

# Lecture notes on Stochastic Variance Reduced Methods.

Robert M. Gower

October 14, 2019

## Abstract

Lecture notes on variance reduction techniques.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Covariates</b>	<b>2</b>
<b>3</b>	<b>The Original Stochastic Variance Reduced (SVRG) method and proof</b>	<b>2</b>
<b>4</b>	<b>Modern version and proof: Free-SVRG</b>	<b>5</b>

## 1 Introduction

Consider the following optimization problem

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

where  $f$  is  $L$ -smooth,  $\lambda$ -strongly convex and  $f_i$  is convex and  $L_i$ -smooth for  $i = 1, \dots, n$ . In other words

$$f(y) + \langle \nabla f(y), w - y \rangle + \frac{\lambda}{2} \|w - y\|_2^2 \leq f(w) \leq f(y) + \langle \nabla f(y), w - y \rangle + \frac{L}{2} \|w - y\|_2^2, \quad (2)$$

and

$$f_i(w) \leq f_i(y) + \langle \nabla f_i(y), w - y \rangle + \frac{L_i}{2} \|w - y\|_2^2, \quad \text{for } i = 1, \dots, n. \quad (3)$$

In last weeks lecture we saw that using Stochastic gradient descent (SGD) to solve (1) can converge much faster in the early iterations as compared to the gradient descent algorithm. But in later iterations the SGD algorithm slows down and thus struggles to reach an accurate solution. Can we get the best of both the fast initial convergence of SGD and the steady linear convergence of gradient descent? Yes we can! The trick is to solve SGD's issues with variance.

What are SGD's issues with variance? Though the stochastic gradient is an unbiased estimator of the gradient, it may have high variance. Indeed, to analyse SGD we had to start by imposing the rather awkward following assumption: That there exists  $B > 0$  such that

$$\mathbb{E}_j[\|\nabla f_j(w^t)\|_2^2] \leq B^2, \text{ for all iterates } w^t \text{ of SGD.}$$

Even with the above assumption, we required decreasing stepsizes to gradually kill off the variance. Yet another glaring issue with SGD is that even if we start the SGD algorithm on the solution  $w^* = w^0$ , the method will not stop. This is because the stochastic gradients are not necessarily zero on the solution, that is  $\nabla f_i(w^*) \neq 0$  is entirely possible. While  $\nabla f(w^*) = 0$ , thus gradient descent will stop once it has reached the solution.

In these notes we set out to describe methods that fix the above issues. Our aim is to have an iterative algorithm of the form

$$w^{t+1} = w^t - \alpha g^t, \tag{4}$$

where  $\alpha > 0$  is a stepsize and  $g^t$  is an estimate of the gradient that satisfies

$$\text{Unbiased:} \quad \mathbf{E}[g^t] = \nabla f(w^t) \tag{5}$$

$$\text{Reducing Variance:} \quad \mathbf{E}[\|g^t\|_2^2] \xrightarrow{w^t \rightarrow w^*} 0. \tag{6}$$

Note that the

$$\mathbf{VAR}[g^t] = \mathbf{E}[\|g^t\|_2^2] - \|\nabla f(w^t)\|_2^2.$$

Consequently if (6) holds, then the variance of  $g^t$  also tends to zero as  $w^t$  tends to  $w^*$ .

Our main tool for building an estimate of the gradient that satisfies the above will be covariates.

## 2 Covariates

Let  $x$  be a random variable. We say that a random variable  $z$  is a covariate of  $x$  if  $\text{cov}[x, z] > 0$ . We can use the covariate  $z$  to build an unbiased estimator of  $x$  that has a small variance. Indeed let

$$x_z = x - z + \mathbf{E}[z],$$

and note that  $\mathbf{E}[x_z] = \mathbf{E}[x]$ . Furthermore

$$\mathbf{VAR}[x_z] = \mathbf{VAR}[x] + \mathbf{VAR}[z] - 2 \text{cov}[x, z].$$

Consequently if  $\text{cov}[x, z]$  is sufficiently large, then  $\mathbf{VAR}[x_z]$  is small. We can build an estimate of the gradient with reduced variance by finding covariates for the stochastic gradient.

