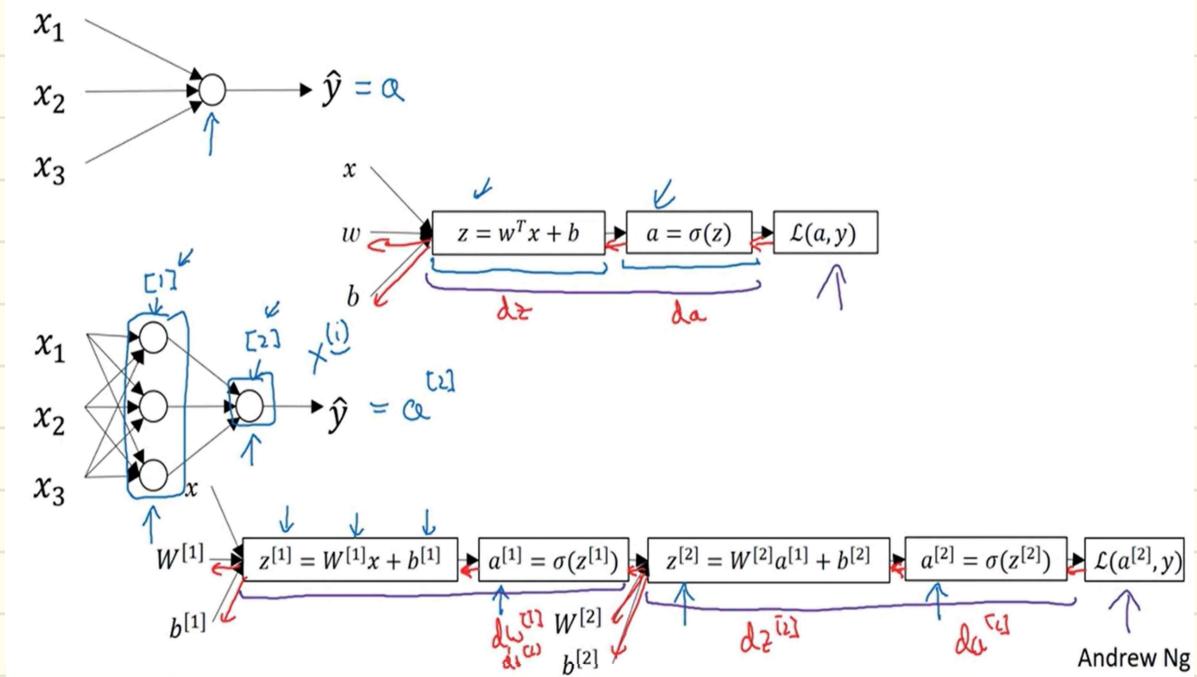


[3-1. Shallow Neural Network]

< Neural Networks Overview >

What is a Neural Network?



