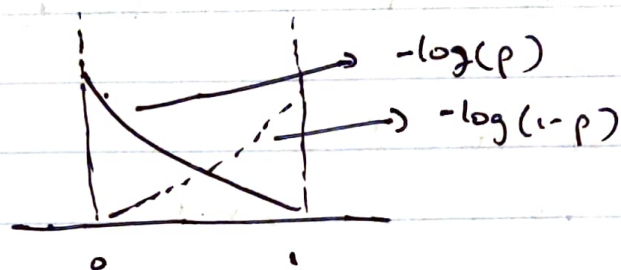


① a) cross entropy loss =  $-y \log(p) - (1-y) \log(1-p)$

$$y=0 \Rightarrow \text{loss} = -\log(1-p)$$

$$y=1 \Rightarrow \text{loss} = -\log(p)$$



- This is better than MSE as we do not have the issue of vanishing gradient.
- Here output is in range of  $[0, 1]$ , classification of input is based on the threshold of the probability of data belonging to particular class.

b) Back propagation cost function:

$$L_{\lambda}(w, b) = \frac{1}{2} \sum_{i=1}^N \text{loss}(h(x_i), y_i) + \frac{\lambda}{2} \sum_{l=1}^{n-1} \|w^{(l)}\|_F^2$$

considering BCE for each layer

$$L_{\lambda} \{w^{(l)}, b^{(l)}\}_{l=1}^{n-1} = -\frac{1}{2} \sum_{j=1}^{s_n} y_j \log(a_j^{(n)}(x^i)) +$$

$$(1-y_j^i) \log(1-a_j^{(n)}(x^i)) + \frac{\lambda}{2} \sum_{l=1}^{L-1} \|w^{(l)}\|_F^2$$

using gradient descent:

i) initialize  $w^{(1)} \in \mathbb{R}^{s_{l+1} \times s_l}$ ,  $b^{(1)} \in \mathbb{R}^{s_{l+1}}$

ii) Until convergence do,

for  $l = 1, 2, \dots, n-1$

$$a) w^{(l)(t+1)} \leftarrow w^{(l)(t)} - \eta \frac{\partial L}{\partial w^{(l)}} \bigg|_{w^{(l)(t)}}$$

$$b) b^{(l)(t+1)} \leftarrow b^{(l)(t)} - \eta \frac{\partial L}{\partial b^{(l)}} \bigg|_{b^{(l)(t)}}$$

$$\therefore \frac{\partial L}{\partial w^{(l)}} = \frac{-1}{2} \sum_{i=1}^N y_j^i \log(a_j^{(n)}(x^i)) + (1-y_j^i) \log(1-a_j^{(n)}(x^i)) + \lambda w^{(l)}$$

$$\frac{\partial L}{\partial b^{(l)}} = \frac{-1}{2} \sum_{i=1}^N y_j^i \log(a_j^{(n)}(x^i)) + (1-y_j^i) \log(1-a_j^{(n)}(x^i))$$

$$\begin{aligned} \Rightarrow 0 &\rightarrow a_1^{(1)}(x) \\ \Rightarrow 0 &\rightarrow a_2^{(1)}(x) \end{aligned} \quad \begin{bmatrix} a_1^{(1)}(x) \\ a_2^{(1)}(x) \end{bmatrix} = h(x^i) \text{ and } \begin{bmatrix} y^i \\ y^i \end{bmatrix} = y^i$$

$$\delta_i^{(n)} \triangleq \frac{\partial J}{\partial z_i^{(n)}} = \frac{\partial}{\partial z_i^{(n)}} \left[ -\frac{1}{2} \sum_j y_j \log(a_j^{(n)}(x)) + (1-y_j) \log(1-a_j^{(n)}(x)) \right]$$

①

$$\delta_i^{(n)} \triangleq \frac{\partial J}{\partial z_i^{(n)}} = \frac{\partial}{\partial z_i^{(n)}} \left[ -\frac{1}{2} \left[ y_i \log(a_i^{(n)}) + (1-y_i) \log(1-a_i^{(n)}) \right] \right]$$

$$= -\frac{1}{2} \left[ \frac{\partial y_i \log(a_i^{(n)})}{\partial z_i^{(n)}} + \frac{\partial (1-y_i) \log(1-a_i^{(n)})}{\partial z_i^{(n)}} \right]$$

$$\cancel{a_i^{(n)}} a_i^{(n)} = f(z_i^{(n)})$$

$$\therefore \delta_i^{(n)} = -\frac{1}{2} \left[ \frac{\partial y_i \log(f(z_i^{(n)}))}{\partial z_i^{(n)}} + \frac{\partial (1-y_i) \log(1-f(z_i^{(n)}))}{\partial z_i^{(n)}} \right]$$

