

DLT Homework 2

20200130 Yujin Kim

#2.1

(Step 1)

$$f: L\text{-smooth, convex} \quad x_{t+1} = x_t - \eta_t \nabla f(x_t) = x_t - \frac{1}{L} \nabla f(x_t).$$

$$\begin{aligned} \|x_{t+1} - x^*\|^2 &= \|x_t - \frac{1}{L} \nabla f(x_t) - x^*\|^2 \\ &= \|x_t - x^*\|^2 - \frac{2}{L} \nabla f(x_t)^T (x_t - x^*) + \frac{1}{L^2} \|\nabla f(x_t)\|^2 \quad \dots \textcircled{1} \end{aligned}$$

(Step 2)

$$\begin{aligned} f(x^*) &\leq f(x - \frac{1}{L} \nabla f(x)) \\ &\leq f(x) + \nabla f(x)^T (-\frac{1}{L} \nabla f(x)) + \frac{L}{2} \left\| -\frac{1}{L} \nabla f(x) \right\|^2 \quad (\text{By smoothness}) \\ &= f(x) - \frac{1}{2L} \|\nabla f(x)\|^2. \quad \dots \textcircled{2} \end{aligned}$$

Descent Lemma

In particular, we have $f(x_{t+1}) \leq f(x_t) - \frac{1}{2L} \|\nabla f(x_t)\|^2 \leq f(x_t)$

and so $f(x_t) \leq f(x_t) \forall t \leq T \quad \dots \textcircled{3}$

(Step 3)

Lemmas for
convex & smooth fns

$g(y) = f(y) - \nabla f(x)^T y$ is L -smooth convex and has minimum at $y=x$.

By $\textcircled{2}$, we have $g(x) \leq g(y) - \frac{1}{2L} \|\nabla g(y)\|^2$

$$\begin{aligned} \Rightarrow g(y) - g(x) &= f(y) - f(x) - \nabla f(x)^T (y-x) \geq \frac{1}{2L} \|\nabla g(y)\|^2 \\ &= \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2 \end{aligned}$$

Hence, $\nabla f(x)^T (y-x) \leq f(y) - f(x) - \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2$

In particular, $\nabla f(x)^T (x^* - x) \leq f(x^*) - f(x) - \frac{1}{2L} \|\nabla f(x)\|^2 \quad \dots \textcircled{4}$

(Step 4)

$$\begin{aligned} \text{Using } \textcircled{4} \text{ in } \textcircled{1}, \quad \|x_{t+1} - x^*\| &\leq \|x_t - x^*\|^2 + \frac{2}{L} (f(x^*) - f(x_t) - \frac{1}{2L} \|\nabla f(x_t)\|^2) \\ &\quad + \frac{1}{L^2} \|\nabla f(x_t)\|^2 \\ &= \|x_t - x^*\|^2 + \frac{2}{L} (f(x^*) - f(x_t)) \\ &\leq \|x_t - x^*\|^2 + \frac{2}{L} (f(x^*) - f(x_T)) \quad (\text{By } \textcircled{3}) \end{aligned}$$

Sum this over $t=0, \dots, T-1$ to obtain $\|x_T - x^*\| \leq \|x_0 - x^*\|^2 + \frac{2T}{L} (f(x^*) - f(x_T))$.

Convergence Rate

$$\begin{aligned} \text{Thus, } f(x_T) - f(x^*) &\leq \frac{L}{2T} (\|x_0 - x^*\|^2 - \|x_T - x^*\|^2) \\ &\leq \frac{L}{2T} \|x_0 - x^*\|^2. \end{aligned}$$

#2.2

(a). By definition of Θ , for $\theta \in \Theta$, $(x_i, f(x_i; \theta))$ is interpolated by a linear model. i.e. $\exists w_0, b_0$ s.t. $f(x_i; \theta) = \langle w_0, x_i \rangle + b_0 \theta_i$.

