

Lecture 10

*Lecturer: Xiangyu Chang**Scribe: Xiangyu Chang**Edited by: Xiangyu Chang*

1 SVRG

1.1 Covariate of Stochastic Gradient

We still consider the finite-sum optimization problem. Let $\{(\mathbf{a}_i, b_i)\}_{i=1}^m$ be a dataset, $\mathcal{F} = \{h_{\mathbf{x}} | h_{\mathbf{x}} : \mathcal{A} \rightarrow \mathcal{B}, \mathbf{x} \in \mathbb{R}^n\}$ be a class of predictor function and ℓ be a loss function. Then the corresponding finite-sum optimization problem is

$$\min_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{m} \sum_i^m f_i(\mathbf{x}), \quad (1)$$

where $f_i(\mathbf{x}) := \ell(b_i, h_{\mathbf{x}}(\mathbf{a}_i))$.

Though the stochastic gradient is an unbiased estimator of the gradient, it may have high variance. Indeed, to analysis SGD we had to start by imposing the assumption about its variance

$$\mathbb{E}_{i_t} [\|\nabla f_{i_t}(\mathbf{x}^t)\|^2] \leq \sigma^2 + \|\nabla f(\mathbf{x}^t)\|^2.$$

Even with the above assumption, we required decreasing step sizes to gradually kill off the variance. Yet another glaring issue with SGD is that even if we start the SGD algorithm on the optimal point $\mathbf{x}^0 = \mathbf{x}^*$, then method will not stop. This is because the stochastic gradients are not necessarily zero on the solution, that is $\nabla f_{i_t}(\mathbf{x}^*) \neq 0$ is entirely possible (see the example in SGD section).

The aim of SVRG is to construct a new “gradient” \mathbf{g}^t such that

- Unbiased: $\mathbb{E}[\mathbf{g}^t] = \nabla f(\mathbf{x}^t)$
- Reducing Variance:

$$\mathbb{E}[\|\mathbf{g}^t\|_2^2] \rightarrow 0, \text{ as } \mathbf{x}^t \rightarrow \mathbf{x}^*.$$

How to reduce the variance of stochastic gradient? The basic idea is to consider an important method in MCMC. That is to construct a covariate variable of $\nabla f_{i_t}(\mathbf{x}^t)$.

Let \mathbf{x} be a random variable. We say that a random variable \mathbf{z} is a covariate of \mathbf{x} if $\text{cov}(\mathbf{x}, \mathbf{z}) > 0$. Then, utilizing \mathbf{z} can build an unbiased estimator of \mathbf{x} that has a small variance. Let

$$\mathbf{x}_{\mathbf{z}} = \mathbf{x} - \mathbf{z} + \mathbb{E}[\mathbf{z}]. \quad (2)$$

And note that $\mathbb{E}[\mathbf{x}_{\mathbf{z}}] = \mathbb{E}[\mathbf{x}]$, and

$$\text{Var}[\mathbf{x}_{\mathbf{z}}] = \text{Var}[\mathbf{x}] + \text{Var}[\mathbf{z}] - 2\text{Cov}(\mathbf{x}, \mathbf{z}). \quad (3)$$

Consequently, if $\text{Cov}(\mathbf{x}, \mathbf{z})$ is sufficiently large, then $\text{Var}[\mathbf{x}_{\mathbf{z}}]$ is small.

1.2 SVRG Algorithm

We can build an estimate of the gradient with reduced variance by finding covariates for the stochastic gradient.

Let $\mathbf{x}^t \in \mathbb{R}^n$ be our current iterate and $\tilde{\mathbf{x}}^k$ be a *reference point*. If \mathbf{x}^t is sufficiently close to $\tilde{\mathbf{x}}^k$ it is reasonable to expect that $\nabla f_{i_t}(\mathbf{x}^t)$ is the covariate of $\nabla f_{i_t}(\tilde{\mathbf{x}}^k)$. Then, i_t is uniformly sampled from $[m]$ and

$$\mathbf{g}^t = \nabla f_{i_t}(\mathbf{x}^t) - \nabla f_{i_t}(\tilde{\mathbf{x}}^k) + \nabla f(\tilde{\mathbf{x}}^k). \quad (4)$$

Let us first refer to the SVRG [1] Algorithm 1.

Algorithm 1 SVRG

Parameters update frequency T and learning rate s

Initialize $\tilde{\mathbf{x}}^0$

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

$\mathbf{x}^0 = \tilde{\mathbf{x}}^{k-1}$

for $t = 1, 2, \dots, T$ **do**

 Randomly pick $i_t \in \{1, \dots, n\}$ and update weight

$$\mathbf{g}^{t-1} = \nabla f_{i_t}(\mathbf{x}^{t-1}) - \nabla f_{i_t}(\tilde{\mathbf{x}}^{k-1}) + \nabla f(\tilde{\mathbf{x}}^{k-1}), \quad (5)$$

$$\mathbf{x}^t = \mathbf{x}^{t-1} - s\mathbf{g}^{t-1}. \quad (6)$$

end for

Last Option: $\tilde{\mathbf{x}}^k = \mathbf{x}^T$;

Average Option: $\tilde{\mathbf{x}}^k = 1/T \sum_t \mathbf{x}^t$;

Random Option : $\tilde{\mathbf{x}}^k = \mathbf{x}^t$ for randomly chosen $t \in \{1, \dots, T\}$.

end for

1.3 Convergence Analysis

We suppose that f is β smooth and α -strongly convex, and f_i is convex and β_i smooth. Before the convergence analysis, we present the following lemmas.

