

Non-convex Finite-Sum Optimization Via SCSG Methods

Authored by:

Michael I. Jordan
Lihua Lei
Cheng Ju
Jianbo Chen

Abstract

We develop a class of algorithms, as variants of the stochastically controlled stochastic gradient (SCSG) methods, for the smooth nonconvex finite-sum optimization problem. Only assuming the smoothness of each component, the complexity of SCSG to reach a stationary point with $\mathbb{E} \|\nabla f(x)\| \leq \epsilon$ is $\mathcal{O}(\min\{\epsilon^{-5/3}, \epsilon^{-1}\}n^{2/3})$, which strictly outperforms the stochastic gradient descent. Moreover, SCSG is never worse than the state-of-the-art methods based on variance reduction and it significantly outperforms them when the target accuracy is low. A similar acceleration is also achieved when the functions satisfy the Polyak-Lojasiewicz condition. Empirical experiments demonstrate that SCSG outperforms stochastic gradient methods on training multi-layers neural networks in terms of both training and validation loss.

1 Paper Body

We study smooth non-convex finite-sum optimization problems of the form
$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(x) \quad (1)$$

where each component $f_i(x)$ is possibly non-convex with Lipschitz gradients. This generic form captures numerous statistical learning problems, ranging from generalized linear models [22] to deep neural networks [19]. In contrast to the convex case, the non-convex case is comparatively under-studied. Early work focused on the asymptotic performance of algorithms [11, 7, 29], with non-asymptotic complexity bounds emerging more recently [24]. In recent years, complexity results have been derived for both gradient methods [13, 2, 8, 9] and stochastic gradient methods [12, 13, 6, 4, 26, 27, 3]. Unlike in the convex case, in the non-convex case one can not expect a gradient-based algorithm to converge to the global minimum if only smoothness is assumed. As a consequence, instead

of measuring functionvalue suboptimality $E f(x) - \inf_x f(x)$ as in the convex case, convergence is generally measured in terms of the squared norm of the gradient; i.e., $E \| \nabla f(x) \|^2$. We summarize the best available rates in Table 1. We also list the rates for Polyak-Łojasiewicz (P-L) functions, which will be defined in Section 2. The accuracy for minimizing P-L functions is measured by $E f(x) - \inf_x f(x)$.

It is also common to use $E \| \nabla f(x) \|^2$ to measure convergence; see, e.g. [2, 8, 9, 3]. Our results can be readily transferred to this alternative measure by using Cauchy-Schwartz inequality, $E \| \nabla f(x) \|^2 \leq E \| \nabla f(x) \|^4$, although not vice versa. The rates under this alternative can be made comparable to ours by replacing ϵ by ϵ^2 . 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

Table 1: Computation complexity of gradient methods and stochastic gradient methods for the finite-sum non-convex optimization problem (1). The second and third columns summarize the rates in the smooth and P-L cases respectively. ϵ is the P-L constant and H is the variance of a stochastic gradient. These quantities are defined in Section 2. The final column gives additional required assumptions beyond smoothness or the P-L condition. The symbol \min denotes a minimum and $O(\cdot)$ the usual Landau big-O notation with logarithmic terms hidden. Smooth Polyak-Łojasiewicz additional cond. Gradient Methods GD Best available

ϵ	O		
n	$n^{2/3}$ [24, 13]		
ϵ	$n^{1/8} O(\epsilon)$ [9]	$n^{1/5} O(\epsilon)$ [9]	
ϵ	$n^{1/8}$		
ϵ			
Stochastic Gradient Methods			
SGD	$O(\epsilon^{1/2})$ [24, 26]		
$2/3$	Best available $O(n^{2/3} + n^{1/3} \epsilon)$ [26, 27]		
1	$n^{2/3} \epsilon^{1/5}$ SCSG $O(\epsilon^{1/5})$		
O			
	[25, 17]		
-			
-			
smooth gradient smooth Hessian			
1	$\epsilon^{1/2}$		
H	$= O(1)$		
	[17]		
$2/3$	$n^{2/3} \epsilon^{1/5} + O(n^{1/3} \epsilon)$ [26, 27]		
ϵ	$(1 - \epsilon/n) + 1 - (1 - \epsilon/n)^{2/3} O(\epsilon^{1/5})$		
H	$= O(1)$		

As in the convex case, gradient methods have better dependence on ϵ in the non-convex case but worse dependence on n . This is due to the requirement of computing a full gradient. Comparing the complexity of SGD and the best achievable rate for stochastic gradient methods, achieved via variance-reduction

methods, the dependence on ϵ is significantly improved in the latter case. However, unless $\epsilon \gg n^{1/2}$, SGD has similar or even better theoretical complexity than gradient methods and existing variance-reduction methods. In practice, it is often the case that n is very large ($10^5 \sim 10^9$) while the target accuracy is moderate ($10^{-1} \sim 10^{-3}$). In this case, SGD has a meaningful advantage over other methods, deriving from the fact that it does not require a full gradient computation. This motivates the following research question: Is there an algorithm that ϵ achieves/beats the theoretical complexity of SGD in the regime of modest target accuracy; ϵ and achieves/beats the theoretical complexity of existing variance-reduction methods in the regime of high target accuracy? The question has been partially answered in the convex case by [21] in their formulation of the stochastically controlled stochastic gradient (SCSG) methods. When the target accuracy is low,

