$ is randomly chosen from $\{x_i^k\}_{i=0,\ldots,\tau-1}^{k=1,\ldots,K}$, it satisfies (48). Moreover, in this case, we have $b = \mathcal{O}(\tau/\epsilon)$ and $s = \mathcal{O}(\tau)$. To find an $\epsilon$-stationary point in expectation, we further set $\tau = \epsilon^{-1/2}$ and $K = \mathcal{O}(\epsilon^{-1/2})$, which implies $\mathbf{E}\big[\|\mathcal{G}(x_{i*}^{k*})\|^2\big] \leq \mathcal{O}(\epsilon)$. Consequently, the sample complexity for the component mappings is

$$KB + K\tau b = \mathcal{O}(\epsilon^{-1/2}) \cdot \mathcal{O}(\epsilon^{-2}) + \mathcal{O}(\epsilon^{-1/2}) \cdot \epsilon^{-1/2} \cdot \mathcal{O}(\epsilon^{-3/2}) = \mathcal{O}(\epsilon^{-5/2}),$$

and the sample complexity for the Jacobians is

$$KS + K\tau s = \mathcal{O}(\epsilon^{-1/2}) \cdot \mathcal{O}(\epsilon^{-1}) + \mathcal{O}(\epsilon^{-1/2}) \cdot \epsilon^{-1/2} \cdot \mathcal{O}(\epsilon^{-1/2}) = \mathcal{O}(\epsilon^{-3/2}).$$

This finishes the proof. □

**Remark 3.** *The sample complexities $\mathcal{O}(\epsilon^{-5/2})$ and $\mathcal{O}(\epsilon^{-3/2})$, for component mappings and their Jacobians respectively, are also obtained using the SARAH estimator in [53], but for a slightly different stationarity measure. Specifically, their results are derived for finding a point $x$ that satisfies $\mathbf{E}[\|\widetilde{\mathcal{G}}(x)\|^2] \leq \epsilon$, where $\widetilde{\mathcal{G}}(x)$ is the approximate gradient mapping defined in (13) and (14). Since $\|\widetilde{\mathcal{G}}(x)\| = 0$ alone may not be a good measure for stationarity, [53] defined a primal-dual stationarity measure which requires additional conditions. In contrast, our results directly guarantee $\mathbf{E}[\|\mathcal{G}(x)\|^2] \leq \epsilon$, where $\mathcal{G}(x)$ is the (exact) gradient mapping defined in (11).*

## 5　The smooth and finite-average case

In this section, we consider problem (2) under the assumption that $f$ is smooth and convex. Specifically, we assume that the component mappings $g_i$ satisfy Assumption 4. For $f$ and $h$, in addition to Assumption 1, we make the following additional assumption.

**Assumption 6.** *The gradient of $f$, denoted as $f'$, is differentiable and $L_f$-Lipschitz continuous.*

Under Assumptions 4 and 6, the composite function $f \circ g$ is smooth and its gradient has a Lipschitz constant

$$L_{f\circ g} \triangleq \ell_f L_g + L_f \ell_g^2. \tag{51}$$

See [61] for a proof of this claim.

Algorithms for solving problem (2) under the above assumptions have been studied in [29, 35, 59, 61]. The best sample complexity is $\mathcal{O}(N + N^{1/2}\epsilon^{-1})$ obtained in [61], using the SARAH/SPIDER estimator for $g'(x)f'(g(x))$, which is the gradient of $f(g(x))$ by the chain rule. In this section, we study an algorithm using the proximal mapping of $f$ instead of the composite gradient. It is no surprising that we can attain the sample complexity here. Despite the same sample complexity in theory, it is often observed in practice that algorithms based on proximal mappings can be more efficient than those based on gradients (e.g., [1, 2, 16]). Therefore, it is very meaningful to establish the sample complexity of proximal-mapping based methods when $f$ is smooth.

We again apply the SARAH/SPIDER estimator to construct $\tilde{g}_i^k$ and $\tilde{J}_i^k$. For $i > 0$, we use (44) and (45), where $\xi$ is interpreted as a random index drawn from $\{1, \dots, N\}$ with replacement. For $i = 0$, we exploit the finite-average structure of $g$ by using the construction in (22) and (23), i.e.,

$$\tilde{g}_0^k = g(x_0^k) \qquad \text{and} \qquad \tilde{J}_0^k = g'(x_0^k). \tag{52}$$

This implies that $\mathbf{E}[\|\tilde{g}_0^k - g(x_0^k)\|^2] = 0$ and $= \mathbf{E}[\|\tilde{J}_0^k - g'(x_0^k)\|^2] = 0$, which can be substituted into Lemma 7 to get the following result.

**Corollary 1.** *Suppose Assumption 5 holds. Let $\tilde{g}_i^k$ and $\tilde{J}_i^k$ be constructed according to (52) for $i = 0$ and (44) and (45) for $i = 1, \dots, \tau - 1$. Then we have for $i = 0, 1, \dots, \tau - 1$,*

$$\mathbf{E}\big[\|\tilde{g}_i^k - g(x_i^k)\|^2\big] \leq \sum_{r=1}^{i} \frac{\ell_g^2}{|\mathcal{B}_r^k|} \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big],$$

$$\mathbf{E}\big[\|\tilde{J}_i^k - g'(x_i^k)\|^2\big] \leq \sum_{r=1}^{i} \frac{L_g^2}{|\mathcal{S}_r^k|} \mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big].$$

Next, we prove a descent property of the algorithm. The additional assumption that $f$ is smooth allows us to derive a tighter descent bound than Lemma 4. In particular, we can replace the term $2\ell_f \|\tilde{g}_i^k - g(x_i^k)\|$ in (27) with $L_f \|\tilde{g}_i^k - g(x_i^k)\|^2$, which leads to reduction of the sample complexity for the component mappings.

**Lemma 8.** *Suppose Assumptions 1, 2 and 6 hold. Then Algorithm 1 has the following descent property:*

$$\Phi(x_{i+1}^k) \leq \Phi(x_i^k) - \frac{M - 2L_{f \circ g}}{2M^2} \big\|\widetilde{\mathcal{G}}(x_i^k)\big\|^2$$
$$+ L_f \big\|\tilde{g}_i^k - g(x_i^k)\big\|^2 + \frac{\ell_f}{2L_g} \big\|\tilde{J}_i^k - g'(x_i^k)\big\|^2, \tag{53}$$

*where $L_{f \circ g}$ is defined in (51).*

*Proof.* We revisit the proof of Lemma 4. In particular, the inequality (29) still holds, i.e.,

$$\Phi(x_{i+1}^k) = \Phi(x_i^k) - \frac{M - \ell_f L_g}{2} \|x_{i+1}^k - x_i^k\|^2 + T_1 + T_2 + T_3. \tag{54}$$

Moreover, we can reuse the bound for $T_2$ in (30). and only need to rebound the terms $T_1$ and $T_3$.

Under Assumption 6, we denote the Hessian of $f$ as $f''$ and it holds that $\|f''(z)\| \le L_f$ for all $z \in \mathbf{R}^m$. For the ease of notation, we denote

$$\Delta_i^k \triangleq \tilde{g}_i^k - g(x_i^k), \qquad z_i^k \triangleq g(x_i^k) + g'(x_i^k)(x_{i+1}^k - x_i^k).$$

Starting with $T_1$, which is defined in (28), we use the second-order Taylor expansion of $f$ to obtain

$$
\begin{aligned}
T_1 &= f(z_i^k) - f(z_i^k + \Delta_i^k) \\
&= f(z_i^k) - \left( f(z_i^k) + \langle f'(z_i^k), \Delta_i^k \rangle + \frac{1}{2}(\Delta_i^k)^T f''(z_i^k + \theta \Delta_i^k)\Delta_i^k \right) \\
&= -\langle f'(z_i^k), \Delta_i^k \rangle - \frac{1}{2}(\Delta_i^k)^T f''(z_i^k + \theta \Delta_i^k)\Delta_i^k,
\end{aligned}
$$

where $\theta \in [0,1]$. Since $f$ is convex and the spectral norm of $f''$ is bounded by $L_f$, we have

$$
\begin{aligned}
T_1 &\le -\langle f'(z_i^k), \Delta_i^k \rangle \\
&\le -\langle f'(g(x_i^k)), \Delta_i^k \rangle + |\langle f'(g(x_i^k)) - f'(z_i^k), \Delta_i^k \rangle| \\
&\le -\langle f'(g(x_i^k)), \Delta_i^k \rangle + L_f \|g(x_i^k) - z_i^k\| \|\Delta_i^k\| \\
&= -\langle f'(g(x_i^k)), \Delta_i^k \rangle + L_f \|g'(x_i^k)(x_{i+1}^k - x_i^k)\| \|\Delta_i^k\|.
\end{aligned}
$$

