





Goal: Empirical Risk Minimization

Consider the optimization problem

$$x^{*} = \underset{x \in \mathbb{R}^{d}}{\operatorname{arg\,min}} \left\{ f(x) := \frac{1}{n} \sum_{i=1}^{n} f_{i}(x) \right\} \quad , \tag{1}$$

where

- f is L-smooth and μ -strongly convex
- each f_i is L_{max} -smooth

Stochastic Variance Reduced Gradient

Algorithm 1 SVRG [4]
Parameters: inner loop size
$$m \ge \frac{L_{\max}}{\mu}$$
, step size α , $p_t := \frac{1}{m}$
Initialization: $w_0 = x_0^m \in \mathbb{R}^d$
for $s = 1, 2, ...$ do
 $x_s^0 = w_{s-1}$
for $t = 0, 1, ..., m - 1$ do
Sample i_t uniformly at random in $\{1, ..., n\}$
 $g_s^t = \nabla f_{i_t}(x_s^t) - \nabla f_{i_t}(w_{s-1}) + \nabla f(w_{s-1})$
 $x_s^{t+1} = x_s^t - \alpha g_s^t$
end for
 $w_s = \sum_{t=0}^{m-1} p_t x_s^t$
end for

Problem: SVRG differs from practice

- Constraint on the size of the loop m
- First iterate reset to the average of past iterates
- No theoretical justification for benefits of mini-batching

Motivations

- Close gap between theory and practice of SVRG
- Offer theoretical convergence guarantees
- Demonstrate benefits from mini-batching

Stochastic Reformulation

Problem (1) can be reformulated as

$$x^* = \underset{x \in \mathbb{R}^d}{\arg \min} \mathbb{E}_{v \sim D} \left[\frac{1}{n} \sum_{i=1}^n v_i f_i(x) \right] =: \mathbb{E}_{v \sim D} \left[f_v(x) \right] , \quad (2)$$
where $\mathbb{E}_{v \sim D} \left[v \right] = \mathbf{1}_n$. To solve (2), we can use SVRG:
 $x_s^{t+1} = x_s^t - \alpha \left(\nabla f_{v_i}(x_s^t) - \nabla f_{v_i}(w_{s-1}) + \nabla f(w_{s-1}) \right) ,$
where $v_t \sim \mathcal{D}$ is sampled at each iteration.
Arbitrary sampling includes all types of sampling.
Example: mini-batching without
replacement
Let $S \subset \{1, \dots, n\}$ be a random set such that
 $\mathbb{P}\left[S = B\right] = 1/{\binom{n}{b}}$ for all $B \subset \{1, \dots, n\}, |B| = b$.
Let $v_i = \begin{cases} n/b & \text{if } i \in S \\ 0 & \text{otherwise} \end{cases}$
Then, $f_v(x) = \frac{1}{b} \sum_{i \in S} f_i(x)$ and $\nabla f_v(x) = \frac{1}{b} \sum_{i \in S} \nabla f_i(x)$.

Towards Closing the Gap between the Theory and Practice of SVRG

Othmane Sebbouh¹

¹LTCI, Télécom Paris, Institut Polytechnique de Paris ²ORFE Department, Princeton University ³INRIA, École Normale Supérieure, PSL Research University

Proposed algorithm: Free-SVRG

Algorithm 2 Free-SVRG (or 1-SVRG [5]) **Parameters:** Free inner loop length m, step size α , $p_t := \left(1 - \alpha \mu\right)^{m-1-t} / \sum_{i=0}^{m-1} (1 - \alpha \mu)^{m-1-i}$ Initialization: $w_0 = x_0^m \in \mathbb{R}^d$ for s = 1, 2, ... do $x_{s}^{0} = x_{s-1}^{m}$ for t = 0, 1, ..., m - 1 do Sample $v_t \sim \mathcal{D}$ $g_s^t = \nabla f_{v_t}(x_s^t) - \nabla f_{v_t}(w_{s-1}) + \nabla f(w_{s-1})$ $x_s^{t+1} = x_s^t - \alpha g_s^t$ end for $w_s = \sum_{t=0}^{m-1} p_t x_s^t$

end for

Solves several issues with SVRG

- Inner iterates (x_s^t) continuously updated (no resetting)
- Free choice of the inner loop size
- Much easier analysis

Algorithm analysis

An essential constant for the analysis:

Lemma: Expected smoothness

Let $v \sim \mathcal{D}$ be a sampling vector. There exists $\mathcal{L} \geq 0$ such that for all $x \in \mathbb{R}^d$,

$$\mathbb{E}_{v \sim D} \left\| \|\nabla f_v(x) - \nabla f_v(x^*)\|_2^2 \right\| \le 2\mathcal{L} \left(f(x) - f(x^*) \right)$$

Example: mini-batching without replacement
$$[1, 2]$$

$$\mathcal{L} = \mathcal{L}(\mathbf{b}) = \frac{1}{\mathbf{b}} \frac{n - \mathbf{b}}{n - 1} L_{\max} + \frac{n \mathbf{b} - 1}{\mathbf{b}} L_{\max} + \frac{n \mathbf{b} - 1}{\mathbf{b}} L \quad .$$

In particular, $\mathcal{L}(\mathbf{1}) = L_{\max}$ and $\mathcal{L}(\mathbf{n}) = L$.

