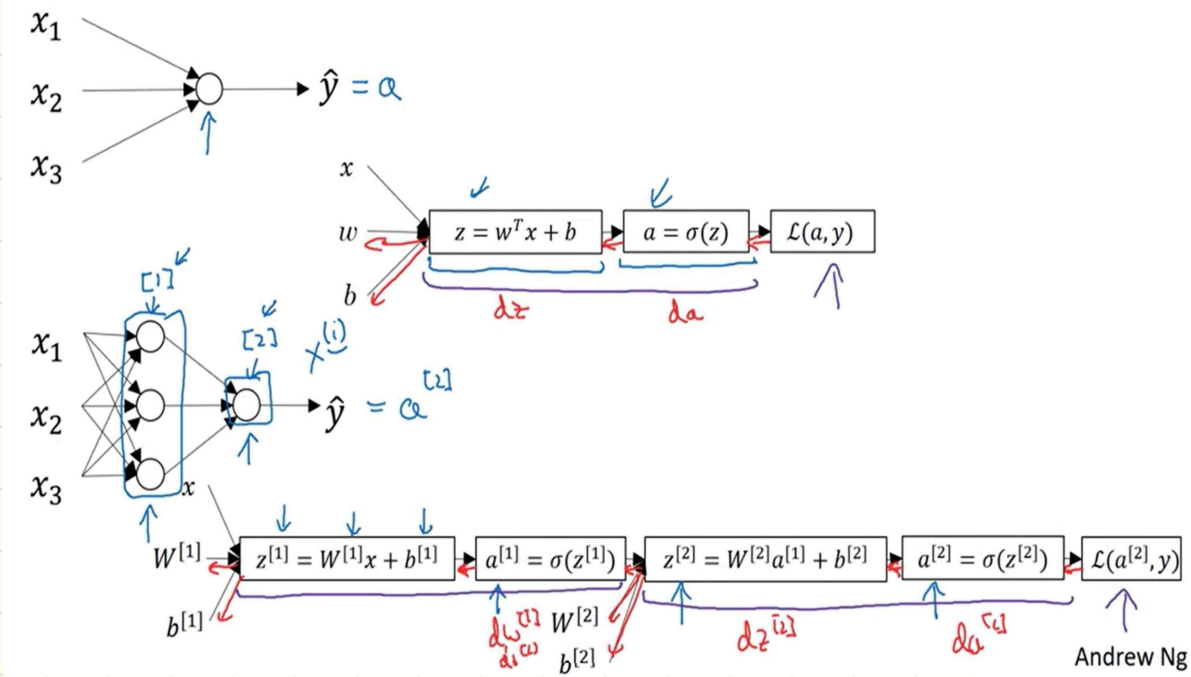


# [7-1. Shallow Neural Network]

## <Neural Networks Overview>

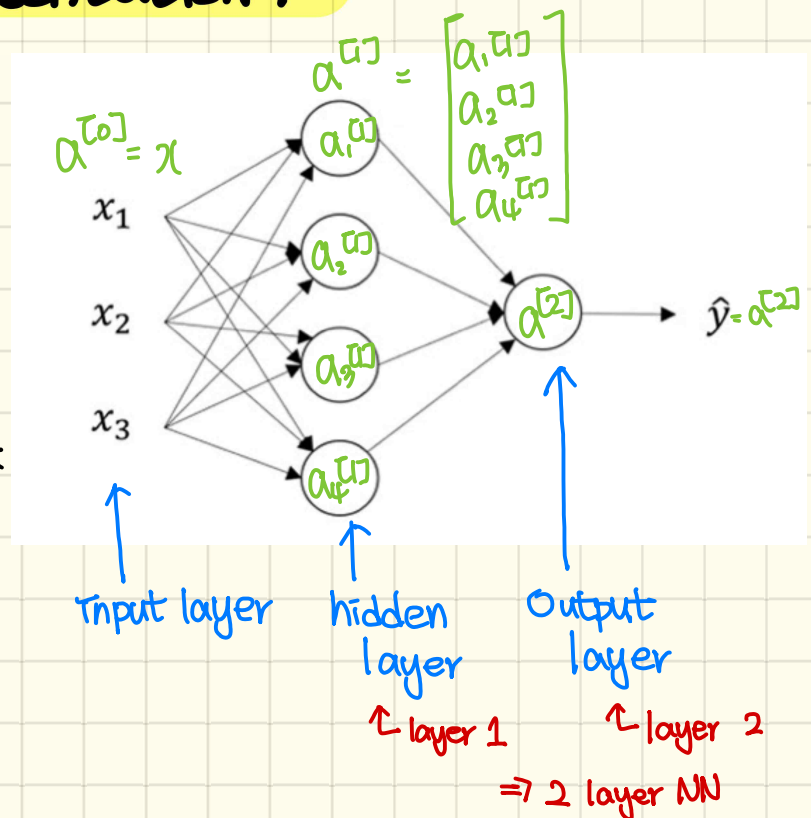
What is a Neural Network?