Andrew Ng

< Neural Network Representation >

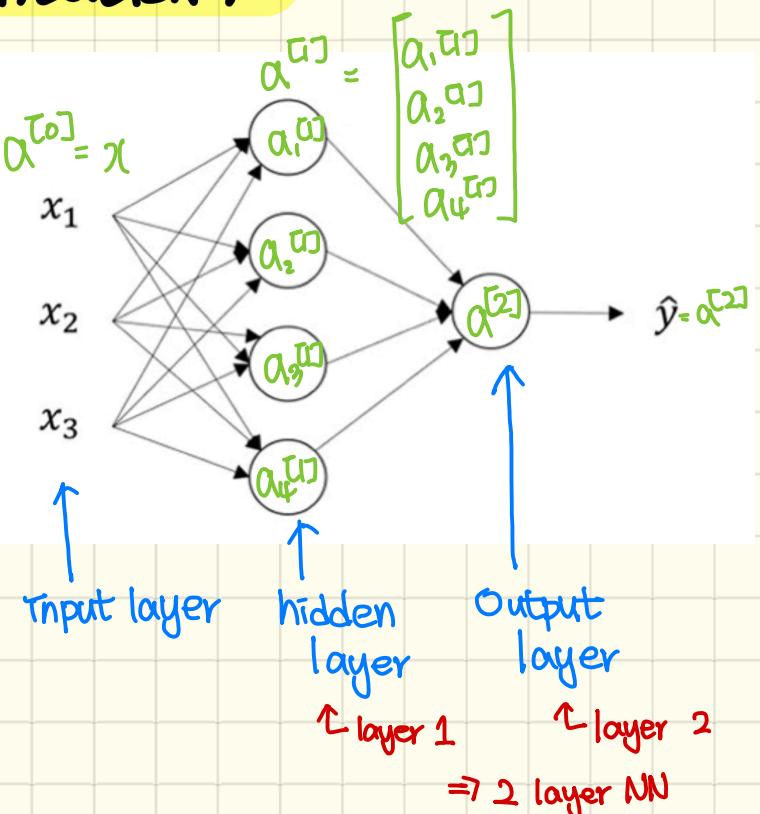
* Input layer (입력층) : $a^{[0]}$

Hidden layer (은닉층) : $a_n^{[1]}$

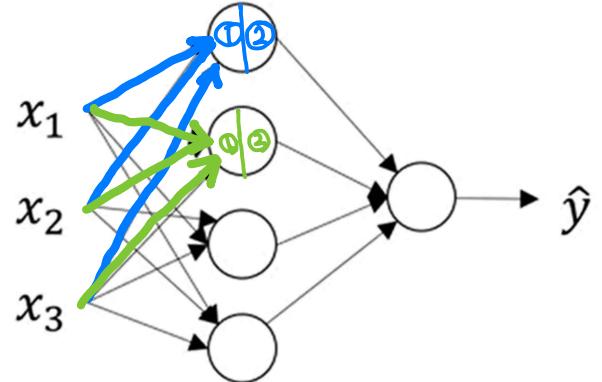
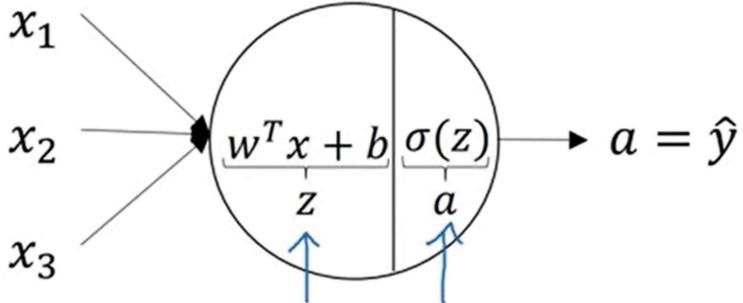
Output layer (출력층)

* 신경망의 층의 개수를 세 때는 입력층은 제외한다

∴ 오른쪽 사진은 2-layer Neural Network



<Computing a Neural Network's Output>

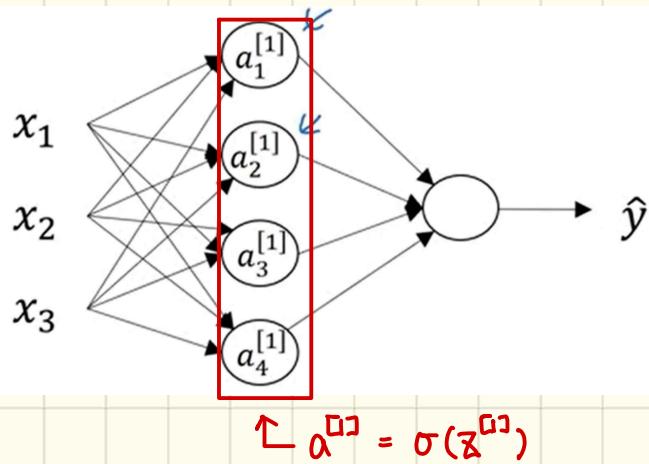


$$z = w^T x + b$$

$$a = \sigma(z)$$

* 파산식 표시

$$\begin{aligned} ① z^{[1]} &= w_1^{[1]T} x + b_1^{[1]} \\ ② a_1^{[1]} &= \sigma(z^{[1]}) \end{aligned}$$

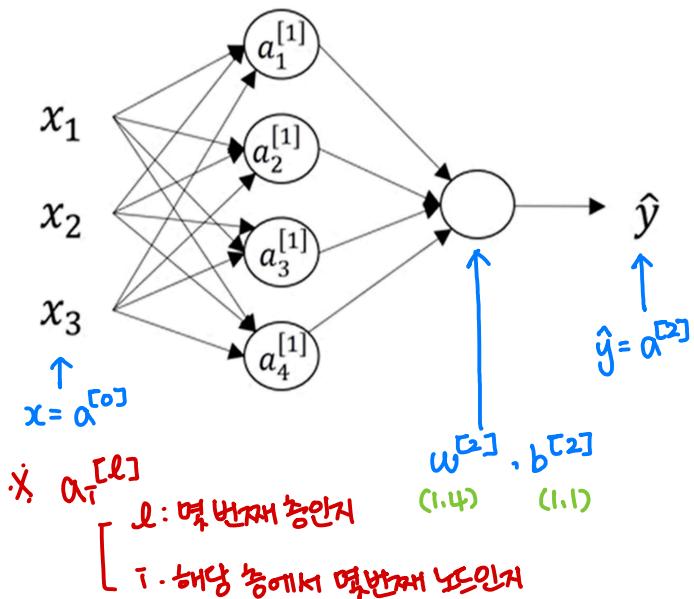


* 연두식 표시

$$\begin{aligned} ① z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]} \\ ② a_2^{[1]} &= \sigma(z_2^{[1]}) \end{aligned}$$

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

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



Given input x :

