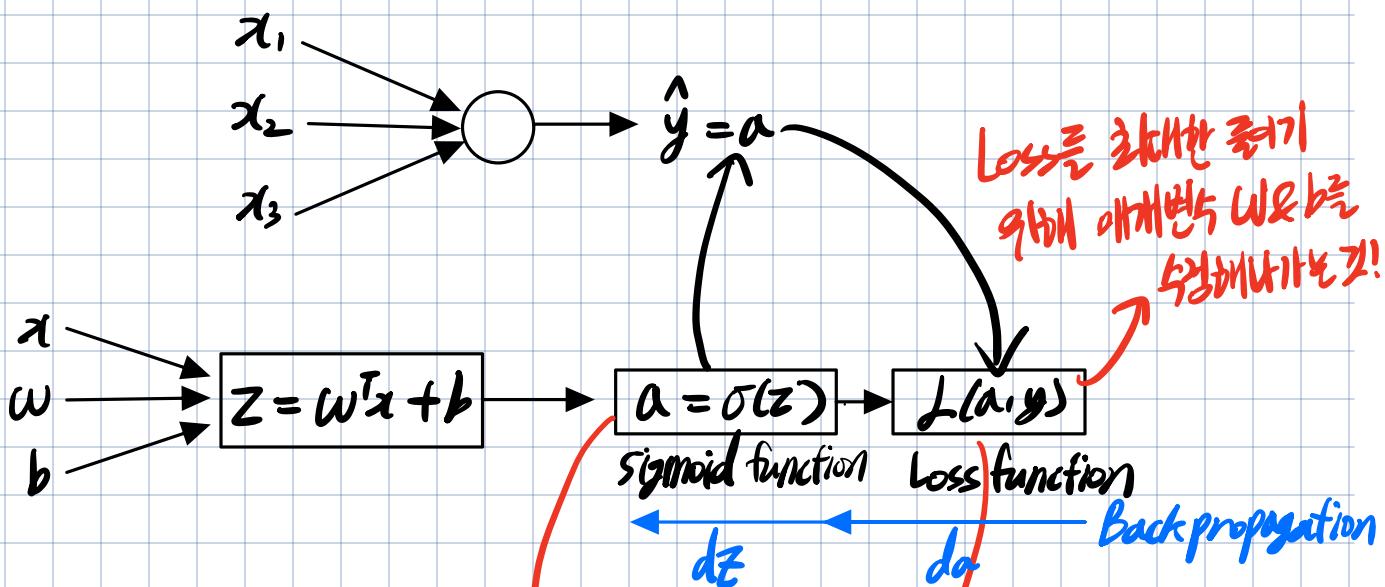


# Week 3. Shallow Neural Networks

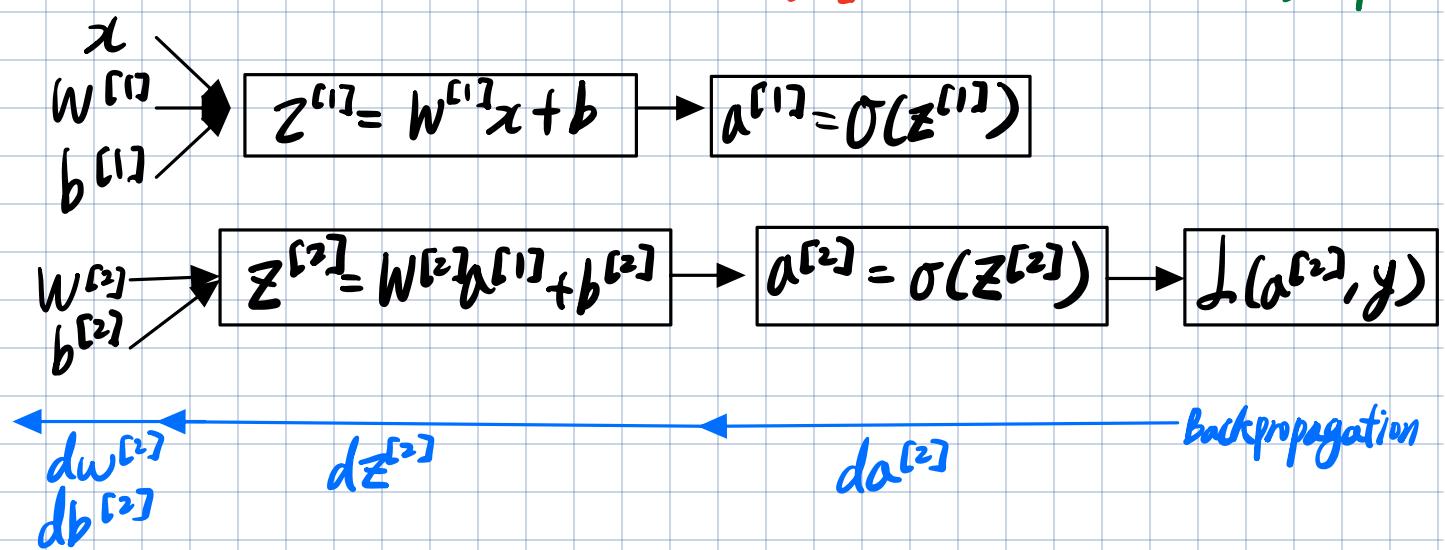
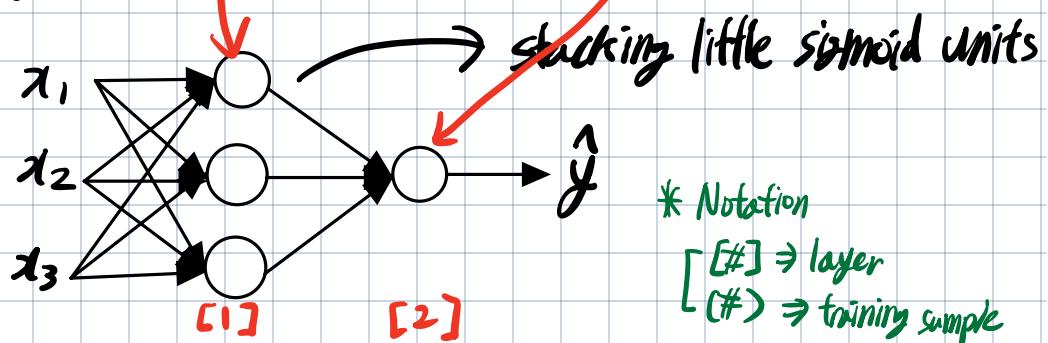
## 1. One hidden layer Neural Network