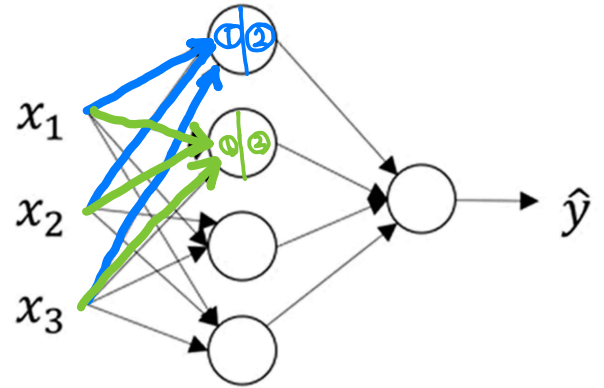
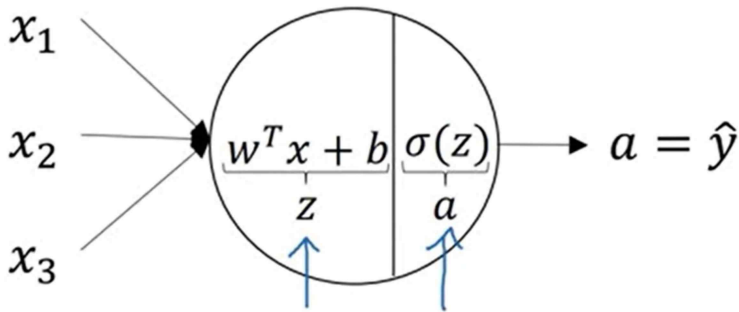
## <Neural Network Representation>

- \* Input layer (입력층) :  $a^{[0]}$
- Hidden layer (은닉층) :  $a_n^{[i]}$
- Output layer (출력층)

\* 신경망의 층의 개수를 셀 때는 입력층은 제외한다  
 $\therefore$  오른쪽 사진은 2-layer Neural Network



# <Computing a Neural Network's Output>



$$z = w^T x + b$$

$$a = \sigma(z)$$

\* 파란색 표시

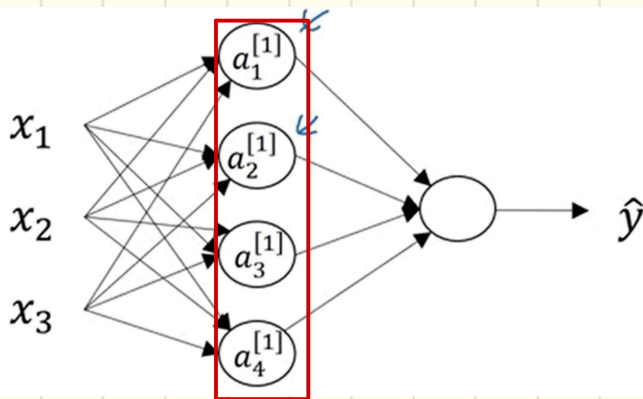
$$\textcircled{1} z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$$

