

①

$$x = \{[0 \ 0]^T, [1 \ 0]^T, [0 \ 1]^T, [1 \ 1]^T\}$$

$$y = f^*(x) = \{0, 1, 1, 0\}$$

$x \in X$

$$J(\theta) = \frac{1}{4} \sum_{x \in X} [f^*(x) - f(x; \theta)]^2$$

for linear function, $f(x; \theta) = f(x; w, b) = x^T w + b$. $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$

$$\therefore J(\theta) = J(w, b) = \frac{1}{4} \sum_{x \in X} \{f^*(x) - f(x; w, b)\}^2 = \frac{1}{4} \sum_{x \in X} \{f^*(x) - x^T w - b\}^2$$

To minimize $J(\theta)$, we find θ s.t. $\nabla_{\theta}(J(\theta)) = 0 = \nabla_{w,b}(J(w,b))$

$$\nabla_{\theta} J = \frac{1}{4} \nabla_w \sum_{x \in X} \{f^*(x) - x^T w - b\}^2$$

$$= \cancel{\frac{1}{4}} \cancel{\nabla_w} \left[\right]$$

$$= \frac{1}{4} \nabla_w \left[\left(0 - [0 \ 0] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - b \right)^2 + \left(1 - [0 \ 1] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - b \right)^2 + \left(1 - [1 \ 0] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - b \right)^2 + \left(0 - [1 \ 1] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} - b \right)^2 \right]$$

$$= \frac{1}{4} \left[(0-b)^2 + (1-w_2-b)^2 + (1-w_1-b)^2 + (0-w_1-w_2-b)^2 \right]$$

$$= \frac{1}{4} \cdot \left[2 \cdot (0-b) \cdot 0 + 2 \cdot (1-w_2-b) \cdot \nabla_w (1-w_2-b) + 2 \cdot (1-w_1-b) \cdot \nabla_w (1-w_1-b) + 2 \cdot (-w_1-w_2-b) \cdot \nabla_w (-w_1-w_2-b) \right]$$

$$= \frac{2}{4} \left[(1-w_2-b) \begin{bmatrix} 0 \\ -1 \end{bmatrix} + (1-w_1-b) \begin{bmatrix} 1 \\ 0 \end{bmatrix} + -(w_1+w_2+b) \begin{bmatrix} -1 \\ -1 \end{bmatrix} \right]$$

$$= \frac{1}{2} \cdot \left[(0) + (1 - w_1 - b) + (w_1 + w_2 + b) \right] \\ (0) + (1 - w_2 - b) + (w_1 + w_2 + b)$$

$$= \frac{1}{2} \cdot \begin{bmatrix} w_2 \\ w_1 \end{bmatrix}$$

Setting to zero, we have, $w_1 = w_2 = 0.$

For bias,

$$\nabla_b J = \frac{1}{4} \nabla_b \sum_{x \in X} \left\{ f^*(x) - x^T w - b \right\}^2$$

$$= \frac{1}{4} \cdot \nabla_b \left[b^2 + (1 - w_2 - b)^2 + (1 - w_1 - b)^2 + (w_1 + w_2 + b)^2 \right]$$

$$= \frac{1}{4} \left[2b + 2(1 - w_2 - b)(-1) + 2(1 - w_1 - b)(-1) + 2(w_1 + w_2 + b) \right]$$

$$= \frac{1}{4} \{ 2b - 2 - 2(1 - w_2 - b) - 2(1 - w_1 - b) + 2(w_1 + w_2 + b) \}$$

$$= \frac{1}{4} [2b - 2 + 2w_2 + 2b \cancel{- 2} + 2 - 2w_1 - 2b + 2w_1 + 2w_2 + 2b]$$

$$= \frac{1}{4} [8b - 4] \quad \dots \text{setting } w_1 = w_2 = 0.$$

$$= 2b - 1$$

Setting $\nabla_b J = 0$, $2b - 1 = 0$

$$\Rightarrow b = \frac{1}{2} \Rightarrow b = 0.5$$

Thus, for XOR problem, $w = [0 \ 0]^T \neq b = 0.5$
minimizes the MSE loss function.

