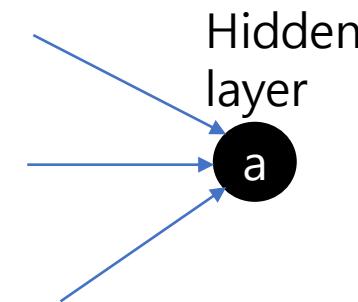


# Activation Function

- Activation function
- $a = g(w^T x + b)$ ,  $z = w^T x + b$

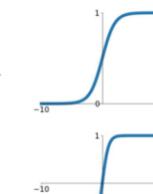
- Sigmoid :  $a = \frac{1}{1+e^{-z}}$
- Tanh :  $a = \frac{\sinh z}{\cosh z} \rightarrow \text{mean value } \sim 0$
- ReLU :  $a = \max(0, z)$



## Activation Functions

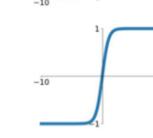
### Sigmoid

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



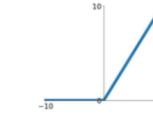
### tanh

$$\tanh(x)$$



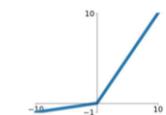
### ReLU

$$\max(0, x)$$



### Leaky ReLU

$$\max(0.1x, x)$$

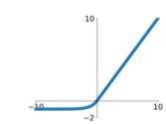


### Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

### ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



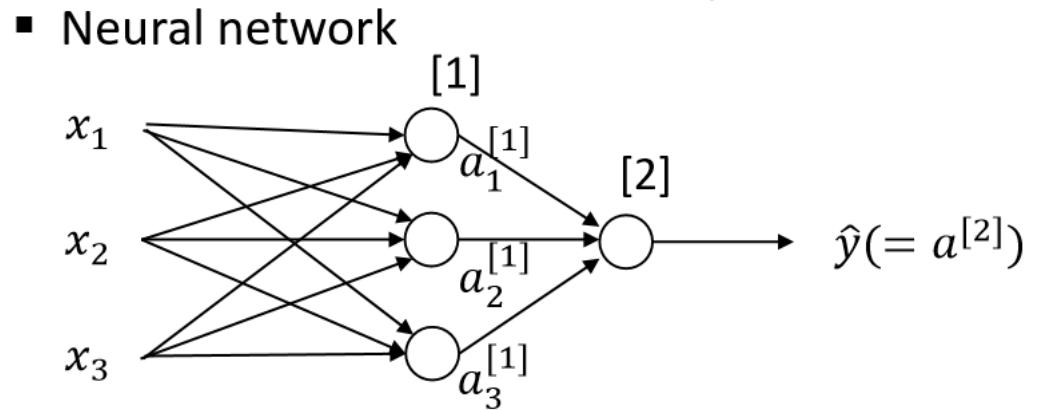
# Why Nonlinear Activation function

- Linear activation : matrix multiplication
- Output depends on only input. -> No meaning of hidden layers
- Assume  $g(z) = z$
- $a^{[1]} = w_1^T x + b_1, a^{[2]} = w_2^T a^{[1]} + b_2$
- $a^{[2]} = w_2^T w_1^T x + w_2^T b_1 + b_2$
- Linear activation is used in regression problem (in output layer)

# Derivatives of activation functions

- Pass
- Derivatives is used in gradient descent for neural network.

# Gradient descent



- Back propagation to give feedback to the neural network
- Parameters :  $w^1, w^2, b^1, b^2$
- Cost function(loss function) :  $J(w^1, w^2, b^1, b^2) = 1/m \sum L(\hat{y}, y)$
- Parameter random initialization is important.

# Gradient descent

$z^1 = W^1x + b^1$	$a^1 = \sigma(z^1)$	$z^2 = W^2a^1 + b^2$	$a^2 = \sigma(z^2)$	$L(a^2, y)$
--------------------	---------------------	----------------------	---------------------	-------------

- Primary input:  $x, W1, b1$
- $w^i = w^i - \alpha \frac{\partial J}{\partial w^i}$
- $b^i = b^i - \alpha \frac{\partial J}{\partial b^i}$
- $\alpha$  is learning rate.
- Assume Loss function :  $-y\log\hat{y} - (1 - y)\log(1 - \hat{y})$

# Gradient descent

- $dz^2 = a^2 - y$
- $dW^2 = dz^2 a^{1T}$
- $db^2 = dz^2 , db^1 = dz^1$
- $dz^1 = W^{2T} dz^2 * g^{1'}(z^1)$
- $dW^1 = dz^1 x^T$

# Weight initialization

- If we initialize the weight 0,...
- $W^1 = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, b^1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$
- Make  $a_1^1 = a_2^1$  and  $dz_1^1 = dz_2^1 \rightarrow$  Full symmetric, a hidden layer's nodes always have same value.
- Random generation is good. But give small value.