Notice that by Assumption 4 we have $\|g'(x_i^k)\| \le \ell_g$, which gives

$$
\begin{aligned}
T_1 &\le -\langle f'(g(x_i^k)), \Delta_i^k \rangle + L_f \ell_g \|x_{i+1}^k - x_i^k\| \|\Delta_i^k\| \\
&\le -\langle f'(g(x_i^k)), \Delta_i^k \rangle + L_f \left( \frac{\ell_g^2}{2} \|x_{i+1}^k - x_i^k\|^2 + \frac{1}{2}\|\Delta_i^k\|^2 \right) \\
&= \frac{L_f}{2}\|\Delta_i^k\|^2 + \frac{L_f \ell_g^2}{2}\|x_{i+1}^k - x_i^k\|^2 - \langle f'(g(x_i^k)), \Delta_i^k \rangle. \qquad (55)
\end{aligned}
$$

For the term $T_3$ in (29), we have for some $\theta \in [0,1]$,

$$
\begin{aligned}
T_3 &= f(g(x_i^k) + \Delta_i^k) - f(g(x_i^k)) \\
&= f(g(x_i^k)) + \langle f'(g(x_i^k)), \Delta_i^k \rangle + \frac{1}{2}(\Delta_i^k)^T f''(g(x_i^k) + \theta \Delta_i^k)\Delta_i^k - f(g(x_i^k)) \\
&\le \frac{L_f}{2}\|\Delta_i^k\|^2 + \langle f'(g(x_i^k)), \Delta_i^k \rangle.
\end{aligned}
$$

Substituting the new bounds on $T_1$ and $T_3$ and the existing bound on $T_2$ in (30) into (54), we obtain

$$
\begin{aligned}
\Phi(x_{i+1}^k) &\le \Phi(x_i^k) - \left( \frac{M}{2} - \ell_f L_g - \frac{1}{2}L_f \ell_g^2 \right) \|x_{i+1}^k - x_i^k\|^2 \\
&\quad + L_f \|\tilde{g}_i^k - g(x_i^k)\|^2 + \frac{\ell_f}{2L_g} \|\tilde{J}_i^k - g'(x_i^k)\|^2.
\end{aligned}
$$

The desired result holds by noting the definitions of $L_{f \circ g}$ and $\widetilde{\mathcal{G}}(x_i^k)$. $\qquad\square$

Parallel to Lemma 2, we have the following result.

**Lemma 9.** *Suppose Assumptions 1, 2 and 6 hold. Let $x_{i+1}^k$ and $\hat{x}_{i+1}^k$ are defined in (13) and (15) respectively. Then we have*

$$\frac{M - L_{f\circ g}}{M^2}\big\|\mathcal{G}(x_i^k)\big\|^2 \;\;\leq\;\; \frac{2M + L_{f\circ g}}{M^2}\big\|\widetilde{\mathcal{G}}(x_i^k)\big\|^2 \tag{56}$$
$$+ 3L_f\big\|\tilde{g}_i^k - g(x_i^k)\big\|^2 + \frac{2\ell_f}{L_g}\big\|\tilde{J}_i^k - g'(x_i^k)\big\|^2.$$

*Proof.* We revisit the proof of Lemma 2, and start with the inequality (20), which is

$$M\big\|\hat{x}_{i+1}^k - x_{i+1}^k\big\|^2 \leq F(x_{i+1}^k; x_i^k) - \widetilde{F}(x_{i+1}^k; x_i^k) + \widetilde{F}(\hat{x}_{i+1}^k; x_i^k) - F(\hat{x}_{i+1}^k; x_i^k).$$

We can establish a tighter bound for the right-hand-side when $f$ is smooth. From the definitions of $F$ and $\widetilde{F}$ in (16) and (17) and the definitions of $T_1$ nd $T_2$ in (28), we have

$$F(x_{i+1}^k; x_i^k) - \widetilde{F}(x_{i+1}^k; x_i^k) \;\; = \;\; f\big(g(x_i^k) + g'(x_i^k)(x_{i+1}^k - x_i^k)\big) - f\big(\tilde{g}_i^k + \tilde{J}_i^k(x_{i+1}^k - x_i^k)\big)$$
$$= \;\; T_1 + T_2$$
$$\leq \;\; \frac{L_f}{2}\big\|\tilde{g}_i^k - g(x_i^k)\big\|^2 + \frac{\ell_f}{2L_g}\big\|\tilde{J}_i^k - g'(x_i^k)\big\|^2 \tag{57}$$
$$+ \frac{\ell_f L_g + L_f \ell_g^2}{2}\big\|x_{i+1}^k - x_i^k\big\|^2 - \big\langle f'\big(g(x_i^k)\big), \Delta_i^k \big\rangle,$$

where the last inequality is due to (30) and (55). Following similar arguments, we can derive

$$\widetilde{F}(\hat{x}_{i+1}^k; x_i^k) - F(\hat{x}_{i+1}^k; x_i^k) \;\; \leq \;\; L_f\big\|\tilde{g}_i^k - g(x_i^k)\big\|^2 + \frac{\ell_f}{2L_g}\big\|\tilde{J}_i^k - g'(x_i^k)\big\|^2 \tag{58}$$
$$+ \frac{\ell_f L_g + L_f \ell_g^2}{2}\big\|\hat{x}_{i+1}^k - x_i^k\big\|^2 + \big\langle f'\big(g(x_i^k)\big), \Delta_i^k \big\rangle.$$

Summing up (57) and (58) and noting the definition of $L_{f\circ g}$, we have

$$M\big\|\hat{x}_{i+1}^k - x_{i+1}^k\big\|^2 \;\; \leq \;\; \frac{3L_f}{2}\big\|\tilde{g}_i^k - g(x_i^k)\big\|^2 + \frac{\ell_f}{L_g}\big\|\tilde{J}_i^k - g'(x_i^k)\big\|^2$$
$$+ \frac{L_{f\circ g}}{2}\big\|\hat{x}_{i+1}^k - x_i^k\big\|^2 + \frac{L_{f\circ g}}{2}\big\|x_{i+1}^k - x_i^k\big\|^2.$$

Combining the above inequality with

$$M\big\|\hat{x}_{i+1}^k - x_i^k\big\|^2 \leq 2M\big\|x_{i+1}^k - x_i^k\big\|^2 + 2M\big\|\hat{x}_{i+1}^k - x_{i+1}^k\big\|^2$$

yields (following similar arguments at the end of proof for Lemma 2)

$$(M - L_{f\circ g})\big\|\hat{x}_{i+1}^k - x_i^k\big\|^2 \;\; \leq \;\; (2M + L_{f\circ g})\big\|x_{i+1}^k - x_i^k\big\|^2$$
$$+ 3L_f\big\|\tilde{g}_i^k - g(x_i^k)\big\|^2 + \frac{2\ell_f}{L_g}\big\|\tilde{J}_i^k - g'(x_i^k)\big\|^2.$$

Finally we obtain the desired result using the definitions of $\mathcal{G}(x_i^k)$ and $\widetilde{\mathcal{G}}(x_i^k)$. $\qquad\square$

The main result of this section is given by the following theorem.