$$\textcircled{2} a_1^{[1]} = \sigma(z_1^{[1]})$$

\* 연두색 표시

$$\textcircled{1} z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$$

$$\textcircled{2} a_2^{[1]} = \sigma(z_2^{[1]})$$

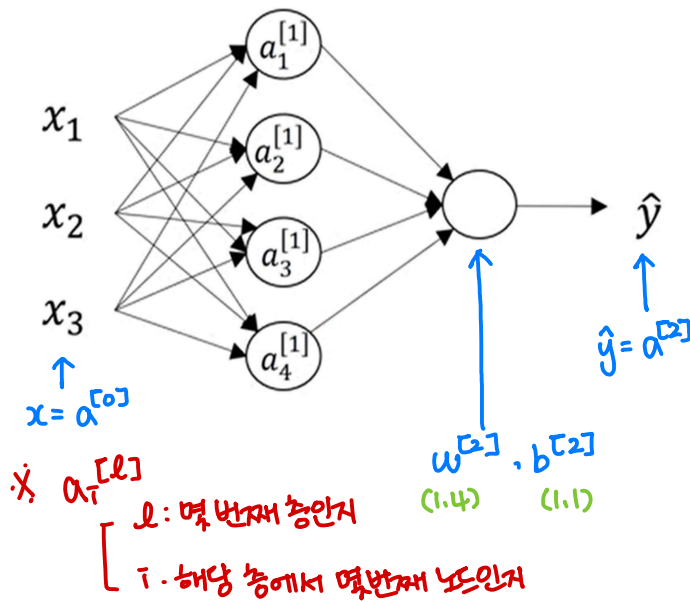


$$\uparrow a^{[1]} = \sigma(z^{[1]})$$

$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, & a_1^{[1]} &= \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, & a_2^{[1]} &= \sigma(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, & a_3^{[1]} &= \sigma(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, & a_4^{[1]} &= \sigma(z_4^{[1]}) \end{aligned}$$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} = \underbrace{\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix}}_{(4,3)} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{(3,1)} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}}_{(4,1)} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix}$$

$\uparrow$   $b^{[1]}$



Given input  $x$ :

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$(4,1) \quad (4,3) \quad (3,1) \quad (4,1)$

$$a^{[1]} = \sigma(z^{[1]})$$

$(4,1) \quad (4,1)$

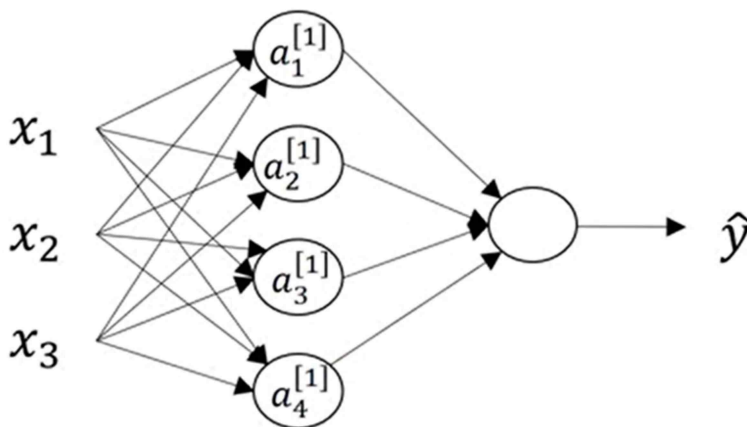
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$(1,1) \quad (1,4) \quad (4,1) \quad (1,1)$

$$a^{[2]} = \sigma(z^{[2]})$$

$(1,1) \quad (1,1)$

## < Vectorizing Across Multiple Examples >



$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

$x \longrightarrow a^{[2]} = \hat{y}$   
 m개의 training set 으로 확장  
 $\Rightarrow x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$   
 $x^{(2)} \longrightarrow a^{[2](2)} = \hat{y}^{(2)}$   
 $\dots$   
 $x^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)}$   
 $\dots a^{[i](j)}$   
 $i$ : 몇 번째 층인지  
 $j$ : 몇 번째 훈련 샘플인지