SCSG has the same $O(\epsilon^{-2})$ rate as SGD but with a much smaller data-dependent constant factor (which does not even require bounded gradients). When the target accuracy is high, SCSG achieves the same rate as the best non-accelerated methods, $O(n\epsilon)$. Despite the gap between this and the optimal rate, SCSG is the first known algorithm that provably achieves the desired performance in both regimes. In this paper, we generalize SCSG to the non-convex setting which, surprisingly, provides a completely affirmative answer to the question raised before. By only assuming smoothness of each component as in almost all other works, SCSG is always $O(\epsilon^{1/3})$ faster than SGD and is never worse than recently developed stochastic gradient methods that achieve the best rate. When $\epsilon \ll n^{-1}$, SCSG is at least $O((n\epsilon)^{2/3})$ faster than the best SVRG-type algorithms. Comparing with the gradient methods, SCSG has a better convergence rate provided $\epsilon \ll n^{6/5}$, which is the common setting in practice. Interestingly, there is a parallel to recent advances in gradient methods; [9] improved the classical $O(\epsilon^{-1})$ rate of gradient descent to $O(\epsilon^{5/6})$; this parallels the improvement of SCSG over SGD from $O(\epsilon^{-2})$ to $O(\epsilon^{5/3})$. Beyond the theoretical advantages of SCSG, we also show that SCSG yields good empirical performance for the training of multi-layer neural networks. It is worth emphasizing that the mechanism by which SCSG achieves acceleration (variance reduction) is qualitatively different from other speed-up

techniques, including momentum [28] and adaptive stepsizes [18]. It will be of interest in future work to explore combinations of these various approaches in the training of deep neural networks. The rest of paper is organized as follows: In Section 2 we discuss our notation and assumptions and we state the basic SCSG algorithm. We present the theoretical convergence analysis in Section 3. Experimental results are presented in Section 4. All the technical proofs are relegated to the Appendices.

2

Notation, Assumptions and Algorithm

We use $\|\cdot\|$ to denote the Euclidean norm and write $\min\{a, b\}$ as $a \wedge b$ for brevity throughout the paper which hides logarithmic terms, will only be used to maximize readability in paper. The notation O , our presentation but will not be used in the formal analysis. We define computation cost using the IFO

framework of [1] which assumes that sampling an index i and P accessing the pair $(f_i(x), \nabla f_i(x))$ incur a unit of cost. For brevity, we write $\nabla f_i(x)$ for $\nabla f_i(x)$. Note that calculating $\nabla f_i(x)$ incurs $\mathcal{O}(1)$ units of computational cost. x is called an ϵ -accurate solution iff $\mathbb{E} \| \nabla f(x) \|^2 \leq \epsilon$. The minimum IFO complexity to reach an ϵ -accurate solution is denoted by $C_{\text{comp}}(\epsilon)$. Recall that a random variable N has a geometric distribution, $N \sim \text{Geom}(\epsilon)$, if N is supported on the non-negative integers \mathbb{N} with $P(N = k) = \epsilon (1 - \epsilon)^k$,

$$P(N = k) = \epsilon (1 - \epsilon)^k, \quad k = 0, 1, \dots$$

An elementary calculation shows that $\mathbb{E} N \sim \text{Geom}(\epsilon) = 1/\epsilon$.

(2)

To formulate our complexity bounds, we define $f^* = \inf_x f(x)$, x^*

$$f^* = f(x^*) \quad \text{if } x^* \in \mathbb{R}^d.$$

Further we define H as an upper bound of the variance of stochastic gradients, i.e. n

$$H^2 = \sup_x \sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x)\|^2.$$

$$\sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x)\|^2 \leq H^2.$$

(3)

The assumption A1 on the smoothness of individual functions will be made throughout this paper. A1 f_i is differentiable with $\| \nabla f_i(x) - \nabla f_i(y) \| \leq L \| x - y \|^k$ for some $L \geq 0$ and all $i \in \{1, \dots, n\}$. As a direct consequence of assumption A1, it holds for any $x, y \in \mathbb{R}^d$ that $\| \nabla f_i(x) - \nabla f_i(y) \| \leq L \| x - y \|^k$.

(4)

In this paper, we also consider the following Polyak-Lojasiewicz (PL) condition [25]. It is weaker than strong convexity as well as other popular conditions that appeared in optimization literature; see [17] for an extensive discussion. A2 $f(x)$ satisfies the P-L condition with $\mu > 0$ if $\| \nabla f(x) \|^2 \geq 2\mu(f(x) - f(x^*))$ where x^* is the global minimum of f .