**Theorem 4.** *Suppose Assumptions 1, 2, 3 and 6 hold for problem (2). In Algorithm 1, let $\tilde{g}_i^k$ and $\tilde{J}_i^k$ be constructed according to (52) for $i = 0$ and (44) and (45) for $i = 1, \ldots, \tau - 1$. If we choose $M \geq 4L_{f \circ g}$ and $\tau = \lceil \sqrt{N} \rceil$, and set the batch sizes $|\mathcal{B}_i^k| = |\mathcal{S}_i^k| = 2\lceil \sqrt{N} \rceil$ for $i = 1, \ldots, \tau - 1$, then*

$$\mathbf{E}\left[\left\|\mathcal{G}(x_{i*}^{k^*})\right\|^2\right] \quad \leq \quad \frac{24M\left(\Phi(x_0^1) - \Phi_*\right)}{K\tau}. \tag{59}$$

*The total sample complexity of reaching an $\epsilon$-stationary point in expectation is $\mathcal{O}(N + \sqrt{N}\epsilon^{-1})$.*

*Proof.* Under the assumption $M \geq 4L_{f \circ g}$, we have $\frac{1}{2} \cdot \frac{M - 2L_{f \circ g}}{2(2M + L_{f \circ g})} \in \left[\frac{1}{18}, \frac{1}{8}\right]$. Multiplying both sides of (56) by $\frac{1}{2} \cdot \frac{M - 2L_{f \circ g}}{2(2M + L_{f \circ g})}$ and adding them to (53), we obtain

$$
\begin{aligned}
\frac{1}{24M}\left\|\mathcal{G}(x_i^k)\right\|^2 \quad \leq \quad & \Phi(x_i^k) - \Phi(x_{i+1}^k) - \frac{M - 2L_{f \circ g}}{4M^2}\left\|\widetilde{\mathcal{G}}(x_i^k)\right\|^2 \\
& + \frac{11}{8}L_f\left\|\tilde{g}_i^k - g(x_i^k)\right\|^2 + \frac{3\ell_f}{4L_g}\left\|\tilde{J}_i^k - g'(x_i^k)\right\|^2.
\end{aligned} \tag{60}
$$

Taking expectation on both sides of the above inequality and applying Corollary 1, we have

$$
\begin{aligned}
\frac{1}{24M}\mathbf{E}\left[\left\|\mathcal{G}(x_i^k)\right\|^2\right] \quad \leq \quad & \mathbf{E}\left[\Phi(x_i^k)\right] - \mathbf{E}\left[\Phi(x_{i+1}^k)\right] - \frac{M - 2L_{f \circ g}}{4M^2}\mathbf{E}\left[\left\|\widetilde{\mathcal{G}}(x_i^k)\right\|^2\right] \\
& + \frac{11}{8}L_f\ell_g^2\sum_{r=1}^{i}\frac{1}{|\mathcal{B}_r^k|}\mathbf{E}\left[\left\|x_r^k - x_{r-1}^k\right\|^2\right] \\
& + \frac{3\ell_f L_g}{4}\sum_{r=1}^{i}\frac{1}{|\mathcal{S}_r^k|}\mathbf{E}\left[\left\|x_r^k - x_{r-1}^k\right\|^2\right].
\end{aligned}
$$

We will use constant batch sizes and let $|\mathcal{B}_i^k| = |\mathcal{S}_i^k| = S$ for all $k = 1, \ldots, K$ and $i = 1, \ldots, \tau - 1$. In addition, we can increase the summation from $\sum_{r=1}^{i}$ to $\sum_{r=1}^{\tau}$, which leads to

$$
\begin{aligned}
\frac{1}{24M}\mathbf{E}\left[\left\|\mathcal{G}(x_i^k)\right\|^2\right] \quad \leq \quad & \mathbf{E}\left[\Phi(x_i^k)\right] - \mathbf{E}\left[\Phi(x_{i+1}^k)\right] - \frac{M - 2L_{f \circ g}}{4M^2}\mathbf{E}\left[\left\|\widetilde{\mathcal{G}}(x_i^k)\right\|^2\right] \\
& + \frac{11}{8}\frac{L_f\ell_g^2}{S}\sum_{r=1}^{\tau}\mathbf{E}\left[\left\|x_r^k - x_{r-1}^k\right\|^2\right] \\
& + \frac{3\ell_f L_g}{4S}\sum_{r=1}^{\tau}\mathbf{E}\left[\left\|x_r^k - x_{r-1}^k\right\|^2\right].
\end{aligned}
$$

Plugging in $\widetilde{\mathcal{G}}(x_i^k) = -M(x_{i+1}^k - x_i^k)$ and noticing that

$$\frac{11}{8}\frac{L_f\ell_g^2}{S} + \frac{3\ell_f L_g}{4S} \leq \frac{3}{2}\frac{\ell_f L_g + L_f\ell_g^2}{S} = \frac{3}{2}\frac{L_{f \circ g}}{S},$$

we obtain

$$
\begin{aligned}
\frac{1}{24M}\mathbf{E}\left[\left\|\mathcal{G}(x_i^k)\right\|^2\right] \quad \leq \quad & \mathbf{E}\left[\Phi(x_i^k)\right] - \mathbf{E}\left[\Phi(x_{i+1}^k)\right] - \frac{M - 2L_{f \circ g}}{4}\mathbf{E}\left[\left\|x_{i+1}^k - x_i^k\right\|^2\right] \\
& + \frac{3}{2}\frac{L_{f \circ g}}{S}\sum_{r=1}^{\tau}\mathbf{E}\left[\left\|x_r^k - x_{r-1}^k\right\|^2\right].
\end{aligned}
$$

26

Summing up the above inequality for $i = 0, \ldots, \tau - 1$ yields

$$\frac{1}{24M} \sum_{i=0}^{\tau-1} \mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \leq \mathbf{E}\big[\Phi(x_0^k)\big] - \mathbf{E}\big[\Phi(x_0^{k+1})\big] - \left(\frac{M}{4} - \frac{2\tau}{S}L_{f \circ g}\right) \sum_{i=0}^{\tau-1} \mathbf{E}\big[\|x_{i+1}^k - x_i^k\|^2\big].$$

The choices of $M \geq 4L_{f \circ g}$, $\tau = \lceil\sqrt{N}\rceil$ and $S = 2\tau$ ensure $\frac{M}{4} - \frac{2\tau}{S}L_{f \circ g} \geq 0$. Therefore

$$\frac{1}{24M} \sum_{i=0}^{\tau-1} \mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \leq \mathbf{E}\big[\Phi(x_0^k)\big] - \mathbf{E}\big[\Phi(x_0^{k+1})\big].$$

Summing up the above inequality for $k = 1, \ldots, K$ and noticing the choice of $x_{i^*}^{k^*}$ in Algorithm 1, we obtain (59).

To get an $\epsilon$-stationary point in expectation, we need to set $K\tau = \mathcal{O}(\epsilon^{-1})$, which implies

$$K = \mathcal{O}(\tau^{-1}\epsilon^{-1}) = \mathcal{O}(N^{-1/2}\epsilon^{-1}).$$

Consequently, the sample complexity for both the component mappings and their Jacobians is

$$KN + K\tau S = \mathcal{O}(N^{-1/2}\epsilon^{-1}) \cdot N + \mathcal{O}(\epsilon^{-1}) \cdot 2N^{1/2} = \mathcal{O}(N + N^{1/2}\epsilon^{-1}).$$

This finishes the proof. $\square$

## 6 The smooth and expectation case

In this section we focus on problem (3) when $f$ is smooth and convex. Specifically, we proceed with Assumptions 1, 5 and 6. Under these assumptions, we still use the SARAH/SPIDER estimators in (43), (44) and (45). Since that the mean-square error bounds bounds on the estimators in Lemma 7 only depends on Assumption 5, they remain valid in this section. We have the following result.