$\Rightarrow$  for  $i=1$  to  $m$ ,

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(z^{[2](i)})$$

for  $i = 1$  to  $m$ :

$$z^{[1]}(i) = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1]}(i) = \sigma(z^{[1]}(i))$$

$$z^{[2]}(i) = W^{[2]}a^{[1]}(i) + b^{[2]}$$

$$a^{[2]}(i) = \sigma(z^{[2]}(i))$$

$$z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = \sigma(z^{[1]})$$

Vectorizing

$$z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(z^{[2]})$$

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}, \quad z^{[1]} = \begin{bmatrix} | & | & | \\ z^{[1]}(1) & z^{[1]}(2) & \dots & z^{[1]}(m) \\ | & | & | \end{bmatrix}$$

$(n \times m)$

$$A^{[1]} = \begin{bmatrix} | & | & | \\ a^{[1]}(1) & a^{[1]}(2) & \dots & a^{[1]}(m) \\ | & | & | \end{bmatrix} \quad \left. \begin{array}{l} \text{hidden units} \\ \text{training unit} \end{array} \right\}$$

## <Explanation for Vectorized Implementation>

$$z^{[1]}(1) = W^{[1]}x^{(1)} + b^{[1]}, \quad z^{[1]}(2) = W^{[1]}x^{(2)} + b^{[1]}, \quad z^{[1]}(3) = W^{[1]}x^{(3)} + b^{[1]}$$

$$W^{[1]} = \begin{bmatrix} \text{wavy} \\ \text{wavy} \\ \text{wavy} \\ \text{wavy} \end{bmatrix}$$

$$W^{[1]}x^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

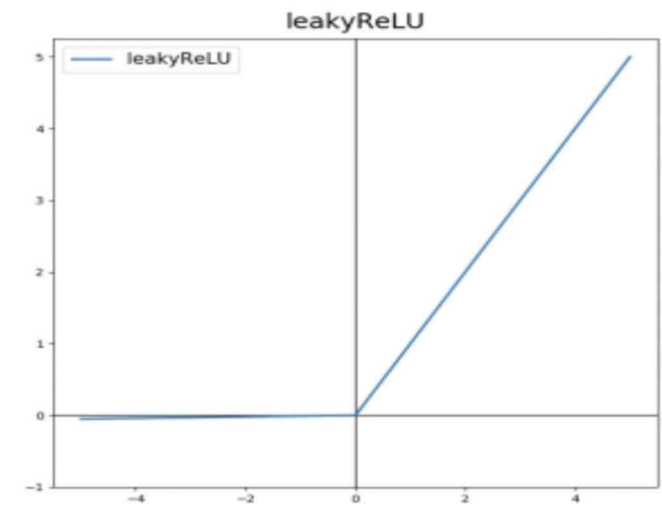
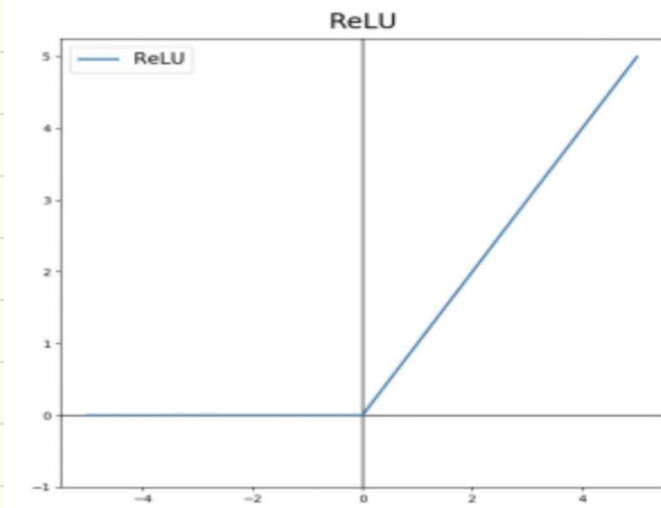
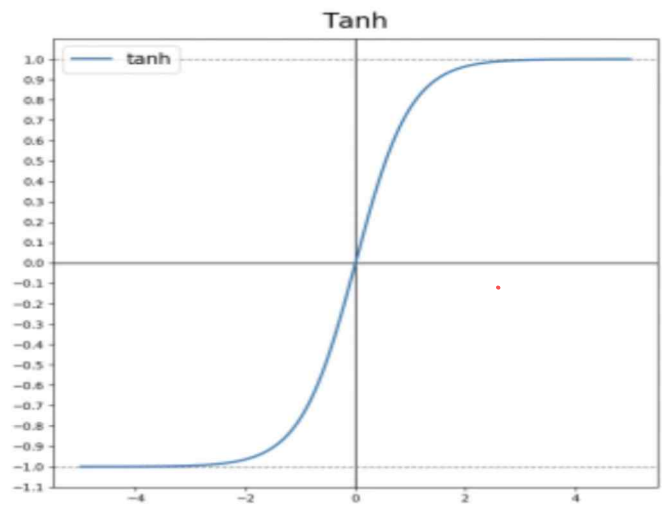
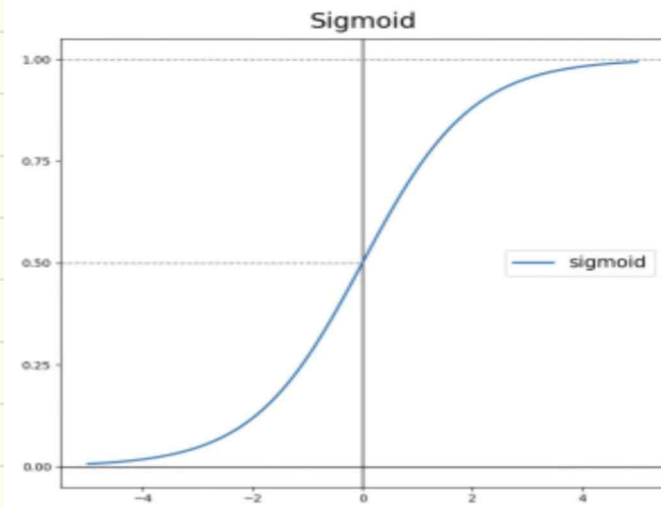
$$W^{[1]}x^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$W^{[1]}x^{(3)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$W^{[1]} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots \\ | & | & | \end{bmatrix} = \begin{bmatrix} \bullet & \bullet & \bullet & \dots \\ \bullet & \bullet & \bullet & \dots \\ \bullet & \bullet & \bullet & \dots \\ \bullet & \bullet & \bullet & \dots \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{[1]}(1) & z^{[1]}(2) & z^{[1]}(3) & \dots \\ | & | & | \end{bmatrix} = z^{[1]}$$

