

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



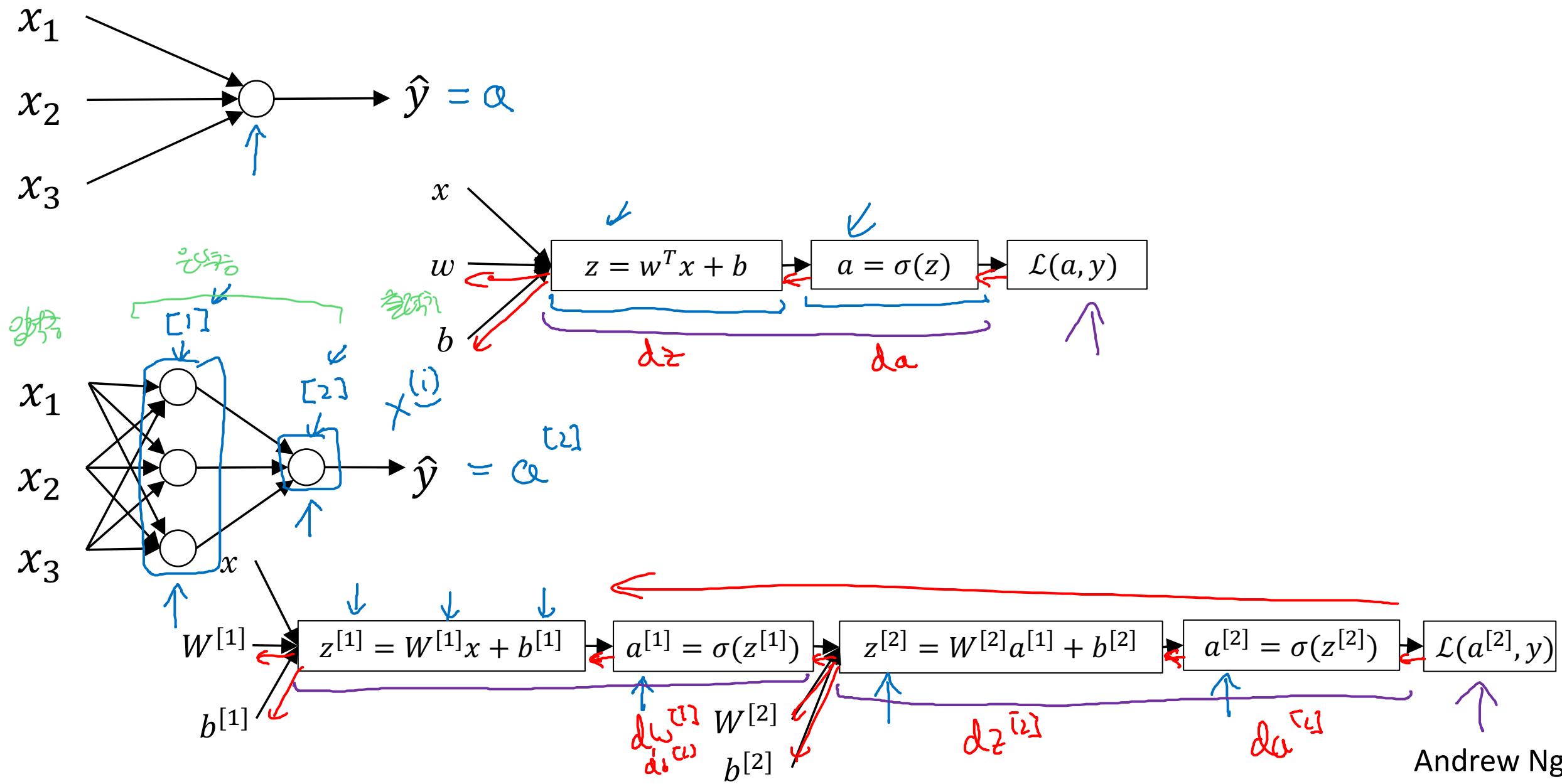
deeplearning.ai

One hidden layer  
Neural Network

---

Neural Networks  
Overview

# What is a Neural Network?





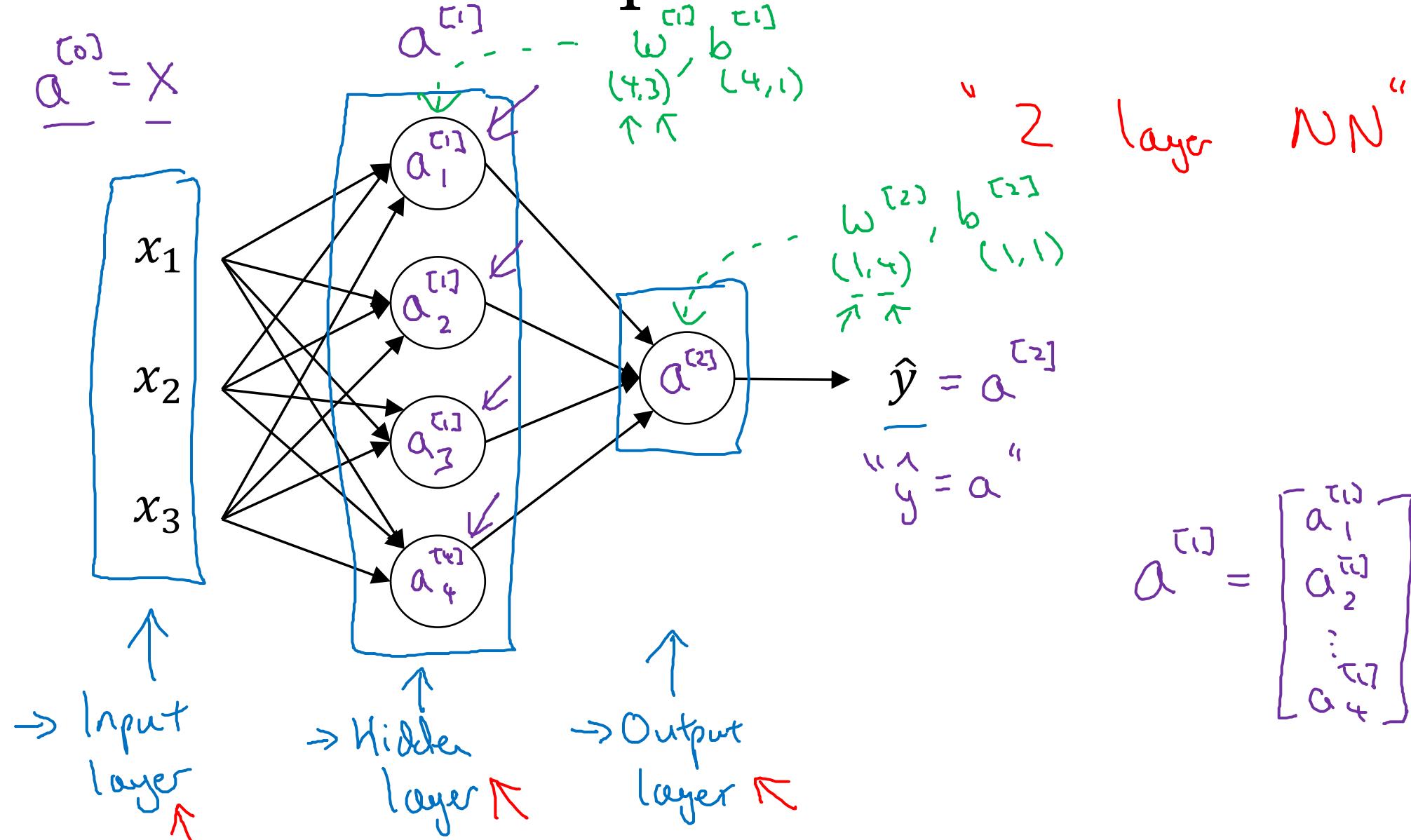
deeplearning.ai

One hidden layer  