Here we allow N to be zero to facilitate the analysis.

3

2.1

Generic form of SCSG methods

The algorithm we propose in this paper is similar to that of [14] except (critically) the number of inner loops is a geometric random variable. This is an essential component in the analysis of SCSG, and, as we will show below, it is key in allowing us to extend the complexity analysis for SCSG to the non-convex case. Moreover, that algorithm that we present here employs a mini-batch procedure in the inner loop and outputs a random sample instead of an average of the iterates. The pseudo-code is shown in Algorithm 1. Algorithm 1 (Mini-Batch) Stochastically Controlled Stochastic Gradient (SCSG) method for smooth non-convex finite-sum objectives Inputs: Number of stages T , initial iterate x_0 , stepsizes $(\eta_j)_{j=1}^T$, batch sizes $(B_j)_{j=1}^T$, mini-batch sizes $(b_j)_{j=1}^T$. Procedure 1: for $j = 1, 2, \dots, T$ do 2: Uniformly sample a batch $I_j \subseteq \{1, \dots, n\}$ with $|I_j| = B_j$; 3: $g_j \leftarrow \frac{1}{B_j} \sum_{i \in I_j} \nabla f_i(x_{j-1})$; 4: $x_0 \leftarrow x_{j-1}$; 5: Generate $N_j \sim \text{Geom}(b_j / (B_j + b_j))$; 6: for $k = 1, 2, \dots, N_j$ do 7:

Randomly pick $I^k_1 \in [n]$ with $\Pr(I^k_1 = b_j) = \frac{1}{B_j}$; (j) (j) 8: $x^k_1 = x^0 + g_j$; (j)

(j)

(j)

9: $x^k_1 = x^0 + g_j$; 10: end for (j) 11: $x^k = x^N$; 12: end for

Output: (Smooth case) Sample x^k from $\{x^j\}_{j=1}^T$ with $P(x^k = x^j) = \frac{B_j}{\sum_{j=1}^T B_j}$; (P-L case) x^k . As seen in the pseudo-code, the SCSSG method consists of multiple epochs. In the j -th epoch, a minibatch of size B_j is drawn uniformly from the data and a sequence of mini-batch SVRG-type updates are implemented, with the total number of updates being randomly generated from a geometric distribution, with mean equal to the batch size. Finally it outputs a random sample from $\{x^j\}_{j=1}^T$. This is the standard way, proposed by [23], as opposed to computing $\arg \min_j \mathbb{E} \|x^j - x^*\|^2$ which requires additional overhead. By (2), the average total cost is $T \sum_{j=1}^T B_j$.

$\sum_{j=1}^T B_j =$

$\sum_{j=1}^T B_j$

$\sum_{j=1}^T B_j$

$\sum_{j=1}^T B_j$

$\sum_{j=1}^T B_j$

$\sum_{j=1}^T B_j = 2 \sum_{j=1}^T B_j$

(5)

Define $T(\epsilon)$ as the minimum number of epochs such that all outputs afterwards are ϵ -accurate solutions, i.e. $T(\epsilon) = \min\{T : \mathbb{E} \|x^T - x^*\|^2 \leq \epsilon\}$ for all $T \geq T(\epsilon)$. Recall the definition of $C_{\text{comp}}(\epsilon)$ at the beginning of this section, the average IFO complexity to reach an ϵ -accurate solution is $T(\epsilon) \sum_{j=1}^T B_j$.

2.2

Parameter settings

The generic form (Algorithm 1) allows for flexibility in both stepsize, η_j , and batch/mini-batch size, (B_j, b_j) . In order to minimize the amount of tuning needed in practice, we provide several default settings which have theoretical support. The settings and the corresponding complexity results are summarized in Table 2. Note that all settings fix $b_j = 1$ since this yields the best rate as will be shown in Section 3. However, in practice a reasonably large mini-batch size b_j might be favorable due to the acceleration that could be achieved by vectorization; see Section 4 for more discussions on this point. 4

Table 2: Parameter settings analyzed in this paper. B_j, b_j Type of Objectives $C_{\text{comp}}(\epsilon)$

$\frac{2}{3} \frac{1}{L} O(1) \eta \frac{1}{n}$ Smooth $O(\frac{5}{3} \frac{1}{n})$

$\frac{3}{1} \frac{n}{2} \frac{2}{3} \frac{5}{3} \frac{1}{j} \frac{1}{n}$ Smooth $O(\frac{1}{n})$

$\frac{1}{1} \frac{1}{1} \frac{1}{1} \frac{1}{n} + \frac{1}{n} \frac{2}{3} O(\frac{1}{n})$ Polyak-Łojasiewicz $O(\frac{1}{n})$

η_j Version 1 Version 2 Version 3

3.3.1

$\frac{1}{2} \frac{2}{3} \frac{1}{1} \frac{2}{3} \frac{2}{3} \frac{1}{2} \frac{2}{3} \frac{2}{3}$

Convergence Analysis One-epoch analysis

First we present the analysis for a single epoch. Given j , we define $e_j = \nabla f_j(x_{j-1}) - \nabla f(x_{j-1})$.

