

LECTURE - 7

Note Title

28-Jun-19

Stochastic (Variance Reduced) Gradient (SVRG)

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x). \quad \text{E.g., } f_i(x) = \frac{1}{2} (\alpha_i^T x - b_i)^2$$

Each f_i is L -Smooth & Convex

$$L = \max_i \|\alpha_i\|^2$$

f is μ -Strongly Convex.

$$\mu = \sigma_{\min} \left(\frac{1}{n} \sum_i \alpha_i \alpha_i^T \right)$$

GD: $O\left(\underbrace{n \cdot d}_{\text{one } \nabla f \text{ Computation}} \cdot \underbrace{\frac{L}{\mu}}_{\text{\# iterations}} \log \frac{1}{\epsilon} \right)$ Computations.

SGD: $O\left(d \cdot \frac{\sigma \cdot \|\alpha_0 - \alpha^*\|}{\epsilon^2}\right)$

Std. dev. in the Stochastic gradients
 one ∇f_i computation. # iterations

In most cases: $\sigma \ll n \cdot \frac{L}{\mu}$

	Dep. on param. prob.	Desired acc.
SGD	Usually better	
GD		Usually better

→ Best of both worlds?

Variance
reduction
step

Pick x_0
Compute $\nabla f(x_0) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_0)$

Do $t=1, \dots, T$

Pick 'i' uniformly at random from $1, \dots, n$

$$x_t = x_{t-1} - \eta [\nabla f_i(x_{t-1}) - \nabla f_i(x_0) + \nabla f(x_0)].$$

Suppose $x_0 = x^*$

$$\begin{aligned} x_1 &= x_0 - \eta [\nabla f_i(x_0) \\ &\quad - \nabla f_i(x_0) \\ &\quad + \nabla f(x_0)] \\ &= x_0 = x^* \end{aligned}$$

Contrast with SGD step: $x_t = x_{t-1} - \eta \nabla f_i(x_{t-1})$.

$$\mathbb{E}_i [\nabla f_i(x_0) - \nabla f(x_0)] = \mathbb{E}_i [\nabla f_i(x_0)] - \nabla f(x_0) = \nabla f(x_0) - \nabla f(x_0) = \underline{\underline{0}}.$$

Recall that the convergence guarantee of SGD depends on the variance in stochastic gradients.

$$\sigma^2 \triangleq \mathbb{E} \left[\underbrace{\| \nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0) - \nabla f(x_t) \|^2}_{\text{Stochastic gradient}} \right]$$

$\underbrace{\| \nabla f(x_0) - \nabla f(x_t) \|^2}_{\text{True gradient}} = \mathbb{E}[\text{stoch. grad.}]$

$$[\because (a+b)^2 \leq 2a^2 + 2b^2]$$

$$\begin{aligned} &\leq 2 \mathbb{E} \left[\| \nabla f_i(x_t) - \nabla f_i(x_0) \|^2 \right] + 2 \mathbb{E} \left[\| \nabla f(x_0) - \nabla f(x_t) \|^2 \right] \\ \text{Jensen's} &\stackrel{\text{ineq.}}{\geq} \underline{4 \mathbb{E} \left[\| \nabla f_i(x_t) - \nabla f_i(x_0) \|^2 \right]} \quad \left\{ \begin{array}{l} \mathbb{E}[\|w\|^2] \geq \|\mathbb{E}[w]\|^2 \\ \text{Jensen's inequality} \end{array} \right\} \end{aligned}$$

$$\rightarrow x^* = \operatorname{argmin}_x f(x)$$

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

$$\mathbb{E} \left[\|\nabla f_i(x_t) - \nabla f_i(x_0)\|^2 \right] \leq 2 \left[\mathbb{E} \left[\|\nabla f_i(x_t) - \nabla f_i(x^*)\|^2 \right] + \mathbb{E} \left[\|\nabla f_i(x_0) - \nabla f_i(x^*)\|^2 \right] \right].$$

$$\mathbb{E} \left[\|\nabla f_i(x) - \nabla f_i(x^*)\|^2 \right] \leq \underline{\frac{2L}{3} [f(x) - f(x^*)]}.$$