Neural Network

---

Neural Network  
Representation

# Neural Network Representation



$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_4^{(1)} \end{bmatrix}$$



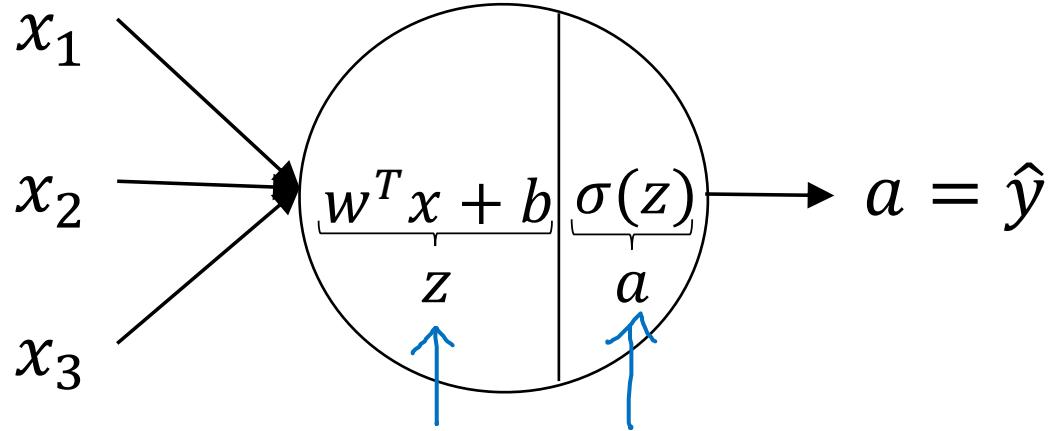
deeplearning.ai

# One hidden layer Neural Network

---

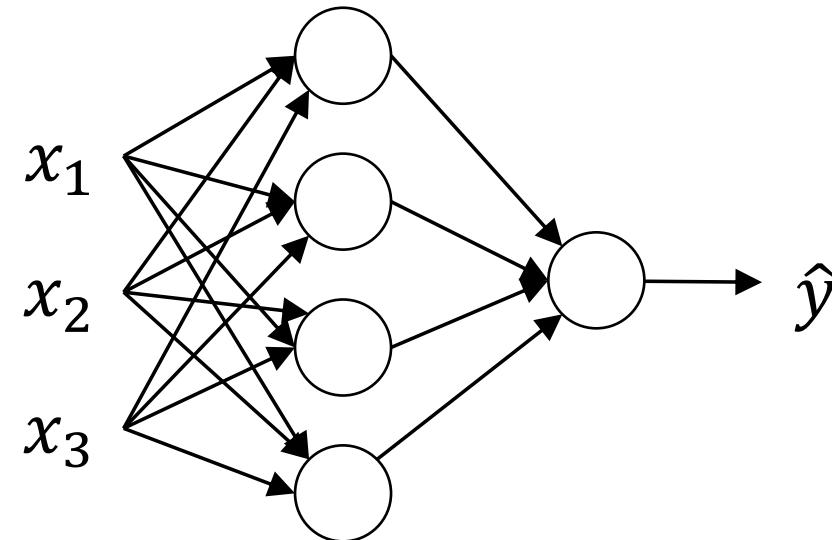