(6)
(j)
(j)

As shown in [14], the gradient update ∇f_k is a biased estimate of the gradient $\nabla f(x_k)$ conditioning on the current random index i_k . Specifically, within the j -th epoch, (j)

(j)
(j)
(j)
(j)

$\mathbb{E}[\nabla f_k] = \nabla f(x_k) + \nabla f_j(x_0) - \nabla f(x_0) = \nabla f(x_k) + e_j$. This reveals the basic qualitative difference between SVRG and SCSG. Most of the novelty in our (j) analysis lies in dealing with the extra term e_j . Unlike [14], we do not assume $\|x_k - x^*\|$ to be bounded since this is invalid in unconstrained problems, even in convex cases. By careful analysis of primal and dual gaps [cf. 5], we find that the stepsize η_j should scale as $2(B_j/b_j)^{1/3}$. Then same phenomenon has also been observed in [26, 27, 4] when $b_j = 1$ and $B_j = n$. 2

Theorem 3.1 Let $\eta_j = 2(B_j/b_j)^{1/3}$. Suppose $\eta \leq 1/6$ and $B_j \leq 9$ for all j , then under Assumption A1, $13b_j \leq 6L(B_j - n) \leq E_k \nabla f(x_j) - \nabla f(x_{j-1}) + \eta H$. (7) $B_j B_j$ The proof is presented in Appendix B. It is not surprising that a large mini-batch size will increase the theoretical complexity as in the analysis of mini-batch SGD. For this reason we restrict most of our subsequent analysis to $b_j \leq 1$. 3.2

Convergence analysis for smooth non-convex objectives

When only assuming smoothness, the output x^T is a random element from $(x_j)_{j=1}^T$. Telescoping (7) over all epochs, we easily obtain the following result. **Theorem 3.2** Under the specifications of Theorem 3.1 and Assumption A1,

$\eta \leq 1/6$ and $B_j \leq 9$ for all j , then $\mathbb{E}[\|x^T - x^*\|^2] \leq \frac{1}{T} \sum_{j=1}^T \frac{B_j}{b_j} \mathbb{E}[\|x_j - x^*\|^2] + \frac{1}{T} \sum_{j=1}^T \frac{B_j}{b_j} \mathbb{E}[\|x_{j-1} - x^*\|^2]$. This theorem covers many existing results. When $B_j = n$ and $b_j = 1$, Theorem 3.2 implies that

$\mathbb{E}[\|x^T - x^*\|^2] = O(1 + \frac{n}{T})$. This yields the same complexity bound $\frac{n}{T} \leq \frac{2}{3} \frac{n}{L} \mathbb{E}[\text{EC}(\eta)] = O(n + \frac{1}{T})$ as SVRG [26]. On the other hand, when $b_j = B_j \leq B$ for some $L \leq H \leq 2B \leq n$, Theorem 3.2 implies that $\mathbb{E}[\|x^T - x^*\|^2] = O(T + B)$. The second term can be made

$\mathbb{E}[\|x^T - x^*\|^2] = O(T + B)$ by setting $B = O(H)$. Under this setting $T = O(\frac{n}{B})$ and $\mathbb{E}[\text{EC}(\eta)] = O(\frac{n}{B})$.

$\mathbb{E}[\|x^T - x^*\|^2] = O(\frac{n}{B})$

This is the same rate as in [26] for SGD. 5

However, both of the above settings are suboptimal since they either set the batch sizes B_j too large or set the mini-batch sizes b_j too large. By Theorem 3.2, SCSG can be regarded as an interpolation between SGD and SVRG. By leveraging these two parameters, SCSG is able to outperform both methods. We start from considering a constant batch/mini-batch size $B_j \leq B$, $b_j \leq 1$.

Similar to SGD and SVRG, B should be at least $O(H)$. In applications like the training of neural networks, the required accuracy is moderate and hence a small batch size suffices. This is particularly important since the gradient can be computed without communication overhead, which is the bottleneck of SVRG-type algorithms. As shown in Corollary 3.3 below, the complexity of SCSG beats both SGD and SVRG. Corollary 3.3 (Constant batch sizes) Set

$$\begin{aligned} b_j &\geq 1, \\ B_j &\geq B = \min\{2H, n, \frac{1}{\epsilon}\} \\ j &\geq \frac{1}{2} \log \frac{1}{\epsilon} \\ 6LB &\geq 3 \end{aligned}$$

Then it holds that

$$EC_{\text{comp}}(\epsilon) = O$$

$$\frac{1}{\epsilon} \log \frac{1}{\epsilon} H^2 L^2 f H + \frac{1}{\epsilon} n \log \frac{1}{\epsilon}$$

