NN 구성

$X_1$

$X_2$

$X_3$

$\hat{y} = a$

↑
Input
layer

$a^{[0]}$

↑
hidden
layer

$a^{[1]}$

↑
output
layer

$a^{[2]}$

$X_1$

$X_2$

$X_3$

$$z_1^{[1]} = w_1^{[1]T} x + b$$
$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]} x + b$$
$$a_2^{[1]} = \sigma(z_2^{[1]})$$

⋮

- 벡터화

- layer 1



$$\begin{cases} z_1^{[1]} = w_1^{[1]T} x + b \\ a_1^{[1]} = \sigma(z_1^{[1]}) \end{cases}$$

$$\begin{cases} z_2^{[1]} = w_2^{[1]} x + b \\ a_2^{[1]} = \sigma(z_2^{[1]}) \end{cases}$$
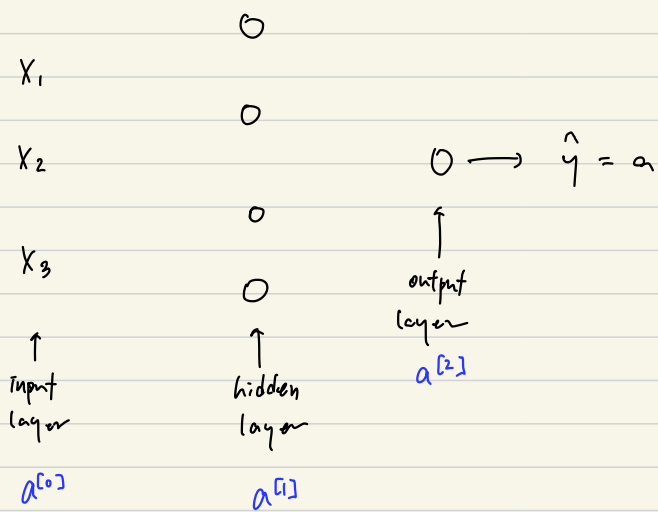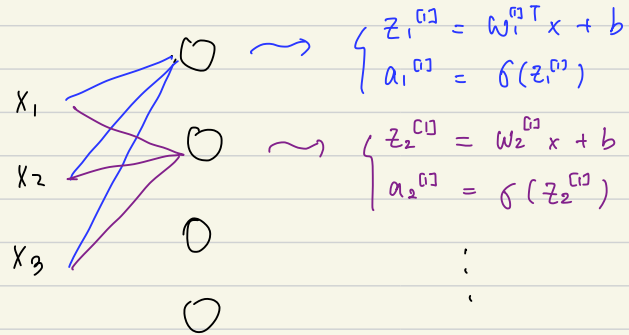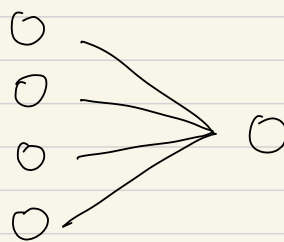
$$z^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ \vdots \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ \vdots \\ z_4^{[1]} \end{bmatrix}$$

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

$$a^{[1]} = \begin{bmatrix} \sigma(z_1^{[1]}) \\ \vdots \\ \sigma(z_4^{[1]}) \end{bmatrix} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix}$$

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

- layer 2



$$W^T = w^{[2]}, \quad b = b^{[2]}$$
$$\quad\quad (1,4) \quad\quad\quad (1,1)$$
$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$
$$a^{[2]} = \sigma(z^{[2]})$$
$$(1,1)$$

- Vectoring across multiple examples

* $a^{[n](i)}$       $n$: layer
                    $i$: example

- 한 개의 샘플에 대한 벡터화

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

$\Rightarrow$

- $m$개의 샘플

for $i = 1$ to $m$,
$$z^{[1](i)} = W^{[1](1)}x^{(i)} + b^{[1]}$$
$$a^{[1](i)} = \sigma(z^{[1](i)})$$
$$z^{[2](i)} = W^{[2](i)}a^{[1](i)} + b^{[2]}$$
$$a^{[2](i)} = \sigma(z^{[2](i)})$$

- $m$개의 samples에 대한 $X$, $Z$, $A$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \cdots & x^{(m)} \end{bmatrix} \quad (n_x, m) \qquad (m: \# \text{ training samples})$$