## Computing a Neural Network's Output

# Neural Network Representation

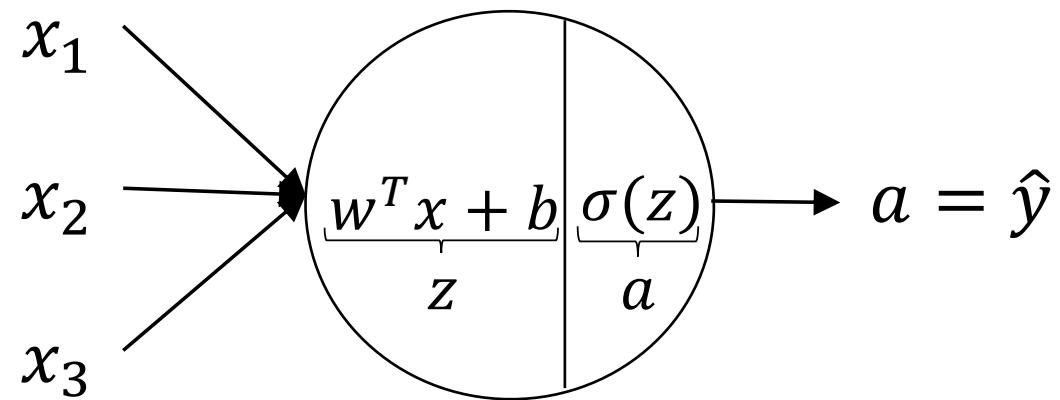


$$z = w^T x + b$$

$$a = \sigma(z)$$

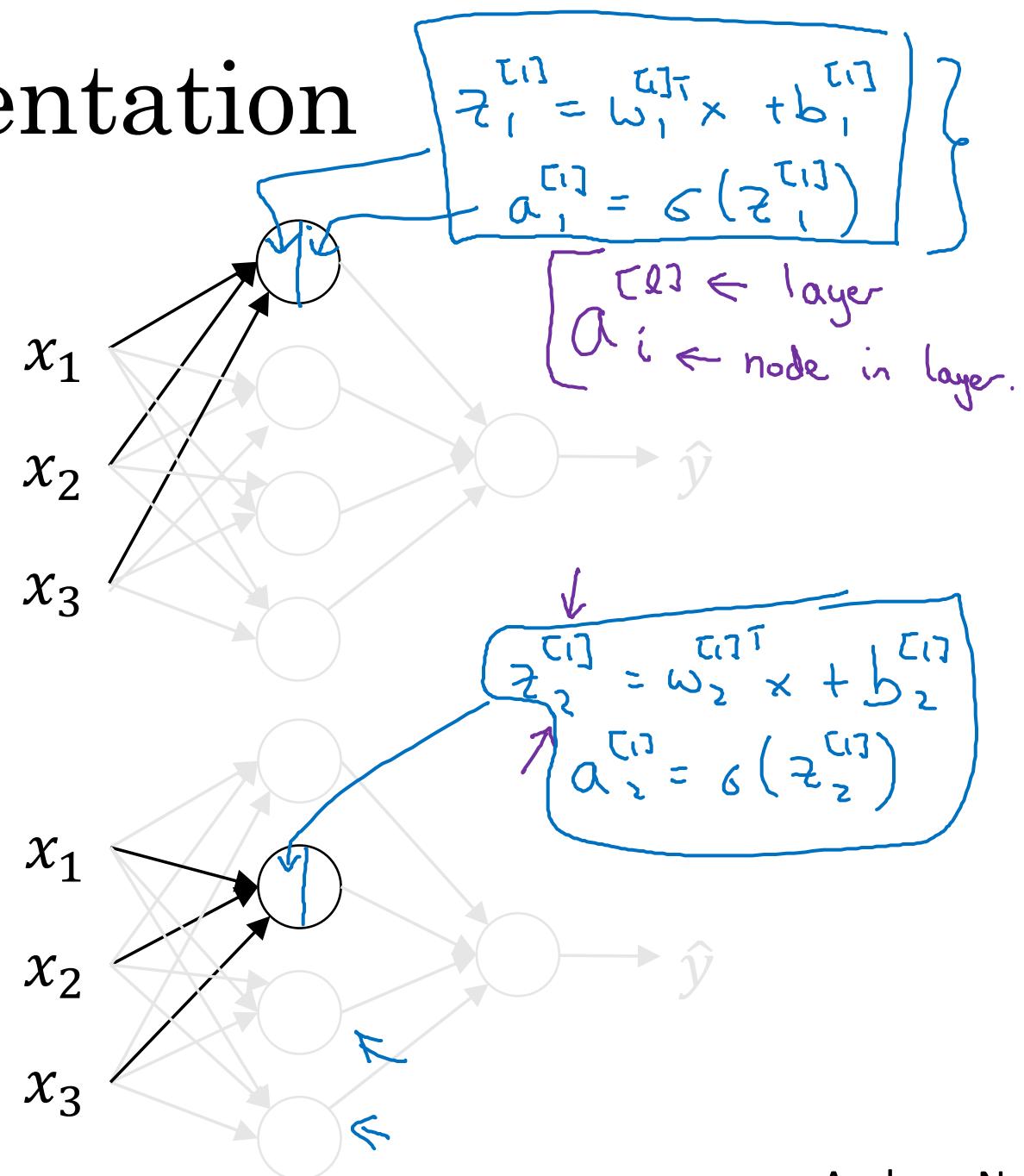


# Neural Network Representation

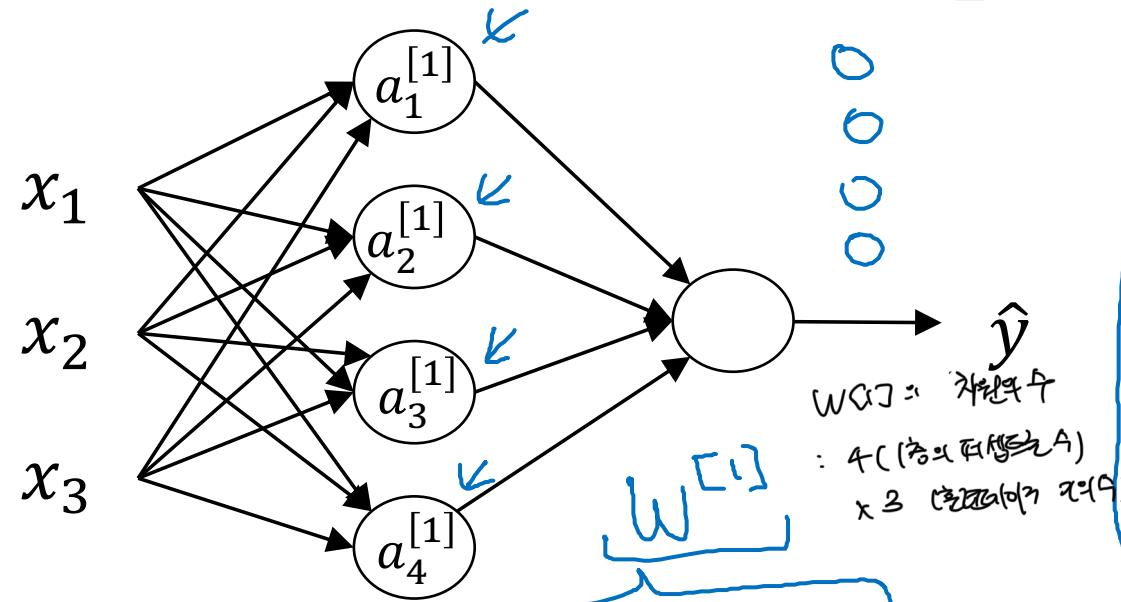


$$z = w^T x + b$$

$$a = \sigma(z)$$



# Neural Network Representation



$w(x) = \text{차원축소} : 4 \times 3 \rightarrow 4$

$$\rightarrow z^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$

$$\rightarrow a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = G(z^{[1]}) \quad \text{where } G(z) = \sigma(z)$$

