

GD & Convergence

2020年9月23日 8:24

Gradient Descent:

$$f(w') = f(w) + \underbrace{\langle \nabla f(w), w' - w \rangle}_{\text{Estimation}} + \underbrace{\frac{\eta}{2} \|w' - w\|^2}_{\text{tail}}$$

- Smoothness assumption:

$$\exists L, |\nabla f(w)| \leq L \text{ for all } w.$$

which means gradient is L -Lipschitz.

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

$$\text{If } w' = w - \eta \nabla f(w)$$

$$\begin{aligned} \Rightarrow f(w') - f(w) &\leq \langle \nabla f(w), -\eta \nabla f(w) \rangle + \frac{\eta^2}{2} \|\nabla f(w)\|^2 \\ &= -\eta \left(1 - \frac{\eta}{2}\right) \|\nabla f(w)\|^2 \end{aligned}$$

$$\text{set } \eta < \frac{2}{L}. \quad \underline{\text{make sure } f(w) \text{ non-increasing}}$$

Convergence Analysis

Suppose f is L -smooth & convex.

$$\text{Goal} \Rightarrow \boxed{\eta \leq \frac{1}{L} \Rightarrow f(w_t) \leq f(w^*) + \frac{\|w_0 - w^*\|^2}{2\eta t}}$$

$$\left. \begin{array}{l} 1. \text{ By } L\text{-smooth. } f(w_{t+1}) \leq f(w_t) - \frac{\eta}{2} \|\nabla f(w_t)\|^2 \quad (\text{proved above}) \\ 2. \text{ By convexity. } f(w_t) \leq f(w^*) + \langle \nabla f(w_t), w_t - w^* \rangle \end{array} \right\}$$

$$\begin{aligned} f(w_{t+1}) &\leq f(w^*) + \langle \nabla f(w_t), w_t - w^* \rangle - \frac{\eta}{2} \|\nabla f(w_t)\|^2 \\ &= f(w^*) - \frac{1}{\eta} \langle w_{t+1} - w_t, w_t - w^* \rangle - \frac{1}{2\eta} \|w_t - w^*\|^2 \\ &\leq f(w^*) + \frac{1}{2\eta} \|w_t - w^*\|^2 - \frac{1}{2\eta} \|w_{t+1} - w^*\|^2 \end{aligned}$$

$$\sum_{i=0}^{t-1} (f(w_{i+1}) - f(w^*)) \leq \frac{1}{2\eta} (\|w_0 - w^*\|^2 - \|w_t - w^*\|^2) \leq \frac{\|w_0 - w^*\|^2}{2\eta}$$

Since $f(w_t)$ non-increasing,

$$f(w_t) - f(w^*) \leq \frac{\|w_0 - w^*\|^2}{2\eta t}$$

GD has two limitations.

- Computing full gradient is slow for big data
- (Theoretical) get stuck at $\nabla f = 0$ points

— (Theoretical) get stuck at $\nabla f = 0$ points

SGD

$$GD: \nabla L(w, X, Y) = \frac{1}{N} \sum_i l(w, x_i, y_i).$$

$$SGD: G_t = \frac{1}{|S|} \sum_{i \in S} l(w, x_i, y_i)$$

S: mini-batch

Analysis of SGD:

f : L -smooth convex function, $\text{Var}(G_t) \leq \sigma^2$

$$\Rightarrow E[f(\bar{w}_t)] \leq f(w^*) + \frac{\|w_0 - w^*\|_2^2}{2\eta t} + \eta \sigma^2$$

$$(\bar{w}_t = \frac{\sum_{i=1}^t w_i}{t})$$

$$\begin{aligned} E[f(w_{t+1})] &\leq f(w_t) + E\langle \nabla f(w_t), w_{t+1} - w_t \rangle + E\left(\frac{L}{2}\|w_{t+1} - w_t\|_2^2\right) \quad (\text{smoothness}) \\ &= f(w_t) - \eta \langle \nabla f(w_t), \nabla f(w_t) \rangle + \frac{\eta^2}{2} E(\|G_t\|_2^2). \\ &= f(w_t) - \eta \|\nabla f(w_t)\|_2^2 + \frac{\eta^2}{2} (\|\nabla f(w_t)\|_2^2 + \text{Var}(G_t)) \quad \text{Var}(S) = \widehat{(S - \bar{S})^2} \\ &\leq f(w_t) - \eta \left(1 - \frac{\eta^2}{2}\right) \|\nabla f(w_t)\|_2^2 + \frac{\eta^2}{2} \sigma^2 \end{aligned}$$