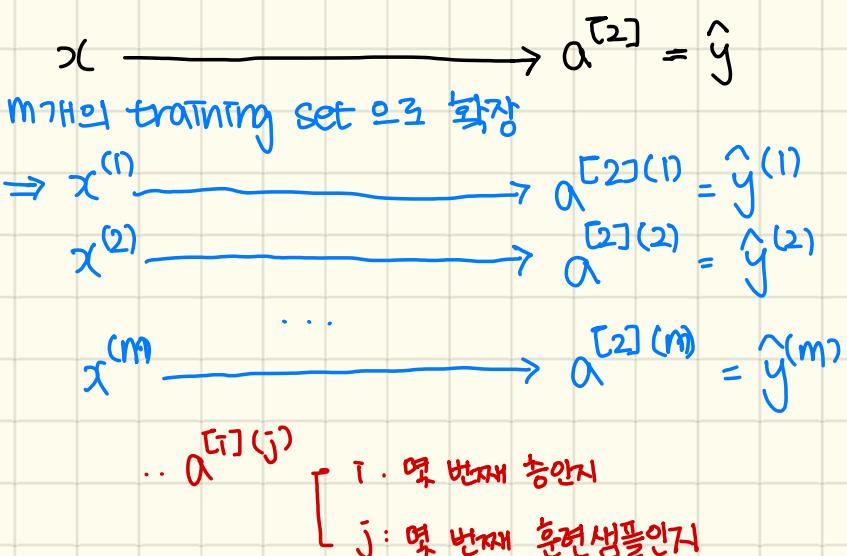
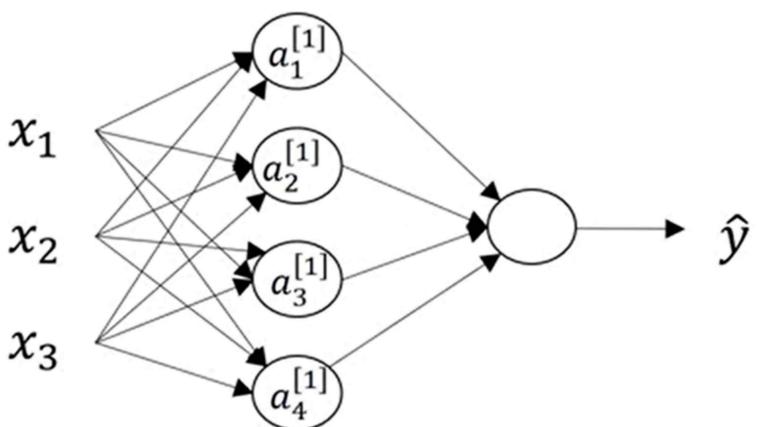
$$z^{[1]} = W^{[1]}x + b^{[1]} \quad (4.1) \quad (4.3) \quad (3.1) \quad (4.1)$$

$$a^{[1]} = \sigma(z^{[1]}) \quad (4.1) \quad (4.1)$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \quad (1.1) \quad (1.4) \quad (4.1) \quad (1.1)$$

$$a^{[2]} = \sigma(z^{[2]}) \quad (1.1) \quad (1.1)$$

<Vectorizing Across Multiple Examples>



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

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

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

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

⇒ for $i=1$ to m ,

$$\begin{aligned} z^{[1](i)} &= W^{[1]}x^{(i)} + b^{[1]} \\ a^{[1](i)} &= \sigma(z^{[1](i)}) \\ z^{[2](i)} &= W^{[2]}a^{[1](i)} + b^{[2]} \\ a^{[2](i)} &= \sigma(z^{[2](i)}) \end{aligned}$$