$(\omega_1^{[1]T} x + b_1^{[1]}) \sigma(z_1^{[1]})$

$(\omega_2^{[1]T} x + b_2^{[1]}) \sigma(z_2^{[1]})$

$(\omega_3^{[1]T} x + b_3^{[1]}) \sigma(z_3^{[1]})$

$(\omega_4^{[1]T} x + b_4^{[1]}) \sigma(z_4^{[1]})$

$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$

$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$

$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}$

$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}$

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

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

$a_3^{[1]} = \sigma(z_3^{[1]})$

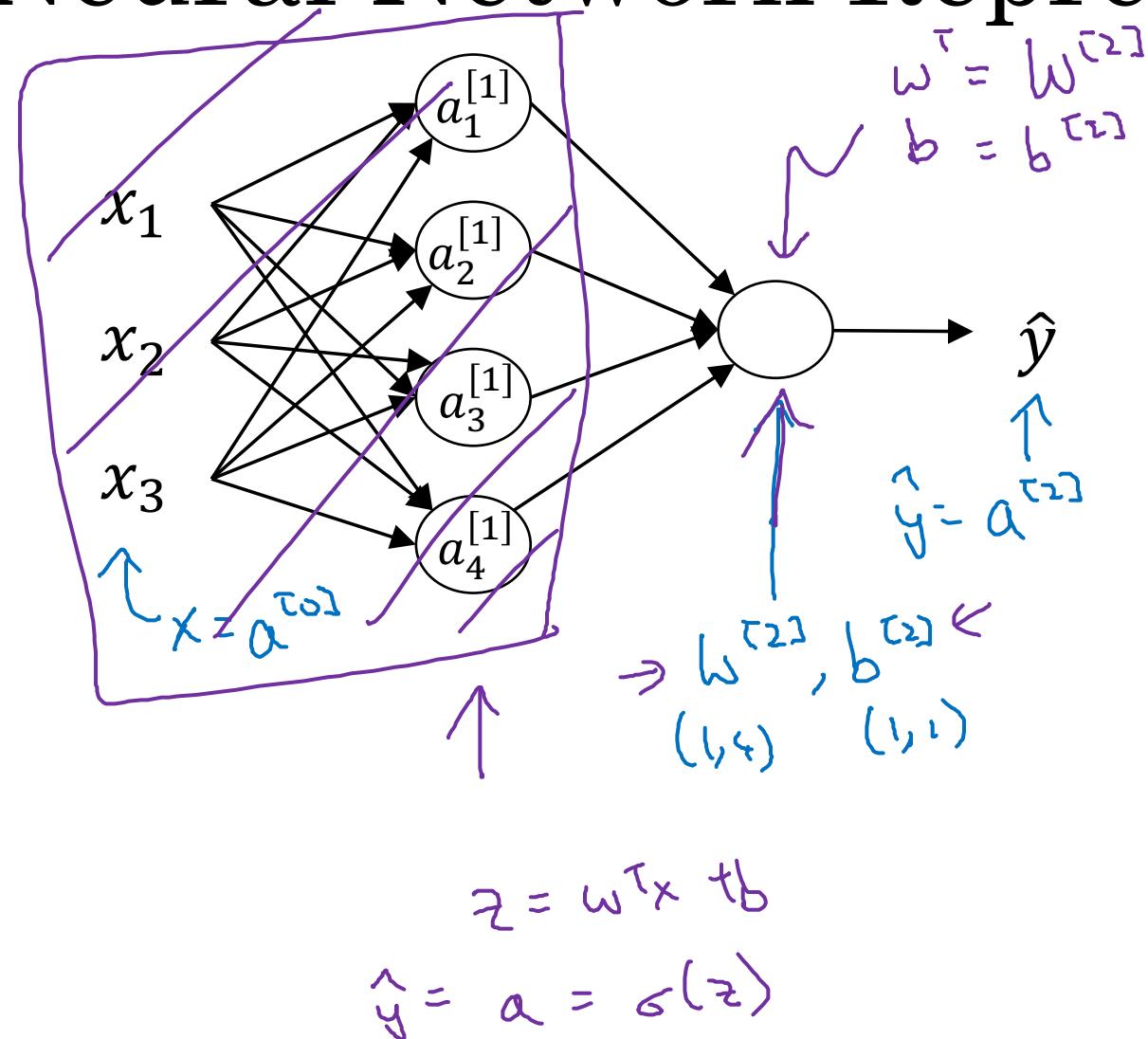
$a_4^{[1]} = \sigma(z_4^{[1]})$

$$= \begin{bmatrix} \rightarrow w_1^{[1]T} x + b_1^{[1]} \\ \rightarrow w_2^{[1]T} x + b_2^{[1]} \\ \rightarrow w_3^{[1]T} x + b_3^{[1]} \\ \rightarrow w_4^{[1]T} x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$\text{여기 차원축소 : } 4 \times 3 \rightarrow 4$