$$\hat{y} = x^T w + b = [x_1 \ x_2] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + 0.5 = [x_1 \ x_2] \begin{bmatrix} 0 \\ 0 \end{bmatrix} + 0.5 = 0.5 + x$$

(2)

$$\hat{y} = \sigma(x^T w + b)$$

$$\text{where, } \sigma = \frac{1}{1 - e^{-x}}.$$

$$\therefore \hat{y} = \frac{1}{1 - e^{-(x^T w + b)}}.$$

Given,

training examples $\rightarrow \{+ve\}$

testing examples $\rightarrow \{+ve, -ve\}$

We know that, $\sigma'(x) = \frac{\partial \sigma(x)}{\partial x} = \sigma(x) \cdot (1 - \sigma(x)) \quad \forall x$

$$\therefore \hat{y} = f(x; w, b) = \sigma(x^T w + b)$$

Consider MSE $\rightarrow \frac{1}{n} \sum_{x \in X} \{f^*(x) - f(x; w, b)\}^2$

where, $f^*(x) \rightarrow \text{desired predictions.}$

Say, we have one example , i.e., $\text{size}(x) = 1$.

$$\therefore \nabla_w J(w) = 0 \quad (\text{to minimize MSE loss})$$

$$= 2 \cdot \left\{ f^*(x) - f(x; w, b) \right\} \cdot \left[\nabla(f^*(x)) - \nabla(f(x; w, b)) \right]$$

$$= -2 \left\{ f^*(x) - f(x; w, b) \right\} \cdot \sigma'(x^T w + b).$$

Now,

$$\sigma'(x^T w + b) = \sigma(x^T w + b) \cdot (1 - \sigma(x^T w + b))$$

$$= \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}} \right)$$

$$= \frac{1 + e^{-z} - 1}{(1 + e^{-z})^2} = \frac{e^{-z}}{(1 + e^{-z})^2}$$

$$= \frac{1}{(1 + e^z)(1 + e^{-z})}$$

$$* = \frac{e^z}{(1 + e^z)^2}$$

$$(z = \mathbf{x}^T \mathbf{w} + b)$$

- (i) Bias term, $b \rightarrow$ converges
because it is scalar & independent of training & testing examples.
- (ii) Weights, $\mathbf{w} \rightarrow$ converge, to minimize log loss. Regardless of values.
- (iii) Training loss \rightarrow converges
the model tries to minimize log loss on training examples.
- (iv) Testing loss \rightarrow depends on testing values.

Q.3

$$\begin{aligned} \text{(i) for } l = K : - \quad \nabla_{w^{(l)}} \hat{y}_K^{(i)} &= \nabla_{w^{(l)}} \frac{e^{z_l}}{\sum_j e^{z_j}} \\ &= \frac{\left(\sum_j e^{z_j} \right) \nabla_{w^{(l)}} e^{z_l} - e^{z_l} \nabla_{w^{(l)}} \left(\sum_j e^{z_j} \right)}{\left(\sum_j e^{z_j} \right)^2} \\ &= \frac{\left(\sum_j e^{z_j} \right) e^{z_l} \nabla_{w^{(l)}} (x^T w^{(l)} + b_K) -}{\left(\sum_j e^{z_j} \right)^2} \\ &\quad \underline{\left(\sum_j e^{z_j} \right)^2} \\ &= \frac{\left(\sum_j e^{z_j} \right) e^{z_l} \nabla_{w^{(l)}} (x^T w^{(l)} + b_K)}{\left(\sum_j e^{z_j} \right)^2} \\ &\quad - \underline{e^{z_l} \nabla_{w^{(l)}} (e^{z_l})} \\ &= \frac{\left(\sum_j e^{z_j} \right) e^{z_l} x^{(i)} - \cancel{e^{z_l} x^{(i)}}}{\left(\sum_j e^{z_j} \right)^2} \end{aligned}$$

only vector part here is $x^{(i)}$:

$$\nabla_{w^{(l)}} \hat{y}_K^{(i)} = x^{(i)} \left(\frac{e^{z_l}}{\sum_j e^{z_j}} \right) \left(\frac{1 - e^{z_l}}{\sum_j e^{z_j}} \right)$$

$$\boxed{\nabla_w^{(l)} \hat{y}_k^{(i)} = X^{(i)} \hat{y}_e^{(i)} (1 - \hat{y}_e^{(i)})}$$

