

1. (a) (1) Sigmoid activation function $y = f(x) = \frac{1}{1+e^{-x}} = (1+e^{-x})^{-1}$

$$\frac{dy}{dx} = -(1+e^{-x})^{-2} \cdot -e^{-x} = e^{-x} (1+e^{-x})^{-2} = \frac{e^{-x}}{(1+e^{-x})^2}$$

$$\lim_{x \rightarrow +\infty} \frac{dy}{dx} = \lim_{x \rightarrow +\infty} \frac{e^{-x}}{(1+e^{-x})^2} = 0 \quad \lim_{x \rightarrow -\infty} \frac{dy}{dx} = \lim_{x \rightarrow -\infty} \frac{e^{-x}}{(1+e^{-x})^2} = 0$$

(2) Tanh activation function $y = f(x) = \frac{e^{2x}-1}{e^{2x}+1}$

$$\frac{dy}{dx} = \frac{2e^{2x}(e^{2x}+1) - (e^{2x}-1)2e^{2x}}{(e^{2x}+1)^2} = \frac{4e^{2x}}{(e^{2x}+1)^2}$$

$$\lim_{x \rightarrow +\infty} \frac{dy}{dx} = \lim_{x \rightarrow +\infty} \frac{4e^{2x}}{(e^{2x}+1)^2} = 0 \quad \lim_{x \rightarrow -\infty} \frac{dy}{dx} = \lim_{x \rightarrow -\infty} \frac{4e^{2x}}{(e^{2x}+1)^2} = 0$$

(3) Leaky ReLU activation function

$$y = f(x) = \max\{\alpha x, x\} = \begin{cases} x, & x \geq 0 \\ \alpha x, & x \leq 0 \end{cases} \text{ for some } \alpha \in (0,1)$$

$$\frac{dy}{dx} = \begin{cases} 1, & x \geq 0 \\ \alpha, & x \leq 0 \end{cases}$$

$$\lim_{x \rightarrow +\infty} \frac{dy}{dx} = 1, \quad \lim_{x \rightarrow -\infty} \frac{dy}{dx} = \alpha \text{ for some } \alpha \in (0,1)$$

(b)

For hidden layer:

$$h = \begin{bmatrix} -0.8 & 0.5 & -1 \\ 1.2 & -0.7 & 0.2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + \begin{bmatrix} -0.4 \\ 0.9 \end{bmatrix} = \begin{bmatrix} -3.2 \\ 1.3 \end{bmatrix}$$

$$\text{ReLU: } \text{ReLU}(h) = \begin{bmatrix} 0 \\ 1.3 \end{bmatrix}$$

For output layer:

$$o = \text{ReLU}\left\{ \begin{bmatrix} 0.6 & 1.1 \end{bmatrix} \begin{bmatrix} 0 \\ 1.3 \end{bmatrix} + \begin{bmatrix} -0.1 \end{bmatrix} \right\} = 1.33$$

Back-propagation

$$\text{loss function: } L = \frac{1}{2}(f(x;\theta) - y)^2, \quad L' = f(x;\theta) - y$$

$$\text{ReLU function: } z(x) = \max(0, x), \quad z'(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x \leq 0 \end{cases}$$

For output layer:

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial z} \cdot \frac{\partial z}{\partial w} = 1 \cdot 1 \cdot \begin{bmatrix} 0 \\ 1.3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1.3 \end{bmatrix}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial z} \cdot \frac{\partial z}{\partial b} = 1 \cdot 1 \cdot 1 = 1$$

For hidden layer:

for the output of hidden layer h :

$$\frac{\partial L}{\partial h} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial z} \cdot \frac{\partial z}{\partial x} = 1 \cdot W^{(2)} = \begin{bmatrix} 0.6 \\ 1.1 \end{bmatrix}$$

$$\begin{aligned} \frac{\partial L}{\partial W_h} &= \frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial z} \cdot \frac{\partial z}{\partial W_h} = \begin{bmatrix} 0.6 \\ 1.1 \end{bmatrix} \odot \begin{bmatrix} 0 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 1.1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 1.1 & 2.2 & 3.3 \end{bmatrix} \end{aligned}$$

$$\frac{\partial L}{\partial b_h} = \frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial z} \cdot \frac{\partial z}{\partial b_h} = \begin{bmatrix} 0.6 \\ 1.1 \end{bmatrix} \odot \begin{bmatrix} 0 \\ 1 \end{bmatrix} \cdot 1 = \begin{bmatrix} 0 \\ 1.1 \end{bmatrix}$$

2. $R(f) := E \{ \rho_{\tau}(Y - f(X)) \}$

$$\begin{aligned} R(f) &= \int_{-\infty}^{+\infty} \rho_{\tau}(Y - f(x)) d\bar{F}_{Y|X}(y) \\ &= \int_{-\infty}^{f(x)} (\tau - 1)(Y - f(x)) d\bar{F}_{Y|X}(y) \\ &\quad + \int_{f(x)}^{+\infty} \tau(Y - f(x)) d\bar{F}_{Y|X}(y) \end{aligned}$$

$$\begin{aligned} \frac{dR(f)}{df} &= \int_{-\infty}^{f(x)} (\tau - 1)(-1) d\bar{F}_{Y|X}(y) \\ &\quad + \int_{f(x)}^{+\infty} \tau(-1) d\bar{F}_{Y|X}(y) \end{aligned}$$

$$= (\tau - 1) \cdot -F_{Y|X}(f(x)) - \tau(1 - F_{Y|X}(f(x)))$$

$$= -\tau F_{Y|X}(f(x)) + \tau F_{Y|X}(f(x)) + F_{Y|X}(f(x)) - \tau = 0$$

$$\Rightarrow F_{Y|X}(f(x)) = \tau, \text{ which means } P(Y \leq f(x)) = \tau.$$

Thus, $f^*(x)$ is the conditional τ -th quantile.

3. At step k , according to lemma 3.1, we have

$$f(\theta^{k+1}) - f(\theta^k) \leq -\frac{1}{2L} \|\nabla f(\theta^k)\|^2$$

Consider $k=1, \dots, T$, then

$$f(\theta^1) - f(\theta^0) \leq -\frac{1}{2L} \|\nabla f(\theta^0)\|^2$$

$$f(\theta^2) - f(\theta^1) \leq -\frac{1}{2L} \|\nabla f(\theta^1)\|^2$$

$$\vdots$$
$$f(\theta^{T+1}) - f(\theta^T) \leq -\frac{1}{2L} \|\nabla f(\theta^T)\|^2$$

Sum them up

$$f(\theta^{T+1}) - f(\theta^0) \leq -\frac{1}{2L} [\|\nabla f(\theta^0)\|^2 + \dots + \|\nabla f(\theta^T)\|^2]$$

$$\leq -\frac{1}{2L} (T \cdot \min_{0 \leq k \leq T} \|\nabla f(\theta^k)\|^2)$$

Consider $\bar{f} \in \mathbb{R}$ s.t. $f(\theta) \geq \bar{f} > -\infty$, then we have

$$\bar{f} - f(\theta^0) \leq -\frac{1}{2L} (T \cdot \min_{0 \leq k \leq T} \|\nabla f(\theta^k)\|^2)$$

$$\frac{2L \{\bar{f} - f(\theta^0)\}}{T} \geq \min_{0 \leq k \leq T} \|\nabla f(\theta^k)\|^2$$