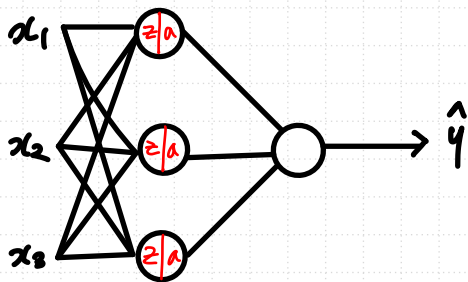


# 활성화 함수



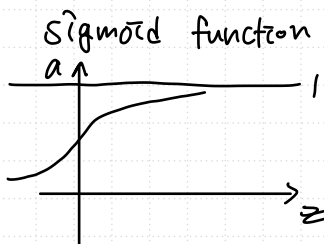
Given  $X$ )  $\begin{bmatrix} w_{01} & w_{11} & w_{21} \\ \vdots & \vdots & \vdots \end{bmatrix} = W^{[1]}$

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

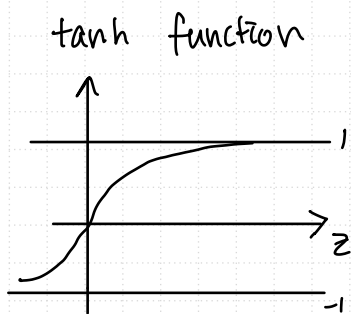
$$a^{[1]} = \sigma(z^{[1]})$$

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

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



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



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

원점을 지나는 바뀐 모양이 달라짐

→  $\sigma \rightarrow \tanh$  in hidden layer?

→ 편향없이 더 다양한 값들이 더 좋다

→ 레이어 중심 0.5 to 0 : 다음 층의 학습이 더 쉬워진다

but 학습능은 아님

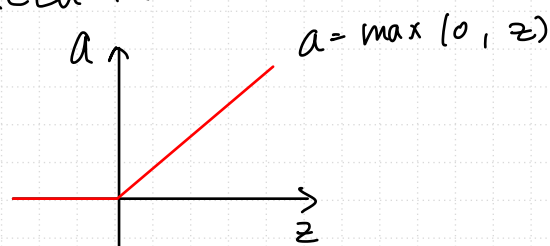
→ 이걸 분류에서는 0과 1 사이로 값이 더 좋기 때문에 시그모이드 함수가 더 좋다.

Sigmoid와 tanh의 단점

→ 근가 굉장히 크거나 작으면 도함수의 값이 거의 0임

→ **경사하강법이 느려짐**

ReLU! (rectified linear unit)



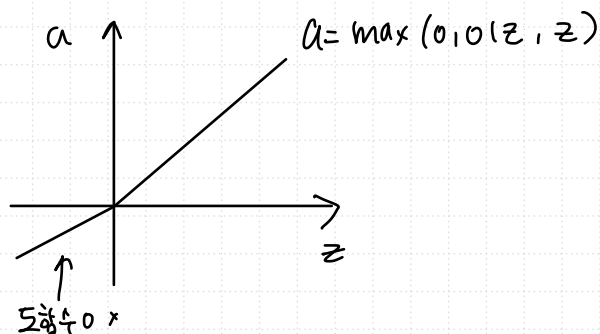
$$z > 0 \quad \frac{da}{dz} = 1$$

$$z < 0 \quad \frac{da}{dz} = 0$$

$z = 0$ ? 엄밀하게 정의 X but 컴퓨터에서 0.000...0이 나오는 거의 항상함

∴ 미분불가능함도 2차

Leaky ReLU?



ReLU와 Leaky ReLU의 장점

→ 대부분의 z에 대해 기울기가 0과 매우 다르다

→ **성능이 훨씬 빠르게 학습됨**

※ 각 근원이 따라 어떤 함수가 좋은지 직접 다 해봐야 한다.

# 왜 비선형 함수를 써야하는가?

if) 활성화 함수 없다면

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

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

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

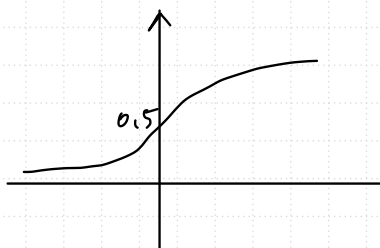
$$= \underbrace{w^{[2]}w^{[1]}}_{w'}x + \underbrace{w^{[2]}b^{[1]} + b^{[2]}}_{b'}$$

$$= w'x + b'$$

층이 얼마나 많은가에, 은닉층이 없는 것과 다름없는 결과가 나온다.

