

Lecture 19:

Stochastic gradient descent

E1 260

(Contd)

- Variance reduction

(SVRG)

- Algorithm

- Convergence

TA 2nd session: Nov. 3rd (Wed.) 18:00 - 19:00

Smooth and Strongly convex functions

ERM: minimize \underline{x} $f(\underline{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\underline{x})$

- f_i is β -smooth
- f is α -strongly convex

Example:

$$\begin{aligned} f(\underline{x}) &= \frac{1}{n} \|A\underline{x} - \underline{y}\|_2^2 \\ &= \frac{1}{n} \sum_{i=1}^n \underbrace{(a_i^T \underline{x} - y_i)^2}_{f_i(\underline{x})} \end{aligned}$$

- SGD needs small step size
- no self tuning ($\tilde{g}_n \not\rightarrow 0$ as $\underline{x} \rightarrow \underline{x}^*$)
- mini-batch reduces variance, but does not yield self tuning

SVRG: Stochastic variance reduced gradient

$$f(\underline{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\underline{x}) ; \tilde{g}_n = \nabla f_I(\underline{x})$$

Key observation:

$$I \sim \text{unif}(1, \dots, n)$$

Any \underline{x} and \underline{z}

$$\tilde{g}_n = \nabla f_I(\underline{x}) - \left(\nabla f_I(\underline{z}) - \nabla f(\underline{z}) \right)$$

is a stochastic gradient.

$$\underline{z} = \underline{x} \rightarrow 0$$

(Recentered gradient)

$$\underline{z} = \underline{x}^* \rightarrow 0$$

• unbiasedness:

$$E_I[\tilde{g}_n] = E[\nabla f_I(\underline{x})] - \left(E_I[\nabla f_I(\underline{z})] - \nabla f(\underline{z}) \right)$$

$$= \nabla f(\underline{x}) - \cancel{\nabla f(\underline{z})} + \cancel{\nabla f(\underline{z})}$$

$$= \nabla f(\underline{x}).$$

• Reducing variance by recentering: \underline{z} is a history point $\underline{x}^{\text{old}}$ and ∇f is the full gradient

The algorithm:

- Operates in epochs:

→ In the k^{th} epoch, take a snapshot of the current iterate $\underline{x}_k^{\text{old}} = \underline{y}_k$ and compute the batch gradient $\nabla f(\underline{x}_k^{\text{old}})$

→ Inner loop:

$$\underline{x}_k^{t+1} = \underline{x}_k^t - \eta \left\{ \nabla f_{i_t}(\underline{x}_k^t) - \left(\nabla f_{i_t}(\underline{x}_k^{\text{old}}) - \nabla f(\underline{x}_k^{\text{old}}) \right) \right\}$$

- Batch gradient is computed **once per epoch** (expensive)

• Inner loop: Requires the same effort as SGD to compute $\nabla f_{i_t}(\underline{x}_t)$

- Take advantage of **both worlds**: batch and SGD

Outer loop:

k^{th} iteration

$$\text{Set } \underline{x}_k = \underline{y}_k$$

Inner loop: for $t = 1, \dots, T$ do

$$I \sim \text{unif}(1, \dots, n)$$

$$\underline{x}_{k,t+1} = \underline{x}_{k,t} - \eta \left\{ \nabla f_I(\underline{x}_{k,t}) \right.$$

$$\left. - \left(\nabla f_I(\underline{y}_k) - \nabla f(\underline{y}_k) \right) \right\}$$

update : $\underline{y}_{k+1} = \frac{1}{T} \sum_{t=1}^T \underline{x}_{k,t}$

Variance reduction lemma:

Let $\{f_i\}$ be β -smooth and $I \sim \text{unif}(1, \dots, n)$

$$\mathbb{E}_I \left[\|\nabla f_I(\underline{x}) - \nabla f_I(\underline{x}^*)\|_2^2 \right] \leq 2\beta \left[f(\underline{x}) - f(\underline{x}^*) \right]$$

diff. goes to zero as $\underline{x} \rightarrow \underline{x}^*$
and nothing about $\nabla f_I(\underline{x})$

Proof:

$$g_i(\underline{x}) = f_i(\underline{x}) - \left[f_i(\underline{x}^*) + \nabla f_i^T(\underline{x}^*) (\underline{x} - \underline{x}^*) \right] \geq 0$$