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

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

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

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

$\underbrace{\hspace{100pt}}_{(n \times m)}$

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

<Explanation for Vectorized Implementation>

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

$$w^{[1]} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}$$

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

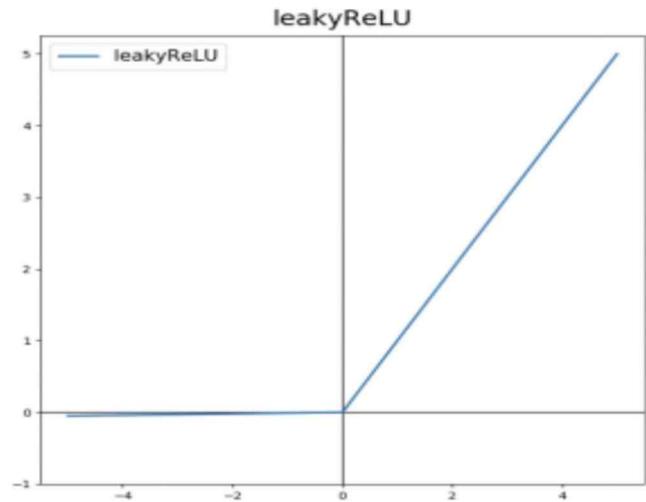
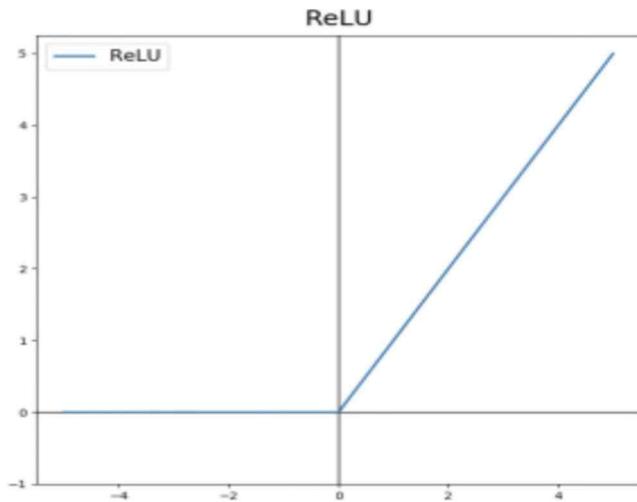
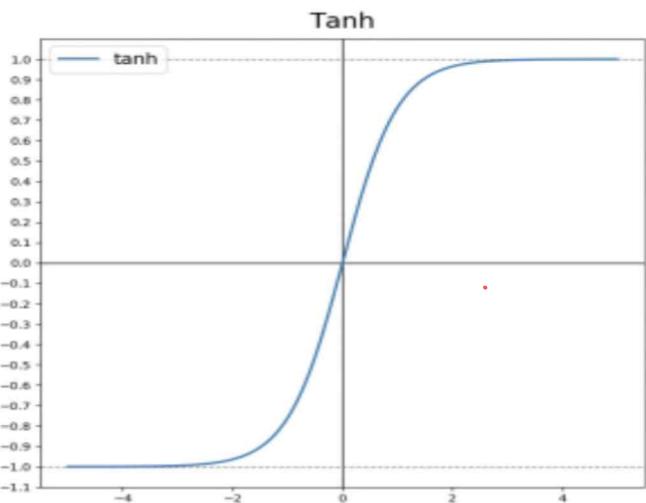
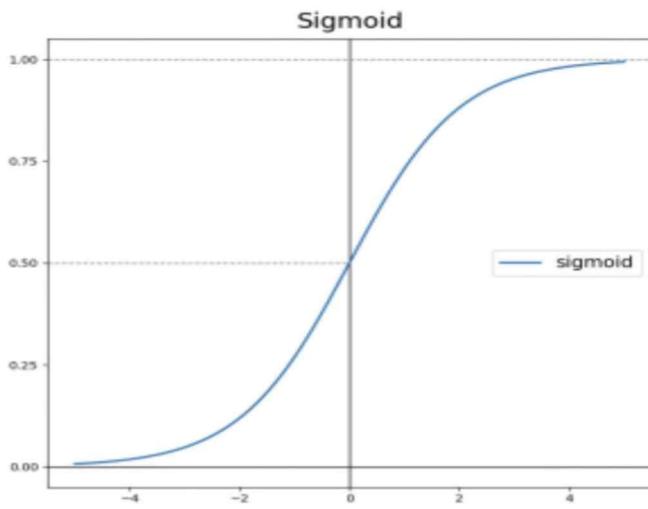
$$w^{[1]}x^{(2)} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

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

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

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

< Activation Functions >



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

tanh의 장점

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

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

- ReLU : $a = \max(0, x)$

ReLU의 장점

- 0보다 큰 활성화수의 미분값이 다른 활성화수의 값보다 크기 때문에 더 효율적이다.

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

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

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

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

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

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

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

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

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

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

$$= \underline{w'x} + b'$$

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

<Derivatives of Activation Functions>

* sigmoid

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

slope of $g(z)$ at z

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

* Tanh

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

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

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

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

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

* ReLU

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

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

↑ 실제로는 정의역에 어려움

* Leaky ReLU

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

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

<Gradient Descent for Neural Networks>

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

- Cost Function : $J(w^{[0]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}, y)$

- Gradient Descent :

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

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

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

<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>

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

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

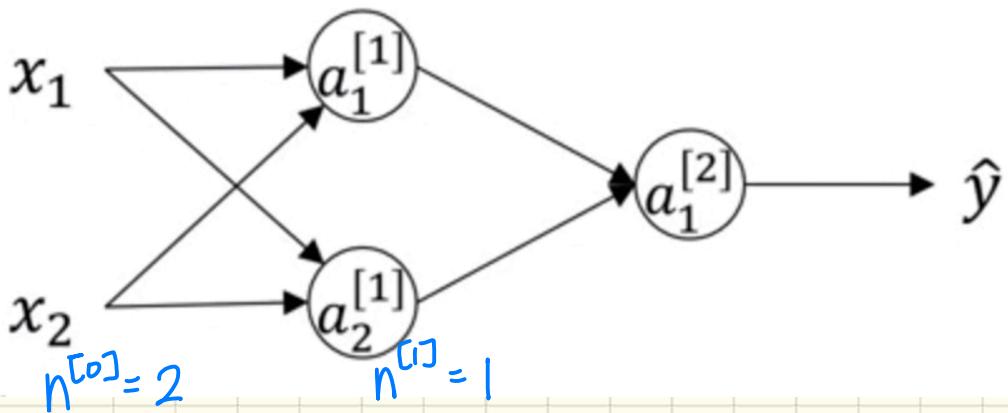
$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}), \text{axis}=1, \text{keepdims=True}$$

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

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

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

<Random Initialization>



If) w 의 초기값을 '0'으로 초기화 한다면?