↪ (Because it is not affected by the nonlinearity of ReLU).

$\theta^* \in \Theta$. For any $\theta \in \Theta$,

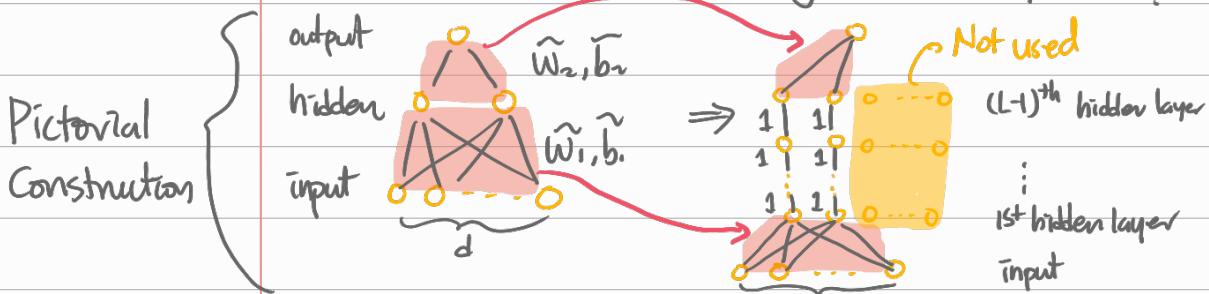
$$\begin{aligned}\hat{R}(\theta^*) &= \frac{1}{n} \sum_{i=1}^n (f(x_i; \theta^*) - y_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\bar{y}_i - y_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\langle w^*, x_i \rangle + b - y_i)^2 \\ &\leq \frac{1}{n} \sum_{i=1}^n (\langle w_0, x_i \rangle + b_0 - y_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (f(x_i; \theta) - y_i)^2 = \hat{R}(\theta)\end{aligned} \quad \dots \textcircled{1}$$

As Θ is finite intersection of open sets $\underbrace{\{\theta \mid \langle w_0, x_i \rangle + b_0 > 0\}}_{\text{open for being inverse image of open set under a continuous map.}}$

Θ is open. $\dots \textcircled{2}$

By $\textcircled{1}, \textcircled{2}$, θ^* is a local minima of \hat{R} .

(b) Let $(\tilde{w}_1, \tilde{b}_1, \tilde{w}_2, \tilde{b}_2)$ be given as the problem.



By the assumption $d \geq 2^{H_L}$, we can let

- | | | |
|--------------------------|---|---|
| Abstract
Construction | $\textcircled{1} \quad (W_1)_{ij} = \begin{cases} (\tilde{w}_1)_{ij} & \text{for } i=1,2 \\ 0 & \text{otherwise} \end{cases}$ | $, (b_1)_i = \begin{cases} (\tilde{b}_1)_i & \text{for } i=1,2 \\ 0 & \text{otherwise} \end{cases}$ |
| | $\textcircled{2} \quad (W_L)_{ij} = \begin{cases} 1 & \text{if } i=j=1 \text{ or } i=j=2 \\ 0 & \text{otherwise} \end{cases}, (b_L)_i = 0 \quad \forall i$ | |
| | $\textcircled{3} \quad (W_L)_{ij} = \begin{cases} (\tilde{w}_2)_{ij} & \text{if } (i,j)=(1,1) \text{ or } (1,2) \\ 0 & \text{otherwise} \end{cases}, b_L = \tilde{b}_2$ | |