**Theorem 5.** *Suppose Assumptions 1, 3, 5 and 6 hold for problem (3). Let the estimates $\tilde{g}_0^k$, $\tilde{J}_0^k$, $\tilde{g}_i^k$ and $\tilde{J}_i^k$ in Algorithm 1 be given in (43)-(45), and we choose $M \geq 4L_{f \circ g}$. For any $\epsilon > 0$, if we choose $\tau = \lceil\epsilon^{-1/2}\rceil$ and the batch sizes as*

$$|\mathcal{B}_0^k| = \left\lceil\frac{11L_f\sigma_g^2}{4\epsilon}\right\rceil, \qquad |\mathcal{S}_0^k| = \left\lceil\frac{3\ell_f^2\sigma_{g'}^2}{2L_g\epsilon}\right\rceil, \qquad |\mathcal{B}_i^k| = |\mathcal{S}_i^k| = 2\left\lceil\epsilon^{-1/2}\right\rceil,$$

*where $i = 1, \ldots, \tau - 1$, then we have*

$$\frac{1}{24M}\mathbf{E}[\|\mathcal{G}(x_{i^*}^{k^*})\|^2] \leq \frac{\Phi(x_0^1) - \Phi(x^*)}{K\tau} + \epsilon. \tag{61}$$

*Consequently, we have $\mathbf{E}\big[\|\mathcal{G}(x_{i^*}^{k^*})\|^2\big] = \mathcal{O}(\epsilon)$ by setting $K = \mathcal{O}(\epsilon^{-1/2})$, and the total sample complexity is $\mathcal{O}(\epsilon^{-3/2})$.*

*Proof.* We choose batch sizes that do not depend on $k$. For the ease of notation, let $|\mathcal{B}_0^k| = B$ and $|\mathcal{S}_0^k| = S$, and $|\mathcal{B}_i^k| = |\mathcal{S}_i^k| = b$ for $i = 1, \ldots, \tau - 1$. Taking expectation of both sizes of (60) and applying Lemma 7, we get

$$
\begin{aligned}
\frac{1}{24M}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \;\leq\; & \mathbf{E}\big[\Phi(x_i^k)\big] - \mathbf{E}\big[\Phi(x_{i+1}^k)\big] - \frac{M - 2L_{f \circ g}}{4M^2}\mathbf{E}\big[\|\widetilde{\mathcal{G}}(x_i^k)\|^2\big] \\
& + \frac{11}{8}\frac{L_f \sigma_g^2}{B} + \frac{3\ell_f \sigma_{g'}^2}{4L_g S} + \frac{11}{8}\frac{L_f \ell_g^2}{b}\sum_{r=1}^{i}\mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big] \\
& + \frac{3\ell_f L_g}{4b}\sum_{r=1}^{i}\mathbf{E}\big[\|x_r^k - x_{r-1}^k\|^2\big].
\end{aligned}
$$

Summing up the above inequality for $i = 0, \ldots, \tau - 1$ and following similar steps as in the proof of Theorem 4, we have

$$
\begin{aligned}
\frac{1}{24M}\sum_{i=0}^{\tau-1}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \;\leq\; & \mathbf{E}\big[\Phi(x_0^k)\big] - \mathbf{E}\big[\Phi(x_0^{k+1})\big] + \tau\left(\frac{11}{8}\frac{L_f \sigma_g^2}{B} + \frac{3\ell_f \sigma_{g'}^2}{4L_g S}\right) \\
& - \left(\frac{M}{4} - \frac{2\tau}{b}L_{f \circ g}\right)\sum_{i=0}^{\tau-1}\mathbf{E}\big[\|x_{i+1}^k - x_i^k\|^2\big].
\end{aligned}
$$

The choices of $M \geq 4L_{f \circ g}$ and $b = 2\tau$ ensure $\frac{M}{4} - \frac{2\tau}{b}L_{f \circ g} \geq 0$, and choices of $B = \frac{11L_f \sigma_g^2}{4\epsilon}$ and $S = \frac{3\ell_f^2 \sigma_{g'}^2}{2L_g \epsilon}$ further ensure the constant term to be less than $\tau\epsilon$. Therefore

$$
\frac{1}{24M}\sum_{i=0}^{\tau-1}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \leq \mathbf{E}\big[\Phi(x_0^k)\big] - \mathbf{E}\big[\Phi(x_0^{k+1})\big] + \tau\epsilon, \tag{62}
$$

which, upon summing over $k = 1, \ldots, K$ and noting the choice of $x_{i^*}^{k^*}$, yields (61). The sample complexities can be calculated as $KB + K\tau b$. $\qquad\square$

In Theorem 5, the choices of $\tau$ and batch sizes all depend on a fixed accuracy $\epsilon$, which can be hard to determine in advance in many situations, and running more iterations will not improve the solution due to the existence of a $\mathcal{O}(\epsilon)$ bias term in (61). Therefore, it would be desirable to develop an algorithm that adaptively chooses the batch sizes to keep improving the accuracy of the solution. Such an adaptive scheme is presented in the following theorem.

**Theorem 6.** *Suppose Assumptions 1, 3, 5 and 6 hold for problem* (3). *Let the estimates $\tilde{g}_0^k$, $\tilde{J}_0^k$, $\tilde{g}_i^k$ and $\tilde{J}_i^k$ in Algorithm 1 be given in* (43)-(45), *and we choose $M \geq 4L_{f \circ g}$. Let $\{\epsilon_k\}_{k=1}^{\infty}$ be a sequence of positive real numbers. If we run each epoch of Algorithm 1 for $\tau_k = \epsilon_k^{-1/2}$ iterations, and set the batch sizes to be*

$$
|\mathcal{B}_0^k| = \left\lceil \frac{11L_f \sigma_g^2}{4\epsilon_k} \right\rceil, \qquad |\mathcal{S}_0^k| = \left\lceil \frac{3\ell_f^2 \sigma_{g'}^2}{2L_g \epsilon_k} \right\rceil, \qquad |\mathcal{B}_i^k| = |\mathcal{S}_i^k| = 2\left\lceil \epsilon_k^{-1/2} \right\rceil,
$$

*where $i = 1, \ldots, \tau - 1$, then we have*

$$
\frac{1}{20M}\mathbf{E}[\|\mathcal{G}(x_{i^*}^{k^*})\|^2] \;\leq\; \frac{\Phi(x_0^1) - \Phi(x^*)}{\sum_{k=1}^{K}\tau_k} + \frac{\sum_{k=1}^{K}\epsilon_k^{1/2}}{\sum_{k=1}^{K}\tau_k} \tag{63}
$$

28

*Specifically, setting $\epsilon_k = k^{-2}$ results in*

$$\frac{1}{20M}\mathbf{E}[\|\mathcal{G}(x_{i*}^{k*})\|^2] = \mathcal{O}\left(\frac{\ln K}{K^2}\right). \tag{64}$$

*Consequently, given any $\epsilon > 0$, we can set $K = \epsilon^{-1/2}$, which leads to an $\mathcal{O}\left(\epsilon \ln \frac{1}{\epsilon}\right)$-stationary solution with total sample complexity of $\mathcal{O}(\epsilon^{-3/2})$.*

*Proof.* Note that the inequality (62) still holds but with a specific set of parameters for each $k$. Specifically, we have

$$\frac{1}{24M}\sum_{i=0}^{\tau_k-1}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \le \mathbf{E}\big[\Phi(x_0^k)\big] - \mathbf{E}\big[\Phi(x_0^{k+1})\big] + \tau_k\epsilon_k, \qquad k = 1, \ldots, K.$$

Since we choose $\tau_k = \epsilon_k^{-1/2}$, it holds that $\tau_k\epsilon_k = \epsilon_k^{1/2}$. Summing this up over $k$ gives

$$\frac{1}{24M}\sum_{k=1}^{K}\sum_{i=0}^{\tau_k-1}\mathbf{E}\big[\|\mathcal{G}(x_i^k)\|^2\big] \le \Phi(x_0^1) - \Phi_* + \sum_{k=1}^{K}\epsilon_k^{1/2}.$$

Because $x_{i*}^{k*}$ is randomly chosen from $\{x_i^k\}_{i=0,\ldots,\tau-1}^{k=1,\ldots,K}$, we conclude (63) holds.

If we choose $\epsilon_k = k^{-2}$, then $\tau_k = \epsilon^{-1/2} = k$ and we have

$$\sum_{k=1}^{K}\tau_k = \frac{1}{2}K(K+1) \quad \text{and} \quad \sum_{k=1}^{K}\epsilon_k^{1/2} = \sum_{k=1}^{k}k^{-1} \le 1 + \int_1^K z^{-1}dz = 1 + \ln K.$$

Substituting the above relationships into (63) yields (64). The total sample complexity for the $g_\xi$'s for running these $K$ epochs will be

$$\sum_{k=1}^{K}\left(|\mathcal{B}_0^k| + \tau_k|\mathcal{B}_1^k|\right) = \mathcal{O}\left(\sum_{k=1}^{K}\epsilon_k^{-1}\right) = \mathcal{O}(K^3).$$

Similarly, the sample complexity for the Jacobians is also $\mathcal{O}(K^3)$. Finally by setting $K = \epsilon^{-1/2}$, we will get an $\mathcal{O}\left(\epsilon \ln \frac{1}{\epsilon}\right)$-stationary solution with total sample complexity of $\mathcal{O}(\epsilon^{-3/2})$. $\qquad\square$

# 7 Numerical experiments

In this section, we present the numerical experiments of our methods and compare with related methods (following the experiment setup in [53]). For the ease of reference, we denote the algorithms in comparison as follows:

- PL: the deterministic prox-linear algorithm described by (4).