(ii) Now, $l \neq k$,

$$\begin{aligned}
 \nabla_w^l \hat{y}_k^{(i)} &= \nabla_w \left(\frac{e^{z_k}}{\sum_j e^{z_j}} \right) = \left(\sum_j e^{z_j} \right) \nabla_w^l e^{z_k} - e^{z_k} \nabla_w^l \left(\sum_j e^{z_j} \right) \\
 &= \left(\sum_j e^{z_j} \right) e^{z_k} (0) - e^{z_k} (\nabla_w^l e^{z_k}) \\
 &= \frac{-e^{z_k}}{\left(\sum_j e^{z_j} \right)^2} e^{z_k} \nabla_w^l (X^T w^l + b_k) \\
 &= -\frac{e^{z_k}}{\left(\sum_j e^{z_j} \right)} \cdot \frac{e^{z_l}}{\left(\sum_j e^{z_j} \right)} \cdot X^{(i)} \\
 &= -\hat{y}_k^{(i)} \cdot \hat{y}_k^{(i)} \cdot X^{(i)} \\
 &= \boxed{\nabla_w^{(l)} \hat{y}_k^{(i)} = -X^{(i)} \hat{y}_k^{(i)} \hat{y}_k^{(i)} \hat{y}_e^{(i)}}
 \end{aligned}$$

$$\begin{aligned}
 \text{(iii)} \quad \nabla_w^{(u)} f_{CE}(w, b) &= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^C y_k^{(i)} \nabla_w^{(i)} \log \hat{y}_k^{(i)} \\
 &= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^C y_k^{(i)} \left(\frac{\nabla_w^{(i)} \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} \right) \\
 &= -\frac{1}{n} \left(\sum_{i=1}^n \left(\sum_{k \neq l} y_k^{(i)} \left(-\frac{\hat{y}_k^{(i)}}{\hat{y}_l^{(i)}} \right) \hat{y}_l^{(i)} \right) \right. \\
 &\quad \left. y_l^{(i)} \left(\frac{x^{(i)} \hat{y}_l^{(i)} (1 - \hat{y}_l^{(i)})}{\hat{y}_l^{(i)}} \right) \right) \\
 &= -\frac{1}{n} \left(\sum_{i=1}^n \left[x^{(i)} \hat{y}_l^{(i)} (1 - \hat{y}_l^{(i)}) - x^{(i)} \hat{y}_l^{(i)} \sum_{k \neq l} y_k^{(i)} \right] \right) \\
 &= -\frac{1}{n} \sum_{i=1}^n \left[x^{(i)} \hat{y}_l^{(i)} - x^{(i)} \hat{y}_l^{(i)} \hat{y}_k^{(i)} - x^{(i)} \hat{y}_k^{(i)} \right] \left(\sum_k y_k^{(i)} - y_l^{(i)} \right) \\
 &= -\frac{1}{n} \sum_{i=1}^n \left[x^{(i)} \hat{y}_l^{(i)} - x^{(i)} \hat{y}_l^{(i)} \right]
 \end{aligned}$$

$$\boxed{\nabla_w^{(u)} f_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n x^{(i)} (y^{(i)} - \hat{y}^{(i)})}$$

$$\begin{aligned}
 \nabla_b f_{CE}(w, b) &= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^C y_k^{(i)} \nabla_b (\log \hat{y}_k^{(i)}) \\
 &= -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^C y_k^{(i)} \left(\frac{\nabla_b \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} \right) \quad \text{--- (1)}
 \end{aligned}$$

$$\text{Now, } \nabla_b \hat{y}_k^{(i)} = \nabla_b \frac{e^{z_k}}{\sum_j e^{z_j}} = \frac{\left(\sum_j e^{z_j} \right) \nabla_b e^{z_k} - e^{z_k} \nabla_b \left(\sum_j e^{z_j} \right)}{\left(\sum_j e^{z_j} \right)^2}$$