## 활성화 함수의 미분

1) Sigmoid

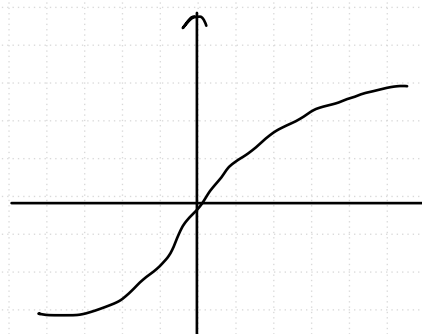


$$\begin{aligned} \frac{dg(z)}{dz} &= \frac{1}{1+e^{-z}} \left( 1 - \frac{1}{1+e^{-z}} \right) \\ &= g(z)(1-g(z)) = a(1-a) \end{aligned}$$

$$\oplus z=0) g'(z) = \frac{1}{4}$$

∴ 이 때  $g(z)$ 의 값을 구해준다면  $g'(z)$ 는 빠르게 구할 수 있다.

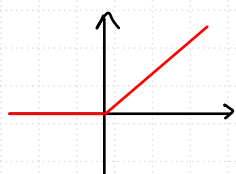
2) tanh



$$\frac{dg(z)}{dz} = 1 - (\tanh(z))^2 = 1 - a^2$$

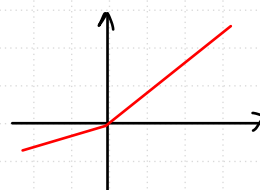
∴ tanh 함수를 미분하면  $\tanh(z)$ 가 있다면 빠르게 구할 수 있다.

3) ReLU



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

4) Leaky-Relu



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

# 신경망 벡터화코딩 경사하강법

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

$n_x = n^{[0]}$ ,  $n^{[1]}$ ,  $n^{[2]} = 1$   
 층의 뉴런 수

Cost func:  $J(\text{params}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}_i, y_i)$   
 (이전 블록)

Gradient descent: 0 이 아닌 수를 초기화해주는 것

Repeat { Until convergence

Compute predicts ( $\hat{y}^{(i)}$ ,  $i=1 \dots m$ )

1 iter  $\left\{ \begin{array}{l} dw^{[0]} = \frac{dJ}{dw^{[0]}}, db^{[0]} = \frac{dJ}{db^{[0]}}, \dots \\ w^{[0]} := w^{[0]} - \alpha dw^{[0]} \\ b^{[0]} := b^{[0]} - \alpha db^{[0]} \\ w^{[1]} := w^{[1]} - \alpha dw^{[1]} \\ b^{[1]} := b^{[1]} - \alpha db^{[1]} \end{array} \right.$

Forward propagation

$z^{[0]} = W^{[0]}x + b^{[0]}$

$A^{[0]} = g^{[0]}(z^{[0]})$

$z^{[1]} = W^{[1]}A^{[0]} + b^{[1]}$

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

이전블록에서  
마지막 활성화함수는 sigmoid

Back propagation

$dz^{[0]} = A^{[0]} - y$  ( $y = [y^{(1)} \dots y^{(m)}]$ )

$dW^{[0]} = \frac{1}{m} dz^{[0]} A^{[0]T}$

$db^{[0]} = \frac{1}{m} \text{np.sum}(dz^{[0]}, \text{axis}=1, \text{keepdims}=\text{True})$

$dz^{[1]} = \underbrace{W^{[1]T}}_{(n^{[1]}, m)} dz^{[0]} \otimes \underbrace{g^{[0]'}(z^{[0]})}_{(n^{[0]}, m)}$   
 element wise product

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

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

$\frac{d\mathcal{L}}{da} = \frac{1}{da} (-y \log a - (1-y) \log(1-a))$

$= -\frac{y}{a} + \frac{1-y}{1-a}$

$\frac{d\mathcal{L}}{dz} = a - y$

$\frac{d\mathcal{L}}{dw} = dz \cdot \sigma$

$\frac{d\mathcal{L}}{db} = dz$

## 랜덤 초기화

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

$$\Rightarrow a_1^{[0]} = a_2^{[0]} \quad \Rightarrow d_{z_1}^{[0]} = d_{z_2}^{[0]} \quad \left. \vphantom{\begin{matrix} \Rightarrow a_1^{[0]} = a_2^{[0]} \\ \Rightarrow d_{z_1}^{[0]} = d_{z_2}^{[0]} \end{matrix}} \right\} \text{가중치의 영향은 받지 않음 서로 다른 뉴런의 값이 똑같다.}$$

: Symmetric

⇒ 같은 것을 제시하기 때문에 뉴런층에 Unit이 하나 있는 것만  
다음 없다

$$W^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

해결책

$W^2$  랜덤 초기화한다

$$W^{[0]} = \text{np.random.randn}((2,2)) \times 0.01$$

$$b^{[0]} = \text{np.zeros}((2,1))$$

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

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

→ 뉴런층의 값에 따라 이 숫자가 달라진다

가중치나 bias에 따라