Lyapunov Convergence Theorem 1

Let
$$\phi_s := \|x_s^m - x^*\|_2^2 + 8\alpha^2 \mathcal{L}S_m(f(w_s) - f(x^*)),$$

where $S_m = \sum_{i=0}^{m-1} (1 - \alpha \mu)^{m-1-i}$. If $\alpha \leq 1/6\mathcal{L}$, then the
iterates of Algorithm 2 converge with
 $\mathbb{E}[\phi_s] \leq \beta^s \phi_0,$ where $\beta = \max\left\{(1 - \alpha \mu)^m, \frac{1}{2}\right\}$.

Total complexity for mini-batching

The **total complexity** of finding an $\epsilon > 0$ approximate solution that satisfies $\mathbb{E}\left| \|x_s^m - x^*\|_2^2 \right| \leq \epsilon \phi_0$ is

 $C_m(\mathbf{b}) := 2\left(\frac{n}{m} + 2\mathbf{b}\right) \max\left\{\frac{3\mathcal{L}(\mathbf{b})}{\mu}, m\right\} \log\left(\frac{1}{\epsilon}\right)$. And for **mini-batching** (dropping the log term): $C_m(\mathbf{b}) := 2\left(\frac{n}{m} + 2\mathbf{b}\right) \max\left\{\frac{3n - \mathbf{b}L_{\max}}{\mathbf{b}n - 1} + \frac{3n\mathbf{b} - 1L}{\mathbf{b}n - 1u}, m\right\}.$





Same total complexity and optimal parameter settings as *Free-SVRG* (up to constants).

We found a **range of values** minimizing the total complexity. If $m \in [\min(n, L_{\max}/\mu), \max(n, L_{\max}/\mu)]$, then

Nidham Gazagnadou¹ Samy Jelassi² Francis Bach³ Robert M. Gower¹

Alternative algorithm: L-SVRG-D

Problem: SVRG requires the strong convexity • SVRG relies on knowing μ

Solution: [3] proposed a **loopless** version of SVRG. **Improvement:** when the variance of the estimate of the gradient is high, decrease the step size.

Algorithm 3 L-SVRG-D (Loopless-SVRG-Decrease) **Parameters:** step size $\alpha, p \in (0, 1]$ Initialization: $w^0 = x^0 \in \mathbb{R}^d, \ \alpha_0 = \alpha$ for k = 0, 1, 2, ... do

Sample $v_k \sim \mathcal{D}$ $g^k = \nabla f_{v_k}(x^k) - \nabla f_{v_k}(w^k) + \nabla f(w^k)$ $x^{k+1} = x^k - \alpha_k g^k$ $(w^{k+1}, \alpha_{k+1}) = \begin{cases} (x^k, \alpha) & \text{with prob. } p \\ (w^k, \sqrt{1-p} \alpha_k) & \text{with prob. } 1-p \end{cases}$

end for

Lyapunov Convergence Theorem 2

Benefits

- **Bigger step size** for the first iterations of the loop, when the **variance is low**
- **Smaller step size** for the last iterations of the loop, when the **variance is high**

How to set the inner loop size?

$$C_m(1) = O\left(\left(n + \frac{L_{\max}}{\mu}\right)\log\left(\frac{1}{\epsilon}\right)\right)$$

 \wedge Includes the practical choice $m = n \wedge$

How to set the mini-batch size?

For any fixed inner loop size m

• the **total complexity** is a **convex function** of b • the step size is an increasing function of b



Figure: The total complexity (left) and the step size (right) as b increases.







Optimal Mini-Batch and Step Sizes for SAGA

International Conference on Machine Learning, 2019.

[2] R. M. Gower, N. Loizou, X. Qian, A. Sailanbayev, E. Shulgin, and P. Richtárik. SGD: General Analysis and Improved Rates.

International Conference on Machine Learning, 2019.

[3] T. Hofmann, A. Lucchi, S. Lacoste-Julien, and B. McWilliams.

Variance Reduced Stochastic Gradient Descent with Neighbors. In Advances in Neural Information Processing Systems, 2015.

[4] R. Johnson and T. Zhang.

Accelerating Stochastic Gradient Descent using Predictive Variance Reduction. In Advances in Neural Information Processing Systems, 2013.

[5] A. Raj and S. U. Stich.

k-SVRG: Variance Reduction for Large Scale Optimization. arXiv:1805.00982, 2018.