$$\therefore \underline{W^{[1]}x^{(i)} = z^{[1]}(i)}$$

# < Activation Functions >



- Sigmoid :  $a = \frac{1}{1+e^{-x}}$

☆☆  $\tan h$  :  $a = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

☆☆☆☆ ReLU :  $a = \max(0, x)$

- leaky ReLU :  $a = \max(0.01x, x)$

$\tan h$ 의 장점

- 값이  $[-1, 1]$  사이에 있고 평균이 0 이므로, 데이터를 이동하는 효과가 있으며, 평균이 0.5인 Sigmoid 보다 더 효율적이다.

ReLU의 장점

- 0보다 큰 활성화함수의 기울기가 다른 활성화함수의 값보다 크기 때문에 더 효율적이다.

# <Why do you need Non-Linear Activation Function?>

$$z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]}) \rightarrow \text{선형 활성화 함수: } g(z) = z \text{ 를 사용하면,}$$

$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]}) \rightarrow z^{[2]}$$

$$a^{[1]} = z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[2]} = z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$$= w^{[2]} \{ w^{[1]}x + b^{[1]} \} + b^{[2]}$$

$$= \underbrace{(w^{[2]}w^{[1]})}_{w'}x + \underbrace{(w^{[2]}b^{[1]} + b^{[2]})}_{b'}$$

$$= w'x + b'$$

∴ 선형 활성화 함수를 사용하면,  $g(g(g(z))) = z$  로 세개의 은닉층을 쌓아도 효과를 얻지 못한다.