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

$\Rightarrow dw$ 를 계산했을 때, 모든 속이 같은 값을 갖게된다.

\therefore 0이 아님 np.random.rand()를 이용하여 랜덤값을 부여해야 함.

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

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

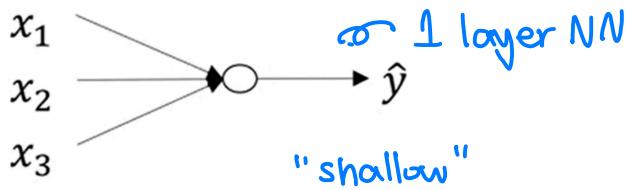
매우 작은 값으로 만들어주기 위함.

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

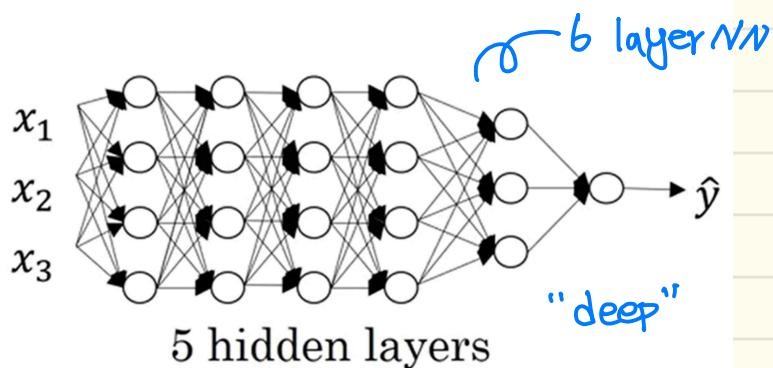
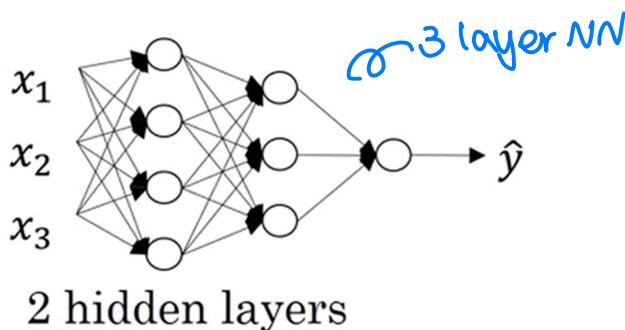
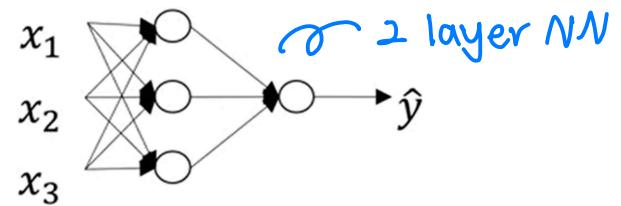
$$b^{[2]} = 0$$

[4-1. Deep Neural Network]

<Deep L-layer Neural Network>

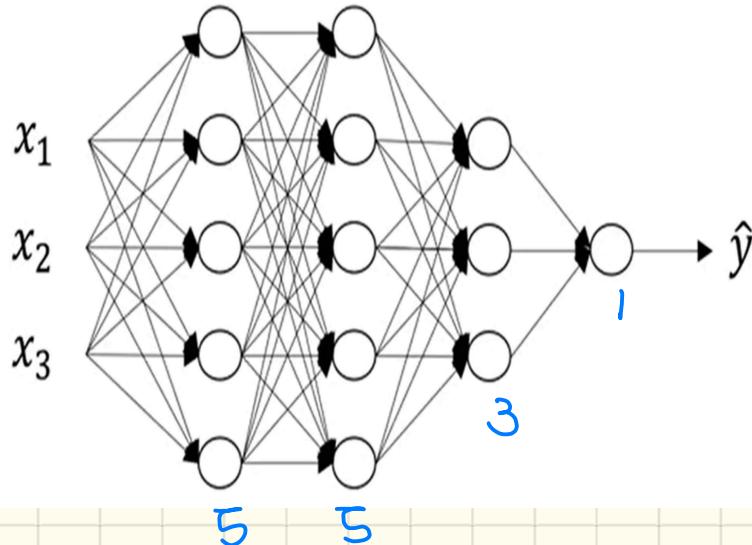


logistic regression



* 얼마나 깊은 신경망을 사용해야 하는지는 예측하기 어렵지만, 하이퍼파라미터를 조정하면서 깊이를 결정해야 한다.