Assume that $L^2 f, H = O(1)$, the above bound can be simplified to

$$\begin{aligned} \frac{1}{\epsilon} \log \frac{1}{\epsilon} n + \frac{1}{\epsilon} n &= O(EC_{\text{comp}}(\epsilon)) = O\left(\frac{1}{\epsilon} \log \frac{1}{\epsilon}\right) \\ 2 \\ n^3 \log \frac{1}{\epsilon} &\geq 3 \end{aligned}$$

When the target accuracy is high, one might consider a sequence of increasing batch sizes. Heuristically, a large batch is wasteful at the early stages when the iterates are inaccurate. Fixing the batch size to be n as in SVRG is obviously suboptimal. Via an involved analysis, we find that $B_j \geq j^2$ gives the best complexity among the class of SCSG algorithms. Corollary 3.4 (Time-varying batch sizes) Set $b_j \geq 1$,

$$\begin{aligned} n &\geq B_j = \min\{2H, n, \frac{1}{\epsilon}\} \\ j &\geq \frac{1}{2} \log \frac{1}{\epsilon} \end{aligned}$$

Then it holds that

$$EC_{\text{comp}}(\epsilon) = O \min$$

$$\begin{aligned} \frac{1}{\epsilon} \log \frac{1}{\epsilon} n^3 H^2 \log \frac{1}{\epsilon} + (H^2 \log \frac{1}{\epsilon} + (L^2 f + H \log n) \log \frac{1}{\epsilon}) \log \frac{1}{\epsilon} \\ \geq (8) \end{aligned}$$

The proofs of both Corollary 3.3 and Corollary 3.4 are presented in Appendix C. To simplify the bound (8), we assume that $L^2 f, H = O(1)$ in order to highlight the dependence on ϵ and n . Then (8) can be simplified to

$$\frac{1}{\epsilon} \log \frac{1}{\epsilon} n^3 \log \frac{1}{\epsilon} + \frac{1}{\epsilon} n^3 \log \frac{1}{\epsilon} + EC_{\text{comp}}(\epsilon) = O\left(\frac{1}{\epsilon} n^3 \log \frac{1}{\epsilon}\right) + O\left(\frac{1}{\epsilon} \log \frac{1}{\epsilon}\right)$$

The log-factor $\log \frac{1}{\epsilon}$ is purely an artifact of our proof. It can be reduced to $\log 2 + \frac{1}{2} \log \frac{1}{\epsilon}$ by setting $B_j \geq j^2 (\log j)^2 + \frac{1}{2}$; see remark 1 in Appendix C. 3.3

$$1 \geq$$

for any

Convergence analysis for P-L objectives

When the component $f_i(x)$ satisfies the P-L condition, it is known that the global minimum can be found efficiently by SGD [17] and SVRG-type algorithms [26, 4]. Similarly, SCSG can also achieve this. As in the last subsection, we start from a generic result to bound $E(f(\bar{x}_T) - f^*)$ and then consider specific settings of the parameters as well as their complexity bounds. 6

1

Theorem 3.5 Let $\eta_j =$

$5Lb_j/3$

1

. Then under the same settings of Theorem 3.2,

$\eta_j B_j/3 + 5Lb_j/3$

$E(f(\bar{x}_T) - f^*) \leq \eta_j T^{-1} \dots \eta_1 T^{-1} f + 6\eta_j H$

$T \sum_{j=1}^T \eta_j T^{-1} \dots \eta_{j+1} T^{-1} I(B_j; n)$

$j=1$

2

.

$\eta_j B_j + 5Lb_j/3$

The proofs and additional discussion are presented in Appendix D. Again, Theorem 3.5 covers existing complexity bounds for both SGD and SVRG. In fact, when $B_j = b_j/B$ as in SGD, via some calculation, we obtain that !

$T \sum_{j=1}^T \eta_j E(f(\bar{x}_T) - f^*) = O(\eta_j f + \dots + L \eta_j B)$

L The second term can be made $O(\eta_j)$ by setting $B = O(H \eta_j)$, in which case $T(\eta_j) = O(\eta_j \log \eta_j)$

$\eta_j f$

). As

$O(LH \eta_j^2)$

a result, the average cost to reach an η_j -accurate solution is $EC_{comp}(\eta_j) =$, which is the same as [17]. On the other hand, when $B_j \leq n$ and $b_j \leq 1$ as in SVRG, Theorem 3.5 implies that ! T

$L \eta_j \eta_j f \cdot E(f(\bar{x}_T) - f^*) = O(1/\eta_j^3 + L$

2/3 This entails that $T(\eta_j) = O((1 + \eta_j^{11/3}) \log 1/\eta_j)$ and hence $EC_{comp}(\eta_j) = O((n + n \eta_j) \log 1/\eta_j)$, which is the same as [26]. By leveraging the batch and mini-batch sizes, we obtain a counterpart of Corollary 3.3 as below. Corollary

3.6 Set

$b_j \leq 1,$

$B_j \leq B = \min$

$12H, n, \eta_j$

$\eta_j \leq \eta_j =$

$1/2$

$6LB/3$

Then it holds that $(EC_{comp}(\eta_j) = O$

!

