

Stochastic Gradient Descent

Gradient Descent \rightarrow linear regression
 Coordinate Descent \rightarrow lasso
 Stochastic G.D.

W

Machine Learning Problems

$$n \gg 1/M$$
$$d \gg 1/M$$

- Given data:

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$

- Learning a model's parameters: $\frac{1}{n} \sum_{i=1}^n \ell_i(w) = \frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top w)^2$

Gradient Descent:

$$w_{t+1} = w_t - \eta \underbrace{\nabla_w \left(\frac{1}{n} \sum_{i=1}^n \ell_i(w) \right)}_{O(d) \times n} \Big|_{w=w_t}$$

Machine Learning Problems

- Given data:

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$

- Learning a model's parameters: $\frac{1}{n} \sum_{i=1}^n \ell_i(w) = \ell(w)$

$$\nabla \ell(w) = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(w)$$

Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla_w \left(\frac{1}{n} \sum_{i=1}^n \ell_i(w) \right) \Big|_{w=w_t}$$

Stochastic Gradient Descent:

old time
Random

$$w_{t+1} = w_t - \eta \nabla_w \ell_{I_t}(w) \Big|_{w=w_t}$$

I_t drawn uniform at random from $\{1, \dots, n\}$

$$\mathbb{E}[\nabla \ell_{I_t}(w)] = \sum_{i=1}^n \underbrace{\mathbb{P}(I_t = i)}_{\frac{1}{n}} \cdot \nabla \ell_i(w) = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(w)$$

Stochastic Gradient Descent

Theorem

Let $w_{t+1} = w_t - \eta \nabla_w \ell_{I_t}(w) \Big|_{w=w_t}$ I_t drawn uniform at random from $\{1, \dots, n\}$ so that

$$\mathbb{E}[\nabla \ell_{I_t}(w)] = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(w) =: \nabla \ell(w)$$

If $\|w_0 - w_*\|_2^2 \leq R$ and $\sup_w \max_i \|\nabla \ell_i(w)\|_2^2 \leq G$ then

$$\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] \leq \frac{R}{2T\eta} + \frac{\eta G}{2} \leq \sqrt{\frac{RG}{T}} \approx \left(\frac{1}{\sqrt{T}}\right)^\eta = \sqrt{\frac{R}{GT}}$$

↑
instead of w_T

$$\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$$

(In practice use last iterate)

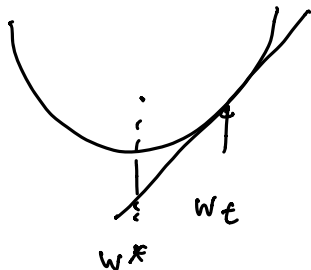
Stochastic Gradient Descent

show: $\mathbb{E}[\ell(w_t) - \ell(w_*)] \leq ?$, $w_{t+1} = w_t - \eta \cdot \nabla \ell_{I_t}(w_t)$

Proof

$$\begin{aligned} \mathbb{E}[\|w_{t+1} - w_*\|_2^2] &= \mathbb{E}[\|w_t - \eta \nabla \ell_{I_t}(w_t) - w_*\|_2^2] \\ &= \underbrace{\mathbb{E}[\|w_t - w_*\|_2^2]} + \underbrace{\eta^2 \mathbb{E}[\|\nabla \ell_{I_t}(w_t)\|_2^2]}_{\leq G} \underbrace{[-2\eta \mathbb{E}[(w_t - w_*)^\top \nabla \ell_{I_t}(w_t)]]}_{\text{blue box}} \\ &\leq \mathbb{E}[\|w_t - w_*\|_2^2] + \eta^2 G + (\ell(w_*) - \ell(w_t)) 2\eta \end{aligned}$$

Convexity $\ell(w)$: $\ell(w_) \geq \ell(w_t) + \nabla \ell_i(w_t)^\top (w_* - w_t)$



$$\rightarrow \ell(w_*) - \ell(w_t) \geq \underbrace{[-\nabla \ell_i(w_t)^\top (w_t - w_*)]}_{\text{blue box}}$$

Stochastic Gradient Descent

Proof

$$\mathbb{E}[\|w_{t+1} - w_*\|_2^2] = \mathbb{E}[\|w_t - \eta \nabla \ell_{I_t}(w_t) - w_*\|_2^2]$$

$$\leq \mathbb{E}[\|w_t - w_*\|_2^2] + \eta^2 G + 2\eta \mathbb{E}[\ell(w_*) - \ell(w_t)]$$

$$\mathbb{E}[\ell(w_t) - \ell(w_*)] \leq \frac{\mathbb{E}[\|w_t - w_*\|_2^2] - \mathbb{E}[\|w_{t+1} - w_*\|_2^2]}{2\eta} + \eta^2 G$$