## 3 The Original Stochastic Variance Reduced (SVRG) method and proof

Let  $w^k \in \mathbb{R}^d$  be our current iterate and let  $\tilde{w}^t \in \mathbb{R}^d$  be a *reference point*. If  $w^k$  is sufficiently close to  $\tilde{w}^t$  it is reasonable to expect that  $\nabla f_i(w^k)$  and  $\nabla f_i(\tilde{w}^t)$  are close (and are thus covariates) for

every  $i = 1, \dots, n$ . Consequently, if  $i \in \{1, \dots, n\}$  is sampled uniformly then

$$g^k = \nabla f_i(w^k) - \nabla f_i(\tilde{w}^t) + \nabla f(\tilde{w}^t), \quad (7)$$

is an unbiased estimate of the gradient with reduced variance.

Before convergence, we need the following three Lemmas.

**Lemma 1** *If  $f$  is an  $L$ -smooth function then*

$$f(y - \frac{1}{L}\nabla f(y)) - f(y) \leq -\frac{1}{2L}\nabla f(y). \quad (8)$$

**Proof:** Setting  $w = y - \frac{1}{L}\nabla f(y)$  in the right hand of (2) gives

$$f(y - \frac{1}{L}\nabla f(y)) - f(y) \leq \langle \nabla f(y), -\frac{1}{L}\nabla f(y) \rangle + \frac{L}{2} \left\| -\frac{1}{L}\nabla f(y) \right\|_2^2 = -\frac{1}{2L}\nabla f(y). \quad \blacksquare$$

**Lemma 2** *If each  $f_i$  is  $L_i$ -smooth then*

$$\mathbf{E} \left[ \left\| \nabla f_i(w) - \nabla f_i(w^*) \right\|_2^2 \right] \leq 2L_{\max}(f(w) - f(w^*)). \quad (9)$$

**Proof:** Let  $g_i(w) = f_i(w) - f_i(w^*) - \langle \nabla f_i(w^*), w - w^* \rangle$  which is  $L_i$ -smooth. By the convexity of  $f_i$  we have that  $g_i(w) \geq 0$  for all  $w$ . From (8) we have that

$$-g_i(w) \stackrel{g_i(w - \frac{1}{L_i}\nabla g_i(w)) \geq 0}{\leq} g_i(w - \frac{1}{L_i}\nabla g_i(w)) - g_i(w) \leq -\frac{1}{2L_i} \left\| \nabla g_i(w) \right\|_2^2 \leq -\frac{1}{2L_{\max}} \left\| \nabla g_i(w) \right\|_2^2.$$

By substituting  $g_i(w) = f_i(w) - f_i(w^*) - \langle \nabla f_i(w^*), w - w^* \rangle$  the above can be re-written as

$$\frac{1}{2L_{\max}} \left\| \nabla f_i(w) - \nabla f_i(w^*) \right\|_2^2 \leq f_i(w) - f_i(w^*) - \langle \nabla f_i(w^*), w - w^* \rangle.$$

Taking expectation with respect to  $i$  and using that  $\frac{1}{n} \sum_{i=1}^n \nabla f_i(w^*) = \nabla f(w^*) = 0$  gives the result.  $\blacksquare$

**Lemma 3** *The second moment of the SVRG gradient estimate is bounded*

$$\mathbf{E} \left[ \left\| g^t \right\|_2^2 \right] \leq 4L_{\max}(f(w^k) - f(w^*)) + 4L_{\max}(f(\tilde{w}^t) - f(w^*)). \quad (10)$$

**Proof:**

$$\begin{aligned} \mathbf{E} \left[ \left\| g^k \right\|_2^2 \right] &\leq \mathbf{E} \left[ \left\| \nabla f_i(w^k) - \nabla f_i(w^*) + \nabla f_i(w^*) - \nabla f_i(\tilde{w}^t) + \nabla f(\tilde{w}^t) \right\|_2^2 \right] \\ &= 2\mathbf{E} \left[ \left\| \nabla f_i(w^k) - \nabla f_i(w^*) \right\|_2^2 \right] + 2\left\| \nabla f_i(w^*) - \nabla f_i(\tilde{w}^t) + \nabla f(\tilde{w}^t) \right\|_2^2 \\ &\leq 2\mathbf{E} \left[ \left\| \nabla f_i(w^k) - \nabla f_i(w^*) \right\|_2^2 \right] + 2\left\| \nabla f_i(w^*) - \nabla f_i(\tilde{w}^t) \right\|_2^2 \\ &\stackrel{(9)}{\leq} 4L_{\max}(f(w^k) - f(w^*)) + 4L_{\max}(f(\tilde{w}^t) - f(w^*)). \end{aligned}$$

Where we used in the first inequality that  $\mathbf{E} \left[ \left\| X - \mathbf{E}[X] \right\|_2^2 \right] \leq \mathbf{E} \left[ \left\| X \right\|_2^2 \right]$  with  $X = \nabla f_i(w^*) - \nabla f_i(\tilde{w}^t) + \nabla f(\tilde{w}^t)$ .