ex) 4 layer NN



$$L = \# \text{layers} = 4$$

$$n^{[l]} = \# \text{ units in layer } l$$

$$n^{[0]} = 3, n^{[1]} = 5, n^{[2]} = 5$$

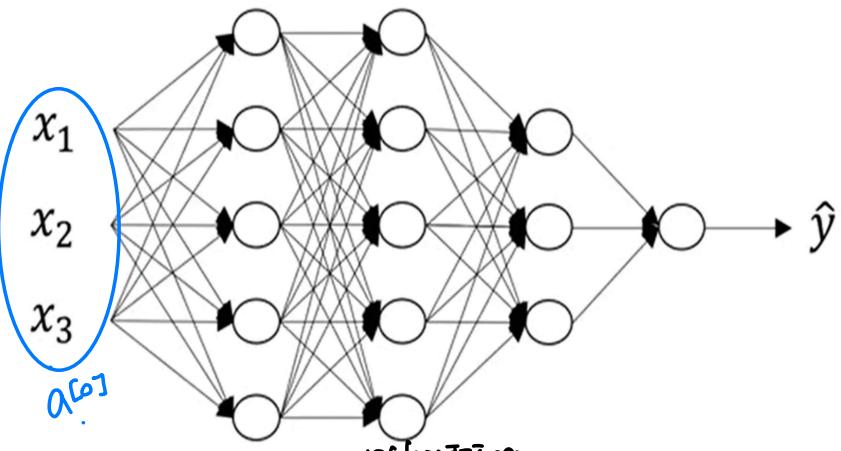
$$n^{[3]} = 3, n^{[4]} = n^{[L]} = 1$$

$$a^{[l]} = g^{[l]}(z^{[l]}) = \text{activation in layer } l$$

$$a^{[0]} = x, a^{[L]} = \hat{y}$$

$$w^{[l]}, b^{[l]}. \text{ weights for } z^{(l)}$$

<Forward Propagation In a Deep Network >



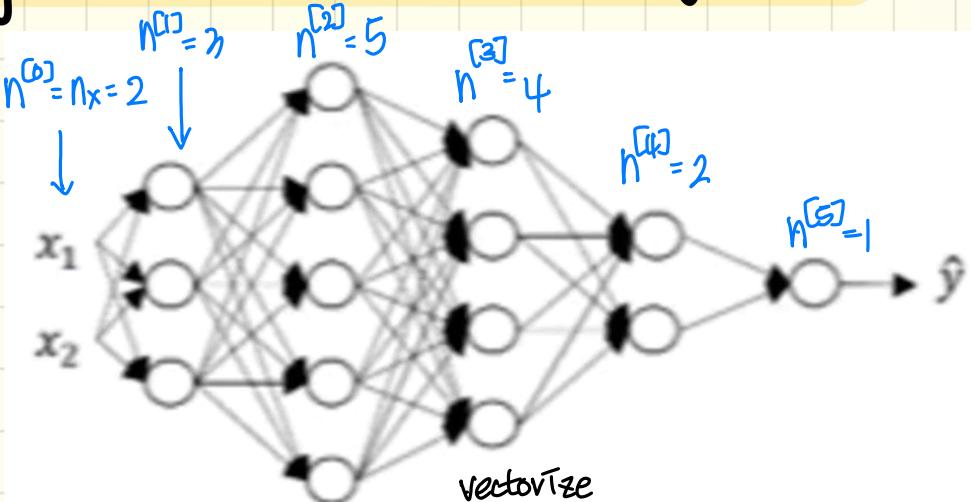
$$\begin{aligned} z^{[1]} &= w^{[1]} x + b^{[1]} \\ a^{[1]} &= g^{[1]}(z^{[1]}) \\ z^{[2]} &= w^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} &= g^{[2]}(z^{[2]}) \\ &\vdots \\ \Rightarrow z^{[l]} &= w^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

vectorizing

$$\begin{aligned} z^{[1]} &= W^{[1]} A^{[0]} + b^{[1]} \\ A^{[1]} &= g^{[1]}(z^{[1]}) \\ z^{[2]} &= W^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(z^{[2]}) \\ &\vdots \\ \hat{y} &= g(z^{[4]}) = A^{[4]} \end{aligned}$$