### • Neural Networks Overview

#### - Logistic Regression Model



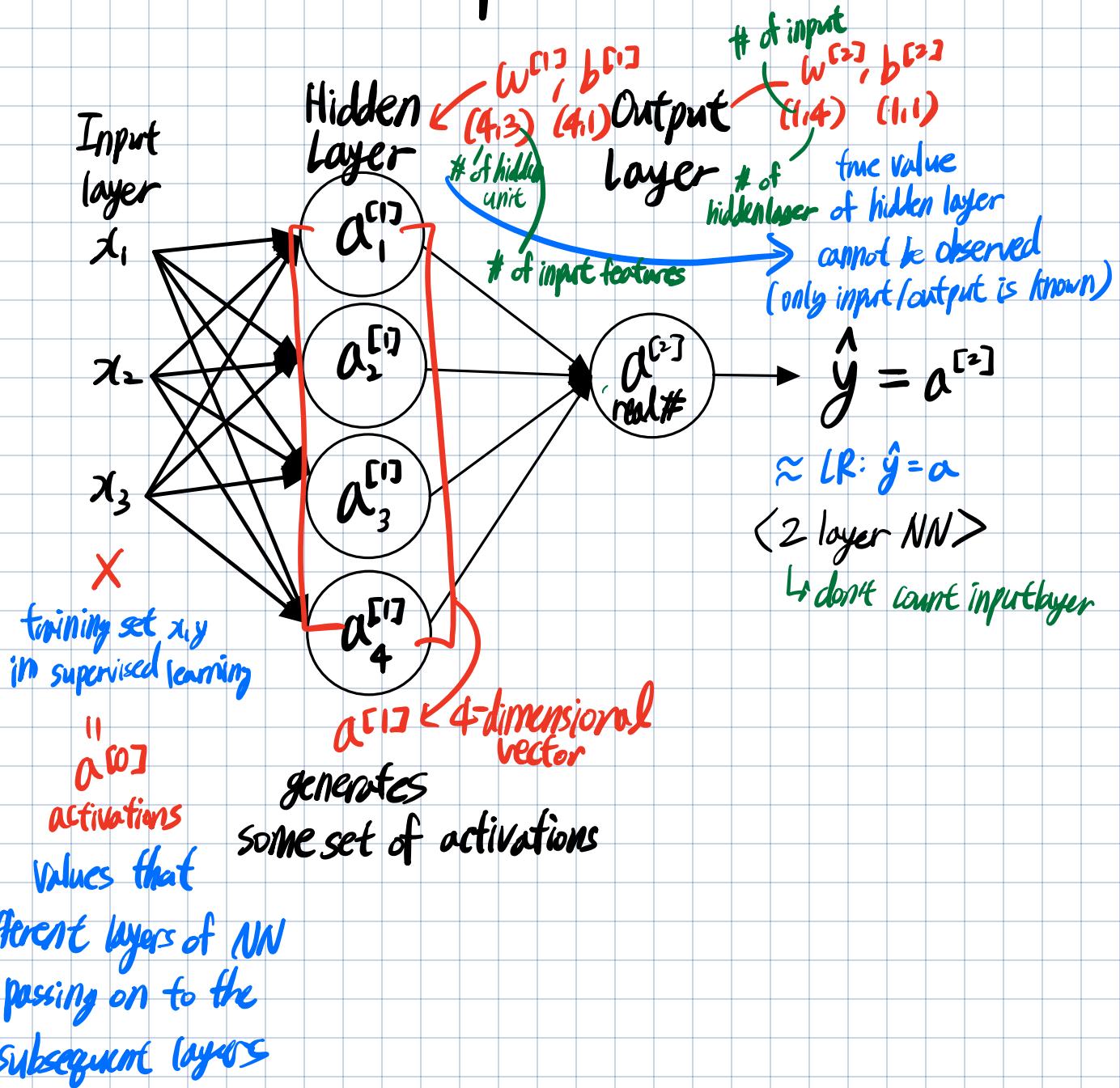
#### - Neural Network



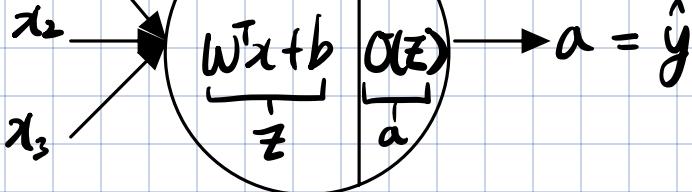
[ Logistic Regression: calculate  $z$  &  $\alpha$

[ Neural Network: calculate  $z$  &  $\alpha$  multiple times  
& finally compute loss function

## • Neural Network presentation



## • Computing a Neural Network's output.



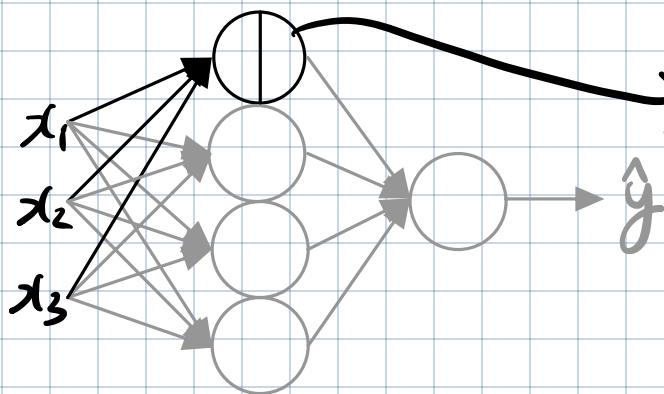
$z = W \cdot x + b$   
 $a = \sigma(z)$

In NN, this computation repeats

- 1st node in hidden layer

\* Notation

$A^m \Rightarrow m^{\text{th}} \text{ layer}$   
 $A_n \Rightarrow n^{\text{th}} \text{ node}$