[$f_i(\underline{x})$ is convex ; $\Rightarrow g_i(\underline{x})$ is convex]

Recall : if h is convex and β -smooth

$$h(\underline{y}) \leq h(\underline{x}) + \nabla h^T(\underline{x}) (\underline{y} - \underline{x}) + \frac{\beta}{2} \|\underline{x} - \underline{y}\|_2^2$$

$$\underline{y} := \underline{x} - \frac{1}{\beta} \nabla h(\underline{x})$$

$$\begin{aligned}
 \Rightarrow h\left(\underline{x} - \frac{1}{\beta} \nabla h(\underline{x})\right) &\leq h(\underline{x}) + \nabla h^T(\underline{x}) \left(-\frac{1}{\beta} \nabla h(\underline{x})\right) \\
 &\quad + \frac{\beta}{2} \left\| -\frac{1}{\beta} \nabla h(\underline{x}) \right\|_2^2 \\
 &= h(\underline{x}) - \frac{1}{2\beta} \left\| \nabla h(\underline{x}) \right\|_2^2
 \end{aligned}$$

Apply to $g_i(\underline{x})$:

$$0 \leq g_i\left(\underline{x} - \frac{1}{\beta} \nabla g_i(\underline{x})\right) \leq g_i(\underline{x}) - \frac{1}{2\beta} \left\| \nabla g_i(\underline{x}) \right\|_2^2$$

$$\Rightarrow -g_i(\underline{x}) \leq -\frac{1}{2\beta} \left\| \nabla g_i(\underline{x}) \right\|_2^2$$

$$\Rightarrow \left\| \nabla g_i(\underline{x}) \right\|_2^2 \leq 2\beta g_i(\underline{x})$$

Substitute for $\nabla g_i(\underline{x}) = \nabla f_i(\underline{x}) - \nabla f_i(\underline{x}^*)$

$$\Rightarrow \|\nabla f_i(\underline{x}) - \nabla f_i(\underline{x}^*)\|_2^2 \leq 2\beta \left[f_i(\underline{x}) - (f_i(\underline{x}^*) - \nabla f_i^T(\underline{x}^*)(\underline{x} - \underline{x}^*)) \right]$$

$$\mathbb{E} \left[\|\nabla f_I(\underline{x}) - \nabla f_I(\underline{x}^*)\|_2^2 \right]$$

$$\leq 2\beta \left[\mathbb{E} \left[f_I(\underline{x}) - f_I(\underline{x}^*) \right] \right]$$

$$+ \mathbb{E} \left[\nabla f_I^T(\underline{x}^*)(\underline{x} - \underline{x}^*) \right]$$

$$\mathbb{E} \left[f_I(\underline{x}) \right] = \sum_{i=1}^m \frac{1}{m} f_i(\underline{x})$$

$$= f(\underline{x})$$

$$= 0$$

$\nabla f(\underline{x}^*)$

$$\mathbb{E} \left[\|\nabla f_I(\underline{x}) - \nabla f_I(\underline{x}^*)\|_2^2 \right] \leq 2\beta \left(f(\underline{x}) - f(\underline{x}^*) \right)$$



Convergence analysis of SVRG:

Let f be α strongly convex and $\{f_i\}$ be β -smooth, then SVRG with a fixed step size $\eta = \frac{1}{10\beta}$ and inner loop size

$$T = 10 \left(\frac{\beta}{\alpha} \right) \quad \left[\frac{\beta}{\alpha} : \text{condition number} \right]$$

Then, after $S+1$ epochs (outer loop)

$$\mathbb{E} \left[f(\underline{y}_{S+1}) - f(\underline{x}^*) \right] \leq 0.9^S \left(f(\underline{y}_1) - f(\underline{x}^*) \right)$$

- Linear convergence
- Taking no. of inner loop iteration as a factor $\left(\frac{\beta}{\alpha} \right)$; convergence does not depend of $\frac{\beta}{\alpha}$; unlike GD.

