# Softmax loss

$N$ = # training example

$m$ = dimention of each training example

$c$ = # class

$$X^{(i)} = \begin{bmatrix} 1 \\ X_1^{(i)} \\ X_2^{(i)} \\ \vdots \end{bmatrix} \in R^{(m+1) \times 1} \quad \left\{ \begin{array}{c} i^{th} \text{ training example with} \\ \text{bias dimension} \end{array} \right\}$$

$$y^{(i)} \in \{1, 2, \cdots, c\} \quad \left\{ \begin{array}{c} \text{class lable associated with} \\ i^{th} \text{ training example} \end{array} \right\}$$

$$X \in R^{N \times (m+1)} \quad \{\text{Design matrix}\}$$

$$X = \begin{bmatrix} X^{(1)^T} \\ X^{(2)^T} \\ \vdots \\ X^{(N)^T} \end{bmatrix}$$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} \quad \{\text{Label vector}\}$$

$$S^{(i)} = W X^{(i)} \in R^{c \times 1} \quad \{\text{Score vector for } i^{th} \text{ train example}\}$$

$$\{ S_j^{(i)} = \{\text{Score given to } i^{th} \text{ class lable}\}$$

where,

$$W \in R^{c \times (m+1)} \quad \{\text{Weight matrix}\} \quad W_j = \{j^{th} \text{ row of } W\}$$

$$S = X W^T \in R^{N \times c} \quad \{\text{Score matrix}\}$$

$\Rightarrow$ In softmax we assume score to be **Unnormalized log probability.**

$\rightarrow$ Let Y be the random variable of lable associated with input exaple random variable X.

$$P(Y=k \mid X=x) = \frac{e^{S_k}}{\sum_j e^{S_j}} \quad \text{where,} \quad S = Wx$$

$\Rightarrow$ For every training example, we want to maximise the assigned probability to the correct class.

i.e. $\rightarrow$ Maximize log likelihood

i.e. $\rightarrow$ Minimize negative log likelihood

$$L_i(W) = -\log\left(P(Y=y^{(i)} \mid X=x^{(i)})\right)$$

$$L(W) = \frac{1}{n} \sum_{i=1}^{n} L_i(W)$$

# ★ Calculating Gradient

$$\nabla_w L(w) = \frac{1}{n} \sum_{i=1}^{n} \boxed{\nabla_w L_i(w)}$$

$$\nabla_w L_i(w) = -\nabla_w \log \left( \frac{e^{S_{y^{(i)}}^{(i)}}}{\sum_j e^{S_j^{(i)}}} \right)$$

$$= -\nabla_w \log \left( \frac{e^{w_{y^{(i)}} x^{(i)}}}{\sum_j e^{w_j x^{(i)}}} \right)$$

$$= -\frac{\sum_j e^{w_j x^{(i)}}}{e^{w_{y^{(i)}} x^{(i)}}} \nabla_w \left( \frac{e^{w_{y^{(i)}} x^{(i)}}}{\sum_j e^{w_j x^{(i)}}} \right)$$

$$\frac{\left( \sum_j e^{w_j x^{(i)}} \right) \nabla_w e^{w_{y^{(i)}} x^{(i)}} - \left( e^{w_{y^{(i)}} x^{(i)}} \right) \sum_j \nabla_w e^{w_j x^{(i)}}}{\left( \sum_j e^{w_j x^{(i)}} \right)^2}$$

③     ①                                              ②

② $\nabla_\omega e^{\omega_j x^{(i)}}$

$= e^{\omega_j x^{(i)}} \cdot \boxed{\nabla_\omega \left( \omega_j x^{(i)} \right)}$

$\left( \nabla_\omega \left( \omega_j x^{(i)} \right) \right)_{a,b} = \frac{\delta}{\delta \omega_{ab}} \left( \omega_j x^{(i)} \right)$

$$= \begin{cases} 0 & \text{if } a \neq j \\ x_b^{(i)} & \text{if } a = j \end{cases}$$

$\left( \nabla_\omega \left( \omega_j x^{(i)} \right) \right)_a = \begin{cases} x^{(i)} & \text{if } a = j \\ 0 & \text{if } a \neq j \end{cases}$

$$\begin{bmatrix} \longleftarrow & 0 & \longrightarrow \\ \longleftarrow & x^{(i)T} & \longrightarrow \\ \longleftarrow & 0 & \longrightarrow \end{bmatrix} \quad \leftarrow j^{th} \text{ row}$$

