

SGD on Least Squares

Problem:

Let $\hat{L}_n(w) = \frac{1}{n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i)^2$, $w \in \mathbb{R}^d$

where $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$ and assume $X^T X$ is invertible for $X = [x_1, \dots, x_n]^T$

The ERM is $\hat{w}_n = (X^T X)^{-1} X^T y$

We run SGD:

$$w_{t+1} = w_t - \eta_t g_t,$$

$$g_t(w_t) = 2(\langle w_t, x_{i_t} \rangle - y_{i_t}) x_{i_t}$$

where each i_t is chosen uniformly at

random from $\{1, \dots, n\}$ and the step sizes satisfy

$$\sum_{t=1}^{\infty} \eta_t = \infty, \quad \sum_{t=1}^{\infty} \eta_t^2 < \infty$$

Assume $\|x_i\| \leq R \in \mathbb{R}$ $\forall i$

Show that $E \|w_t - \hat{w}_n\|^2 \rightarrow 0$ as $t \rightarrow \infty$

Solution:

The gradient:

$$\nabla \hat{L}_n(w) = \frac{2}{n} X^T (Xw - y) = Hw - \frac{2}{n} X^T y$$

Hessian:

$$H = \nabla^2 \hat{L}_n(w) = \frac{2}{n} \underbrace{X^T X}_{\geq 0}$$

$\Rightarrow H$ is positive definite ≥ 0

Let $\mu = \lambda_{\min}(H) > 0$, $L = \lambda_{\max}(H)$

Then \hat{L}_n is μ -strongly convex and has L -Lipschitz gradient

$$\langle \nabla \hat{L}_n(w) - \nabla \hat{L}_n(\hat{w}_n), w - \hat{w}_n \rangle$$

$$= \langle H(w - \hat{w}_n), w - \hat{w}_n \rangle$$

$$= (w - \hat{w}_n)^T H (w - \hat{w}_n) \geq 0$$

$$\geq \mu \|w - \hat{w}_n\|^2 \quad (1)$$

$$\| \nabla \hat{L}_n(w) - \underbrace{\nabla \hat{L}_n(\hat{w}_n)}_D \| = \| H(w - \hat{w}_n) \|$$

$$\leq L \|w - \hat{w}_n\| \quad (2)$$

(2)

For any w

$$E[g_+(w)] | w = w]$$

$$= \frac{1}{n} \sum_{i=1}^n 2(\langle w, x_i \rangle - y_i) x_i = \nabla \hat{L}_n(w)$$

$\Rightarrow g_+$ is an unbiased estimator of
the true gradient

Let $r_i = \langle \hat{w}_n, x_i \rangle - y_i$ be the residuals
at the ERM

$$\begin{aligned} \|g_+(w) - g_+(\hat{w}_n)\| &= \|2\langle w - \hat{w}_n, x_i \rangle\| \\ &= 2 \|x_i^\top (w - \hat{w}_n)\| \\ &\leq 2 \|x_i\| \|w - \hat{w}_n\| \|x_i\| \\ &\leq 2R^2 \|w - \hat{w}_n\| \end{aligned}$$

Using $\|a+b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$

$$\underline{E \|g_+(w)\|^2} = E \underbrace{\|g_+(w) - g_+(\hat{w}_n)\|}_a^2 + \underbrace{\|g_+(\hat{w}_n)\|}_b^2$$

$$\leq 2E \|g_+(w) - g_+(\hat{w}_n)\|^2 + 2E \|g_+(\hat{w}_n)\|^2$$
$$4R^4 \|w - \hat{w}_n\|^2$$