- S-PL: the mini-batch stochastic prox-linear algorithm in Algorithm 2. We note that S-PL coincides with SGN method in the concurrent work [53].

- SVR-PL: Algorithm 1 where the SVRG estimator is applied and augmented with 1st-order correction technique.

- Sarah-PL: Algorithm 1 using the SARAH estimator. Specifically, when $f$ is nonsmooth, Sarah-PL overlaps with SGN2 [53].

- When $f$ is smooth, we also compare with the CIVR method [61] and the N-Spider method [60] for stochastic composite optimization.

## 7.1 Nonsmooth nonlinear systems

In this experiment, we solve the following nonsmooth problem:

$$\underset{x \in \mathbf{R}^n}{\text{minimize}} \ \Phi(x) := \Big\| \frac{1}{N} \sum_{j=1}^N g_j(x) \Big\|_1 + \beta \|x\|_1,$$

where we want to find a sparse $x$ s.t. $\frac{1}{N} \sum_{j=1}^N g_j(x)$ is close to 0. Let $a_j \in \mathbf{R}^n$ be the $j$-th data point and $b_j \in \{-1, +1\}$ be the corresponding label. The function $g_j : \mathbf{R}^n \mapsto \mathbf{R}^4$ is defined as

$$g_j(x) = \begin{bmatrix} 1 - \tanh(z_j) \\ \left(1 - \frac{1}{1+e^{-z_j}}\right)^2 \\ \log\left(1 + e^{-z_j}\right) - \log\left(1 + e^{-z_j - 1}\right) \\ \log\left(1 + (z_j - 1)^2\right) \end{bmatrix} \quad \text{with} \quad z_j = b_j \cdot a_j^T x, \tag{65}$$

where each row of $g_j$ corresponds a certain type of binary classification loss, which can be viewed as a mixture of multiple models. We test our methods with the ijcnn1 dataset[1] and the MNIST dataset[2]. For ijcnn1, we randomly extract $N = 10000$ data points. For MNIST, we extract $N = 10000$ data points of two digits (Figure 1 shows the plots for "1" and "9"). Specifically, each data point $a_j$ in the ijcnn1 dataset is 22 dimensional, we set $\beta = 0$ for ijcnn1 dataset, meaning that we do not require the solution to be sparse. For the MNIST dataset, each $a_j$ are 784 dimensional where most entries are 0. In this case, we set $\beta = N^{-1}$ as the sparsity penalty parameter.

In the experiment, we test PL, S-PL/SGN, SVR-PL and Sarah-PL/SGN2 algorithms. For SVR-PL and Sarah-PL, we estimate $g$ and $g'$ with the mini-batch sizes suggested in Remark 1. Specifically, for SVR-PL, we choose $|\mathcal{S}_i^k| = |\mathcal{B}_i^k| = \lceil cN^{4/5} \rceil$. For Sarah-PL, we choose $\epsilon = 10^{-2}$, therefore we choose the large batches to be $|\mathcal{S}_0^k| = |\mathcal{B}_0^k| = \epsilon^{-2} = N$. For this finite sum problem we slightly revise the Sarah-PL such that $\tilde{g}_0^k = g(x_0^k)$ and $\tilde{J}_0^k = g'(x_0^k)$ and we set $|\mathcal{S}_i^k| = |\mathcal{B}_i^k| = \lceil c\epsilon^{-3/2} \rceil$ for $i > 0$. For both SVR-PL and Sarah-PL, $c$ is set to be $c = 0.1$ after tuning from the set $\{0.01, 0.05, 0.1, 0.5, 1, 2\}$. For S-PL, the batch size is set to be 500. For all methods, we select the best performing $M$ from the discrete range $\{1, 5, 10, 20, 40, 60, 80, 100\}$. For ijcnn1 dataset, $M = 1$ works best for all methods; For MNIST dataset, $M = 40$ works best for all methods. All methods start from the initial solution $x = 0$.

The results are shown in Figure 1, where each curve is plotted by averaging 5 rounds of running an algorithm. We can see that in terms of sample complexity, all stochastic methods significantly outperforms the deterministic PL algorithm. Among the stochastic methods, SVR-PL and Sarah-PL perform better than S-PL (mini-batch only). Sarah-PL performs the best for this particular experiment, benefiting from using $|\mathcal{S}_0^k| = |\mathcal{B}_0^k| = N$ in the finite-sum setting, even though we do not have theory to support its advantage.

---

[1] https://www.csie.ntu.edu.tw/ cjlin/libsvmtools/datasets/binary.html
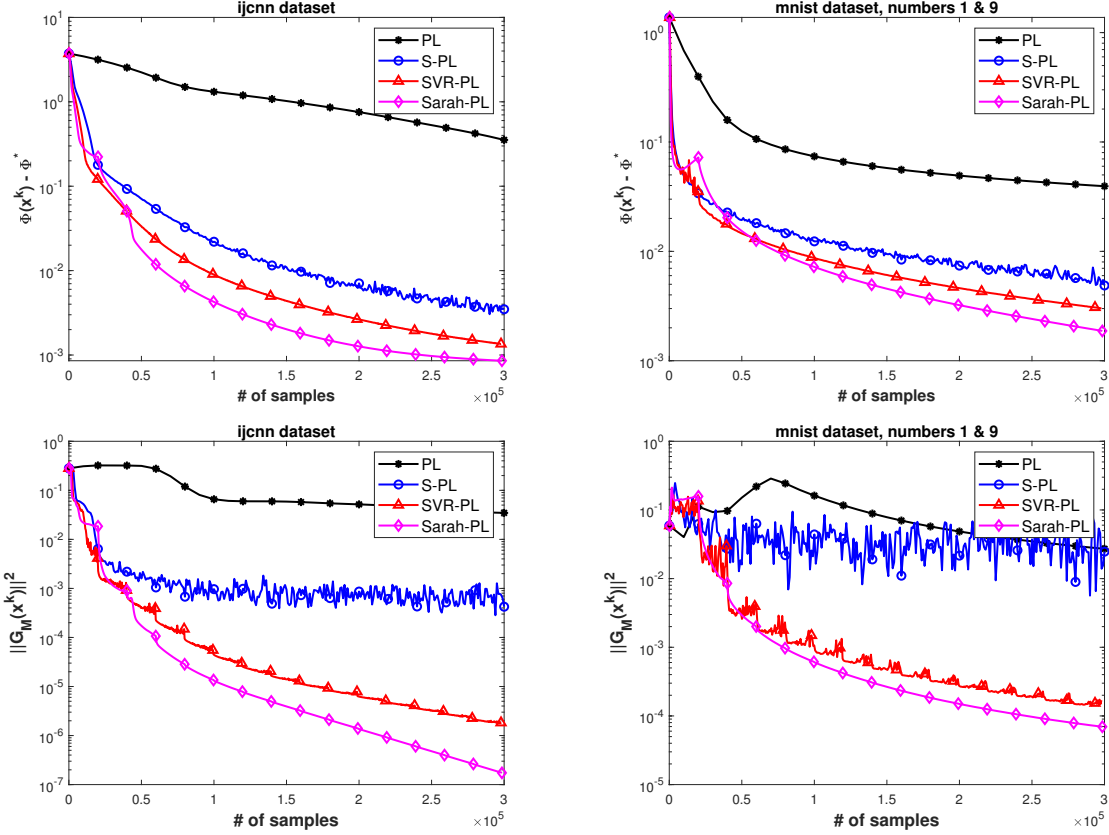[2] http://yann.lecun.com/exdb/mnist/

Figure 1: Comparison of the deterministic and stochastic prox-linear methods for nonsmooth composite optimization: the left column is for the ijcnn1 dataset and the right column is for the MNIST dataset (digits "1" and "9"). The first row shows the decrease of objective gap versus number of samples (where $\Phi^*$ is approximated by collecting the lowest value after running all algorithms for a much longer time). The second row shows squared norm of the (exact) gradient mapping (computed off-line using the full dataset).