$\sum_j \nabla_\omega e^{\omega_j x^{(i)}} = \begin{bmatrix} \leftarrow e^{\omega_1 x^{(i)}} x^{(i)T} \rightarrow \\ 0 \end{bmatrix} + \begin{bmatrix} \leftarrow e^{\omega_2 x^{(i)}} x^{(i)T} \rightarrow \\ 0 \\ 0 \end{bmatrix}$

$+ \cdots + \begin{bmatrix} 0 \\ 0 \\ \leftarrow e^{\omega_c x^{(i)}} x^{(i)T} \rightarrow \end{bmatrix}$

$$\sum_j \nabla_\omega e^{\omega_j x^{(i)}} = \begin{bmatrix} \leftarrow e^{\omega_1 x^{(i)}} X^{(i)T} \longrightarrow \\ \leftarrow e^{\omega_2 x^{(i)}} X^{(i)T} \longrightarrow \\ \vdots \\ \leftarrow e^{\omega_c x^{(i)}} X^{(i)T} \longrightarrow \end{bmatrix} //$$

① $\quad \nabla_\omega e^{\omega_{y^{(i)}} x^{(i)}}$

$$= e^{\omega_{y^{(i)}} x^{(i)}} \cdot \nabla \omega_{y^{(i)}} x^{(i)}$$

$$= \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \leftarrow e^{\omega_{y^{(i)}} x^{(i)}} X^{(i)T} \longrightarrow \\ 0 \\ \vdots \\ 0 \end{bmatrix} \longrightarrow y^{(i)th} \text{ grow}$$

②ⓝ $\begin{matrix} \\ y^{(i)th} \text{ row} \end{matrix}$ $\begin{bmatrix} 0 \\ 0 \\ \vdots \\ \left( e^{\omega_{y^{(i)}} x^{(i)}} \sum_i e^{\omega_j x^{(i)}} \right) X^{(i)T} \\ 0 \\ \vdots \\ 0 \end{bmatrix} - \begin{bmatrix} e^{(\omega_1 + \omega_{y^{(i)}}) x^{(i)}} X^{(i)T} \\ \vdots \\ e^{2\omega_{y^{(i)}} x^{(i)}} X^{(i)T} \\ \vdots \end{bmatrix}$

$$= \begin{bmatrix} -e^{(\omega_1 + \omega_{y^{(i)}}) x^{(i)}} X^{(i)T} \\ \vdots \\ e^{\omega_{y^{(i)}} x^{(i)}} \left( \sum_{j \neq y^{(i)}} e^{\omega_j x^{(i)}} \right) X^{(i)T} \\ \vdots \end{bmatrix}$$

$$\nabla_\omega L_i(\omega) = \frac{1}{e^{\omega_{y^{(i)}} x^{(i)}} \left(\sum_j e^{\omega_j x^{(i)}}\right)} \begin{bmatrix} e^{(\omega_1 + \omega_{y^{(i)}}) x^{(i)}} \, x^{(i)T} \\ \vdots \\ -e^{\omega_{y^{(i)}} x^{(i)}} \left(\sum_{j \neq y^{(i)}} e^{\omega_j x^{(i)}}\right) x^{(i)T} \\ \vdots \\ e^{(\omega_c + \omega_{y^{(i)}}) x^{(i)}} \, x^{(i)T} \end{bmatrix}$$

$$\nabla_\omega L_i(\omega) = \frac{1}{\sum_j e^{\omega_j x^{(i)}}} \begin{bmatrix} e^{\omega_1 x^{(i)}} \, x^{(i)T} \\ e^{\omega_2 x^{(i)}} \, x^{(i)t} \\ \vdots \\ -\left(\sum_{j \neq y^{(i)}} e^{\omega_j x^{(i)}}\right) x^{(i)T} \\ \vdots \\ e^{\omega_c x^{(i)}} \, x^{(i)T} \end{bmatrix}$$

$$\nabla_\omega L_i(\omega) = \begin{bmatrix} P(Y=1 \mid X = x^{(i)}) \, x^{(i)T} \\ P(Y=2 \mid X = x^{(i)}) \, x^{(i)T} \\ \vdots \\ -1 + P(Y = y^{(i)} \mid X = x^{(i)}) \, x^{(i)T} \\ \vdots \\ P(Y = c \mid X = x^{(i)}) \, x^{(i)T} \end{bmatrix}$$