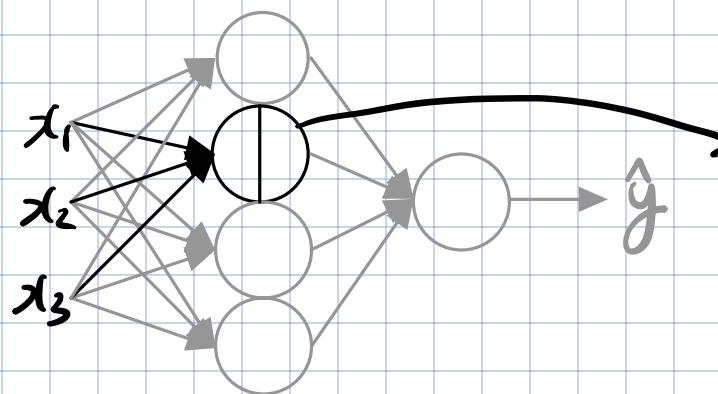


2 steps of computation

$$L: z_1^{[1]} = W_1^T x + b_1^{[1]}$$

$$R: a_1^{[1]} = \sigma(z_1^{[1]})$$

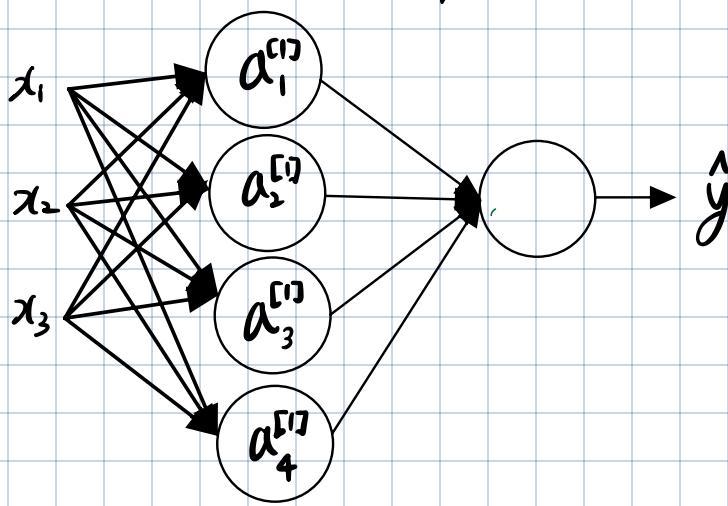
- 2nd node in hidden layer



$$L: z_2^{[1]} = W_2^T x + b_2^{[1]}$$

$$R: a_2^{[1]} = \sigma(z_2^{[1]})$$

- Vectorization of  $Z, W$



$$z_1^{[1]} = W_1^{[1]T} x + b_1^{[1]}, a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = W_2^{[1]T} x + b_2^{[1]}, a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = W_3^{[1]T} x + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = W_4^{[1]T} x + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]})$$

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

→ stack row vectors as a matrix

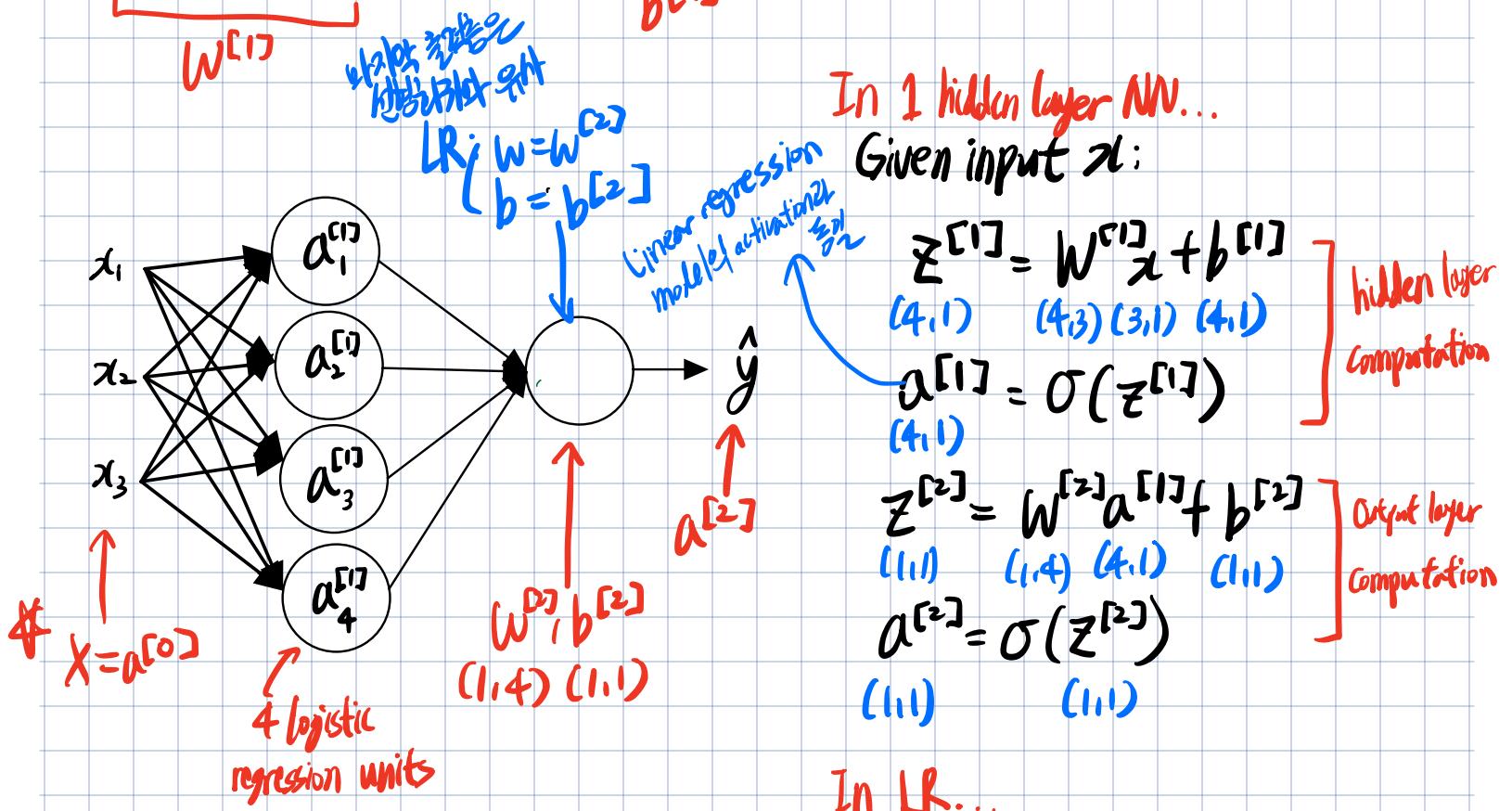
$$\left[ \begin{array}{c} -W_1^{[1]T} - \\ \vdots \end{array} \right] \times \left[ \begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right] + \left[ \begin{array}{c} b_1^{[1]} \\ \vdots \end{array} \right] = \left[ \begin{array}{c} W_1^{[1]T} x + b_1^{[1]} \\ \vdots \\ W_4^{[1]T} x + b_4^{[1]} \end{array} \right] = \left[ \begin{array}{c} z_1^{[1]} \\ \vdots \\ z_4^{[1]} \end{array} \right]$$

$$\begin{bmatrix} - & w_2^{[1]} & - \\ - & w_3^{[1]T} & - \\ - & w_4^{[1]T} & - \end{bmatrix} \begin{bmatrix} x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_2 \\ z_3 \\ z_4 \end{bmatrix}$$

