

Nonconvex Variance Reduced Optimization with Arbitrary Sampling

Samuel Horváth¹ Peter Richtárik^{1, 2, 3}

¹ KAUST, ²University of Edinburgh, ³Moscow Institute of Physics and Technology

The Problem

$$\min_{x \in \mathbb{R}^d} f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \quad (1)$$

- f_i is L_i -smooth but **non-convex**
- n is **big**

Arbitrary Sampling

- Sampling**: a random set-valued mapping S with values being subsets of $[n] := \{1, 2, \dots, n\}$. A sampling is used to generate minibatches in each iteration.
- Probability matrix** associated with sampling S : $P_{ij} \stackrel{\text{def}}{=} \text{Prob}(\{i, j\} \subseteq S)$
- Probability vector** associated with sampling S : $p = (p_1, \dots, p_n)$, $p_i \stackrel{\text{def}}{=} \text{Prob}(i \in S)$
- Minibatch size**: $b = \mathbb{E}[|S|]$ (expected size of S)
- Proper sampling**: Sampling for which $p_i > 0$ for all $i \in [n]$
- “**Arbitrary sampling**” = any proper sampling

Main Contributions

- We develop **arbitrary sampling variants** of 3 popular variance-reduced methods for solving the non-convex problem (1): **SVRG** [1], **SAGA** [2], **SARAH** [3].
- Our rates for $b = 1$: **up to $n \times$ better** (depending on $\{L_i\}$). Improvements even in the case when $L_i = L_j$ for all i, j (for **SVRG** & **SAGA**).
- Our rates for $b \geq 1$: **Linear or superlinear speedup in minibatch size b** . That is, # of iterations needed to output a solution of a given accuracy drops by a factor equal or greater to b .
- We design **importance sampling & approximate importance sampling for minibatches**, which vastly outperform standard uniform minibatch strategies in practice.

Key Lemma

Let $\zeta_1, \zeta_2, \dots, \zeta_n$ be vectors in \mathbb{R}^d and let $\bar{\zeta} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \zeta_i$ be their average. Let S be a proper sampling. Let $v = (v_1, \dots, v_n) > 0$ be such that

$$P - pp^\top \preceq \text{Diag}(p_1 v_1, p_2 v_2, \dots, p_n v_n). \quad (2)$$

Then

$$\mathbb{E} \left[\left\| \sum_{i \in S} \frac{\zeta_i}{n p_i} - \bar{\zeta} \right\|^2 \right] \leq \frac{1}{n^2} \sum_{i=1}^n \frac{v_i}{p_i} \|\zeta_i\|^2.$$

Whenever (2) holds, it must be the case that

$$v_i \geq 1 - p_i.$$

Optimal Sampling & Superlinear Speedup

- Under our analysis, the **independent sampling S^*** defined by

$$p_i \stackrel{\text{def}}{=} \begin{cases} (b+k-n) \frac{L_i}{\sum_{j=1}^k L_j}, & \text{if } i \leq k \\ 1, & \text{if } i > k \end{cases},$$

- is **optimal**, where k is the largest integer satisfying $0 < b+k-n \leq \sum_{i=1}^k L_i$.

- All 3 methods enjoy **superlinear speed in b** up to the minibatch size

$$b_{\max} := \max\{b \mid bL_n \leq \sum_{i=1}^n L_i\}.$$

Stochastic Gradient Evaluations to Achieve $\mathbb{E}[\|\nabla f(x)\|^2] \leq \epsilon$

Alg	Uniform sampling	Arbitrary sampling [NEW]	S^* [NEW]
SVRG	$\max\left\{n, \frac{(1+4/3)L_{\max}c_1 n^{2/3}}{\epsilon}\right\}$ [1]	$\max\left\{n, \frac{(1+4\alpha/3)\bar{L}c_1 n^{2/3}}{\epsilon}\right\}$	$\max\left\{n, \frac{(1+4(n-b)/3n)\bar{L}c_1 n^{2/3}}{\epsilon}\right\}$
SAGA	$n + \frac{2L_{\max}c_2 n^{2/3}}{\epsilon}$ [2]	$n + \frac{(1+\alpha)\bar{L}c_2 n^{2/3}}{\epsilon}$	$n + \frac{(1+n-b)\bar{L}c_2 n^{2/3}}{\epsilon}$
SARAH	$n + \frac{n-b}{n-1} \frac{L_{\max}^2 c_3}{\epsilon^2}$ [3]	$n + \frac{\alpha \bar{L}^2 c_3}{\epsilon^2}$	$n + \frac{n-b}{n} \frac{\bar{L}^2 c_3}{\epsilon^2}$

Constants: $L_{\max} = \max_i L_i$ $\bar{L} = \frac{1}{n} \sum_i L_i$ $c_1, c_2, c_3 = \text{universal constants}$ $\alpha := \frac{b}{L^2 n^2} \sum_{i=1}^n \frac{v_i L_i^2}{p_i}$

Samplings

- Uniform S^u** : Every subset of $[n]$ of size b (minibatch size) is chosen with the same probability: $1/\binom{n}{b}$
- Independent S^*** : For each $i \in [n]$ we independently flip a coin, and with probability p_i include element i into S .
- Approximate Independent S^a** : Fix some $k \in [n]$ and let $a = \lceil k \max_{i \leq k} p_i \rceil$. We now sample a single set S' of cardinality a using the uniform minibatch sampling S^u . Subsequently, we apply an independent sampling S^* to select elements of S' , with selection probabilities $p'_i = k p_i / a$. The resulting random set is S^a .