Next we prove the convergence of the original SVRG method in the following theorem.

---

**Algorithm 1** Original SVRG

---

- 1: **Parameters** number of inner iterations  $m$  and learning rate  $\alpha$ .
  - 2: Choose  $\tilde{w}_0$ .
  - 3: **for**  $t = 1, 2, \dots$  **do**
  - 4:      $\tilde{w} = \tilde{w}_{t-1}$ .
  - 5:     Compute the full gradient  $\nabla f(\tilde{w})$ .
  - 6:     Set  $w_0 = \tilde{w}$ .
  - 7:     **for**  $k = 0, \dots, m - 1$  **do**
  - 8:         Sample  $i_k$  from  $\{1, \dots, n\}$ .
  - 9:          $g_k(w^k) = \nabla f_{i_k}(w^k) - \nabla f_{i_k}(\tilde{w}) + \nabla f(\tilde{w})$ .
  - 10:         Update  $w^{k+1} = w^k - \alpha g_k(w^k)$ .
  - 11:     Choose the following reference point  $\tilde{w}_s$ , according to the options below.
  - 12:     **Option Last:**  $\tilde{w}_t = w_m$ .
  - 13:     **Option Average:** Choose  $\tilde{w}_t$  such that  $\tilde{w}_t = \frac{1}{m} \sum_{i=0}^{m-1} w^i$
- 

**Theorem 4** Consider the iterates of Algorithm 1. If we choose the stepsize  $\alpha = 1/10L_{\max}$  and the number of inner iterations as  $m = \lambda/L_{\max}$  then the SVRG method (7) converges according to

$$\mathbb{E}[f(\tilde{w}^t)] - f(w^*) \leq 0.9^t(f(\tilde{w}^0) - f(w^*)). \quad (11)$$

**Proof:** First note that

$$\begin{aligned} \mathbf{E}_j \left[ \|w^{k+1} - w^*\|_2^2 \right] &= \|w^k - w^*\|_2^2 - 2\alpha \left\langle \nabla f(w^k), w^k - w^* \right\rangle + \mathbf{E}_j \left[ \|g^k\|_2^2 \right] \\ &\leq \|w^k - w^*\|_2^2 - 2\alpha(f(w^k) - f(w^*)) + \mathbf{E}_j \left[ \|g^k\|_2^2 \right] \\ &\stackrel{(10)}{\leq} \|w^k - w^*\|_2^2 - 2\alpha(1 - 2\alpha L_{\max})(f(w^k) - f(w^*)) + 4\alpha L_{\max}(f(\tilde{w}^t) - f(w^*)). \end{aligned}$$

Taking total expectation, summing up over  $k = 0 \dots m - 1$  and using telescopic cancellation we have that

$$\begin{aligned} \mathbb{E} \left[ \|w^m - w^*\|_2^2 \right] &\leq \mathbb{E} \left[ \|w^0 - w^*\|_2^2 \right] - 2\alpha(1 - 2\alpha L_{\max})\mathbb{E} \left[ \sum_{k=0}^{m-1} (f(w^k) - f(w^*)) \right] \\ &\quad + 4m\alpha^2 L_{\max} \mathbb{E} \left[ f(\tilde{w}^t) - f(w^*) \right]. \end{aligned}$$

Using that  $w^0 = \tilde{w}^t$ , strong convexity  $f(\tilde{w}^t) - f(w^*) \geq \frac{\lambda}{2} \|\tilde{w}^t - w^*\|_2^2$  and re-arranging we have that

$$\begin{aligned} 2\alpha(1 - 2\alpha L_{\max})\mathbb{E} \left[ \sum_{k=0}^{m-1} (f(w^k) - f(w^*)) \right] &\leq \mathbb{E} \left[ \|\tilde{w}^t - w^*\|_2^2 \right] - \mathbb{E} \left[ \|w^m - w^*\|_2^2 \right] \\ &\quad + 4m\alpha^2 L_{\max} \mathbb{E} \left[ f(\tilde{w}^t) - f(w^*) \right] \\ &\leq \left( 4m\alpha^2 L_{\max} + \frac{2}{\lambda} \right) \mathbb{E} \left[ f(\tilde{w}^t) - f(w^*) \right] \end{aligned}$$