## 7.2 Smooth nonlinear systems

In this experiment, we solve the following smooth problem:

$$\underset{x \in \mathbf{R}^n}{\text{minimize}} \ \ \Phi(x) := \left\| \frac{1}{N} \sum_{j=1}^{N} g_j(x) \right\|^2,$$

where $g_j : \mathbf{R}^n \mapsto \mathbf{R}^4$ is defined by (65). We compare Sarah-PL, CIVR [61], and the N-Spider algorithm [60]. For all three methods, the batch sizes are set to be $\lceil \sqrt{N} \rceil$. For N-Spider, we set $\epsilon_i^k = \frac{10}{1+k}$ for ijcnn1 and $\epsilon_i^k = \frac{10^{-2}}{1+k}$ for MNIST after some tuning. For both Sarah-PL and CIVR, their parameter $M$ or step size $\eta = M^{-1}$ are chosen from the set $\{0.1, 0.5, 1, 5, 10, 20, 40, 60\}$. For ijcnn1, Sarah-PL works best with $M = 0.1$ and CIVR works best with $\eta = 0.5^{-1}$. For MNIST, Sarah-PL works best with $M = 10$ and CIVR works best with $\eta = 20^{-1}$.

The results are shown in Figure 2, where each curve is plotted by averaging 5 rounds of running
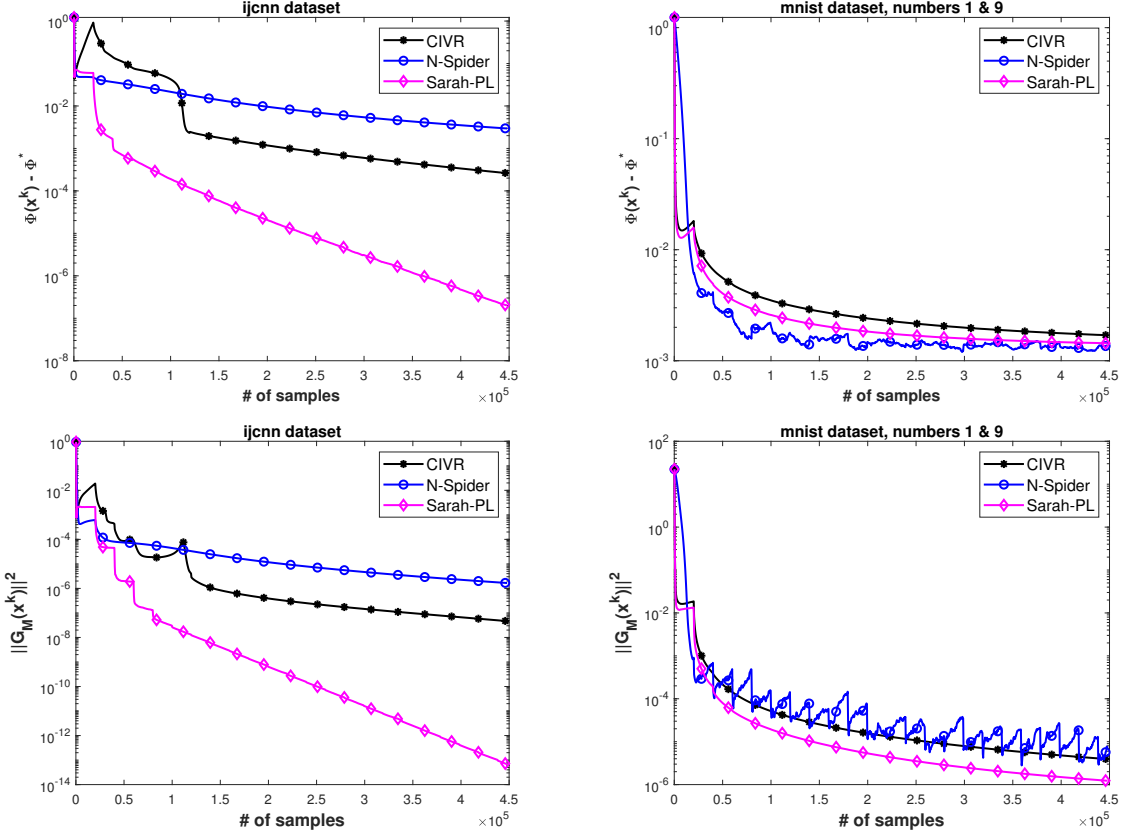
Figure 2: Comparison of stochastic variance-reduction methods for smooth composite optimization: the left column is for the ijcnn1 dataset and the right column is for the MNIST dataset (digits "1" and "9"). The first row shows objective gap versus number of samples (where $\Phi^*$ is approximated by collecting the lowest value after running all algorithms for a much longer time). The second row shows squared norm of the (exact) gradient mapping (computed off-line using the full dataset).

an algorithm. For ijcnn1, Sarah-PL significantly outperforms CIVR and N-Spider, demonstrating the potential advantage of prox-linear algorithms over chain-rule based methods (both with variance reduction). For MNIST, all three methods performs similarly.

## 7.3 Constrained stochastic optimization through penalty method

We consider a risk-sensitive portfolio optimization problem. Let $r_i \in \mathbf{R}^d$ be the vector of expected reward of $d$ stocks at time period $i$, for $i = 1, 2, ..., N$. The problem of maximizing the expected total reward across $N$ periods, with a constraint on the conditional value at risk (CVaR) is formulated as [49, 33]

$$\underset{x \in \Delta_d, \tau \in \mathbf{R}}{\text{maximize}} \ \left( \frac{1}{N} \sum_{i=1}^{N} r_i \right)^T x \qquad \text{s.t.} \qquad \tau + \frac{1}{\beta N} \sum_{i=1}^{N} \max\{-r_i^T x - \tau, 0\} \leq 0,$$
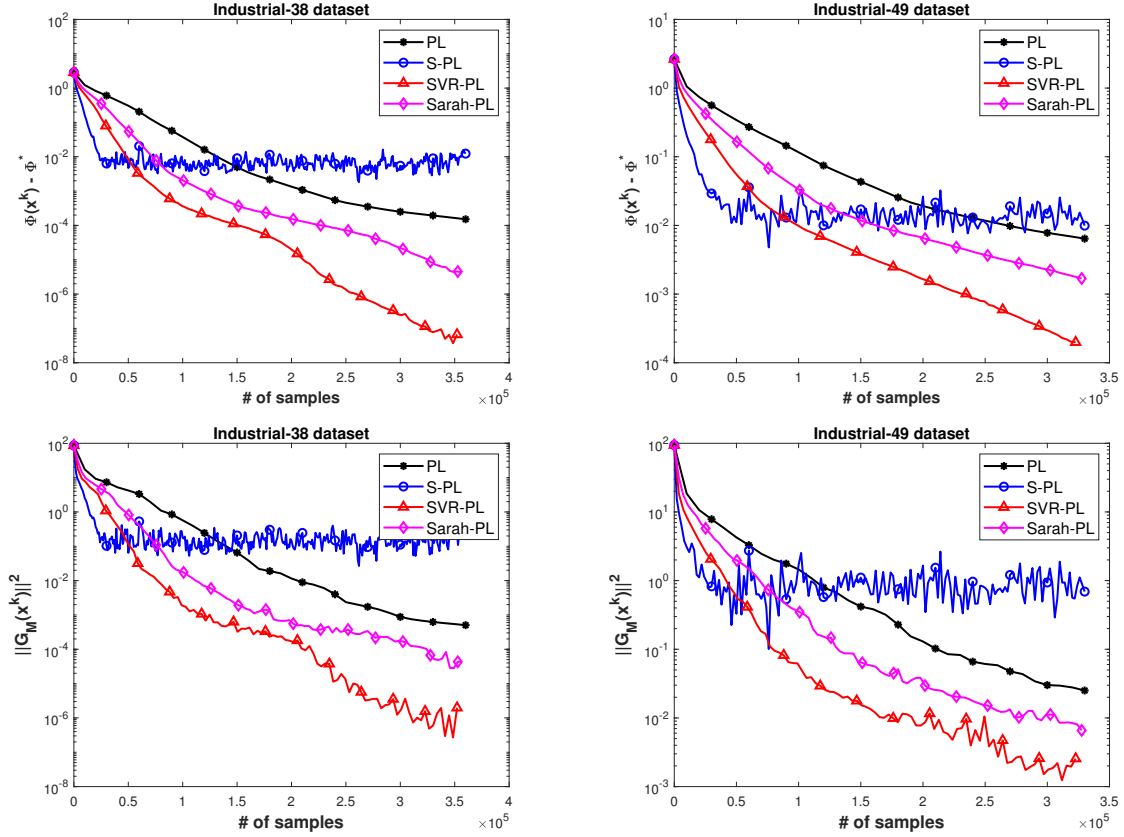
32