$32) 1/H \eta_j H \eta_j f/n + \eta_j \cdot \log \eta_j \eta_j \eta_j$

Recall the results from Table 1, SCSG is $O(\sqrt{1 + \frac{1}{\alpha}})$ faster than SGD and is never worse than SVRG. When both α and β are moderate, the acceleration of SCSG over SVRG is significant. Unlike the smooth case, we do not find any possible choice of setting that can achieve a better rate than Corollary 3.6.

4

Experiments

We evaluate SCSG and mini-batch SGD on the MNIST dataset with (1) a three-layer fully-connected neural network with 512 neurons in each layer (FCN for short) and (2) a standard convolutional neural network LeNet [20] (CNN for short), which has two convolutional layers with 32 and 64 filters of size 5×5 respectively, followed by two fully-connected layers with output size 1024 and 10. Max pooling is applied after each convolutional layer. The MNIST dataset of handwritten digits has 50,000 training examples and 10,000 test examples. The digits have been size-normalized and centered in a fixed-size image. Each image is 28 pixels by 28 pixels. All experiments were carried out on an Amazon p2.xlarge node with a NVIDIA GK210 GPU with algorithms implemented in TensorFlow 1.0. Due to the memory issues, sampling a chunk of data is costly. We avoid this by modifying the inner loop: instead of sampling mini-batches from the whole dataset, we split the batch B_j into B_j/b_j mini-batches and run SVRG-type updates sequentially on each. Despite the theoretical advantage of setting $b_j = 1$, we consider practical settings $b_j \geq 1$ to take advantage of the acceleration obtained

by vectorization. We initialized parameters by TensorFlow's default Xavier uniform initializer. In all experiments below, we show the results corresponding to the best-tuned stepsizes. We consider three algorithms: (1) SGD with a fixed batch size $B \in \{512, 1024\}$; (2) SCSG with a fixed batch size $B \in \{512, 1024\}$ and a fixed mini-batch size $b = 32$; (3) SCSG with time-varying batch sizes $B_j = d_j^{3/2} \sqrt{n}$ and $b_j = dB_j/32e$. To be clear, given T epochs, the IFO complexity PT of the three algorithms are $T B$, $2T B$ and $2 \sum_{j=1}^T B_j$, respectively. We run each algorithm with 20 passes of data. It is worth mentioning that the largest batch size in Algorithm 3 is $d \sqrt{2751.5 n} = 4561$, which is relatively small compared to the sample size 50000. We plot in Figure 1 the training and the validation loss against the IFO complexity (i.e., the number of passes of data) for fair comparison. In all cases, both versions of SCSG outperform SGD, especially in terms of training loss. SCSG with time-varying batch sizes always has the best performance and it is more stable than SCSG with a fixed batch size. For the latter, the acceleration is more significant after increasing the batch size to 1024. Both versions of SCSG provide strong evidence that variance reduction can be achieved efficiently without evaluating the full gradient.

4

2

6

8

10

#grad / n

12

10-1
 4
 2
 6
 8
 10
 #grad / n
 12
 14
 0
 2
 4
 6
 8
 10
 #grad / n
 12
 14
 100
 10-1
 0
 2
 4
 6
 8
 10
 #grad / n
 12
 14
 10-1
 10-2 0
 2
 4
 6
 8
 10
 #grad / n
 12
 14
 100
 10-1 0
 2
 4
 6
 8
 10

#grad / n
 12
 FCN
 Training Log-Loss
 Training Log-Loss
 10-2
 14
 100
 0
 10-1
 FCN
 SGD (B = 512) SCSG (B = 512, b = 32) SCSG (B = j1.5, B/b = 32)
 100
 14
 Validation Log-Loss
 Validation Log-Loss
 0
 SGD (B = 1024) SCSG (B = 1024, b = 32) SCSG (B = j1.5, B/b = 32)
 100
 Validation Log-Loss
 10-2
 Training Log-Loss
 10-1
 CNN
 Validation Log-Loss
 Training Log-Loss
 CNN
 SGD (B = 512) SCSG (B = 512, b = 32) SCSG (B = j1.5, B/b = 32)
 100
 SGD (B = 1024) SCSG (B = 1024, b = 32) SCSG (B = j1.5, B/b = 32)
 100
 10-1
 10-2 0
 2
 4
 #grad / n
 6
 8
 10
 12
 14
 2
 4
 #grad / n
 6
 8

10
12
14
100
10-1 0

Figure 1: Comparison between two versions of SCSG and mini-batch SGD of training loss (top row) and validation loss (bottom row) against the number of IFO calls. The loss is plotted on a log-scale. Each column represents an experiment with the setup printed on the top.