Re-arranging again and using Jensen's inequality we have

$$\begin{aligned}
\mathbb{E}[f(\sum_{k=0}^{m-1} \frac{w^k}{m})] - f(w^*) &\leq \frac{1}{m} \mathbb{E}[\sum_{k=0}^{m-1} f(w^k)] - f(w^*) \\
&\leq \frac{4m\alpha^2 L_{\max} + 2\lambda^{-1}}{2\alpha(1-2\alpha L_{\max})m} \mathbb{E}[f(\tilde{w}^t) - f(w^*)] \\
&= \left( \frac{2\alpha L_{\max}}{1-2\alpha L_{\max}} + \frac{1}{\lambda\alpha(1-2\alpha L_{\max})m} \right) \mathbb{E}[f(\tilde{w}^t) - f(w^*)]
\end{aligned}$$

It now remains to substitute  $\alpha = 1/10L_{\max}$  and  $m = 20L_{\max}/\mu$  to see that

$$\frac{2\alpha L_{\max}}{1-2\alpha L_{\max}} + \frac{1}{\lambda\alpha(1-2\alpha L_{\max})m} = \frac{2/10}{1-2/10} + \frac{1}{2(1-2/10)} = \frac{2}{8} + \frac{5}{8} = \frac{7}{8}. \quad \blacksquare$$

## 4 Modern version and proof: Free-SVRG

The original SVRG method in Algorithm 1 tends not to work well because the inner iterates are always being reset to the reference point (line 6) and because the number of inner iterates tends to be too big ( $m = L/\mu \gg 1$ ). Rather, in practice it seems that not resetting the inner iterates and using  $m = n$  tends to work better. Here we present a version of SVRG that does just that, see Algorithm 2.

To declutter the notation, we define for a given step size  $\gamma > 0$ :

$$S_m \stackrel{\text{def}}{=} \sum_{i=0}^{m-1} (1-\gamma\mu)^{m-1-i} \quad \text{and} \quad \alpha_t \stackrel{\text{def}}{=} \frac{(1-\gamma\mu)^{m-1-t}}{S_m}, \quad \text{for } t = 0, \dots, m-1. \quad (12)$$

---

### Algorithm 2 *Free-SVRG*

---

**Parameters** inner-loop length  $m$ , step size  $\gamma$ , and  $\alpha_t$  defined in (12)

**Initialization**  $w_0 = x_0^m \in \mathbb{R}^d$

**for**  $s = 1, 2, \dots$  **do**

$$x_s^0 = x_{s-1}^m$$

**for**  $t = 0, 1, \dots, m-1$  **do**

Sample  $i \sim \frac{1}{n}$

$$g_s^t = \nabla f_i(x_s^t) - \nabla f_i(w_{s-1}) + \nabla f(w_{s-1})$$

$$x_s^{t+1} = x_s^t - \gamma g_s^t$$

$$w_s = \sum_{t=0}^{m-1} \alpha_t x_s^t$$


---

**Theorem 5** Consider the setting of Algorithm 2 and the following Lyapunov function

$$\phi_s \stackrel{\text{def}}{=} \|x_s^m - x^*\|_2^2 + C(f(w_s) - f(x^*)) \quad \text{where} \quad C \stackrel{\text{def}}{=} 8\gamma^2 L_{\max} S_m. \quad (13)$$

Let  $f_i$  be  $L_i$ -smooth and  $f$  be  $\mu$ -strongly convex with

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|x - y\|_2^2. \quad (14)$$

If  $\gamma \leq \frac{1}{6L_{\max}}$ , then

$$\mathbf{E}[\phi_s] \leq \beta^s \phi_0, \quad \text{where } \beta = \max\left\{(1 - \gamma\mu)^m, \frac{1}{2}\right\}. \quad (15)$$