Lemma 1 *Let f is a β smooth function then*

$$f(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|^2. \quad (7)$$

Lemma 2 *If each f_{i_t} is β_{i_t} smooth and convex, then*

$$\mathbb{E}[\|\nabla f_{i_t}(\mathbf{x}) - \nabla f_{i_t}(\mathbf{x}^*)\|^2] \leq 2\beta_{\max}(f(\mathbf{x}) - f(\mathbf{x}^*)). \quad (8)$$

Proof 1 *Let $h_{i_t}(\mathbf{x}) = f_{i_t}(\mathbf{x}) - f_{i_t}(\mathbf{x}^*) - \langle \nabla f_{i_t}(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \geq 0$ due to the convexity. By Lemma 2, we have that*

$$-h_{i_t}(\mathbf{x}) \leq h_{i_t}(\mathbf{x} - \frac{1}{\beta_{i_t}} \nabla h_{i_t}(\mathbf{x})) - h_{i_t}(\mathbf{x}) \leq -\frac{1}{2\beta_{i_t}} \|\nabla h_{i_t}(\mathbf{x})\|^2 \leq -\frac{1}{2\beta_{\max}} \|\nabla h_{i_t}(\mathbf{x})\|^2. \quad (9)$$

By substituting h_{i_t} , then

$$\|\nabla f_{i_t}(\mathbf{x}) - \nabla f_{i_t}(\mathbf{x}^*)\|^2 \leq 2\beta_{\max}(f_{i_t}(\mathbf{x}) - f_{i_t}(\mathbf{x}^*) - \langle \nabla f_{i_t}(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle). \quad (10)$$

Then taking expectation with respect to i_t and using that $\mathbb{E}[\nabla f_{i_t}(\mathbf{x}^)] = \nabla f(\mathbf{x}^*) = 0$, we can finish the prove.*

Now we bound the variance of stochastic gradient.

Lemma 3 *The second moment of SVRG gradient is bounded as*

$$\mathbb{E}_{i_t}[\|\mathbf{g}^t\|^2] \leq 4\beta_{\max}(f(\mathbf{x}^t) - f(\mathbf{x}^*)) + 4\beta_{\max}(f(\tilde{\mathbf{x}}^k) - f(\mathbf{x}^*)). \quad (11)$$

Proof 2

$$\mathbb{E}_{i_t}\|\mathbf{g}^t\|^2 = \mathbb{E}_{i_t}[\|\nabla f_{i_t}(\mathbf{x}^t) - \nabla f_{i_t}(\mathbf{x}^*) + \nabla f_{i_t}(\mathbf{x}^*) - \nabla f_{i_t}(\tilde{\mathbf{x}}^k) + \nabla f(\tilde{\mathbf{x}}^k)\|^2] \quad (12)$$

$$\leq 2\mathbb{E}_{i_t}\|\nabla f_{i_t}(\mathbf{x}^t) - \nabla f_{i_t}(\mathbf{x}^*)\|^2 + 2\mathbb{E}_{i_t}\|\nabla f_{i_t}(\mathbf{x}^*) - \nabla f_{i_t}(\tilde{\mathbf{x}}^k) + \nabla f(\tilde{\mathbf{x}}^k)\|^2 \quad (13)$$

$$\leq 2\mathbb{E}_{i_t}\|\nabla f_{i_t}(\mathbf{x}^t) - \nabla f_{i_t}(\mathbf{x}^*)\|^2 + 2\mathbb{E}_{i_t}\|\nabla f_{i_t}(\mathbf{x}^*) - \nabla f_{i_t}(\tilde{\mathbf{x}}^k)\|^2 \quad (14)$$

$$\leq 4\beta_{\max}(f(\mathbf{x}^t) - f(\mathbf{x}^*)) + 4\beta_{\max}(f(\tilde{\mathbf{x}}^k) - f(\mathbf{x}^*)). \quad (15)$$

Based on these lemmas, we can prove the convergence of SVRG.

Theorem 1 *The sequence $\{\tilde{\mathbf{x}}_k\}$ in Algorithm 1 has the following property*

$$\mathbb{E}[f(\tilde{\mathbf{x}}^k) - f(\mathbf{x}^*)] \leq \left[\frac{1}{\alpha s(1-2s\beta)T} + \frac{2s\beta}{1-2s\beta} \right] \mathbb{E}[f(\tilde{\mathbf{x}}^{k-1}) - f(\mathbf{x}^*)], \quad (16)$$

where $\beta = \beta_{\max}$.