Figure 3: Comparison of different prox-linear algorithms for constrained stochastic optimization through penalty formulation: the left column is for the Industrial-38 dataset and the right column is for the Industrial-49 dataset. The first row shows objective gap versus number of samples and the second row shows squared norm of the (exact) gradient mapping.

where $\Delta_d := \left\{ x : x \geq 0, \ \sum_{j=1}^d x_j = 1 \right\}$ is the probability simplex. Using the exact penalty method (Section 1.1), this problem can be reformulated as

$$\operatorname*{minimize}_{x \in \Delta_d, \tau \in \mathbf{R}} \ -\left( \frac{1}{N} \sum_{i=1}^N r_i \right)^T x + \rho \cdot \max\left\{ 0, \tau + \frac{1}{\beta N} \sum_{i=1}^N \max\{-r_i^T x - \tau, 0\} \right\}.$$

Following the suggestion of [53], the nonsmooth term $\beta^{-1} \cdot \max\{-r_i^T x - \tau, 0\}$ is smoothed as $\frac{1}{2\beta} \left( \sqrt{(r_i^T x + \tau)^2 + \gamma^2} - r_i^T x - \tau - \gamma \right)$.

In this experiment, we test different methods on the Industrial-38 and the Industrial-49 dataset[3]. From each dataset, $N = 10000$ data points are extracted for the experiment. The parameters in the problem formulation are set to be $\beta = 10^{-1}$, $\rho = 5$, and $\gamma = 10^{-3}$. For algorithmic parameters, their tuning process is the same as that described in Section 7.1. The following results are obtained. In Industrial-38 dataset, $M = 40, 60, 40, 60$ works best for PL, S-PL,SVR-PL and

---

[3]http://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data_library.html

Sarah-PL respectively; For S-PL, the batch size is chosen to be 1000; For both SVR-PL and Sarah-PL, the batch sizes are the same as those in Section 7.1 with $c = 2$.

Figure 3 shows the results, again averaged over 5 runs of each algorithm. In this experiment, SVR-PL performs the best. It is worth noting that although S-PL has fast convergence in the initial stage, it stagnates at a relatively high error floor. SVR-PL and Sarah-PL reach higher accuracy due to their advanced variance-reduction schemes.

# 8 Discussions

In this paper, we have mostly relied on the SARAH/SPIDER estimators for variance reduction, except that for the nonsmooth and finite-average case (Section 3) we used a modified SVRG estimator with first-order correction. If we use the SVRG type of estimators for other cases, then the resulting sample complexities are suboptimal. More specifically, when $f$ is smooth, we have derived sample complexity of $\mathcal{O}(N + N^{2/3}\epsilon^{-1})$ and $\mathcal{O}(\epsilon^{-5/3})$ for the cases of $g$ being a finite average and a general expectation respectively. They are inferior compared to the $\mathcal{O}(N + \sqrt{N}\epsilon^{-1})$ and $\mathcal{O}(\epsilon^{-3/2})$ bounds using the SARAH/SPIDER estimators obtained in Sections 5 and 6.

The sample complexities of our methods for smooth $f$ are the same as the stochastic gradient descent type of methods that use the chain-rule to construct gradient estimators [61, 60]. However, it is often observed that algorithms based on proximal mappings can be more efficient than those based on gradients in practice (see e.g., [1, 2, 16]). Here we shed more light from a theoretical perspective. Consider the least squares problem of minimizing $F(x) := \frac{1}{2}\|g(x)\|^2$, where $g(x) = \frac{1}{N}\sum_{i=1}^{N} g_i(x)$. Given any SARAH/SPIDER variance reduced estimator $\tilde{g}_i^k$ and $\tilde{J}_i^k$, our proxi-linear scheme construct the update as

$$x_{i+1}^k = x_i^k - \left(M \cdot I + [\tilde{J}_i^k]^T \tilde{J}_i^k\right)^{-1} [\tilde{J}_i^k]^T \tilde{g}_i^k,$$

which is a *damped Gauss-Newton* iteration. Note that if $g(x^*) \approx 0$, then $\nabla^2 F(x^*) \approx [g'(x^*)]^T g'(x^*)$ (see [39]). This indicates that the Gauss-Newton matrix $[\tilde{J}_i^k]^T \tilde{J}_i^k$ becomes a better approximation of the Hessian $\nabla^2 F(x_i^k)$ as $x_i^k$ moves closer to $x^*$. Therefore, prox-linear based methods (Gauss-Newton especially) can take advantage of the second-order information whenever possible, while chain-rule based gradient methods cannot.

It is worth noting that both SVRG and SARAH/SPIDER schemes need a large sample batch at the beginning of each epoch, and slightly smaller sample batches in later iterations. However, under many circumstances it is more preferable if constant small batches are taken in each iteration. Recently, a STOchastic Recursive Momentum (STORM) variance reduction scheme that takes one sample per iteration has been proposed to solve smooth stochastic programming problem [14], and has been extended to a distributionally robust optimization (DRO) problem of form (3) with $f$ being smooth [45]. An optimal $\mathcal{O}(\epsilon^{-3/2})$ sample complexity is achieved in these works. However, we were not able to extend the STORM technique to problems with nonsmooth $f$. Deriving an algorithm with (constant) small mini-batch sizes for problems (2) and (3) with nonsmooth $f$ remains open.

## Acknowledgments

# References

[1] Hilal Asi and John C. Duchi. The importance of better models in stochastic optimization. *Proceedings of the National Academy of Sciences*, 116(46):22924–22930, 2019.

[2] Hilal Asi and John C. Duchi. Stochastic (approximate) proximal point methods: Convergence, optimality, and adaptivity. *SIAM Journal on Optimization*, 29(3):2257–2290, 2019.

[3] Yu Bai, John Duchi, and Song Mei. Proximal algorithms for constrained composite optimization, with applications to solving low-rank sdps. *arXiv preprint arXiv:1903.00184*, 2019.

[4] Dimitri P Bertsekas. Approximation procedures based on the method of multipliers. *Journal of Optimization Theory and Applications*, 23(4):487–510, 1977.

[5] Dimitri P Bertsekas. Nonlinear programming. *Journal of the Operational Research Society*, 48(3):334–334, 1997.

[6] Jose Blanchet, Donald Goldfarb, Garud Iyengar, Fengpei Li, and Chaoxu Zhou. Unbiased simulation for optimizing stochastic function compositions. *arXiv preprint arXiv:1711.07564*, 2017.

[7] James V Burke. Descent methods for composite nondifferentiable optimization problems. *Mathematical Programming*, 33(3):260–279, 1985.

[8] James V Burke and Abraham Engle. Line search methods for convex-composite optimization. *arXiv preprint arXiv:1806.05218*, 2018.

[9] James V Burke and Abraham Engle. Strong metric (sub) regularity of KKT mappings for piecewise linear-quadratic convex-composite optimization. *arXiv preprint arXiv:1805.01073*, 2018.

[10] James V Burke and Michael C Ferris. A Gauss-Newton method for convex composite optimization. *Mathematical Programming*, 71(2):179–194, 1995.

[11] Coralia Cartis, Nicholas IM Gould, and Philippe L Toint. On the evaluation complexity of composite function minimization with applications to nonconvex nonlinear programming. *SIAM Journal on Optimization*, 21(4):1721–1739, 2011.

[12] Vasileios Charisopoulos, Damek Davis, Mateo Díaz, and Dmitriy Drusvyatskiy. Composite optimization for robust blind deconvolution. *arXiv preprint arXiv:1901.01624*, 2019.