$b^{[1]}$

4x3 matrix

$W^{[1]}$



## • Vectorizing across multiple examples

<One training sample>

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

< $m$  training samples>

$$\begin{aligned} x &\longrightarrow a^{[2]} = \hat{y} \\ x^{(1)} &\longrightarrow a^{[2](1)} = \hat{y}^{(1)} \\ x^{(2)} &\longrightarrow a^{[2](2)} = \hat{y}^{(2)} \\ &\vdots \\ x^{(m)} &\longrightarrow a^{[2](m)} = \hat{y}^{(m)} \end{aligned}$$

## [ Unvectorized ]

~~for~~  $i = 1$  to  $M$ ,  $x^i = a^{[0]i}$

$$\begin{aligned} z^{[1]}(i) &= w^{[1]} \lambda^{(i)} + b^{[1]} \\ \lambda^{[1]}(i) &= \sigma(-f^{[1]}(i)) \end{aligned}$$

$$x = 20 \text{ (km)}$$

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

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

가) training samples  
혹은 features 가 7212  
개다!

$A^{[2]}, Z^{[2]}$

政治小史

stinkim을 풍해

matrix로 나타낼 수 있음

# [Vectorized]

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}, \quad X \in \mathbb{R}^{n \times m}$$

↑ horizontally to  
from matrix

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

Add  $X$

different input features

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

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

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

stack vector  $x$

horizontally to form Matrix

different input features

$$Z^{[1]} = \begin{bmatrix} Z^{[1]}(1) & Z^{[1]}(2) & \dots & Z^{[1]}(m) \end{bmatrix}$$

activation at  
1st training sample  
of first hidden unit (node)

hidden units  
node

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

- Explanation for vectorized implementation

## - Justification for vectorized implementation

$$\mathcal{Z}^{[1]}(1) = W^{[1]}\chi^{(1)} + b^{[1]}$$

for simplification

$$\mathcal{Z}^{[1][2]} = W^{[1]}_{2L} \mathcal{Z}^{[2]} + b^{[1]}_0$$

$$z^{[1][3]} = w^{[1]} z^{[1^3]} + b^{[0]}$$

$$W[1] = \begin{bmatrix} \parallel \\ \parallel \\ \parallel \end{bmatrix}$$

$$\omega^{[1]} x^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\omega \begin{bmatrix} 1 \\ x \end{bmatrix}^{(2)} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$W^{[1]} x^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$W^{[1]} \begin{bmatrix} 1 & 1 & 1 & \dots \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots \\ 1 & 1 & 1 & \dots \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & z^{[1](1)} & z^{[1](2)} & z^{[1](3)} & \dots \\ 1 & b^{(1)} & b^{(2)} & b^{(3)} & \dots \end{bmatrix} = Z^{[1]}$$

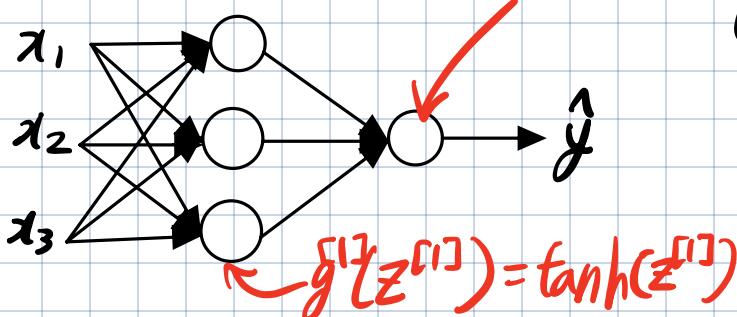
X  
training set

broadcasting

single train sample:  $Z^{[1](i)} = W^{[1]} x^{(i)}$  일 때, training sample인  $x^{(i)}$ 를 열로 쪼갤 수 있다면,  $Z^{[1]}$ 는  $Z$  vector가 열로 쪼갠 matrix로 만들어짐 +  $b^{(i)}$ 는 broadcasting을 통해 각각의 같은 정수.  
∴ forward propagation 1) training sample을 열로 쪼개서  $Z$ ,  $a$  형태로 열로 쪼갠다.

## • Activation function

### - Sigmoid function



Given  $z$ :

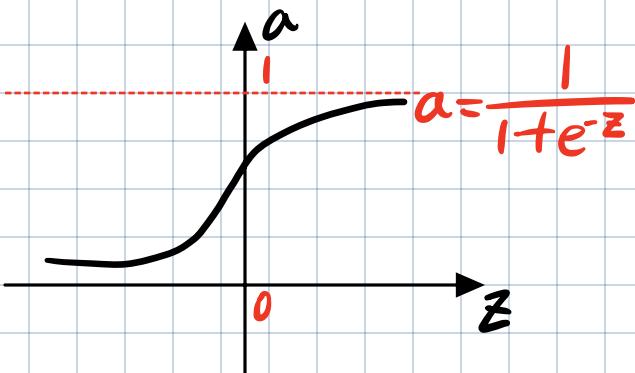
$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$a^{[1]} = \sigma(Z^{[1]}) \rightarrow g(Z^{[1]})$$

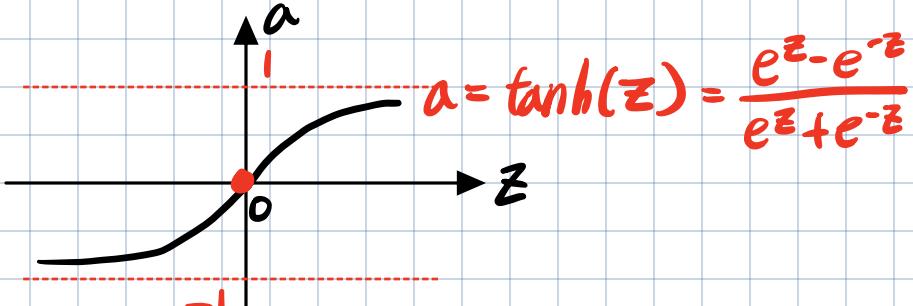
$$Z^{[2]} = W^{[2]} X + b^{[2]}$$

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

Sigmoid가 아닌 다른 활성화 함수가 올 수도 있다.