By convexity $f(w) \leq f(w^*) + \langle \nabla f(w), w - w^* \rangle$

$$\begin{aligned} E[f(w_{t+1})] &\leq f(w^*) + \langle \nabla f(w_t), w_t - w^* \rangle - \frac{\eta}{2} \|\nabla f(w_t)\|_2^2 + \frac{\eta}{2} \sigma^2 \\ &\leq f(w^*) + E(\underbrace{\langle G_t, w_t - w^* \rangle - \frac{\eta}{2} \|G_t\|_2^2}_{= \frac{1}{2\eta} (\|w_t - w^*\|_2^2 - \|w_{t+1} - w^*\|_2^2)}) + \eta \sigma^2 \end{aligned}$$

↓

(Bregman divergence)

$$\text{Hence, } E(f(w_{t+1})) \leq f(w^*) + \frac{1}{2\eta} E(\|w_t - w^*\|_2^2 - \|w_{t+1} - w^*\|_2^2) + \eta \sigma^2$$

$$\begin{aligned} \sum_{i=0}^t [E(f(w_{i+1})) - f(w^*)] &\leq \frac{1}{2\eta} (\|w_0 - w^*\|_2^2 - E\|w_t - w^*\|_2^2) + t\eta \sigma^2 \\ &\leq \frac{1}{2\eta} \|w_0 - w^*\|_2^2 + t\eta \sigma^2 \end{aligned}$$

$$f(\bar{w}_t) \leq \frac{\sum_{i=1}^t f(w_i)}{t} \quad (\text{we do not have } E f(w_i) \text{ non-increasing}).$$

$$\Rightarrow E f(\bar{w}_t) \leq f(w^*) + \frac{\|w_0 - w^*\|_2^2}{2\eta t} + \eta \sigma^2$$

If f is strongly convex and smooth,

GD $\Rightarrow (1-p)^t$ linear convergence

SGD $\Rightarrow \frac{1}{t}$ convergence \Rightarrow noise helps escaping local optimum

SVRG (Stochastic variance reduction gradient)

Compute full gradient for once in stats

SVRG (Stochastic variance reduction gradient)

Compute full gradient for every m steps

$$X_0 \Rightarrow \text{full gradient } g_1$$

$$\Rightarrow SG \text{ (with same data of } X_t) : g_2$$

$$X_t : SG \quad g_3$$

$$\Rightarrow \text{unbiased } \underbrace{g_1 + g_3 - g_2}_{\downarrow}$$

$$\mathbb{E}(g_1 + g_3 - g_2) = \mathbb{E}(g_3) \text{ an estimation for GD}(X_t)$$

but reduce variance

SVRG algorithm

For $s=1, 2, \dots$

$$\tilde{w} = \tilde{w}_{s-1} \Leftarrow \text{place we need to compute full gradient.}$$

$$\tilde{u} = \frac{1}{N} \sum_{i=1}^N \nabla l_i(\tilde{w}) = \nabla f(\tilde{w}) \quad (\text{full Gradient})$$

For $t=1, 2, \dots, m$ # SG.