} for $l=1$ to 4

<Getting your Matrix Dimensions Right>

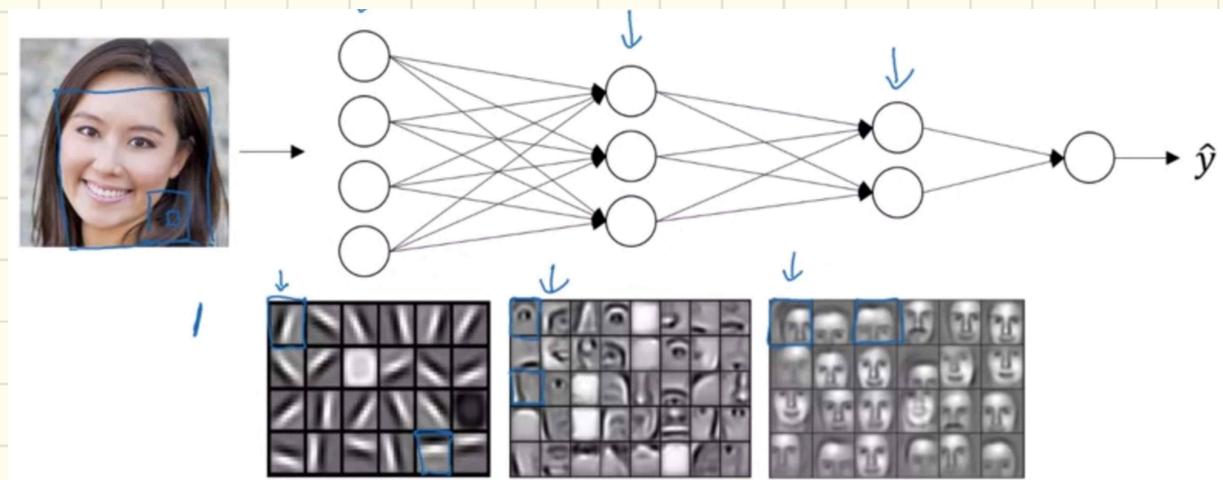


vectorize

$$\begin{aligned} z^{[1]} &= W^{[1]} x + b^{[1]} \Rightarrow z^{[l]} = (n^{[l]}, 1) \\ (3,1) \quad (3,2) \quad (2,1) \quad (3,1) & \quad W^{[l]} = (n^{[l]}, n^{[l+1]}) \\ z^{[2]} &= W^{[2]} a^{[1]} + b^{[2]} \\ (5,1) \quad (5,2) \quad (3,1) \quad (3,1) & \quad b^{[l]} = (n^{[l]}, 1) \\ dw^{[l]} &= (n^{[l]}, n^{[l+1]}) \\ db^{[l]} &= (n^{[l]}, 1) \end{aligned}$$

$$\begin{aligned} z^{[1]} &= W^{[1]} x + b^{[1]} \\ (3,m) \quad (3,2) \quad (2,m) \quad (3,1) & \quad \Rightarrow z^{[l]}, A^{[l]}, (n^{[l]}, m) \\ \quad & \quad dz^{[l]}, dA^{[l]}, (n^{[l]}, m) \end{aligned}$$

<Why Deep Representations?>



직관 1) 네트워크가 깊어질수록 더 많은 특징을 잡아낼 수 있다. 낮은 층에서는 간단한 특징을 찾고 깊은 층에서는 낮은 층에서 텅자란 간단한 것들을 통합해서 더 복잡한 특성을 찾아낼 수 있다.

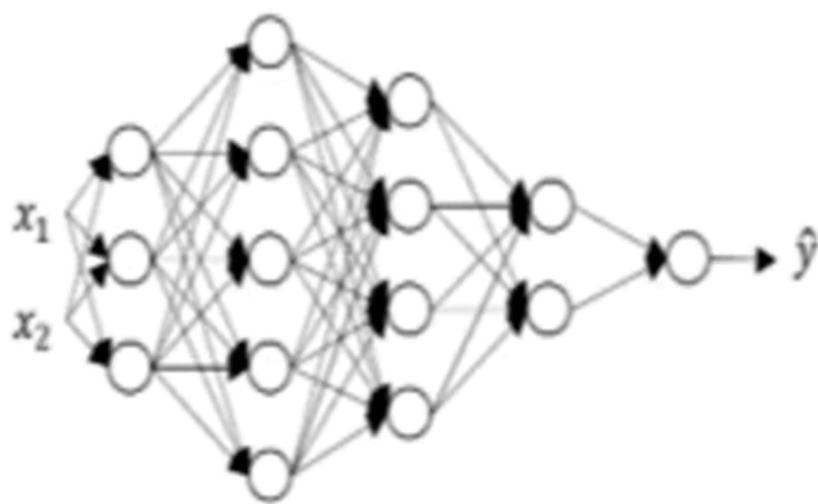
Circuit theory and deep learning

Informally: There are functions you can compute with a “small” L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

직관 2) 순환이론에 따르면. 얕은 네트워크보다 깊은 네트워크에서 더 저산하기 쉬운 행위가 있다.

