

Stochastic Gradient Descent:

(60)

Parameters: $\eta > 0$, $T > 0$

Initialization: $w^{(1)} = 0$

Steps: $t = 1, \dots, T$

Choose $v^{(t)}$ from a distribution

so that $\mathbb{E}[v^{(t)} | w^{(t)}] \in \partial f(w^{(t)})$.

$$w^{(t+1)} = w^{(t)} - \eta v^{(t)}$$

Output: $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$

Theorem: Let $B, \rho > 0$. Let

f be a convex ^{ρ -Lipschitz} function and let

$w^* \in \text{argmin}_{w: \|w\| \leq B} f(w)$. Run SGD

for T steps with $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$. Then

$$\mathbb{E}[f(\bar{w})] - f(w^*) \leq \frac{B\rho}{\sqrt{T}}$$

Note: This is the same result we had before except that now the statement is about $\mathbb{E}[f(\bar{w})]$ instead of $f(\bar{w})$.

Assumo che $\omega \neq 1$

$$\|V(H)\| \leq \rho.$$

Proof: The proof proceeds along the same lines as before.

$$E [f(\bar{w}) - f(w^*)] \stackrel{\text{convexity}}{\leq}$$

$$E \left[\frac{1}{T} \sum_{t=1}^T f(w^{(t)}) - f(w^*) \right] \rightarrow$$

Note that the expectation is over all the choices of $v^{(t)}$, $t=1, \dots, T-1$. We have

$$E \left[\frac{1}{T} \sum_{t=1}^T \langle w^{(t)} - w^*, v^{(t)} \rangle \right] =$$

$$= \frac{1}{T} \sum_{t=1}^T E_{v_1^t} \left[\langle w^{(t)} - w^*, v^{(t)} \rangle \right] =$$

$$= \frac{1}{T} \sum_{t=1}^T E_{v_1^{t-1}} E_{v_1^t} \left[\langle w^{(t)} - w^*, v^{(t)} \rangle \mid v_1^{t-1} \right]$$

$$= \frac{1}{T} \sum_{t=1}^T E_{v_1^{t-1}} \left[\langle w^{(t)} - w^*, \underbrace{E_{v_1^t} [v^{(t)} \mid v_1^{t-1}]}_{\in \partial f(w^{(t)})} \rangle \right]$$

$\underbrace{E_{v_1^t} [v^{(t)} \mid v_1^{t-1}]}_{\in \partial f(w^{(t)})}$

$$\leq \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\Delta \langle w^{(H)} - w^*, \nabla f(w^{(H)}) \rangle \right] \quad (62)$$

$$\leq \frac{B\rho}{\sqrt{T}} \quad \text{identical proof as before.}$$

Note: There is of course still the variance!

Note: X, Y rvs; $\tilde{\mathcal{F}}$ filtration
 Y is $\tilde{\mathcal{F}}$ -measurable

$$\mathbb{E} [XY | \tilde{\mathcal{F}}] = Y \mathbb{E} [X | \tilde{\mathcal{F}}]$$

Lemma: Let $V = \operatorname{argmin}_{x \in H} \|x - w\|^2$

where H is closed and convex.

Then $\forall u \in H$

$$\|V - u\|^2 \leq \|w - u\|^2.$$

See exercise.

The previous proof works more or less as is because of the

Adding a projection step.

(64)

We require w^* to be in a ball of radius B . But we cannot guarantee that \bar{w} fulfills this.

We can add a projection step.

$$w^{(t+\frac{1}{2})} = w^{(t)} - \eta \nabla f(w^{(t)})$$

$$w^{(t+1)} = \underset{w \in \mathcal{B}}{\operatorname{argmin}} \|w - w^{(t+\frac{1}{2})}\|$$

We get the same convergence guarantees in this way, but now \bar{w} is also in the same ball.

Variable step size.

(65)

$$\eta_t = \frac{B}{\rho \sqrt{t}}$$

After averaging.

Output average over last αT steps.

Strongly convex functions.

$$f(u) \geq f(w) + \langle u-w, \nabla f(w) \rangle + \frac{\xi}{2} \|w-u\|^2.$$

$$f(\alpha w + \bar{\alpha} u) \leq \alpha f(w) + \bar{\alpha} f(u) - \frac{\xi}{2} \alpha \bar{\alpha} \|w-u\|^2.$$

Learning with SGD

(66)

$$\mathcal{L}_{\mathcal{D}}(w) = \mathbb{E}_{z \sim \mathcal{D}} [\ell(w, z)]$$

Note:

$$\begin{aligned} \nabla \mathcal{L}_{\mathcal{D}}(w^{(t)}) &= \nabla \mathbb{E}_{z \sim \mathcal{D}} [\ell(w^{(t)}, z)] \\ &= \mathbb{E}_{z \sim \mathcal{D}} [\nabla \ell(w^{(t)}, z)] \end{aligned}$$

Therefore, as long as we pick a new sample for each step of the SGD algorithm we do get an unbiased gradient!

hidden
One-layer NNs and SGD: (67)

$$\hat{f}(x; \theta) = \frac{1}{N} \sum_{i=1}^N \sigma_i(x; \theta_i)$$

N number of neurons in hidden layer.

$x \in \mathbb{R}^d$; d neurons in input layer

$$\sigma_i(x; \theta_i) = a_i \sigma(\langle w_i, x \rangle + b_i)$$

$\sigma: \mathbb{R} \rightarrow \mathbb{R}$; e.g.: $\sigma(x) = \frac{1}{1 + e^{-x}}$

sigmoid \nearrow derivative 1 at $x=0$

$$\sigma(x) = \max\{0, x\} \swarrow \text{ReLU}$$