CNN
scsg (B=j1.5, B/b=16) sgd (B=j1.5)
10 0 Validation Log Loss
10 0 Training Log Loss
CNN
scsg (B=j1.5, B/b=16) sgd (B=j1.5)
10 -1
10 -1
0
50
100 150 Wall Clock Time (in second)
200
0
50
100 150 Wall Clock Time (in second)
200

Figure 2: Comparison between SCSG and mini-batch SGD of training loss and validation loss with a CNN loss, against wall clock time. The loss is plotted on a log-scale. Given 2B IFO calls, SGD implements updates on two fresh batches while SCSG replaces the second batch by a sequence of variance reduced updates. Thus, Figure 1 shows that the gain due to variance reduction is significant when the batch size is fixed. To further explore this, we compare SCSG with time-varying batch sizes to SGD with the same sequence of batch sizes. The results corresponding to the best-tuned constant stepsizes are plotted in Figure 3a. It is clear that the benefit from variance reduction is more significant when using time-varying batch sizes. We also compare the performance of SGD with that of SCSG with time-varying batch sizes against wall clock time, when both algorithms are implemented in TensorFlow and run on a Amazon p2.xlarge node with a NVIDIA GK210 GPU. Due to the cost of computing variance reduction terms in SCSG, each update of SCSG is slower per iteration compared to SGD. However, SCSG makes faster progress 8

in terms of both training loss and validation loss compared to SCD in wall clock time. The results are shown in Figure 2.

10-2 0
2
4
6

8
 10
 #grad / n
 12
 14
 SGD SCSG
 10-1
 10-2 0
 2
 4
 6
 8
 10
 #grad / n
 12
 CNN
 14
 B/b = 2.0 B/b = 5.0 B/b = 10.0 B/b = 16.0 B/b = 32.0
 100
 10-1
 10-2
 0
 2
 4
 6
 8
 10
 #grad / n
 12
 14
 FCN Training Log-Loss
 10-1
 FCN 100
 Training Log-Loss
 SGD SCSG
 Training Log-Loss
 Training Log-Loss
 CNN 100
 B/b = 2.0 B/b = 5.0 B/b = 10.0 B/b = 16.0 B/b = 32.0
 100
 10-1
 0
 2
 4
 6
 8

10

#grad / n

12

14

(b) SCSG with different B_j/b_j

(a) SCSG and SGD with increasing batch sizes

Finally, we examine the effect of B_j/b_j , namely the number of mini-batches within an iteration, since it affects the efficiency in practice where the computation time is not proportional to the batch size. Figure 3b shows the results for SCSG with $B_j = d_j 3/2 \cdot n$ and $dB_j/b_j \in \{2, 5, 10, 16, 32\}$. In general, larger B_j/b_j yields better performance. It would be interesting to explore the tradeoff between computation efficiency and this ratio on different platforms.

5

Conclusion and Discussion

We have presented the SCSG method for smooth, non-convex, finite-sum optimization problems. SCSG is the first algorithm that achieves a uniformly better rate than SGD and is never worse than SVRG-type algorithms. When the target accuracy is low, SCSG significantly outperforms the SVRG-type algorithms. Unlike various other variants of SVRG, SCSG is clean in terms of both implementation and analysis. Empirically, SCSG outperforms SGD in the training of multi-layer neural networks. Although we only consider the finite-sum objective in this paper, it is straightforward to extend SCSG to the general stochastic optimization problems where the objective can be written as $E[F(x; \cdot)]$: at the beginning of j -th epoch a batch of i.i.d. sample (x_1, \dots, x_{B_j}) is drawn from the distribution F and $B_j = 1 \times \dots \times f(x_{j-1}; \cdot)$ (see line 3 of Algorithm 1); $g_j = \frac{1}{B_j} \sum_{i=1}^{B_j} \nabla f(x_i; \cdot)$ at the k -th step, a fresh sample (x_1, \dots, x_{b_j}) is drawn from the distribution F and (j)

$\frac{1}{B_j} \sum_{i=1}^{B_j} \nabla f(x_i; \cdot) =$

$\frac{1}{B_j} \sum_{i=1}^{B_j} \nabla f(x_i; \cdot) = \frac{1}{B_j} \sum_{i=1}^{B_j} \nabla f(x_0; \cdot) + \frac{1}{B_j} \sum_{i=1}^{B_j} \nabla f(x_i; \cdot) - \frac{1}{B_j} \sum_{i=1}^{B_j} \nabla f(x_0; \cdot)$

(see line 8 of Algorithm 1).

Our proof directly carries over to this case, by simply suppressing the term $I(B_j \cdot n)$, and yields $\frac{1}{5/3}$ for smooth non-convex objectives and the bound $O(\frac{1}{5/3} \cdot \frac{1}{2/3})$ for the bound $O(P-L)$ objectives. These bounds are obtained simply by setting $n = \frac{1}{5/3}$ in our convergence analysis. Compared to momentum-based methods [28] and methods with adaptive stepsizes [10, 18], the mechanism whereby SCSG achieves acceleration is qualitatively different: while momentum aims at balancing primal and dual gaps [5], adaptive stepsizes aim at balancing the scale of each coordinate, and variance reduction aims at removing the noise. We believe that an algorithm that combines these three techniques is worthy of further study, especially in the training of deep neural networks where the target accuracy is modest.