- Randomly pick mini-batch S

$$\cdot w_t = w_{t-1} - \eta \left[\tilde{u} + \sum_{i \in S} \nabla l_i(w_{t-1}) - \sum_{i \in S} \nabla l_i(\tilde{w}) \right]$$

$$\text{Option 1: } \tilde{w}_s = w_m$$

$$\text{Option 2: } \tilde{w}_s = w_i \text{ for } i \text{ randomly chosen from } [m].$$

Analysis: (Assume L -smooth & μ -strongly convex).

$$\mathbb{E} \left\| \tilde{u} + \sum_{i \in S} \nabla l_i(w_{t-1}) - \sum_{i \in S} \nabla l_i(\tilde{w}) \right\|$$

$$\leq 2 \mathbb{E} \left\| \sum \nabla l_i(w_{t-1}) - \sum \nabla l_i(w^*) \right\|^2 + 2 \mathbb{E} \left\| \sum_{i \in S} \nabla l_i(\tilde{w}) - \nabla l_i(w^*) - \nabla f(\tilde{w}) \right\|^2$$

$$= 2 \mathbb{E} \left\| \sum \nabla l_i(w_{t-1}) - \sum \nabla l_i(w^*) \right\|^2 + 2 \mathbb{E} \left\| \sum_{i \in S} \nabla l_i(\tilde{w}) - \nabla l_i(w^*) - \mathbb{E} \left[\sum \nabla l_i(\tilde{w}) - \sum \nabla l_i(w^*) \right] \right\|^2$$

$$\leq 2 \mathbb{E} \left\| \sum \nabla l_i(w_{t-1}) - \nabla l_i(w^*) \right\|^2 + 2 \mathbb{E} \left\| \sum \nabla l_i(\tilde{w}) - \nabla l_i(w^*) \right\|^2$$

$$\leq 4L (f(w_{t-1}) - f(w^*) + f(\tilde{w}) - f(w^*))$$

$$\mathbb{E} \|w_t - w^*\|^2$$

$$= \mathbb{E} \|w_{t-1} - w^*\|^2 - 2\eta \langle w_{t-1} - w^*, \mathbb{E}(w_t) \rangle + \eta^2 \mathbb{E}(\|w_t\|^2)$$

$$\leq \mathbb{E} \|w_{t-1} - w^*\|^2 - 2\eta \langle w_{t-1} - w^*, \nabla f(w_{t-1}) \rangle + 4L\eta^2 (f(w_{t-1}) - f(w^*) + f(\tilde{w}) - f(w^*))$$

$$\leq \mathbb{E} \|w_{t-1} - w^*\|^2 - 2\eta (f(w_{t-1}) - f(w^*)) + 4L\eta^2 (f(w_{t-1}) - f(w^*) + f(\tilde{w}) - f(w^*))$$

$$= \mathbb{E} \|w_{t-1} - w^*\|^2 - 2\eta (1-2L\eta) (f(w_{t-1}) - f(w^*)) + 4L\eta^2 (f(\tilde{w}) - f(w^*))$$

$$\Rightarrow \mathbb{E} \|w_m - w^*\|^2 \leq \mathbb{E} \|w_0 - w^*\|^2 - 2\eta (1-2L\eta) \mathbb{E} \left(\sum_{t=1}^m f(w_{t-1}) - f(w^*) \right) + 4Lm\eta^2 \mathbb{E} (f(\tilde{w}) - f(w^*))$$

$$= \mathbb{E} \|\tilde{w} - w^*\|^2 - 2\eta (1-2L\eta) m \mathbb{E} (f(\tilde{w}) - f(w^*)) + 4Lm\eta^2 \mathbb{E} (f(\tilde{w}) - f(w^*))$$

$$\leq \frac{2}{\mu} \mathbb{E} (f(\tilde{w}) - f(w^*)) - 2\eta (1-2L\eta) m \mathbb{E} (f(\tilde{w}) - f(w^*)) + 4Lm\eta^2 \mathbb{E} (f(\tilde{w}) - f(w^*))$$

$$\Rightarrow \mathbb{E}[f(\tilde{w}_s) - f(w^*)] \leq \left[\frac{1}{\mu \gamma^{(1-2L\gamma)m}} + \frac{2L\gamma}{1-2L\gamma} \right] \mathbb{E}[f(\tilde{w}_{s-1}) - f(w^*)]$$

\Rightarrow Linear convergence rate