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[\ell(w_t) - \ell(w_*)] \leq \frac{1}{2\eta T} \left\{ \mathbb{E}[\|w_0 - w_*\|_2^2] - \mathbb{E}[\|w_T - w_*\|_2^2] \right\} + \eta^2 G$$

Stochastic Gradient Descent

Proof

$$\begin{aligned}
 \mathbb{E}[\|w_{t+1} - w_*\|_2^2] &= \mathbb{E}[\|w_t - \eta \nabla \ell_{I_t}(w_t) - w_*\|_2^2] \\
 &= \mathbb{E}[\|w_t - w_*\|_2^2] - 2\eta \mathbb{E}[\nabla \ell_{I_t}(w_t)^T (w_t - w_*)] + \eta^2 \mathbb{E}[\|\nabla \ell_{I_t}(w_t)\|_2^2] \\
 &\leq \mathbb{E}[\|w_t - w_*\|_2^2] - 2\eta \mathbb{E}[\ell(w_t) - \ell(w_*)] + \eta^2 G
 \end{aligned}$$

$$\mathbb{E}[\nabla \ell_{I_t}(w_t)^T (w_t - w_*)] = \mathbb{E}[\mathbb{E}[\nabla \ell_{I_t}(w_t)^T (w_t - w_*) | I_1, w_1, \dots, I_{t-1}, w_{t-1}]]$$

$$\begin{aligned}
 \mathbb{E}[\ell(w_t)] &= \mathbb{E}[\nabla \ell(w_t)^T (w_t - w_*)] \\
 &\geq \mathbb{E}[\ell(w_t) - \ell(w_*)] \\
 \sum_{t=1}^T \mathbb{E}[\ell(w_t) - \ell(w_*)] &\leq \frac{1}{2\eta} (\mathbb{E}[\|w_1 - w_*\|_2^2] - \mathbb{E}[\|w_{T+1} - w_*\|_2^2] + T\eta^2 G) \\
 &\leq \frac{R}{2\eta} + \frac{T\eta G}{2}
 \end{aligned}$$

Stochastic Gradient Descent

Proof

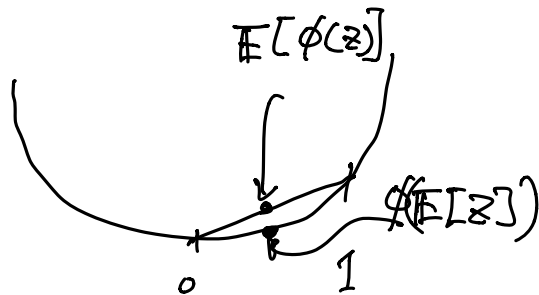
Jensen's inequality:

For any random $Z \in \mathbb{R}^d$ and convex function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$, $\phi(\mathbb{E}[Z]) \leq \mathbb{E}[\phi(Z)]$

$$Z \in \{0, 1\}$$

$$\frac{1}{2}, \frac{1}{2}$$

$$\phi(z) = z^2$$



$$\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] \leq \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\ell(w_t) - \ell(w_*)]$$

\uparrow
 $\{w_1, \dots, w_T\}$
 v.p $\frac{1}{T} \dots \frac{1}{T}$

$$\leq \frac{R}{2\eta T} + \frac{\eta G^2}{2}$$

\nwarrow
 $\mathbb{E}[\|w_0 - w_*\|^2]$

$$\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$$

Stochastic Gradient Descent

$$\ell(w) = \underbrace{\frac{1}{n} \sum_{i=1}^n \ell_i(w)}_{\text{S.G.D.}} + \underbrace{\|w\|_2^2}_{\text{G.D.}}$$

SGD for $t=1 \dots T$.

Proof

Jensen's inequality:

For any random $Z \in \mathbb{R}^d$ and convex function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$, $\phi(\mathbb{E}[Z]) \leq \mathbb{E}[\phi(Z)]$

$$-\eta \left(\underbrace{\nabla \ell(w)}_{\text{S.G.D.}} + \underbrace{\nabla R(w)}_{\text{G.D.}} \right) \mathbb{E}[\cdot] = \nabla \ell(w)$$

$$\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] \leq \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\ell(w_t) - \ell(w_*)]$$

$$\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$$

$$\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] \leq \frac{R}{2T\eta} + \frac{\eta G}{2} \leq \sqrt{\frac{RG}{T}}$$

$$\eta = \sqrt{\frac{R}{GT}}$$

η = step size

T = how many iterations

Mini-batch SGD

Instead of one iterate, average B stochastic gradient together

Advantages:

- Smaller variance
- Parallelization

$$-\eta \cdot \frac{1}{B} \sum_{j=1}^B \nabla \ell_{I_{t,j}}(\omega)$$