$\times 1 \text{ 차원축소는 } 4 \times 1 \text{ 차원축소는 } 1 \times 1$

# Neural Network Representation learning



Given input  $x$ :

$$\rightarrow z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

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

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

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



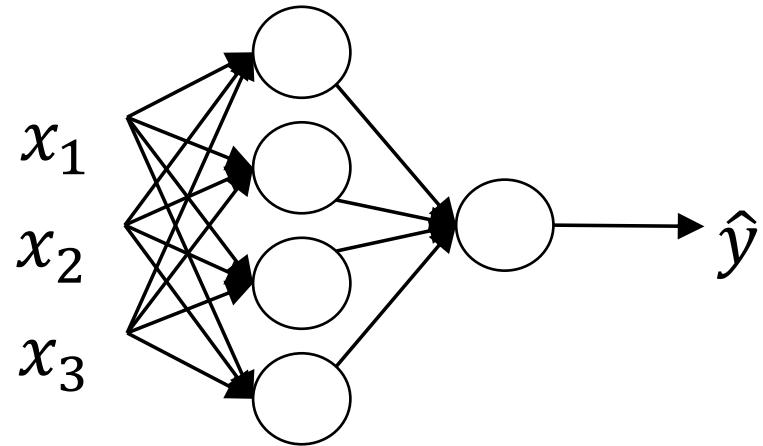
deeplearning.ai

# One hidden layer Neural Network

---

Vectorizing across  
multiple examples

# Vectorizing across multiple examples



$x \rightarrow a^{[2]} = y$   
 $x^{(1)} \rightarrow a^{[2](1)} = y^{(1)}$   
 $x^{(2)} \rightarrow a^{[2](2)} = y^{(2)}$   
 $\vdots$   
 $x^{(n)} \rightarrow a^{[2](m)} = y^{(m)}$