#2.3

$$(a) f_w(x) = \frac{1}{\sqrt{m}} \sum_{j=1}^m s_j \sigma(w_j^T x)$$

$$f_{w_0 + w\Sigma}(x) = \frac{1}{\sqrt{m}} \sum_{j=1}^m s_j \sigma((w_j^0 + \sum_{i \neq j} w_i)^T x)$$

where $W_0 = [w_1^0 \dots w_m^0]$, $W = [w_1 \dots w_m]$

$$\sigma(w_j^0 + \sum_{i \neq j} w_i)^T x \approx \sigma(w_j^0 x) + \sigma'(w_j^0 x) (\sum_{i \neq j} w_i^T x) + \frac{\sigma''(w_j^0 x)}{2} (\sum_{i \neq j} w_i^T x)^2$$

$$\text{Thus, } f_{w\Sigma}(x) = \frac{1}{\sqrt{m}} \sum_{j=1}^m s_j \sigma'(w_j^0 x) (\sum_{i \neq j} w_i^T x)$$

$$= \frac{1}{\sqrt{m}} \underbrace{\langle s \sigma'(W_0^T x), \sum W^T x \rangle}_{\text{diag}\{s_1, \dots, s_m\}}$$

$$\sigma' \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix} = \begin{bmatrix} \sigma'(v_1) \\ \vdots \\ \sigma'(v_m) \end{bmatrix}$$

$$\text{diag}\{s_1, \dots, s_m\}$$

$$\text{and } f_{w\Sigma}^Q(x) = \frac{1}{2\sqrt{m}} \sum_{j=1}^m s_j \sigma''(w_j^0 x) (\sum_{i \neq j} w_i^T x)^2$$

$$= \frac{1}{2\sqrt{m}} (\sum W^T x)^T H_0 (\sum W^T x)$$

$$\text{, for } H_0 = \text{diag}\{s_1 \sigma''(w_1^0 x), \dots, s_m \sigma''(w_m^0 x)\}$$

$$= \text{diag}\{s \sigma''(W_0^T x)\}$$

$$\text{diag}\{s_1, \dots, s_m\}, \quad \sigma'' \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix} = \begin{bmatrix} \sigma''(v_1) \\ \vdots \\ \sigma''(v_m) \end{bmatrix}$$

Note that $f_{w\Sigma}^L(x)$ is linear to each $(W\Sigma)_{ij}$

and $f_{w\Sigma}^Q(x)$ is quadratic to each $(W\Sigma)_{ij}$.

$$(b) \sum_{ij} \sim \text{Unif}\{\pm 1\}. \text{ Then, } \mathbb{E}[\sum_{ij}] = 0. \Rightarrow \mathbb{E}[f_{w\Sigma}^L(x)] = 0.$$

$$\sum_{ij}^2 \equiv 1 \Rightarrow \sum = I \Rightarrow f_{w\Sigma}^Q(x) \equiv f_w^Q(x).$$

$$(c) \left(\mathbb{E}_\Sigma [|f_{w\Sigma}^L(x)|] \right)^2 \leq \mathbb{E}_\Sigma [(f_{w\Sigma}^L(x))^2]$$

$$= \mathbb{E}_\Sigma \left[\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^m s_i s_j \sigma'(w_i^0 x) \sigma'(w_j^0 x) \sum_{ii} \sum_{jj} (w_i^T x)(w_j^T x) \right]$$

$$= \mathbb{E} \left[\frac{1}{m} \sum_{j=1}^m s_j^2 \sigma'(w_j^0 x)^2 \sum_{ii}^2 (w_i^T x)^2 \right] = O(\frac{1}{m} \times m \times m^{-1/2})$$

$$\hookrightarrow O(1)$$

$$\hookrightarrow \mathbb{E} \sum_{ii} \sum_{jj} = 0 \text{ for } i \neq j$$

$$\text{Thus, } \mathbb{E}_\Sigma [|f_{w\Sigma}^L(x)|] = O(m^{-1/2})$$

(d) In comparison,

$$\mathbb{E}[f_{W\Sigma}^Q(x)] = f_w^Q(x) = \frac{1}{2\sqrt{m}} \sum_{j=1}^m s_j \sigma''(w_j^T x)(w_j^T x)^2 \\ = O\left(\frac{1}{\sqrt{m}} \times m \times (m^{-1/4})^2\right) = O(1).$$

Thus, f^Q dominates over f^L , and hence NTK regime may not be useful this case.