Proof 3 *By conditioning on \mathbf{x}^{t-1} , we have $\mathbb{E}\mathbf{g}^{t-1} = \nabla f(\mathbf{x}^{t-1})$ and this leads to*

$$\mathbb{E}\|\mathbf{x}^t - \mathbf{x}^*\|^2 = \|\mathbf{x}^{t-1} - \mathbf{x}^*\|^2 - 2s(\mathbf{x}^{t-1} - \mathbf{x}^*)^\top \mathbb{E}\mathbf{g}^{t-1} + s^2\mathbb{E}\|\mathbf{g}^{t-1}\|^2 \quad (17)$$

$$\leq \|\mathbf{x}^{t-1} - \mathbf{x}^*\|^2 - 2s(\mathbf{x}^{t-1} - \mathbf{x}^*)^\top \nabla f(\mathbf{x}^{t-1}) + 4\beta s^2[f(\mathbf{x}^{t-1}) - f(\mathbf{x}^*) + f(\tilde{\mathbf{x}}_k) - f(\mathbf{x}^*)] \quad (18)$$

$$\leq \|\mathbf{x}^{t-1} - \mathbf{x}^*\|^2 - 2s(f(\mathbf{x}^{t-1}) - f(\mathbf{x}^*)) + 4\beta s^2[f(\mathbf{x}^{t-1}) - f(\mathbf{x}^*) + f(\tilde{\mathbf{x}}_k) - f(\mathbf{x}^*)] \quad (19)$$

$$= \|\mathbf{x}^{t-1} - \mathbf{x}^*\|^2 - 2s(1-2s\beta)[f(\mathbf{x}^{t-1}) - f(\mathbf{x}^*)] + 4\beta s^2[f(\tilde{\mathbf{x}}_k) - f(\mathbf{x}^*)]. \quad (20)$$

Take total expectation, summing up over $t = 1, \dots, T$, then

$$\mathbb{E}\|\mathbf{x}^T - \mathbf{x}^*\|^2 \leq \mathbb{E}\|\mathbf{x}^0 - \mathbf{x}^*\|^2 - 2s(1-2s\beta)\mathbb{E}\left[\sum_{t=0}^{T-1}(f(\mathbf{x}^t) - f(\mathbf{x}^*))\right] + 4\beta s^2T[f(\tilde{\mathbf{x}}_k) - f(\mathbf{x}^*)]. \quad (21)$$

Using that $\mathbf{x}^0 = \tilde{\mathbf{x}}^k$, strong convexity says $f(\tilde{\mathbf{x}}^k) - f(\mathbf{x}^*) \geq \frac{\alpha}{2}\|\tilde{\mathbf{x}}^k - \mathbf{x}^*\|^2$ and rearranging the formulation, we have

$$2s(1-2s\beta)\mathbb{E}\left[\sum_{t=0}^{T-1}(f(\mathbf{x}^t) - f(\mathbf{x}^*))\right] \leq \mathbb{E}\|\mathbf{x}^0 - \mathbf{x}^*\|^2 - \mathbb{E}\|\mathbf{x}^T - \mathbf{x}^*\|^2 + 4\beta s^2T[f(\tilde{\mathbf{x}}_k) - f(\mathbf{x}^*)] \quad (22)$$

$$\leq (4T\beta s^2 + 2/\alpha)[f(\tilde{\mathbf{x}}_k) - f(\mathbf{x}^*)]. \quad (23)$$

Thus,

$$\mathbb{E}[f(\tilde{\mathbf{x}}^{k+1})] - f^* = \mathbb{E}\left[f\left(\frac{1}{T}\sum_{t=0}^{T-1}\mathbf{x}^t\right)\right] - f^* \quad (24)$$

$$\leq \frac{1}{T}\mathbb{E}\left[\sum_{t=0}^{T-1}f(\mathbf{x}^t)\right] - f^* \quad (25)$$

$$\leq \frac{4T\beta s^2 + 2\alpha^{-1}}{T2s(1-2s\beta)}\mathbb{E}[f(\tilde{\mathbf{x}}_k) - f(\mathbf{x}^*)] \quad (26)$$

$$= \left[\frac{1}{\alpha s(1-2s\beta)T} + \frac{2s\beta}{1-2s\beta} \right] \mathbb{E}[f(\tilde{\mathbf{x}}^k) - f(\mathbf{x}^*)]. \quad (27)$$

Moreover, we can substitute $s = 1/10\beta$ and $T = 20\beta/\alpha$, then

$$\left[\frac{1}{\alpha s(1 - 2s\beta)T} + \frac{2s\beta}{1 - 2s\beta} \right] = 7/8 < 0.9$$

.

Finally, we can obtain

$$\mathbb{E}[f(\tilde{\mathbf{x}}^k)] - f^* \leq (0.9)^k \mathbb{E}[f(\tilde{\mathbf{x}}^0) - f(\mathbf{x}^*)]. \quad (28)$$

References

- [1] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems*, 26, 2013.