Proof!

$$\mathbb{E} [f(\underline{y}_{s+1}) - f(\underline{x}^*)] \leq 0.9 (f(\underline{y}_s) - f(\underline{x}^*))$$

Recall:

$$\underline{y}_{s+1} = \frac{1}{T} \sum_{t=1}^T \underline{x}_t ; \quad \underline{x}_t \text{ is from the } s^{\text{th}} \text{ epoch.}$$

(\underline{x}_t^s)

Similarly, let us use \underline{y}
(instead of \underline{y}_{s+1})

$$\| \underline{x}_{t+1} - \underline{x}^* \|_2^2 = \| \underline{x}_t - \eta \left[\nabla f_{i_t}(\underline{x}_t) - (\nabla f_{i_t}(\underline{y}) - \nabla f(\underline{y})) \right] \|_2^2$$

$$= \textcircled{A} \| \underline{x}_t - \underline{x}^* \|_2^2 + \textcircled{B} \eta^2 \| \underline{u}_t \|^2 - \textcircled{C} 2\eta \underline{u}_t^\top (\underline{x}_t - \underline{x}^*)$$

$$\underline{u}_t = \nabla f_{i_t}(\underline{x}_t) - \nabla f_{i_t}(\underline{y}) + \nabla f(\underline{y})$$

- Recall (3) from SGD is the variance term was not going to zero. So we took small η

$$\textcircled{*} \mathbb{E}_{\mathcal{I}} [\|\underline{u}_t\|_2^2] = \mathbb{E}_{\mathcal{I}} \left[\|\nabla f_{i_t}(\underline{x}_t) - \nabla f_{i_t}(\underline{y}) + \nabla f(\underline{y}) - \nabla f_{i_t}(\underline{x}^*) + \nabla f_{i_t}(\underline{x}^*)\|_2^2 \right]$$

$$(a+b)^2 \leq 2a^2 + 2b^2$$

$$\leq 2 \mathbb{E}_{\mathcal{I}} \left[\|\nabla f_{i_t}(\underline{x}_t) - \nabla f_{i_t}(\underline{x}^*)\|_2^2 \right]$$

$$+ 2 \mathbb{E}_{\mathcal{I}} \left[\|\nabla f_{i_t}(\underline{y}) - \nabla f(\underline{y}) - \nabla f_{i_t}(\underline{x}^*)\|_2^2 \right]$$

$$\mathbb{E} \left[\nabla f_{i_t}(\underline{y}) - \nabla f(\underline{y}) - \nabla f_{i_t}(\underline{x}^*) \right] = 0 \quad \text{and} \quad \mathbb{E} \left[\|\underline{z} - \mathbb{E}(\underline{z})\|_2^2 \right] \leq \mathbb{E} \left[\|\underline{z}\|_2^2 \right]$$

$$\textcircled{*} \mathbb{E}_{\mathcal{I}} \left[\|\underline{u}_t\|_2^2 \right] \leq 2 \mathbb{E}_{\mathcal{I}} \left[\|\nabla f_{i_t}(\underline{x}_t) - \nabla f_{i_t}(\underline{x}^*)\|_2^2 \right] + 2 \mathbb{E}_{\mathcal{I}} \left[\|\nabla f_{i_t}(\underline{y}) - \nabla f_{i_t}(\underline{x}^*)\|_2^2 \right]$$

Recall Variance reduction lemma:

$$\mathbb{E}_{\mathcal{I}} \left[\left\| \nabla f_{\mathcal{I}}(\underline{x}) - \nabla f_{\mathcal{I}}(\underline{x}^*) \right\|_2^2 \right] \leq 2\beta \left[f(\underline{x}) - f(\underline{x}^*) \right]$$

$$\textcircled{*} \quad \mathbb{E}_{\mathcal{I}} \left[\left\| \underline{u}_t \right\|_2^2 \right] \leq 4\beta \left[f(\underline{x}_t) - f(\underline{x}^*) + f(\underline{y}) - f(\underline{x}^*) \right]$$

Now, let work with \textcircled{C} :

$$2\eta \underline{u}_t^\top (\underline{x}_t - \underline{x}^*)$$

$$\begin{aligned} \mathbb{E} \left[2\eta \underline{u}_t^\top (\underline{x}_t - \underline{x}^*) \right] &= 2\eta \mathbb{E} \left[\underline{u}_t \right]^\top (\underline{x}_t - \underline{x}^*) \\ &= 2\eta \nabla f^\top(\underline{x}) (\underline{x}_t - \underline{x}^*) \\ &\geq 2\eta \left(f(\underline{x}_t) - f(\underline{x}^*) \right) \end{aligned}$$