### - tanh function



almost always work better

than sigmoid in hidden layer

the mean of activation  $\approx 0$

$\Rightarrow$  편향의 초기값:  $0.5 \rightarrow 0$

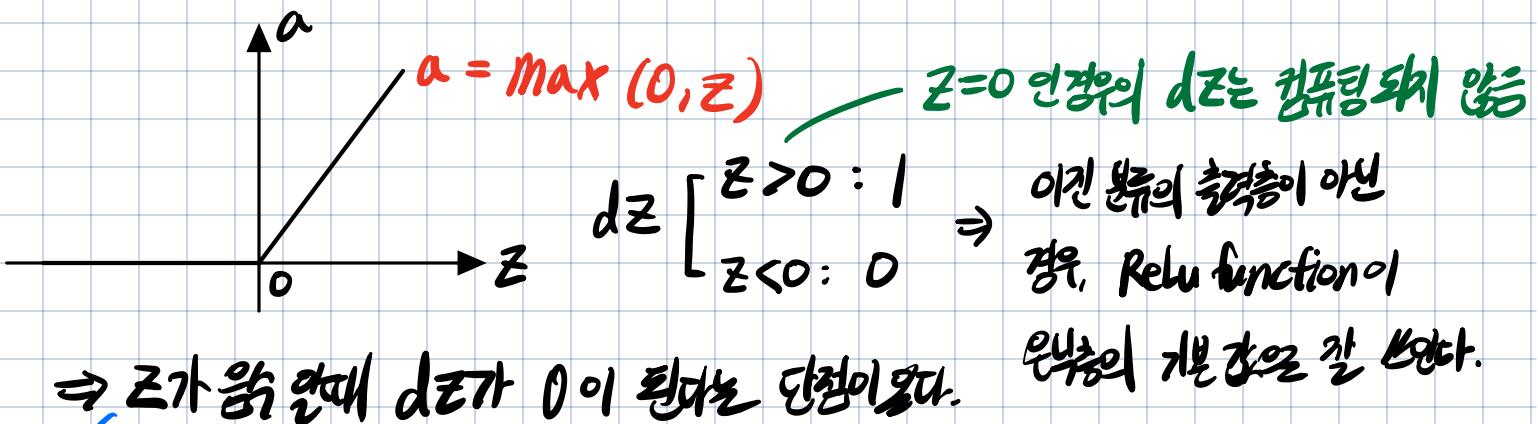
다중 출력의 확률이 합이 1

→ but output layer에서는 사용X

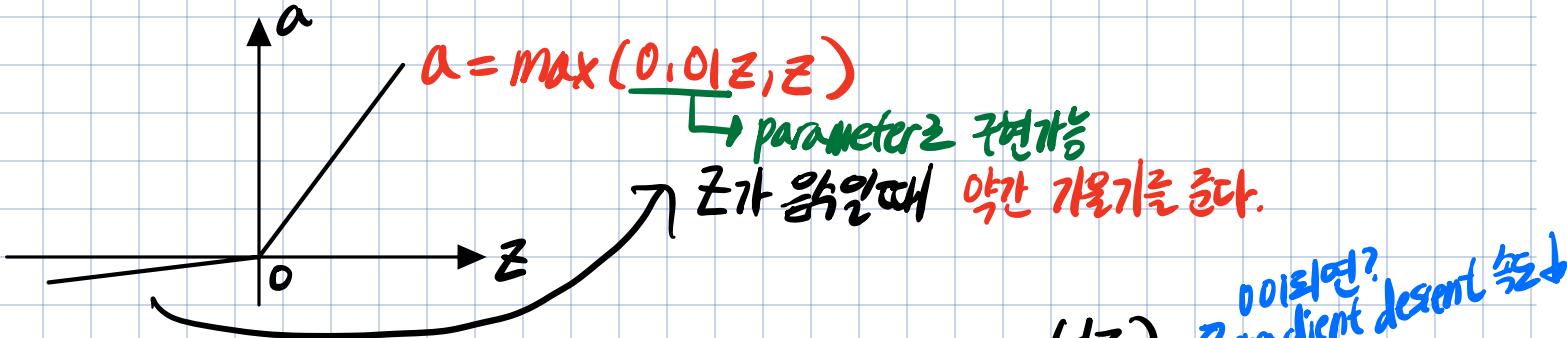
∴  $y \in \{0, 1\}$  이면  $-1 \leq \hat{y} \leq 1$  보다  $0 \leq \hat{y} \leq 1$  을 출력하는 것의 더 좋기 때문

∴ Sigmoid activation을 쓰는 한 가지 예 : **binary classification**

- Sigmoid & tanh 모두  $z$ 가 긍정적이나 부정적일 때 derivative(slope)는 0이거나 가까워진다  $\Rightarrow$  slow down gradient descent
- ReLU (Rectified Linear unit) function



- Leaky ReLU (ReLU보다 성능이 좋지만 잘 안쓰임)

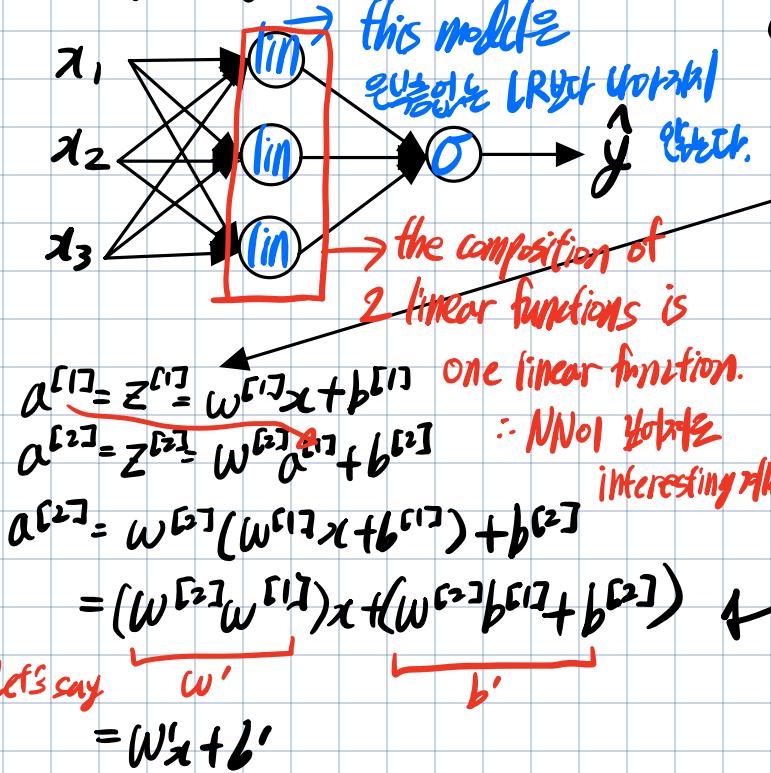


- PReLU (Leaky ReLU 훈련 대체물)의  $z$ 에 대해 기울기가 0과 매우 다른  $(dz)$
- NN의 학습 속도 ↑

- 신경망을 고를 때에는 여러 option을 시도해보고 선택!