[13] Rixon Crane and Fred Roosta. DINGO: Distributed newton-type method for gradient-norm optimization. In *Advances in Neural Information Processing Systems 32*, pages 9498–9508. Curran Associates, Inc., 2019.

[14] Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex sgd. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

[15] Christoph Dann, Gerhard Neumann, and Jan Peters. Policy evaluation with temporal differences: a survey and comparison. *Journal of Machine Learning Research*, 15(1):809–883, 2014.

[16] Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.

[17] Damek Davis and Benjamin Grimmer. Proximally guided stochastic subgradient method for nonsmooth, nonconvex problems. *SIAM Journal on Optimization*, 29(3):1908–1930, 2019.

[18] Ofer Dekel, Ran Gilad-Bachrach, Ohad Shamir, and Lin Xiao. Optimal distributed online prediction using mini-batches. *Journal of Machine Learning Research*, 13:165–202, 2012.

[19] Dmitriy Drusvyatskiy. The proximal point method revisited. *SIAG/OPT Views and News (A Forum for the SIAM Activity Group on Optimization)*, 26(1):1–8, 2017.

[20] Dmitriy Drusvyatskiy and Adrian S Lewis. Error bounds, quadratic growth, and linear convergence of proximal methods. *Mathematics of Operations Research*, 2018.

[21] Dmitriy Drusvyatskiy and Courtney Paquette. Efficiency of minimizing compositions of convex functions and smooth maps. *Mathematical Programming*, 178:503–558, 2019.

[22] John C Duchi and Feng Ruan. Stochastic methods for composite and weakly convex optimization problems. *SIAM Journal on Optimization*, 28(4):3229–3259, 2018.

[23] John C Duchi and Feng Ruan. Solving (most) of a set of quadratic equalities: Composite optimization for robust phase retrieval. *Information and Inference: A Journal of the IMA*, 8(3):471–529, 2019.

[24] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In *Advances in Neural Information Processing Systems*, pages 687–697, 2018.

[25] Roger Fletcher and G Alistair Watson. First and second order conditions for a class of non-differentiable optimization problems. *Mathematical Programming*, 18(1):291–307, 1980.

[26] Saeed Ghadimi, Andrzej Ruszczyński, and Mengdi Wang. A single time-scale stochastic approximation method for nested stochastic optimization. Preprint, arXiv:1812.01094, 2018.

[27] Tamir Hazan, Shoham Sabach, and Sergey Voldman. Stochastic proximal linear method for structured non-convex optimization. *Optimization Methods and Software*, 35:921–937, 2020.

[28] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Convex analysis and minimization algorithms I: Fundamentals*, volume 305. Springer science & business media, 2013.

[29] Zhouyuan Huo, Bin Gu, Ji Jiu, and Heng Huang. Accelerated method for stochastic composition optimization with nonsmooth regularization. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*, pages 3287–3294, 2018.

[30] A. N. Iusem, A. Jofré, R. I. Oliveira, and P. Phompson. Extragradient method with variance reduction for stochastic variational inequalities. *SIAM Journal on Optimization*, 27(2):686–724, 2017.

[31] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.

[32] J. Koshal, A. Nedić, and U. B. Shanbhag. Regularized iterative stochastic approximation methods for stochastic variational inequality problems. *IEEE Transactions on Automatic Control*, 58(3):594–609, 2013.

[33] Guanghui Lan and Zhiqiang Zhou. Algorithms for stochastic optimization with function or expectation constraints. *Computational Optimization and Applications*, pages 1–38, 2020.

[34] Adrian S Lewis and Stephen J Wright. A proximal method for composite minimization. *Mathematical Programming*, 158(1-2):501–546, 2016.

[35] Xiangru Lian, Mengdi Wang, and Ji Liu. Finite-sum composition optimization via variance reduced gradient descent. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 1159–1167, 2017.

[36] Yu Nesterov. Modified gauss–newton scheme with worst case guarantees for global performance. *Optimisation Methods and Software*, 22(3):469–483, 2007.

[37] Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2613–2621. JMLR. org, 2017.

[38] Lam M Nguyen, Marten van Dijk, Dzung T Phan, Phuong Ha Nguyen, Tsui-Wei Weng, and Jayant R Kalagnanam. Optimal finite-sum smooth non-convex optimization with SARAH. *arXiv preprint arXiv:1901.07648*, 2019.

[39] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, 2nd edition, 2006.

[40] Peter Ochs, Jalal Fadili, and Thomas Brox. Non-smooth non-convex Bregman minimization: Unification and new algorithms. *Journal of Optimization Theory and Applications*, 181:244–278, 2019.

[41] James M. Ortega and Werner C. Rheinboldt. *Iterative Solution of Nonlinear Equations in Several Variables*. Academic Press, New York, 1970.

[42] Nhan H Pham, Lam M Nguyen, Dzung T Phan, and Quoc Tran-Dinh. ProxSARAH: An efficient algorithmic framework for stochastic composite nonconvex optimization. *arXiv preprint arXiv:1902.05679*, 2019.

[43] BT Poljak. On the Bertsekas' method for minimization of composite functions. In *International Symposium on Systems Optimization and Analysis*, pages 179–186. Springer, 1979.

[44] Boris T. Polyak. *Introduction to Optimization*. Optimization Software, Inc., 1987.

[45] Qi Qi, Zhishuai Guo, Yi Xu, Rong Jin, and Tianbao Yang. A practical online method for distributionally deep robust optimization. *arXiv preprint arXiv:2006.10138*, 2020.

[46] Hassan Rafique, Mingrui Liu, Qihang Lin, and Tianbao Yang. Non-convex min-max optimization: Provable algorithms and applications in machine learning. *arXiv:1810.02060*, 2018.

[47] Hamed Rahimian and Sanjay Mehrotra. Distributionally robust optimization: a review. arXiv preprint, arXiv:1908.05659, 2019.

[48] R. Tyrrell Rockafellar. Coherent approaches to risk in optimization under uncertainty. *IN-FORMS TutORials in Operations Research*, 2007.

[49] R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of risk*, 2:21–42, 2000.

[50] Fred Roosta, Yang Liu, Peng Xu, and Michael W Mahoney. Newton-MR: Newton's method without smoothness or convexity. *arXiv preprint arXiv:1810.00303*, 2018.

[51] Andrzej Ruszczyński. Advances in risk-averse optimization. *INFORMS TutORials in Operation Research*, 2013.

[52] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 1998.

[53] Quoc Tran-Dinh, Nhan Pham, and Lam Nguyen. Stochastic gauss-newton algorithms for nonconvex compositional optimization. In *International Conference on Machine Learning*, pages 9572–9582. PMLR, 2020.

[54] Mengdi Wang, Ethan X Fang, and Han Liu. Stochastic compositional gradient descent: algorithms for minimizing compositions of expected-value functions. *Mathematical Programming*, 161(1-2):419–449, 2017.

[55] Mengdi Wang, Ji Liu, and Ethan Fang. Accelerating stochastic composition optimization. In *Advances in Neural Information Processing Systems*, pages 1714–1722, 2016.

[56] Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh. SpiderBoost and momentum: Faster variance reduction algorithms. In *Advances in Neural Information Processing Systems 32*, pages 2406–2416. Curran Associates, Inc., 2019.

[57] Lin Xiao and Tong Zhang. A proximal stochastic gradient method with progressive variance reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.

[58] Yue Yu and Longbo Huang. Fast stochastic variance reduced ADMM for stochastic composition optimization. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 3364–3370, 2017.

[59] Junyu Zhang and Lin Xiao. A composite randomized incremental gradient method. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7454–7462, Long Beach, California, USA, June 2019.

[60] Junyu Zhang and Lin Xiao. Multi-level composite stochastic optimization via nested variance reduction. arXiv preprint arXiv:1908:11468, 2019.

[61] Junyu Zhang and Lin Xiao. A stochastic composite gradient method with incremental variance reduction. In *Advances in Neural Information Processing Systems 32*, pages 9078–9088. Curran Associates, Inc., 2019.

[62] Dongruo Zhou, Pan Xu, and Quanquan Gu. Stochastic variance-reduced cubic regularized Newton methods. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5990–5999. PMLR, July 2018.