$$\delta_i = -\frac{1}{2} \left[ y_i \frac{f'(z_i^{(n)})}{f(z_i^{(n)})} + (1-y_i) \frac{-f'(z_i^{(n)})}{1-f(z_i^{(n)})} \right] \quad \text{--- ②}$$

diff ② w.r.t  $z^{(n-1)}$

$$\begin{aligned}\therefore \delta_j^{(n-1)} &\triangleq \frac{\partial J}{\partial z_j^{(n-1)}} = \sum_{i=1}^n \frac{\partial J}{\partial z_i^{(n)}} \times \frac{\partial z_i^{(n)}}{\partial z_j^{(n-1)}} \\ &= \delta_j^{(n)} \times \frac{\partial z_i^{(n)}}{\partial z_j^{(n-1)}}\end{aligned}$$

$$z_i^{(n)} = w_{i1}^{(n-1)} \times a_{11}^{(n-1)} + \dots + w_{ij}^{(n-1)} a_{ij}^{(n-1)} + \dots$$

$w_{is}^{(n-1)} \cdot a_{sn-1}^{(n-1)}$

$$a_{ij}^{(n-1)} = f(z_i^{(n-1)})$$

$$\begin{aligned}\therefore \frac{\partial z_i^{(n)}}{\partial z_j^{(n-1)}} &= w_{ij}^{(n-1)} \times \frac{\partial a_{ij}^{(n-1)}}{\partial z_j} \\ &= w_{ij}^{(n-1)} \times \frac{\partial f(z_j^{(n-1)})}{\partial z_j}\end{aligned}$$

$$\frac{\partial z_i^{(n)}}{\partial z_j^{(n-1)}} = w_{ij}^{(n-1)} \times f'(z_j^{(n-1)})$$



$$s_j^{(n-1)} = \frac{\partial J}{\partial z_j^{(n-1)}} = \sum_{i=1}^{s_n} s_i^{(n)} \times w_{ij}^{(n-1)} \times f'(z_j^{(n-1)})$$

— (3)

generalize derivatives for layer  $l$ ,

$$\sum_{i=1}^{s_l} s_i^{(l)} \cdot w_{ij}^{(l-1)} \cdot f'(z_j^{(l-1)}) = s_j^{(l-1)} = \frac{\partial J}{\partial z_j^{(l-1)}}$$

diff.  $J$  w.r.t  $w_{ij}$  &  $b_i$

$$\frac{\partial J}{\partial w_{ij}^{(l)}} = \frac{\partial J}{\partial z_i^{(l+1)}} \times \frac{\partial z_i^{(l+1)}}{\partial w_{ij}^{(l)}}$$

$$z_i^{(l+1)} = \dots + w_{ij}^{(l)} a_j^{(l)} + \dots + b_i^{(l)} \quad \& \text{ using } \textcircled{1} \& \textcircled{2}$$

$$\frac{\partial J}{\partial w_{ij}} = s_i^{(l+1)} \cdot a_j^{(l)}$$

similarly

$$\frac{\partial J}{\partial b_i^{(l)}} = \frac{\partial J}{\partial z_i^{(l+1)}} \times \frac{\partial z_i^{(l+1)}}{\partial b_i^{(l)}}$$

$$\frac{\partial J}{\partial b_i^{(l)}} = s_i^{(l+1)} \times 1$$

$$\frac{\partial J}{\partial \omega^{(l)}} = \begin{bmatrix} \frac{\partial J}{\partial \omega^{(l)}_{i_1}} & \frac{\partial J}{\partial \omega^{(l)}_{i_2}} & \dots \end{bmatrix} = \begin{bmatrix} \delta^{(l+1)}_{i_1} a_1^{(l)} & \delta^{(l+1)}_{i_2} a_2^{(l)} & \dots \end{bmatrix}$$

$$= \begin{bmatrix} \delta^{(l+1)}_{i_1} \\ \delta^{(l+1)}_{i_2} \\ \vdots \end{bmatrix} [a_1^{(l)} \ a_2^{(l)} \ \dots] = \delta^{(l+1)}_{i_1} \cdot a^{(l)T}$$

$$\frac{\partial J}{\partial b^{(l)}} = \delta^{(l+1)}_{i_1} = \begin{bmatrix} \delta^{(l+1)}_{i_1} \\ \delta^{(l+1)}_{i_2} \\ \vdots \end{bmatrix} = \delta^{(l+1)}$$

c) First derivative of sigmoid function:

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

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

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

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

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

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

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

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

$$\sigma'(z) = \sigma(z) \times (1 - \sigma(z))$$

- (2) a) function  $a_1$  is at  $z_1$   
function  $a_2$  is at  $z_2$

$$z_1 = x_1 w_3 + x_2 w_5 + w_1$$

$$z_2 = x_1 w_4 + x_2 w_6 + w_2$$

Activation

$$\begin{aligned} a_1 &= c z_1 \\ &= c (x_1 w_3 + x_2 w_5 + w_1) \end{aligned}$$

$$\begin{aligned} a_2 &= c z_2 \\ &= c (x_1 w_4 + x_2 w_6 + w_2) \end{aligned}$$

$$\begin{aligned} z_3 &= w_7 + w_8 a_1 + w_9 a_2 \\ &= w_7 + w_8 c [x_1 w_3 + x_2 w_5 + w_1] + w_9 c [x_1 w_4 + x_2 w_6 + w_2] \end{aligned}$$

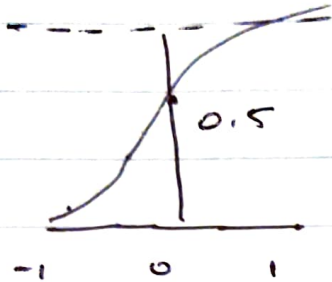
$$\begin{aligned} &= x_1 [c w_3 w_8 + c w_4 w_9] + x_2 [c w_5 w_8 + c w_6 w_9] \\ &\quad + [w_7 + c w_1 w_8 + c w_2 w_9] \end{aligned}$$

$$a_3 = P(y=1 | z, w) = \frac{1}{1 + e^{-z_3}}$$

$$= \frac{1}{1 + e^{-\left( x_1 [c w_3 w_8 + c w_4 w_9] + x_2 [c w_5 w_8 + c w_6 w_9] + [w_7 + c w_1 w_8 + c w_2 w_9] \right)}}$$



Thus we get a sigmoid function ( $\frac{1}{1+e^x}$ )



→ sigmoid function, this gives us final classification boundary.

$$b) z_3 = x_1 [c\omega_3\omega_8 + c\omega_4\omega_9]$$

$$+ x_2 [c\omega_5\omega_8 + c\omega_6\omega_9]$$

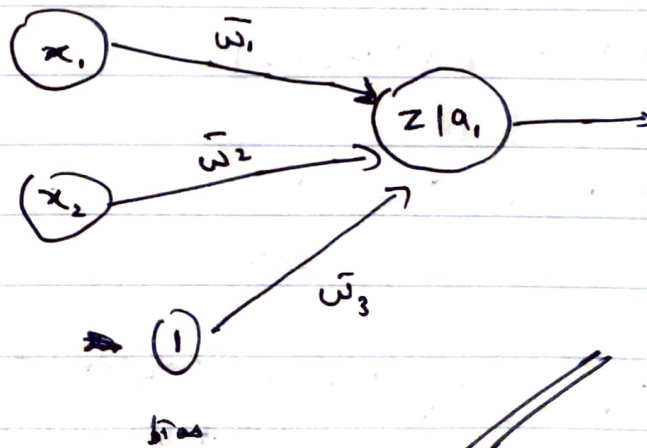
$$+ [\omega_7 + c\omega_1\omega_8 + c\omega_2\omega_9]$$

weights

$$\bar{\omega}_1 = c\omega_3\omega_8 + c\omega_4\omega_9$$

$$\bar{\omega}_2 = c\omega_5\omega_8 + c\omega_6\omega_9$$

$$\bar{\omega}_3 = c\omega_1\omega_8 + c\omega_2\omega_9 + \omega_7$$



c) In such a case, ~~the~~ we get linear activation of the inputs. If we remove hidden layer we can still represent the multilayered neural network.  
TRUE!



③ Assigning linear activation functions for layer 1 :

$$L(z_1) = c(x_1 w_1 + x_2 w_3)$$

$$\text{from } z_1 = x_1 w_1 + x_2 w_3$$

$$L(z_2) = c(x_1 w_2 + x_2 w_4)$$

$$\text{from } z_2 = x_1 w_2 + x_2 w_4$$

$$\begin{aligned} \text{Therefore we get } z_2 &= w_5 [c(x_1 w_1 + x_2 w_3)] \\ &\quad + w_6 [c(x_1 w_2 + x_2 w_4)] \\ &= x_1 (c(w_1 w_5 + w_2 w_6)) + x_2 (c(w_3 w_5 + w_4 w_6)) \end{aligned}$$

Assign sigmoid function for layer 2

$$\begin{aligned} s(z^{(2)}) &= \frac{1}{1 + \exp(-x_1 (c(w_1 w_5 + w_2 w_6)) - x_2 (c(w_3 w_5 + w_4 w_6)))} \\ &= \frac{1}{1 + \exp(x_1 \beta_1 + x_2 \beta_2)} \end{aligned}$$

where

$$\beta_1 = -c(w_1 w_5 + w_2 w_6)$$

$$\beta_2 = -c(w_3 w_5 + w_4 w_6)$$