Now, the term $\nabla_b e^{z_k} = e^{z_k} \begin{bmatrix} \frac{\partial}{\partial b_1} \\ \frac{\partial}{\partial b_2} \\ \vdots \\ \frac{\partial}{\partial b_K} \\ \frac{\partial}{\partial b_C} \end{bmatrix} (x^T w^{(K)} + b_k)$

$$\nabla_b e^{z_k} = e^{z_k} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \rightarrow K^{\text{th}} \text{ position}$$

$$\therefore \nabla_b \hat{y}_k^{(i)} = \underbrace{\left(\sum_j e^{z_j} \right) \nabla_b e^{z_k} - e^{z_k} \left(\sum_j \nabla_b e^{z_j} \right)}_{\left(\sum_j e^{z_j} \right)^2}$$

$$\therefore \nabla_b \hat{y}_k^{(i)} = \left(\frac{\sum_j e^{z_j}}{\sum_j e^{z_j}} \right) \cdot e^{z_k} \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} - \frac{e^{z_k}}{\left(\sum_j e^{z_j} \right)^2} \sum_j \frac{e^{z_j}}{0} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\therefore \nabla_b \hat{y}_k^{(i)} = \frac{e^{z_k}}{\left(\sum_j e^{z_j} \right)} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} - \frac{e^{z_k}}{\left(\sum_j e^{z_j} \right)^2} \begin{bmatrix} e^{z_1} \\ e^{z_2} \\ \vdots \\ e^{z_C} \end{bmatrix}$$

$$= \hat{y}_k^{(i)} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} - \hat{y}_k^{(i)} \begin{bmatrix} e^{z_1} / (e^{z_1})^2 \\ e^{z_2} / (e^{z_2})^2 \\ \vdots \\ e^{z_C} / (e^{z_C})^2 \end{bmatrix}$$

$$\therefore \frac{\nabla_b \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \xrightarrow{k^{\text{th}} \text{ position}} - \begin{bmatrix} \hat{y}_1^{(i)} \\ \hat{y}_2^{(i)} \\ \vdots \\ \hat{y}_c^{(i)} \end{bmatrix}$$

$$\frac{y_k^{(i)} \nabla_b \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} - y_k^{(i)} \begin{bmatrix} \hat{y}_1^{(i)} \\ \hat{y}_2^{(i)} \\ \vdots \\ \hat{y}_c^{(i)} \end{bmatrix}$$

Kth position

$$\therefore \sum_{k=1}^c y_k^{(i)} \left(\frac{\nabla_b \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} \right) = \begin{bmatrix} y_1^{(i)} \\ y_2^{(i)} \\ \vdots \\ y_c^{(i)} \end{bmatrix} - \left(\sum_{k=1}^c y_k^{(i)} \right) \begin{bmatrix} \hat{y}_1^{(i)} \\ \hat{y}_2^{(i)} \\ \vdots \\ \hat{y}_c^{(i)} \end{bmatrix}$$

$$\therefore \sum_{k=1}^c y_k^{(i)} \left(\frac{\nabla_b \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} \right) = \vec{y}^{(i)} - (1) \vec{\hat{y}}^{(i)}$$

$$\therefore -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c y_k^{(i)} \left(\frac{\nabla_b \hat{y}_k^{(i)}}{\hat{y}_k^{(i)}} \right) = -\frac{1}{n} \sum_{i=1}^n (\vec{y}^{(i)} - \vec{\hat{y}}^{(i)})$$

from ①

$$\therefore \nabla_b f_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n (\vec{y}^{(i)} - \vec{\hat{y}}^{(i)})$$