$a^{[2](i)}$  example i  
layer 2

$\left\{ \begin{array}{l} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right.$

for  $i = 1$  to  $m$ ,

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

# Vectorizing across multiple examples

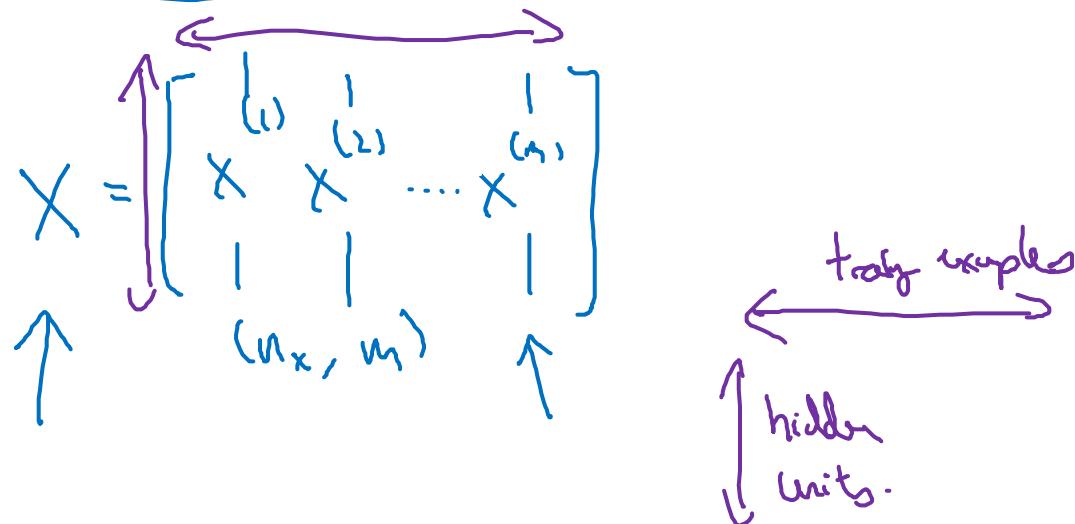
for  $i = 1$  to  $m$ :