Goal: Bound

Approach:

$$g_i(x) = \underbrace{f_i(x)}_{\substack{\downarrow \\ L\text{-Smooth} \\ \text{Convex}}} - \underbrace{\langle \nabla f_i(x^*), x \rangle}_{\substack{\downarrow \\ \text{Linear}}}$$

$$\nabla g_i(x^*) = \nabla f_i(x^*) - \nabla f_i(x^*) = 0$$

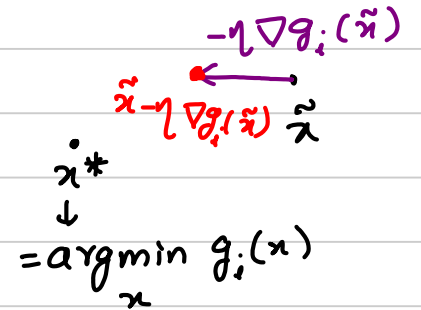
$$\Rightarrow g_i(x^*) \stackrel{\text{Conv. } g_i}{\leq} g_i(\tilde{x} - \eta \nabla g_i(\tilde{x}))$$

$$g_i \stackrel{L\text{-Smooth}}{\leq} g_i(\tilde{x}) - \underbrace{\eta \left(2 - \frac{\eta L}{2}\right) \|\nabla g_i(\tilde{x})\|^2}_{\text{subtraction}}$$

$$\eta = \frac{1}{L} \longrightarrow \leq g_i(\tilde{x}) - \frac{3}{2L} \|\nabla g_i(\tilde{x})\|^2$$

$$\|\nabla g_i(\tilde{x})\|^2 \leq [g_i(\tilde{x}) - g_i(x^*)] \frac{2L}{3}$$

$$\|\nabla f_i(\tilde{x}) - \nabla f_i(x^*)\|^2 \leq [g_i(\tilde{x}) - g_i(x^*)] \frac{2L}{3}$$



Taking an average,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\tilde{x}) - \nabla f_i(x^*)\|^2 &\leq \frac{1}{n} \cdot \frac{2L}{3} \sum_{i=1}^n \left[\underbrace{g_i(\tilde{x}) - g_i(x^*)}_{f_i(\tilde{x}) - \langle \nabla f_i(x^*), \tilde{x} \rangle} \right] \\ &= \frac{2L}{3n} \sum_{i=1}^n f_i(\tilde{x}) - f_i(x^*) - \langle \nabla f_i(x^*), \tilde{x} - x^* \rangle \\ &= \frac{2L}{3} \left\{ \frac{1}{n} \sum_{i=1}^n f_i(\tilde{x}) - \frac{1}{n} \sum_{i=1}^n f_i(x^*) \right. \\ &\quad \left. - \left\langle \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^*), \tilde{x} - x^* \right\rangle \right\} \end{aligned}$$

$$= \frac{2L}{3} \left\{ f(\tilde{x}) - f(x^*) - \underbrace{\langle \nabla f(x^*), \tilde{x} - x^* \rangle}_0 \right\}$$

Lemma: $\mathbb{E} \left[\|\nabla f_i(\tilde{x}) - \nabla f_i(x^*)\|^2 \right] \leq \frac{2L}{3} \left[f(\tilde{x}) - f(x^*) \right].$

Furthermore,

$$\sigma^2 \triangleq \mathbb{E} \left[\|\nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0) - \nabla f(x_t)\|^2 \right]$$

[2nd 3rd page computation] $\leq 8 \left\{ \mathbb{E} \left[\|\nabla f_i(x_t) - \nabla f_i(x^*)\|^2 \right] + \mathbb{E} \left[\|\nabla f_i(x_0) - \nabla f_i(x^*)\|^2 \right] \right\}.$

$$\leq \frac{16L}{3} \left[f(x_t) - f(x^*) + f(x_0) - f(x^*) \right].$$

$$\text{SVRG} \quad : \quad x_{t+1} = x_t - \eta \left[\underbrace{\nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0)}_{\text{stochastic gradient}} \right]$$

$$\mathbb{E} \left[\|x_{t+1} - x^*\|^2 \right] = \mathbb{E} \left[\|x_t - x^*\|^2 \right] - 2\eta \left\langle \underbrace{\mathbb{E} \left[\nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0) \right]}_{\text{stochastic gradient}}, x_t - x^* \right\rangle + \eta^2 \mathbb{E} \left[\underbrace{\| \nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0) \|^2}_{\text{variance}} \right].$$

$$= \mathbb{E} \left[\|x_t - x^*\|^2 \right] - 2\eta \left\langle \underbrace{\nabla f(x_t)}_{\text{true gradient}}, x_t - x^* \right\rangle + \eta^2 \underbrace{\| \nabla f(x_t) \|^2}_{\text{true gradient norm}} + \eta^2 \underbrace{\sigma^2}_{\text{variance}}$$