Acknowledgments The authors thank Zeyuan Allen-Zhu, Chi Jin, Nilesh Tripuraneni and anonymous reviewers for helpful discussions.

2 References

- [1] Alekh Agarwal and Leon Bottou. A lower bound for the optimization of finite sums. ArXiv e-prints abs/1410.0723, 2014. 9
- [2] Naman Agarwal, Zeyuan Allen-Zhu, Brian Bullins, Elad Hazan, and Tengyu Ma. Finding approximate local minima for nonconvex optimization in linear time. arXiv preprint arXiv:1611.01146, 2016.
- [3] Zeyuan Allen-Zhu. Natasha: Faster stochastic non-convex optimization via strongly non-convex parameter. arXiv preprint arXiv:1702.00763, 2017.
- [4] Zeyuan Allen-Zhu and Elad Hazan. Variance reduction for faster non-convex optimization. ArXiv e-prints abs/1603.05643, 2016.
- [5] Zeyuan Allen-Zhu and Lorenzo Orecchia. Linear coupling: An ultimate unification of gradient and mirror descent. arXiv preprint arXiv:1407.1537, 2014.
- [6] Zeyuan Allen-Zhu and Yang Yuan. Improved SVRG for non-strongly-convex or sum-of-nonconvex objectives. ArXiv e-prints, abs/1506.01972, 2015.
- [7] Dimitri P Bertsekas. A new class of incremental gradient methods for least squares problems. SIAM Journal on Optimization, 7(4):913?926, 1997.
- [8] Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Accelerated methods for non-convex optimization. arXiv preprint arXiv:1611.00756, 2016.
- [9] Yair Carmon, Oliver Hinder, John C Duchi, and Aaron Sidford. "convex until proven guilty": Dimension-free acceleration of gradient descent on non-convex functions. arXiv preprint arXiv:1705.02766, 2017.
- [10] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. Journal of Machine Learning Research, 12(Jul):2121?2159, 2011.
- [11] Alexei A Gaivoronski. Convergence properties of backpropagation for neural nets via theory of stochastic gradient methods. part 1. Optimization methods and Software, 4(2):117?134, 1994.
- [12] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. SIAM Journal on Optimization, 23(4):2341?2368, 2013.
- [13] Saeed Ghadimi and Guanghui Lan. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. Mathematical Programming, 156(1-2):59?99, 2016.
- [14] Reza Harikandeh, Mohamed Osama Ahmed, Alim Virani, Mark Schmidt, Jakub Konecny, and Scott Sallinen. Stop wasting my gradients: Practical SVRG. In Advances in Neural Information Processing Systems, pages 2242?2250, 2015.
- [15] Matthew D Hoffman, David M Blei, Chong Wang, and John William Paisley. Stochastic variational inference. Journal of Machine Learning Research, 14(1):1303?1347, 2013.
- [16] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In Advances in Neural Information Processing Systems, pages 315?323, 2013.
- [17] Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal gradient methods under the polyak-?ojasiewicz condition. In Joint European Conference on Machine Learning and Knowledge Discovery in Databases, pages 795?811. Springer, 2016.
- [18] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- [19] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. Nature, 521(7553):436?444, 2015.
- [20] Yann LeCun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-

based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. [21] Lihua Lei and Michael I Jordan. Less than a single pass: Stochastically controlled stochastic gradient method. *arXiv preprint arXiv:1609.03261*, 2016. [22] Peter McCullagh and John A Nelder. *Generalized Linear Models*. CRC Press, 1989. 10

[23] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009. [24] Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*. Kluwer Academic Publishers, Massachusetts, 2004. [25] Boris Teodorovich Polyak. Gradient methods for minimizing functionals. *Zhurnal Vychislitel’noi Matematiki i Matematicheskoi Fiziki*, 3(4):643–653, 1963. [26] Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. *arXiv preprint arXiv:1603.06160*, 2016. [27] Sashank J Reddi, Suvrit Sra, Barnabas Poczos, and Alex Smola. Fast incremental method for nonconvex optimization. *arXiv preprint arXiv:1603.06159*, 2016. [28] Ilya Sutskever, James Martens, George E Dahl, and Geoffrey E Hinton. On the importance of initialization and momentum in deep learning. *ICML (3)*, 28:1139–1147, 2013. [29] Paul Tseng. An incremental gradient (-projection) method with momentum term and adaptive stepsize rule. *SIAM Journal on Optimization*, 8(2):506–531, 1998. [30] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and R in Machine Learning, 1(1–2):1–305, 2008. *variational inference. Foundations and Trends*