$C :=$

$$\leq \underbrace{8R^4 \|w - \hat{w}_n\|^2}_{< \infty} + \boxed{\frac{2}{n} \sum_{i=1}^n 4\gamma_i^2 R^2}$$

bound
on $\|x_i\|^2$
 γ^2 residual def

$E[A, B]$

$E[E[A], E[B]]$

Expand one SGD step

$$\begin{aligned} \|w_{t+1} - \hat{w}_n\|^2 &= \|w_t - \hat{w}_n - \eta_t g_t\|^2 \\ &= \|w_t - \hat{w}_n\|^2 \cancel{- 2\eta_t g_t^\top w_t} \\ &\quad - 2\eta_t \langle w_t - \hat{w}_n, g_t \rangle + \eta_t^2 \|g_t\|^2 \\ \Rightarrow E\|w_{t+1} - \hat{w}_n\|^2 &= E\|w_t - \hat{w}_n\|^2 \\ &\quad - 2\eta_t E\langle w_t - \hat{w}_n, g_t \rangle + \eta_t^2 E\|g_t\|^2 \end{aligned}$$

$$E\langle w_t - \hat{w}_n, \underbrace{E[g_t]}_{\|}\rangle$$

$$\begin{aligned} \text{By } \stackrel{(1)}{\cancel{2}}, E\langle w_t - \hat{w}_n, g_t(w_t) \rangle &= \nabla \hat{L}_n(w_t)^\top - \nabla \hat{L}_n(w_n) \\ &\geq \mu \|w - \hat{w}_n\|^2 \end{aligned}$$

$$E\|w_{t+1} - \hat{w}_n\|^2 = \underbrace{(1 - 2\mu\eta_t + \eta_t^2 \cdot 8R^4)}_{+ C\eta_t^2} E\|w_t - \hat{w}_n\|^2$$

We know $\eta_t \rightarrow 0$ as $t \rightarrow \infty$

$\exists N$ s.t. $\forall t \geq N$

$$\underline{2\mu\eta_t - 8R^4\eta_t^2 \geq \mu\eta_N^2} \quad (\eta_N \leq \frac{\mu}{8R^4})$$

$$\exists E \|w_{t+1} - \hat{w}_n\|^2 \leq (1 - \mu \eta_t^*) \underbrace{E \|w_t - \hat{w}_n\|^2}_{a_t} + C \eta_t^2$$

$$\forall t \geq N \quad a_{t+1} \leq (1 - \mu \eta_t^*)^2 a_t + C \eta_t^2$$

$$a_t \leq (1 - \mu \eta_{t-1}^*)^2 a_{t-1} + C \eta_{t-1}^2$$

$$a_{t+1} \leq (1 - \mu \eta_t^*)^2 (1 - \mu \eta_{t-1}^*)^2 a_{t-1}$$

$$+ \frac{(1 - \mu \eta_t^*)^2}{\Delta} C \eta_{t-1}^2 + C \eta_t^2$$

$\leq \dots$

$$\leq \boxed{a_N} \prod_{k=N}^t (1 - \mu \eta_k^*)^2 +$$

$$C \sum_{k=N}^t \eta_k^2 \prod_{j=k+1}^t (1 - \mu \eta_j^*)$$

$$\text{where } \prod_{j=t+1}^{\infty} (1 - \mu \eta_j^*) = 1$$

$$M = \max_{i=1, \dots, N} \cancel{a_1, a_2, \dots, a_N} \xrightarrow{\rightarrow 0 \text{ as } t \rightarrow \infty} a_i$$

$$a_{t+1} \leq M \left[\prod_{k=N}^t (1 - \mu \eta_k^*) + \cancel{C \sum_{k=N}^t \eta_k^2 \prod_{j=k+1}^t (1 - \mu \eta_j^*)} \right] \xrightarrow{\rightarrow 0}$$

$$\log(1-x) \leq -x \quad \rightarrow 0$$

$$\begin{aligned} (\log \prod (1 - \mu \eta_k^*)) &= \sum \log (1 - \mu \eta_k^*) \\ &\leq \sum -\mu \eta_k^* \leq -\infty \end{aligned}$$

$$\Rightarrow E \|w_{t+1} - \hat{w}_n\|^2 \rightarrow 0 \text{ as } t \rightarrow \infty$$