[By L -smoothness of $f(\cdot)$, $\| \nabla f(x_t) \|^2 \leq L [f(x_t) - f(x^*)]$]

$$\leq \mathbb{E} \left[\|x_t - x^*\|^2 \right] - 2\eta \underbrace{[f(x_t) - f(x^*)]}_{\text{true gradient}} + \eta^2 L \underbrace{(f(x_t) - f(x^*))}_{\text{true gradient}}$$

$$+ \eta^2 \cdot \frac{16L}{3} \left[f(x_t) - f(x^*) + f(x_0) - f(x^*) \right]$$

$$\mathbb{E}[\|x_{t+1} - x^*\|^2] \leq \mathbb{E}[\|x_t - x^*\|^2] - \eta \left(2 - \frac{16\eta L}{3}\right) (f(x_t) - f(x^*)) \\ + \frac{16\eta^2 L}{3} (f(x_0) - f(x^*)).$$

Telescoping and averaging,

$$\frac{1}{T+1} \sum_{t=0}^T \left[\mathbb{E}[f(x_t)] - f(x^*) \right] \leq \frac{\sqrt{\|x_0 - x^*\|^2}}{\eta \left(2 - \frac{16\eta L}{3}\right) (T+1)} + \frac{16\eta L}{3 \left(2 - \frac{16\eta L}{3}\right)} [f(x_0) - f(x^*)].$$

$$f \text{ is } \mu\text{-str. cvx.} \quad \therefore \quad f(x_0) - f(x^*) \geq \frac{\mu}{2} \underbrace{\|x_0 - x^*\|^2}.$$

$$\text{Av. exp. Suboptimality} \leq \left[\underbrace{\frac{2}{\mu \eta (2 - \frac{199L}{3})^{T+1}}}_{\frac{L}{\mu} = K} + \underbrace{\frac{169L}{3 (2 - \frac{199L}{3})}}_{\leq \frac{1}{10}} \right] \underbrace{\{f(x_0) - f(x^*)\}}_{\text{Initial Suboptimality}}$$

$$\eta = \frac{1}{100 \cdot L}$$

$$\approx \left[\frac{2K}{T+1} + \frac{1}{10} \right] [\text{Initial Suboptimality}]$$

$$T > 6 \cdot K \Rightarrow \underline{\text{Av. exp. Suboptimality} \leq \frac{1}{2} \text{ Initial Subopt.}}$$

Pick x_0 .

For epochs $= 1, \dots, \log(\frac{1}{\epsilon})$.

SVRG epoch $\left[\begin{array}{l} \text{Compute } \nabla f(x_0) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_0) \leftarrow O(nd) \\ \text{For } t=0, \dots, T=5k \\ \quad \text{--- pick } i \text{ u.a.r. from } 1, \dots, n. \\ \quad \text{--- } x_{t+1} = x_t - \eta \left[\nabla f_i(x_t) - \nabla f_i(x_0) + \nabla f(x_0) \right] \end{array} \right. \left. \begin{array}{l} O(kd) \end{array} \right.$

Reset $x_0 = x_1, \dots, x_{T+1}$ uniformly at random.

$$O\left(\underbrace{(nd + kd)}_{\text{epoch cost}} \underbrace{\log \frac{1}{\epsilon}}_{\text{number of epochs}}\right)$$

$$f_i(x) = (a_i^T x - b_i)^2$$

$$O\left((nd + kd) \cdot \log \frac{1}{\epsilon} \right)$$

\downarrow
 $\max_i \|a_i\|^2$
 μ

Importance Sampling

i w.p. $p_i \propto \|a_i\|^2$.

$$f(x) = \sum_{i=1}^n p_i \underbrace{\left(\frac{1}{np_i} (a_i^T x - b_i)^2 \right)}$$

$$\text{Smoothness} \left(\frac{1}{np_i} (a_i^T x - b_i)^2 \right) = \frac{\sum_{j=1}^n \|a_j\|^2}{n}$$

$$p_i = \frac{\|a_i\|^2}{\sum_{j=1}^n \|a_j\|^2}$$

Uniform Sampling

$$k_{us} = \frac{\max_i \|a_i\|^2}{\mu}$$

Weighted Sampling

$$k_{ws} = \frac{\sum_{j=1}^n \|a_j\|^2}{n\mu}$$

$$\underline{k_{ws} \leq k_{us}}$$