○○이 되면?  
gradient descent 속도↑

## • Why do you need non-linear activation functions?



## < Forward Propagation >

Given  $x$ :

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

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

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

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

linear activation function  
identity activation function

Let's say

interesting  $\neq f(x) \Rightarrow y$ 는  $x$ 에 대한 선형 함수가 아님

If you use linear/identity activation functions,  
 $\Rightarrow$  the NN will output linear function of the  
input

↳ 1000 hidden layer가 있는 deep NN은 다를 때  
linear function을 쓰거나 activation  
function을 쓰지 않는다면 그다지  
많은 NN은 훈련이나 학습을 계산하기  
때문에 익숙이 없는 경우 다른 있다.

## - $g(z) = z$ regression problem

대부분 machine learning은 ( $z, y = \text{실수 값}$ )

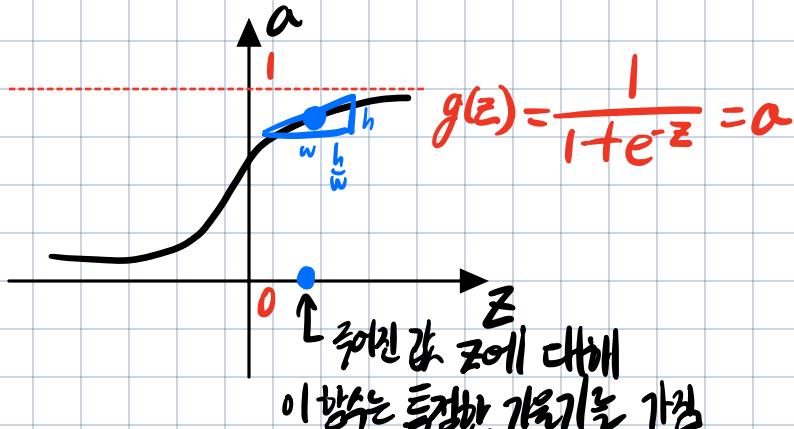
할 때 or 대체로 output layer는 NN

e.g. house pricing: 익숙한 linear ok,  
and yet often need to use Relu instead

## • Derivatives of activation functions

- NN의 back propagation을 계산하기 위해서는 각각의 학습률 계산을 구해야 함

- Sigmoid activation function



$$\frac{d}{dz} g(z)$$

= slope of  $g(x)$  of  $x$   
 $= g'(z)$

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

$$= g(z)(1-g(z))$$

$$= a(1-a)$$

e.g. when  $z$  is large

$$z=10, g(z) \approx 1$$

$$\frac{d}{dz}g(z) \approx 1(1-1) \approx 0$$

when  $z$  is small

$$z=-10, g(z) \approx 0$$

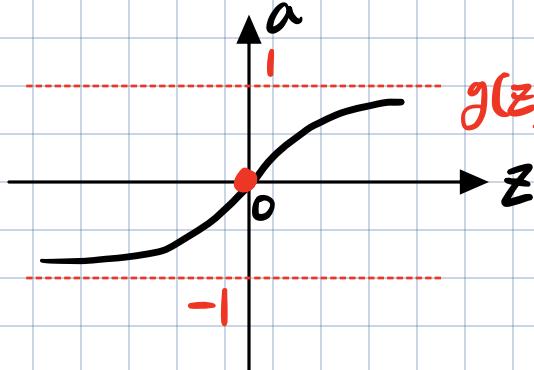
$$\frac{d}{dz}g(z) \approx 0 \cdot (1-0) \approx 0$$

when  $z$  is zero

$$z=0, g(z)=\frac{1}{2}$$

$$\frac{d}{dz}g(z)=\frac{1}{2}(1-\frac{1}{2})=\frac{1}{4}$$

### - Tanh activation function



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

$$\frac{d}{dz}g(z) = \text{slope of } g(x) \text{ at } x = g'(z) = 1 - (\tanh(z))^2 = 1 - a^2$$

e.g. when  $z$  is large

$$z=10, g(z) \approx 1$$

$$g'(z) = 1 - 1^2 \approx 0$$

when  $z$  is small

$$z=-10, g(z) \approx -1$$

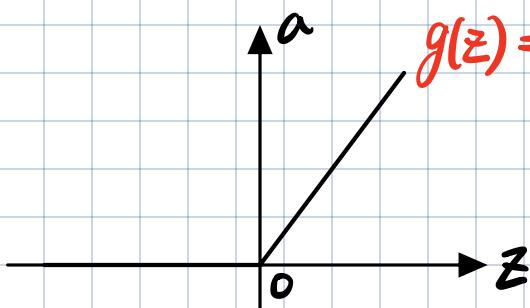
$$g'(z) \approx 0$$

when  $z$  is zero

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

$$g'(z)=1$$

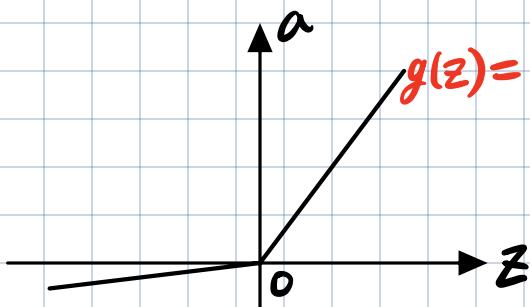
### - ReLU and Leaky ReLU



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

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

ReLU의 미분값은  
(undefined if  $z=0$ )



$$g(z) = \max(0, 0.01z, z)$$

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

Leaky ReLU의 미분값은  
(undefined if  $z=0$ )

## • Gradient descent for Neural Networks

- Equations for implementing Gradient Descent  $\Rightarrow$  back propagation
- Parameters of single layer Neural Network :

$$\begin{array}{cccc} w^{[1]} & b^{[1]} & w^{[2]} & b^{[2]} \\ (n^{[1]}, n^{[0]}) & (n^{[1]}, 1) & (n^{[2]}, n^{[1]}) & (n^{[2]}, 1) \end{array}$$

$n_x = n^{[0]}, n^{[1]}, n^{[2]} = 1 \Rightarrow$  Single hidden unit  
 (input) (hidden) (output)

- Cost function:  $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$
- let's say  
 you're doing  
 binary classification

$a^{[2]}$  true

- Use Gradient descent to train algorithm

① Initialize parameters randomly other than 0

② Repeat Gradient descent to compute predicts ( $\hat{y}^{[i]}, i=1 \sim m$ )

③ Compute derivatives ( $w^{[i]}, b^{[i]}$  of the cost function of  $\sum \ell$ )

$$dw^{[i]} = \frac{dJ}{dw}, db^{[i]} = \frac{dJ}{db}, \dots$$

$$w^{[1]} := w^{[1]} - \alpha dw^{[1]} \quad \dots \text{반복되는 계산을 위한 저작법}$$

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

how to compute partial derivative term?

