

CS292F Lecture 8 stochastic subgradient method

Alg: for $k=1, \dots, K$

$$x^{(k+1)} = x^{(k)} - t_k \hat{g}_k$$

Assumption ① $E[\hat{g}_k | x_k] = g_k$ when g_k is a subgradient of f at x_k

② $E[\|\hat{g}_k - g_k\|^2 | x_k] \leq \underbrace{6^2}_{\text{hidden dimension}}$

③ f is G -Lipschitz, $\|g_k\| \leq G$

$x \rightarrow x^+$

$$\|x^+ - x^*\|^2 = \|x - t\hat{g} - x^*\|^2$$

$$= \|x - x^*\|^2 - 2t \hat{g}^T (x - x^*) + t^2 \|\hat{g}\|^2$$

take $E[\cdot | x]$

$$E[\|x^+ - x^*\|^2 | x] = \|x - x^*\|^2 - 2t \underbrace{E[\hat{g} | x]^T}_{\bar{g}^T} (x - x^*) + t^2 E[\|\hat{g}\|^2 | x]$$

Apply ① $\Rightarrow \|x - x^*\|^2 - 2t \bar{g}^T (x - x^*) + t^2 (\|g\|^2 + E[\|\hat{g} - g\|^2 | x])$

② $\Rightarrow \|x - x^*\|^2 - 2t \bar{g}^T (x - x^*) + t^2 (\|g\|^2 + G^2)$
 $\leq \|x - x^*\|^2 - 2t (f(x) - f^*) + t^2 (G^2 + G^2)$

take full $E[\cdot]$

$$E[\|x^+ - x^*\|^2] \leq E[\|x - x^*\|^2] - 2t (E[f(x)] - f^*) + t^2 (G^2 + G^2)$$

Telescope $k=1, 2, \dots, k$

$$E[\|x^{(2)} - x^*\|^2] \leq E[\|x^{(1)} - x^*\|^2] - 2t_1[E[f(x^{(1)})] - f^*] + t_1^2(C^2 + G^2)$$

$$E[\|x^{(3)} - x^*\|^2] \leq E[\|x^{(2)} - x^*\|^2] - 2t_2[E[f(x^{(2)})] - f^*] + t_2^2(C^2 + G^2)$$

$$E[\|x^{(k+1)} - x^*\|^2] \leq E[\|x^{(k)} - x^*\|^2] - 2t_k[E[f(x^{(k)})] - f^*] + t_k^2(C^2 + G^2)$$

Add up

$$2 \sum_{i=1}^k t_i (E[f(x^{(i)})] - f^*) \leq E[\|x^{(1)} - x^*\|^2] - E[\|x^{(k+1)} - x^*\|^2] + \sum_{i=1}^k t_i^2 (C^2 + G^2)$$

Strategy do $\min_i a_i \sum b_i \leq \sum a_i b_i$ or $\min_i b_i \sum a_i \leq \sum a_i b_i$

1. $2 \left(\sum_{i=1}^k t_i \right) \cdot \min_i (E[f(x^{(i)})] - f^*) \leq$

$$E[\|x^{(1)} - x^*\|^2] + \left(\sum_{i=1}^k t_i^2 \right) (C^2 + G^2)$$

$$2 \left(\sum_{i=1}^k t_i \right) \rightarrow \infty$$

Strategy 2:

$$2 \min_i t_i \cdot \sum_{i=1}^k (E[f(x^{(i)})] - f^*) \leq \|x^{(1)} - x^*\|^2 + \sum_{i=1}^k t_i^2 (G^2 + \sigma^2)$$

Divide both sides by $2K \cdot \min_i t_i$

$$\frac{1}{2K} \sum_{i=1}^k (E[f(x^{(i)})] - f^*) \leq$$

Average suboptimality
through out the updates

$$t_i = \frac{1}{\sqrt{K}}$$

$$\frac{\|x^{(1)} - x^*\|^2 + \sum_{i=1}^k t_i^2 (G^2 + \sigma^2)}{2 \cdot K \cdot \min_i t_i}$$

$$= \frac{\|x^{(1)} - x^*\|^2 + G^2 + \sigma^2}{2 \sqrt{K}}$$

only difference from
non-stochastic
subgradient

Nonconvex case: assumption. f is L -smooth

$$E[\hat{g}|x] = \nabla f(x), \quad E[\|\hat{g} - \nabla f(x)\|^2|x] \leq G^2$$

By L -smoothness

$$f(x^+) \leq f(x) + \langle x^+ - x, \nabla f(x) \rangle + \frac{L}{2} \|x^+ - x\|^2$$

$$\begin{aligned} E[f(x^+)|x] &\leq f(x) - t \langle E[\hat{g}|x], \nabla f(x) \rangle + \frac{L}{2} t^2 E[\|\hat{g}\|^2|x] \\ &= f(x) - t \|\nabla f(x)\|^2 + \frac{L}{2} t^2 [\|\nabla f(x)\|^2 + G^2] \end{aligned}$$

$$t \leq \frac{1}{L}$$

$$\begin{aligned} &= f(x) - \underbrace{\left(t - \frac{L t^2}{2}\right)}_{\frac{t}{2}} \|\nabla f(x)\|^2 + \frac{L}{2} t^2 G^2 \\ &\Rightarrow f(x) - \underbrace{\left(t - \frac{L}{2} t\right)}_{\frac{t}{2}} \|\nabla f(x)\|^2 + \frac{L}{2} t^2 G^2 \end{aligned}$$

new for SGD.

Take $E[\cdot]$

$$E[f(x^+)] \leq \underbrace{E[f(x)]} - \frac{t}{2} E[\|\nabla f(x)\|^2] + \frac{L}{2} t^2 G^2$$

Telescope:

$$\sum_{i=1}^k \frac{t_i}{2} E[\|\nabla f(x^{(i)})\|^2] \leq E[f(x^{(1)})] - E[f(x^{(k)})] + \underbrace{\frac{LG^2}{2} \sum_{i=1}^k t_i^2}_{}$$

$$t_i = \frac{1}{\sqrt{k}}$$

$$\frac{1}{2\sqrt{k}} \sum_{i=1}^k E[\|\nabla f(x^{(i)})\|^2] \leq E[f(x^{(1)})] - f^* + \frac{LG^2}{2} \cdot 1$$

Divide both sides by $\frac{\sqrt{k}}{2}$

$$\frac{1}{k} \sum_{i=1}^k E[\|\nabla f(x^{(i)})\|^2] \leq \frac{2(E[f(x^{(1)})] - f^* + \frac{LG^2}{2})}{\sqrt{k}} \quad \square$$