SVRG with Arbitrary Sampling

Algorithm 1: SVRG (x^0, m, T, η, S)

$\tilde{x}^0 = x_m^0 = x^0$, $M = \lceil T/m \rceil$;

for $s = 0$ **to** $M - 1$ **do**

$x_0^{s+1} = x_m^s$; $g^{s+1} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x}^s)$

for $t = 0$ **to** $m - 1$ **do**

Draw a random subset (minibatch) $S_t \sim S$

$v_t^{s+1} = \sum_{i \in S_t} \frac{1}{n p_i} (\nabla f_i(x_t^{s+1}) - \nabla f_i(\tilde{x}^s)) + g^{s+1}$

$x_{t+1}^{s+1} = x_t^{s+1} - \eta v_t^{s+1}$

end

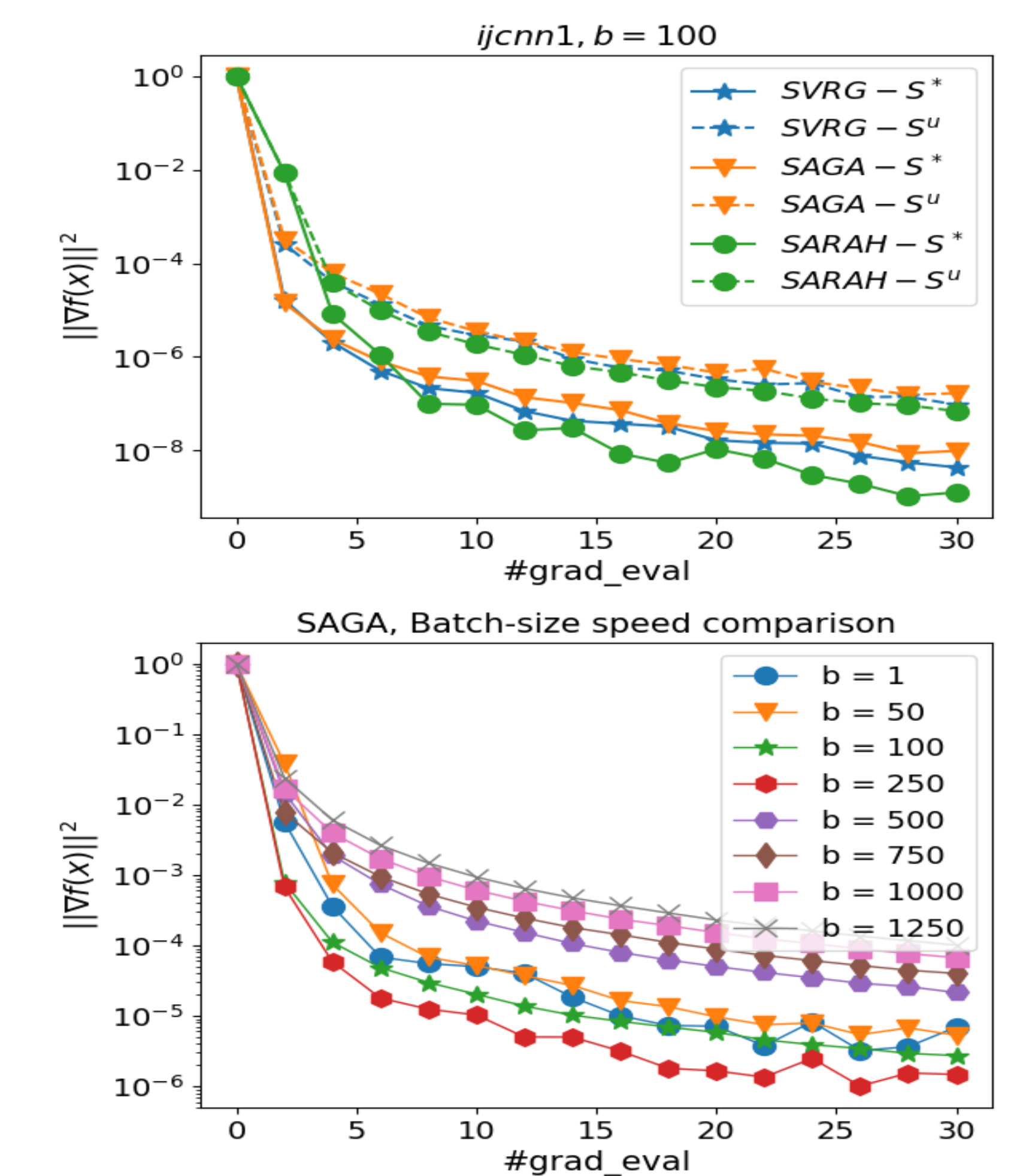
$\tilde{x}^{s+1} = x_m^{s+1}$

end

Output: Iterate x_a chosen uniformly random

from $\{\{x_t^{s+1}\}_{t=0}^m\}_{s=0}^M$

Numerical Results



References

- Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabás Póczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. In *The 33th International Conference on Machine Learning*, pages 314–323, 2016.
- Sashank J Reddi, Suvrit Sra, Barnabás Póczos, and Alex Smola. Fast incremental method for smooth nonconvex optimization. In *Decision and Control (CDC), 2016 IEEE 55th Conference on*, pages 1971–1977. IEEE, 2016.
- Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Stochastic recursive gradient algorithm for nonconvex optimization. *arXiv:1705.07261*, 2017.