## <Derivatives of Activation Functions>

\* Sigmoid

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

slope of  $g(z)$  at  $z$

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

\* Tanh

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

$$g'(z) = 1 - (g(z))^2$$

$$\rightarrow z=10, \tanh(z) \approx 1 \rightarrow g'(z) \approx 0$$

$$z=-10, \tanh(z) \approx -1 \rightarrow g'(z) \approx 0$$

$$z=0, \tanh(z)=0 \rightarrow g'(z)=1$$

\* ReLU

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & (\text{if } z < 0) \\ 1 & (\text{if } z \geq 0) \end{cases}$$

↪ 0 이하로는 정의하지 않음

\* Leaky ReLU

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & (\text{if } z < 0) \\ 1 & (\text{if } z \geq 0) \end{cases}$$



# < Gradient Descent for Neural Networks >

- Parameters :  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$   
 $(n^{[1]}, n^{[0]}) \quad (n^{[1]}, 1) \quad (n^{[2]}, n^{[1]}) \quad (n^{[2]}, 1) \quad n_x = n^{[0]}, n^{[2]} = 1$

- Cost Function :  $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{[i]}, y^{[i]})$   
 $= a^{[2]}$

- Gradient Descent :

(반복 계산)  $dw^{[1]} = \frac{\partial J}{\partial w^{[1]}}, db^{[1]} = \frac{\partial J}{\partial b^{[1]}}, \dots$

$$W^{[1]} = W^{[1]} - \alpha dw^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha db^{[1]}$$

< Forward propagation >

$$z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]}) = \sigma(z^{[2]})$$

< Back Propagation >

$$dz^{[2]} = A^{[2]} - Y, Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims}=\text{True})$$

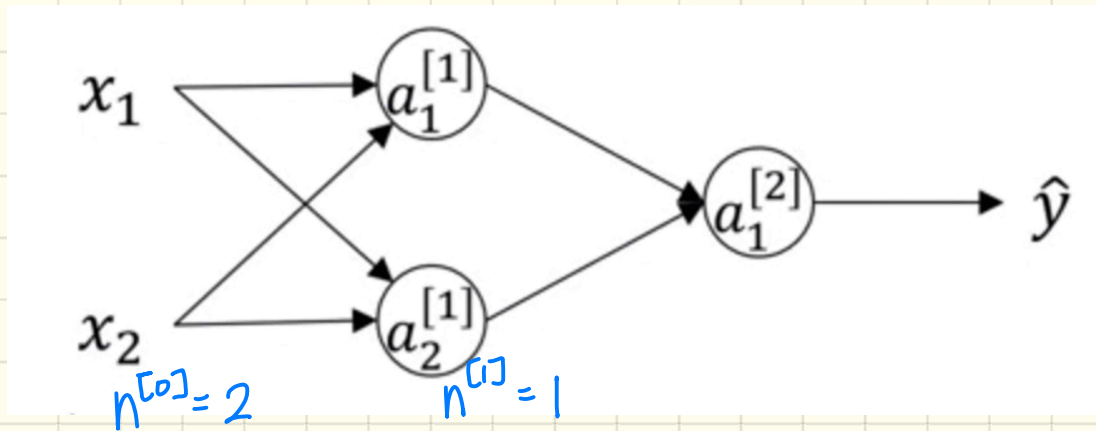
$$dz^{[1]} = \underbrace{W^{[2]T}}_{(n^{[2]}, m)} dz^{[2]} * \underbrace{g^{[1]'}(z^{[1]})}_{(n^{[1]}, m)}$$

:  $(n^{[1]}, 1)$ 로 계산 X  
 $(n^{[1]}, 1)$ 로 계산 O

$$dw^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims}=\text{True})$$

# < Random Initialization >



If)  $w$ 의 초기값을 '0'으로 초기화 한다면?

$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \rightarrow a_1^{[1]} = a_2^{[1]}, \quad dz_1^{[1]} = dz_2^{[1]}$$

$\Rightarrow dw$ 를 계산했을 때, 모든 층이 같은 값을 갖게 된다.

$\therefore$  0이 아니라 `np.random.rand()`를 이용하여 랜덤값을 부여해주어야 함.

$$\Rightarrow W^{[1]} = np.random.randn(2,2) * 0.01$$

$$b^{[1]} = np.zeros((2,1))$$

매우 작은 값으로 만들어주기 위함.

$$W^{[2]} = np.random.randn(1,2) * 0.01$$

$$b^{[2]} = 0$$