## • Formulas for computing derivatives

- Forward propagation:

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

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

activation function

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

- Back propagation

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

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

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

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

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

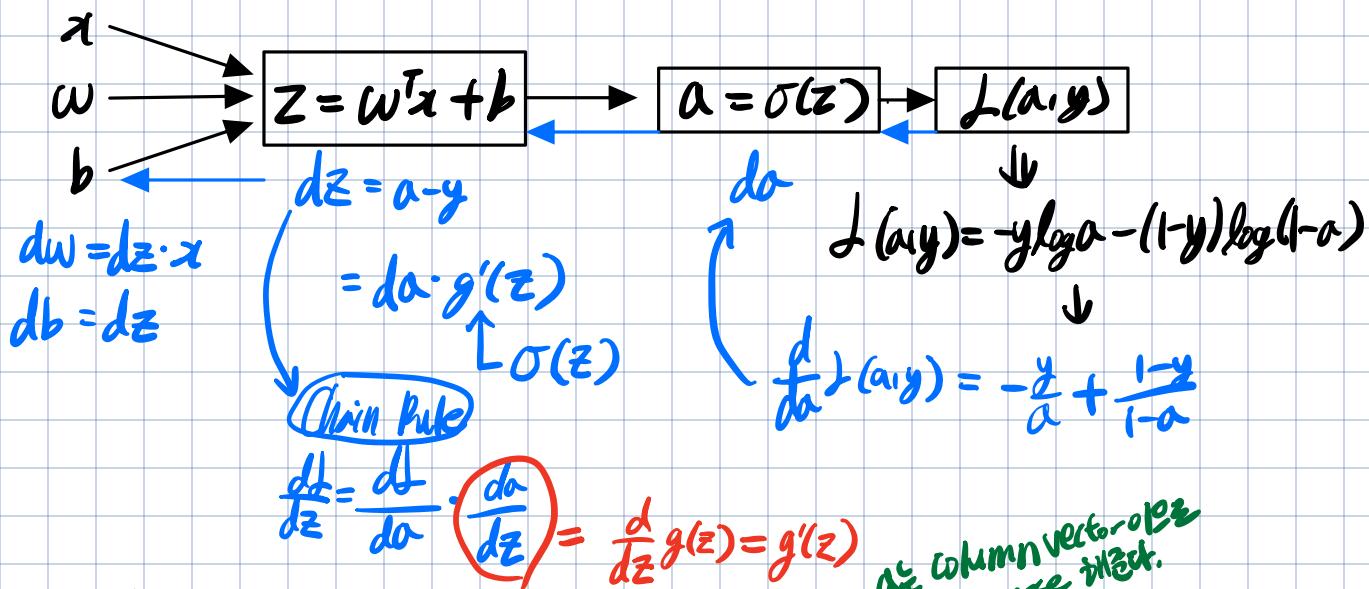
$(n^{[1]}, m)$  element-wise product  $\rightarrow$  행렬끼리의 곱셈  $(n^{[1]}, m)$  matrix

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

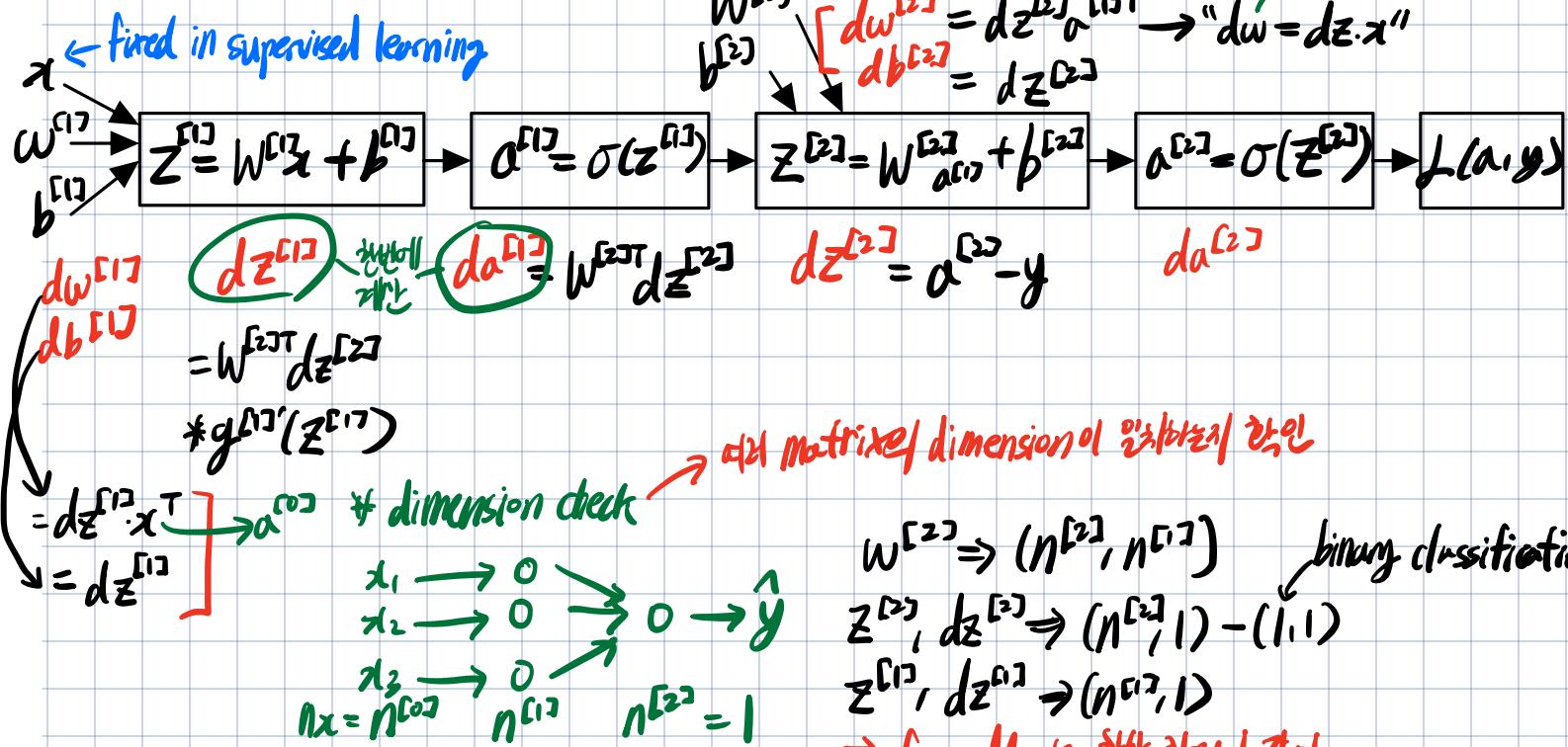
$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True})$$