#2.4

(Step 1)

Expectation
Recursion

$$\mathbb{E}[\theta_t] = \mathbb{E}[\theta_{t-1}] - \gamma_t \mathbb{E}[(\langle x_{i(t)}, \theta_{t-1} \rangle - y_{i(t)}) x_{i(t)}]$$

$$\begin{aligned}\mathbb{E}[(\langle x_{i(t)}, \theta_{t-1} \rangle - y_{i(t)}) x_{i(t)} | \theta_{t-1}] &= \frac{1}{n} \sum_{i=1}^n (\langle x_i, \theta_{t-1} \rangle - y_i) x_i \\ &= \frac{1}{n} X^T (X \theta_{t-1} - y)\end{aligned}$$

$$\begin{aligned}\text{Thus, } \mathbb{E}[(\langle x_{i(t)}, \theta_{t-1} \rangle - y_{i(t)}) x_{i(t)}] &= \mathbb{E}_{\theta_{t-1}} \left[\frac{1}{n} X^T (X \theta_{t-1} - y) \right] \\ &= \frac{1}{n} X^T (X \mathbb{E}[\theta_{t-1}] - y).\end{aligned}$$

$$\begin{aligned}\text{Hence, } \mathbb{E}[\theta_t] &= \mathbb{E}[\theta_{t-1}] - \frac{1}{\|x\|^2} X^T (X \mathbb{E}[\theta_{t-1}] - y) \\ &= (I - \frac{1}{\|x\|^2} X^T X) \mathbb{E}[\theta_{t-1}] + \frac{1}{\|x\|^2} X^T y. \quad \dots \textcircled{1}\end{aligned}$$

(Step 2)

Convergence

By letting $\rho_t = \mathbb{E}[\theta_t]$, $A = (I - \frac{1}{\|x\|^2} X^T X)$, $b = \frac{1}{\|x\|^2} X^T y$,

we have $\rho_0 = 0$, $\rho_{t+1} = A\rho_t + b$, for $0 \leq A < I$. $\dots \textcircled{2}$

$$\Rightarrow \rho_t = A^t \rho_0 + \sum_{i=0}^{t-1} A^i b = \sum_{i=0}^{t-1} A^i b. = \left(\sum_{i=0}^{t-1} A^i \right) \frac{1}{\|x\|^2} X^T y \quad \dots \textcircled{3}$$

\rightarrow convergence in elementwise sense, by $\textcircled{2}$.

$$(I - A) \left(\sum_{i=0}^{t-1} A^i \right) = I - A^t \rightarrow I \text{ as } t \rightarrow \infty. \text{ Hence, } \left(\sum_{i=0}^{t-1} A^i \right) \frac{X^T X}{\|x\|^2} \rightarrow I \text{ as } t \rightarrow \infty$$

Denote $B_t = \sum_{i=0}^{t-1} A^i$. Then, $B_t X^T X \rightarrow \|x\|^2 I$ as $t \rightarrow \infty$.

$$\Rightarrow B_t X^T X X^T \rightarrow \|x\|^2 X^T \text{ as } t \rightarrow \infty$$

$$\Rightarrow B_t X^T \rightarrow \|x\|^2 X^T (X X^T)^{-1} \text{ as } t \rightarrow \infty$$

$$\Rightarrow B_t b = B_t \frac{X^T}{\|x\|^2} y = \frac{1}{\|x\|^2} (B_t X^T) y$$

$$\rightarrow X^T (X X^T)^{-1} y. \text{ asy } t \rightarrow \infty \quad \dots \textcircled{4}$$

Thus, $\textcircled{4}$ with $\textcircled{3}$ gives

$$\mathbb{E}[\theta_t] = \rho_t = B_t b \rightarrow X^T (X X^T)^{-1} y \text{ as } t \rightarrow \infty.$$