# training samples

$$Z^{[1]} = \begin{bmatrix} z^{[1](1)} & z^{[1](2)} & \cdots & z^{[1](m)} \end{bmatrix} \quad \updownarrow \text{\# hidden units}$$

$$A^{[1]} = \begin{bmatrix} a^{[1](1)} & a^{[1](2)} & \cdots & a^{[1](m)} \end{bmatrix}$$

$\Rightarrow$
$$Z^{[1]} = W^{[1]}X + b^{[1]}$$
$$A^{[1]} = \sigma(Z^{[1]})$$
$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$
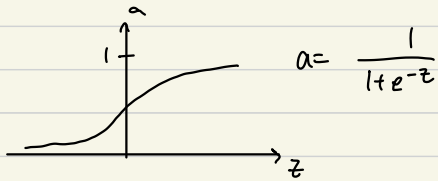$$A^{[2]} = \sigma(Z^{[2]})$$

- Vectorized implementation

$$Z^{[1]} = \left[ W^{[1]} x^{(1)} + b^{[1]} \quad W^{[1]} x^{(2)} + b^{(1)} \quad \cdots \quad W^{[1]} x^{(m)} + b^{[1]} \right]$$

$$= W^{[1]} \left[ x^{(1)} \; x^{(2)} \cdots \; x^{(m)} \right] + b^{[1]}$$

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

## ⟨ Activation functions ⟩

- Sigmoid function



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

⟵ Never use sigmoid except output layer in binary classification!

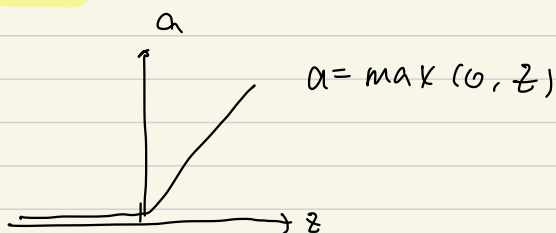- tanh



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

- 은닉층 : tanh > sigmoid   ⟵ 값이 ±1 사이라서 평균이 0에 가깝기 때문.

- 출력층 : sigmoid   ⟵ 값이 0~1 사이라서 이진분류에 유리함

- sigmoid, tanh 단점

  - z가 매우 크거나 작은 때 기울기가 0에 가까워지면서 경사하강법 느려짐

- ==ReLu==



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

- 단점 : z < 0 인때 기울기가 0
  ↓
  Leaky ReLu : z < 0 인때도 약간의 기울기를 줌

- 장점 : 대부분의 z 에서 기울기가 0보다 큼.

⇒ 이진분류 출력층 → sigmoid
   다른 경우 → ReLu

< 신경망 네트워크 경사하강법 >

Parameters :  $W^{[1]}$ , $b^{[1]}$ , $W^{[2]}$ , $b^{[2]}$
dim : $(n^{[1]}, n^{[0]})$, $(n^{[1]}, 1)$ , $(n^{[2]}, n^{[1]})$ , $(n^{[2]}, 1)$

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

Gradient descent :

Repeat {
$\qquad$ Compute Params ($\hat{y}^{(i)}$, $i = 1, \cdots, m$)
$\qquad$ $dW^{[1]} = \dfrac{dJ}{dW^{[1]}}$ ,  $db^{[1]} = \dfrac{dJ}{db^{[1]}}$ ,  $\cdots$  # Params 의 도함수 구함

$\qquad$ $W^{[1]} := W^{[1]} - \alpha \, dW^{[1]}$    # update params
$\qquad$ $b^{[1]} := b^{[1]} - \alpha \, dW^{[1]}$

$\qquad\qquad\qquad \vdots$

} 수렴할 때까지 반복

&lt;Forward propagation&gt;

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

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

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

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

&lt;Backward propagation&gt;

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$(n^{[2]}, ) \rightarrow (n^{[2]}, 1)$

$$db^{[2]} = \frac{1}{m} np.sum (dZ^{[2]}, \; axis=1, \; keepdims=True)$$
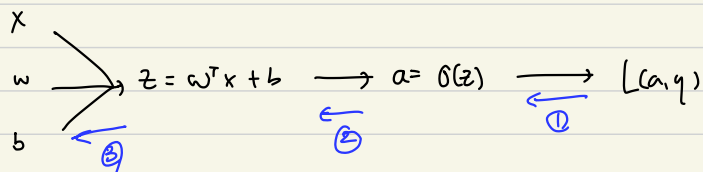
$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'} (Z^{[1]})$$