<Building Blocks of Deep Neural Networks>

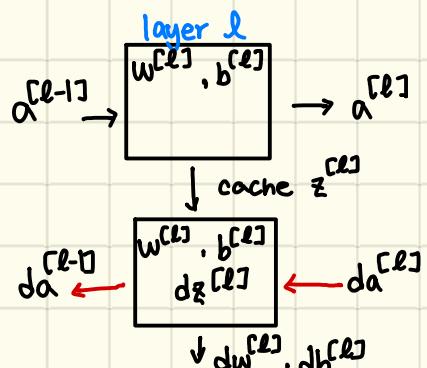
* Forward and Backward functions



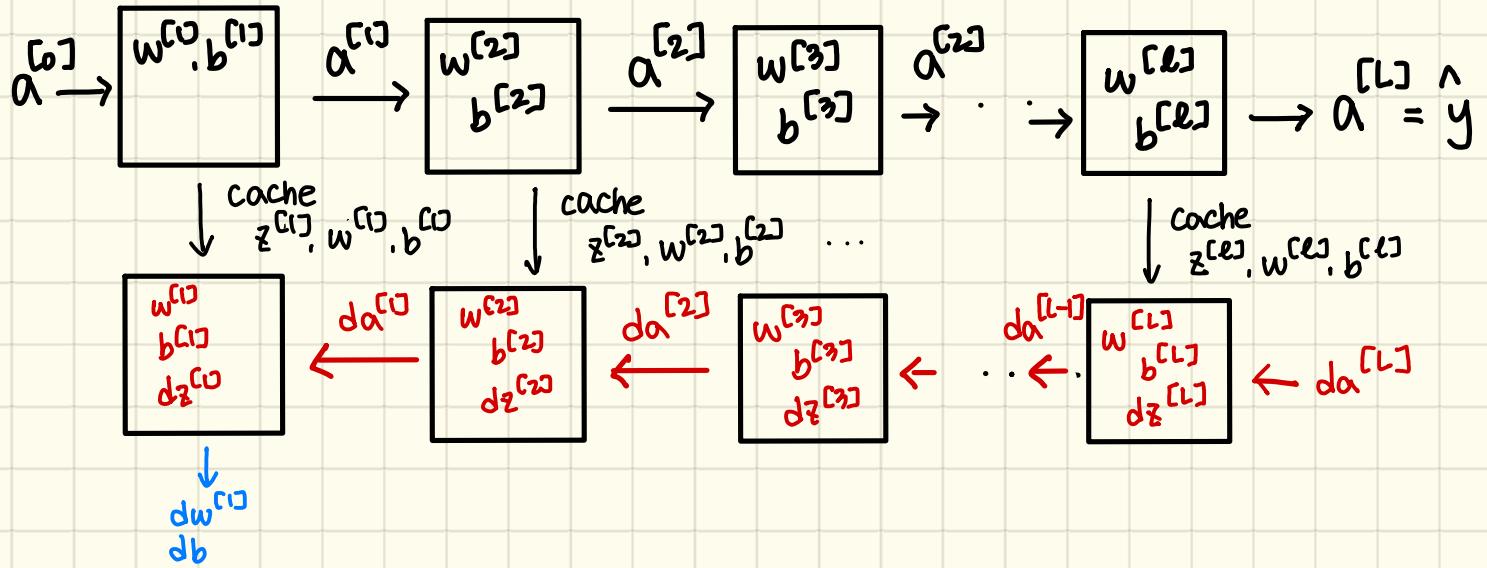
Layer l : $W^{[l]}, b^{[l]}$

Forward · Input $a^{[l-1]}$, Output $a^{[l]}$
cache $z^{[l]}$

Backward · Input $da^{[l]}$
Output $da^{[l-1]}, dw^{[l]}, db^{[l]}$



회상해보면,



<Forward and Backward Propagation>

* Forward propagation for layer l

- Input : $\alpha^{[l-1]}$
 - Output : $a^{[l]}$, cache ($z^{[l]}$) \oplus coding 관점에서 cache $w^{[l]}, b^{[l]}$
- $$z^{[l]} = w^{[l]} \cdot \alpha^{[l-1]} + b^{[l]} \quad \xrightarrow{\text{vectorize}} \quad z^{[l]} = W^{[l]} \cdot A^{[l-1]} + b^{[l]}$$
- $$a^{[l]} = g^{[l]}(z^{[l]}) \quad \quad \quad A^{[l]} = g^{[l]}(z^{[l]})$$

* Backward propagation for layer l

- Input : $da^{[l]}$
- Output : $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$dz^{[l]} = da^{[l]} * g^{[l]}'(z^{[l]})$$

$$dW^{[l]} = dz^{[l]} \cdot \alpha^{[l-1]T}$$

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

$$da^{[l]} = W^{[l]T} \cdot dz^{[l]}$$

$$\rightarrow dz^{[l]} = W^{[l+1]T} \cdot dz^{[l+1]} * g^{[l]}'(z^{[l]})$$

$$dz^{[l]} = dA^{[l]} * g^{[l]}'(z^{[l]})$$

$$dW^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$$

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

$$dA^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$$

< Parameters VS Hyperparameters >

* What are hyperparameters?

- parameters : $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots$
- Hyperparameters : learning rate α
 - # Iterations
 - # hidden layer L
 - # hidden units $n^{[1]}, n^{[2]}, \dots$
 - Choice of activation functions
 - ...