$X \in \mathbb{R}^{m \times 2}$

Input layer     Hidden layer 1     Hidden layer 2     out layer

**Feedforward :**

$$Z_1 = X \cdot W_1 + \vec{b_1}$$
$$A_1 = g_1(Z_1) \qquad ; \text{where } g_1 = \text{activation function}$$

$$Z_2 = A_1 W_2 + \vec{b_2}$$
$$A_2 = g_2(Z_2) \qquad ; \text{where } g_2 = \text{activation function}$$

$$Z_3 = A_2 \cdot W_3 + \vec{b_3}$$
$$A_3 = g_3(Z_3) = \hat{y} \qquad ; \text{where } g_3 = \text{activation function}$$

$g_1 / g_2 \rightarrow$ ReLU activation

$g_3 \rightarrow$ Sigmoid activation

**Loss :** Cross-entropy, $\mathcal{L} = \frac{-1}{m} \sum_{i=0}^{m} \left( y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right)$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

## Activation functions:

$$\text{ReLU}(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{d}{dx} \text{ReLU}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

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

$$\frac{d}{dx} \sigma(x) = \frac{d}{dx}\left(\frac{1}{1+e^{-x}}\right) = -(1+e^{-x})^{-2} \cdot (-e^{-x})$$

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

$$= \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}}$$

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

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

$$= \sigma(x) \cdot (1 - \sigma(x))$$

## Backpropagation:

$$\frac{\partial L}{\partial z_3} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_3} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \sigma(z_3)}{\partial z_3}$$

$$= \left(\frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}}\right) \cdot \left(\hat{y}(1-\hat{y})\right)$$

$$= \frac{-y(1-\hat{y}) + (1-y)\hat{y}}{\hat{y}(1-\hat{y})} \cdot \hat{y}(1-\hat{y})$$

$$= \hat{y} + y \cdot y + \hat{y} - y\hat{y}$$

$$= \hat{y} - y$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_3} = \frac{\partial L}{\partial z_3} \cdot \frac{\partial (A_2 \cdot w_3 + b_3)}{\partial w_3}$$

$$\therefore \frac{\partial L}{\partial w_3} = \hat{y} - y \cdot A_3$$

$$\frac{\partial L}{\partial b_3} = \frac{\partial L}{\partial z_3} \cdot \frac{\partial z_3}{\partial b_3} = \hat{y} - y \qquad ; \frac{\partial z_3}{\partial b_3} = 1$$

$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2} = \frac{\partial L}{\partial z_3} \cdot \frac{\partial z_3}{\partial A_2} \cdot \frac{\partial A_2}{\partial z_2}$$

$$= \hat{y} - y \cdot w_3 \cdot g_2'(z_2)$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot A_1$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2} = \frac{\partial L}{\partial z_2} \cdot 1$$

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial A_1} \cdot \frac{\partial A_1}{\partial z_1} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial A_1} \cdot \frac{\partial A_1}{\partial z_1}$$

$$= \frac{\partial L}{\partial z_2} \cdot w_2 \cdot g_1'(z_1)$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1} = \frac{\partial L}{\partial z_1} \cdot X$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1} = \frac{\partial L}{\partial z_1} \cdot 1$$

Update parameters: $\theta := \theta - \alpha \frac{\partial L}{\partial \theta}$ ; where $\alpha$ = learning rate

$$w_3 = w_3 - \alpha \frac{\partial L}{\partial w_3} \qquad w_2 = w_2 - \alpha \frac{\partial L}{\partial w_2} \qquad w_1 = w_1 - \alpha \frac{\partial L}{\partial w_1}$$

$$b_3 = b_3 - \alpha \frac{\partial L}{\partial b_3} \qquad b_2 = b_2 - \alpha \frac{\partial L}{\partial b_2} \qquad b_1 = b_1 - \alpha \frac{\partial L}{\partial b_1}$$