**Proof:**

$$\begin{aligned} \mathbf{E}_t [\|x_s^{t+1} - x^*\|_2^2] &= \mathbf{E}_t [\|x_s^t - x^* - \gamma g_s^t\|_2^2] \\ &= \|x_s^t - x^*\|_2^2 - 2\gamma \mathbf{E}_t [g_s^t]^\top (x_s^t - x^*) + \gamma^2 \mathbf{E}_t [\|g_s^t\|_2^2] \\ &\stackrel{(10)}{\leq} \|x_s^t - x^*\|_2^2 - 2\gamma \nabla f(x_s^t)^\top (x_s^t - x^*) \\ &\quad + 2\gamma^2 [2L_{\max}(f(x_s^t) - f(x^*)) + 2L_{\max}(f(w_{s-1}) - f(x^*))] \\ &\stackrel{(14)}{\leq} (1 - \gamma\mu) \|x_s^t - x^*\|_2^2 - 2\gamma(1 - 2\gamma L_{\max})(f(x_s^t) - f(x^*)) \\ &\quad + 4\gamma^2 L_{\max}(f(w_{s-1}) - f(x^*)). \end{aligned} \quad (16)$$

Note that since  $\gamma \leq \frac{1}{6L_{\max}}$  and  $L_{\max} \geq 0$ , we have that  $\gamma \leq \frac{1}{2\mu}$ , and consequently  $(1 - \gamma\mu) > 0$ . Thus by iterating (16) over  $t = 0, \dots, m-1$  and taking the expectation, since  $x_s^0 = x_{s-1}^m$ , we obtain

$$\begin{aligned} \mathbf{E} [\|x_s^m - x^*\|_2^2] &\leq (1 - \gamma\mu)^m \mathbf{E} [\|x_{s-1}^m - x^*\|_2^2] \\ &\quad - 2\gamma(1 - 2\gamma L_{\max}) \sum_{t=0}^{m-1} (1 - \gamma\mu)^{m-1-t} \mathbf{E} [f(x_s^t) - f(x^*)] \\ &\quad + 4\gamma^2 L_{\max} \mathbf{E} [f(w_{s-1}) - f(x^*)] \sum_{t=0}^{m-1} (1 - \gamma\mu)^{m-1-t} \\ &\stackrel{(12)+(13)}{=} (1 - \gamma\mu)^m \mathbf{E} [\|x_{s-1}^m - x^*\|_2^2] \\ &\quad - 2\gamma(1 - 2\gamma L_{\max}) S_m \sum_{t=0}^{m-1} \alpha_t \mathbf{E} [f(x_s^t) - f(x^*)] \\ &\quad + \frac{C}{2} \mathbf{E} [f(w_{s-1}) - f(x^*)]. \end{aligned} \quad (17)$$

Since  $f$  is convex, we have by Jensen's inequality that

$$f(w_s) - f(x^*) = f\left(\sum_{t=0}^{m-1} \alpha_t x_s^t\right) - f(x^*) \leq \sum_{t=0}^{m-1} \alpha_t (f(x_s^t) - f(x^*)). \quad (18)$$

Consequently,

$$\mathbf{E} [f(w_s) - f(x^*)] \stackrel{(13)+(18)}{\leq} \sum_{t=0}^{m-1} \alpha_t \mathbf{E} [(f(x_s^t) - f(x^*))]. \quad (19)$$

As a result,

$$\begin{aligned} \mathbf{E}[\phi_s] &\stackrel{(17)+(19)}{\leq} (1 - \gamma\mu)^m \mathbf{E}[\|x_{s-1}^m - x^*\|_2^2] + \frac{C}{2} \mathbf{E}[f(w_{s-1}) - f(x^*)] \\ &\quad - 2\gamma(1 - 6\gamma L_{\max}) S_m \sum_{t=0}^{m-1} \alpha_t \mathbf{E}[(f(x_s^t) - f(x^*))]. \end{aligned}$$

Since  $\gamma \leq \frac{1}{6L_{\max}}$ , the above implies

$$\begin{aligned} \mathbf{E}[\phi_s] &\leq (1 - \gamma\mu)^m \mathbf{E}[\|x_{s-1}^m - x^*\|_2^2] + \frac{C}{2} \mathbf{E}[f(w_{s-1}) - f(x^*)] \\ &\leq \beta \mathbf{E}[\phi_{s-1}], \end{aligned}$$

where  $\beta = \max\{(1 - \gamma\mu)^m, \frac{1}{2}\}$ .

Moreover, if we set  $w_s = x_s^t$  with probability  $\alpha_t$ , for  $t = 0, \dots, m - 1$ , the result would still hold. Indeed (18) would hold with equality and the rest of the proof would follow verbatim.  $\blacksquare$