$\uparrow$  np.sum은 축별 합을  
reduce로 재배열하면  
안쓰는 듯.

## • Back propagation intuition



— 온전한 기계학적 학습의 단계를 전경망에서 거친다거나 한다.



## - Summary of gradient descent

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

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

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

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

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

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

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

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

## Vectorized implementation

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

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T} \quad \text{truth} \quad \text{각 샘플 } (n^{[2]}, n^{[1]}) \text{의 } \frac{1}{m} \text{ 평균 차이 } X$$

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

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

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

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

## Vectorized implementation (multiple training samples)

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

stack Z's columns

$$Z^{[1]} = \begin{bmatrix} | & | & | \\ Z^{[1](1)} & Z^{[1](2)} & \dots & Z^{[1](n)} \\ | & | & | \end{bmatrix}$$

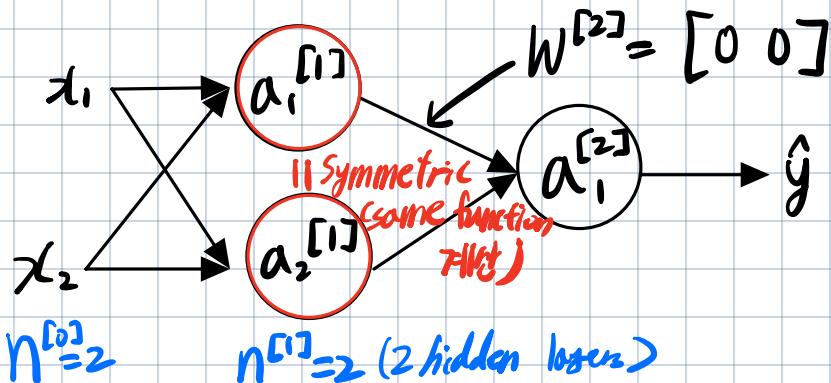
$$a^{[1]} = g^{[1]}(Z^{[1]})$$

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

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

### • Random Initialization

- What happens if you initialize weights to zero?



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

사실  $b$ 는 0으로 해도 ok

$\Rightarrow$  어떤 상수의 경우에도  $a_1^{[1]} = a_2^{[1]}$ , 두 개의 헬름이 같은 것으로 계산된다.

backpropagation 규칙은  $dZ_i^{[l]} = dZ_2^{[l]}$  두 유닛이 같은 값으로 초기화

→ 가증치의 불쾌감이 항상 같음.

let's say  $d\omega = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$  every row same value  $\xrightarrow{\text{weight update}} w^{[1]} = \omega^{[1]} - \alpha d\omega \rightarrow w^{[1]} = \begin{bmatrix} \dots \\ \dots \end{bmatrix}$

첫 번째 // 열이 두 번째  
열과 같아진다.

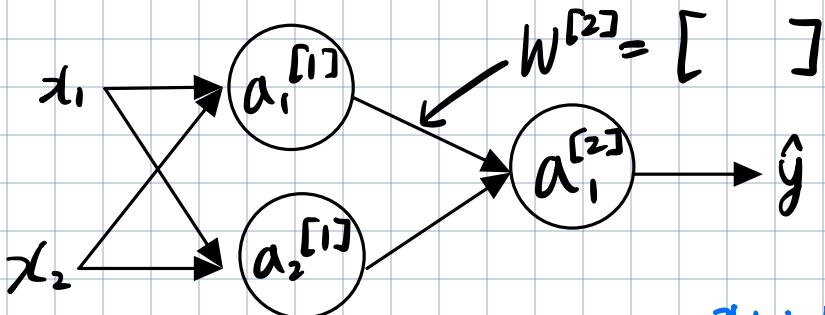
∴ W값은 모두 0으로 초기화  $\rightarrow$  두 유닛이 같은 학습률 계산하는 것으로 시작

→ 두 익우식이 축적우보에 항상 같은 역할을 주는 것이 반복됨

∴ 두 유도형상 같은 항수를 가지게 되므로 맥유보이 하나인 꼴, 다른 항수는

제작하기 위한 각각 다른 유닛이 필요하기 때문에 반드시 일련의 조건을 만족해야 한다.

### - Random initialization



→ 광장화 작은 임의의 수를 허락으로 만들기 위해

①  $w_{1,1} = np.random.randn(2,2) * 0.01$  매개 변수 초기화의 NN에서는 0.01이 일반화를 선택할 수 있음

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

$\hookrightarrow b$  does not have

## Symmetry breaking problem

$\therefore$  W/T randomly set  $\Sigma(\frac{1}{2}, \frac{1}{2})$  퍼센트

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

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

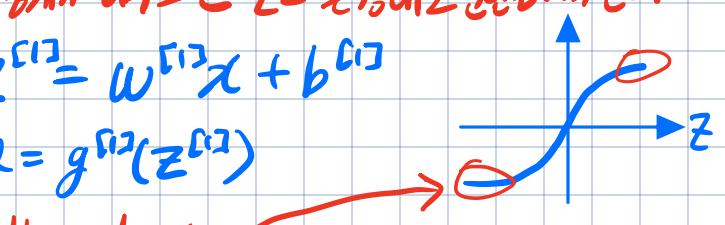
↳ 가중치의 초기값은 매우 작게 설정해야 함

## tanh / sigmoid activation

가장치가 너무 크면 군도 콘 상태로 혼연을 시작한다.

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

$$\alpha = q^{(1)}(\bar{z}^{(1)})$$



Small gradient →  
→ slow learning