↳ element-wise product

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

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

&lt;역전파 도함수 계산식 유도 과정&gt;

X
w
b

$\xrightarrow{}$ $z = \omega^T x + b$ $\xrightarrow{}$ $a = \sigma(z)$ $\xrightarrow{}$ $L(a,y)$

③          ②          ①

① $da = \dfrac{d}{da} L(a,y) = -y\log a - (1-y)\log(1-a) = -\dfrac{y}{a} + \dfrac{1-y}{1-a}$

② $dz = \dfrac{dL}{dz} = \dfrac{da}{dz} \cdot \dfrac{dL}{da}$

$= \dfrac{d}{dz} g(z) \cdot da \qquad (a = g(z) = \sigma(z))$

$= da \cdot g'(z)$

$= a - y$

③ $dw = dz \cdot x$

$db = dz$

$X$

$W^{[1]}$ $\longrightarrow$ $Z^{[1]} = W^{[1]}X + b^{[1]}$ $\longrightarrow$ $a^{[1]} = \sigma(Z^{[1]})$

$b^{[1]}$ $\underset{④}{\longleftarrow}$ $\underset{③}{\longleftarrow}$

$a^{[1]}$

$W^{[2]}$ $\longrightarrow$ $Z^{[2]} = W^{[2]}X + b^{[2]}$ $\longrightarrow$ $a^{[2]} = \sigma(Z^{[2]})$ $\longrightarrow$ $L(a^{[2]}, y)$

$b^{[2]}$ $\underset{②}{\longleftarrow}$ $\underset{①}{\longleftarrow}$

① $dZ^{[2]} = a^{[2]} - y$

② $dw^{[2]} = dZ^{[2]} a^{[1]T}$

$db^{[2]} = dZ^{[2]}$

③ $dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$

$*$ dim: $\quad W^{[2]} \quad (n^{[2]}, n^{[1]}) = (1, n^{[1]})$

$Z^{[2]}, dZ^{[2]} \quad (n^{[2]}, 1) = (1, 1)$

$Z^{[1]}, dZ^{[1]} \quad (n^{[1]}, 1)$

$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$

$(n^{[1]}, 1) \quad (n^{[1]}, n^{[2]}) (n^{[2]}, 1) \quad * (n^{[1]}, 1)$

④ $dw^{[1]} = dZ^{[1]} \cdot X^T = a^{[0]T}$

$db^{[1]} = dZ^{[1]}$

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]}a^{[1]^T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]^T}dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]}x^T$$

$$db^{[1]} = dz^{[1]}$$

- 벡터화

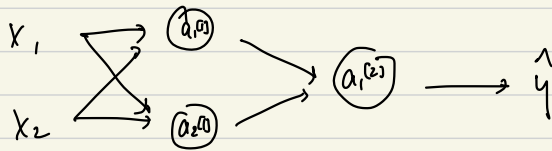$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

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

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

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

<랜덤 초기화>

$x_1$ ⟶ $a_1^{[1]}$
$x_2$ ⟶ $a_2^{[1]}$ ⟶ $a_1^{[2]}$ ⟶ $\hat{y}$

- 파라미터를 0으로 초기화하는 경우

$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \qquad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]} \qquad dz_1^{[1]} = dz_1^{[2]}$$

⟹ $a_1^{[1]}$, $a_2^{[1]}$ hidden unit이 항상 같은 계산을 함.
⟹ hidden unit이 하나인 것과 같음.
⟹ 랜덤하게 파라미터를 초기화해야함

- Random initialization

$W^{[1]} = np.random.rand ((2,2)) * 0.01$
$b^{[1]} = np.zeros ((2,1))$
$W^{[2]} = np.random.rand ((1,2)) * 0.01$
$b^{[2]} = 0$

⟹ 작은 수로 초기화하는 이유?
⟹ $W$가 커지면 $Z$도 커짐
⟹ sigmoid, tanh 에서는 $Z$가 커질수록 기울기가 0에 가까워지기 때문에 학습 속도가 느려짐.