(Convexity)

Combining everything:

$$\begin{aligned}\mathbb{E}_{\underline{\Gamma}} [\|\underline{x}_{t+1} - \underline{x}^*\|] &\leq \|\underline{x}_t - \underline{x}^*\|_2^2 \\ &+ 4\beta\eta^2 \left[f(\underline{x}_t) - f(\underline{x}^*) + f(\underline{y}) - f(\underline{x}^*) \right] \\ &- 2\eta \left[f(\underline{x}_t) - f(\underline{x}^*) \right] \\ &\leq \|\underline{x}_t - \underline{x}^*\|_2^2 \\ &- 2\eta(1 - 2\beta\eta) \left[f(\underline{x}_t) - f(\underline{x}^*) \right] \\ &+ 4\eta^2\beta \left[f(\underline{y}) - f(\underline{x}^*) \right]\end{aligned}$$

Iterating this inequality:

$$\mathbb{E}_I \left[\|\underline{x}_{t+1} - \underline{x}^*\|_2^2 \right] \leq \|\underline{x}_1 - \underline{x}^*\|_2^2 - 2\eta (1 - 2\beta\eta) \mathbb{E} \left[\sum_{k=1}^T f(\underline{x}_k) - f(\underline{x}^*) \right] + 4\eta^2 \beta \cdot \mathbb{E} \left[f(y) - f(\underline{x}^*) \right]$$

- we have $\underline{x}_1 = y$ (our initialization)
- f is α strongly convex

$$\|y - \underline{x}^*\|_2^2 \leq \frac{2}{\alpha} [f(y) - f(\underline{x}^*)]$$

$$2\eta (1 - 2\beta\eta) \left(\mathbb{E} f\left(\frac{1}{T} \sum_k \underline{x}_k\right) - f(\underline{x}^*) \right) \leq \left(\frac{2}{\alpha} + 4\eta^2 \beta T \right) \cdot \frac{1}{T} [f(y) - f(\underline{x}^*)]$$

$$\Rightarrow \mathbb{E} f(\underline{y}_{s+1}) - f(\underline{x}^*) \leq \frac{\left(\frac{2}{\alpha} + 4\eta^2\beta\tau\right) \cdot \frac{1}{\tau} \left[f(\underline{y}_s) - f(\underline{x}^*)\right]}{2\eta(1-2\beta\eta)}$$

$$= 0.9 \left[f(\underline{y}_s) - f(\underline{x}^*) \right]$$

with

$$\eta = \frac{1}{10\beta} \quad \text{and} \quad \tau = 10 \cdot \left(\frac{\beta}{\alpha}\right)$$

GD vs. SGD vs. SVRG

$$\text{GD: } x_{t+1} = x_t - \eta \nabla f(x_t) \cdot [n \text{ grad. comp.}]$$

$$\text{SGD: } x_{t+1} = x_t - \eta \nabla f_{i_t}(x_t) \approx [1 \text{ grad. comp.}]$$

$$\text{SVRG: } x_{t+1} = x_t - \eta \left[\nabla f_{i_t}(x_t) - \left(\nabla f_{i_t}(y) - \nabla f(y) \right) \right] \approx \begin{cases} n + T \\ \text{grad comp.} \end{cases}$$
$$\approx n + 10 \cdot \left(\frac{\beta_0}{\alpha} \right)$$

↑
depends on the
condn. numbers

GD vs. SGD vs. SVRG: Σ -accuracy

GD: $O\left(\frac{B}{\alpha} \log\left(\frac{1}{\epsilon}\right)\right)$ iterations

Momentum $O\left(\sqrt{\frac{B}{\alpha}} \log\left(\frac{1}{\epsilon}\right)\right)$

So $O\left(m \cdot \frac{B}{\alpha} \log\left(\frac{1}{\epsilon}\right)\right)$ grad. computation

SGD: $O\left(\frac{1}{\alpha \epsilon}\right)$ iterations
 $= O\left(\frac{1}{\alpha \epsilon}\right)$ grad. computation

SVRG: $O\left(\log\left(\frac{1}{\epsilon}\right)\right)$ iteration
 $O\left((m + \frac{B}{\alpha}) \log\left(\frac{1}{\epsilon}\right)\right)$ grad. computation