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

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

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

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

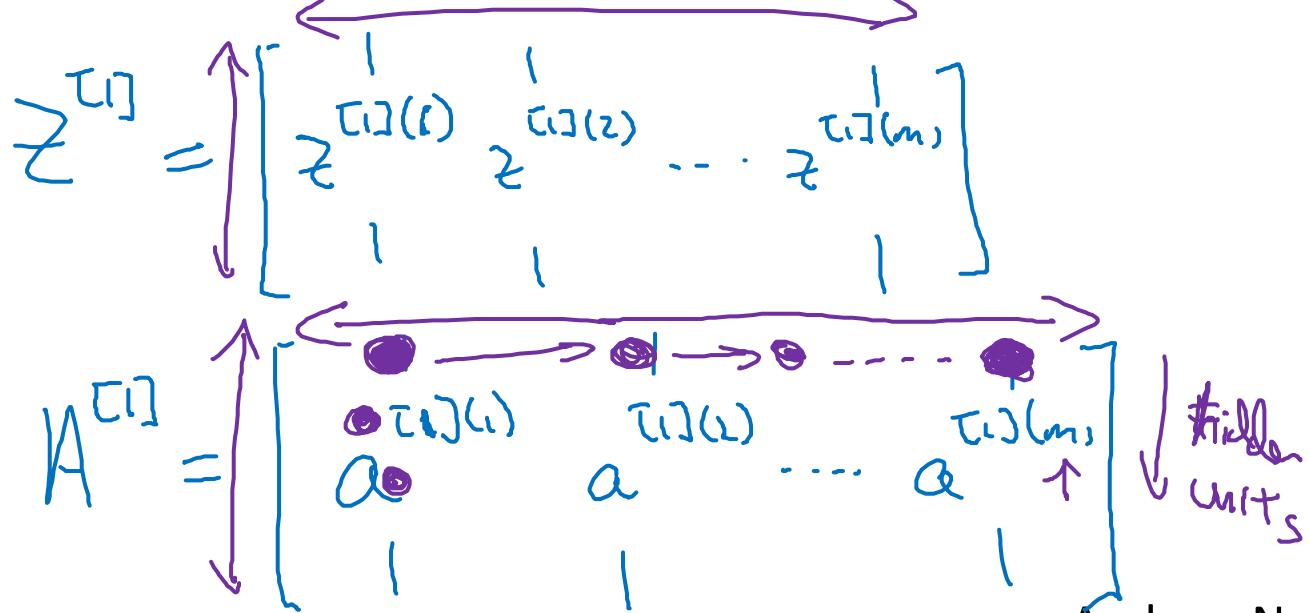


$$z^{[1]} = W^{[1]}X + b^{[1]}$$

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

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

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





deeplearning.ai

# One hidden layer Neural Network

---

## Explanation for vectorized implementation

# Justification for vectorized implementation

$$z^{(i)(1)} = w^{(1)} x^{(1)} + b^{(1)},$$

$$z^{(i)(2)} = w^{(2)} x^{(2)} + b^{(2)},$$

$$z^{(i)(3)} = w^{(3)} x^{(3)} + b^{(3)}$$

$$w^{(1)} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

$$w^{(1)} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

제한 조건: \$w^{(1)}\$의 행렬 \$x^{(1)}\$의 열 \$m\$  
제한 조건: \$w^{(1)}\$의 행렬 \$x^{(1)}\$의 차원은 \$m\$  
 $\Rightarrow x^{(1)}$는 $m \times 1$인 ($3 \times 1$) 행렬$

$$w^{(2)} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{(3)} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

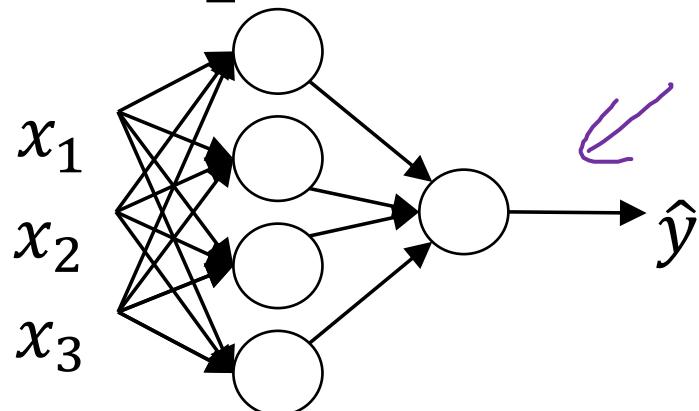
$$w^{(1)} \left[ \begin{array}{c|c|c} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{array} \right] = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} z^{(1)(1)} \\ z^{(1)(2)} \\ z^{(1)(3)} \\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

$$z^{(i)} = w^{(i)} X + b^{(i)}$$

$\Sigma$

$$w^{(i)} x^{(i)} = z^{(i)(i)}$$

# Recap of vectorizing across multiple examples



$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}$$

A vertical stack of input vectors  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ . A purple arrow points from the text "vectorizing across multiple examples" down to this equation.

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

A vertical stack of hidden layer activations  $a^{[1](1)}, a^{[1](2)}, \dots, a^{[1](m)}$ . A purple arrow points from the text "vectorizing across multiple examples" down to this equation.

```

for i = 1 to m
    z[1](i) = W[1]x(i) + b[1]
    → a[1](i) = σ(z[1](i))
    → z[2](i) = W[2]a[1](i) + b[2]
    → a[2](i) = σ(z[2](i))

```

Handwritten annotations:  $x = a^{[1]}$  and  $x^{(i)} = a^{[1](i)}$  are written next to the equations for  $a^{[1]}$  and  $a^{[2]}$  respectively. Brackets group the first two equations and the last two equations.

```

Z[1] = W[1]X + b[1] ← wT,1A[1]+bT,1
A[1] = σ(Z[1])
Z[2] = W[2]A[1] + b[2]
A[2] = σ(Z[2])

```

Handwritten annotations:  $A^{[1]}$  is circled in blue. Brackets group the first two equations and the last two equations.



deeplearning.ai

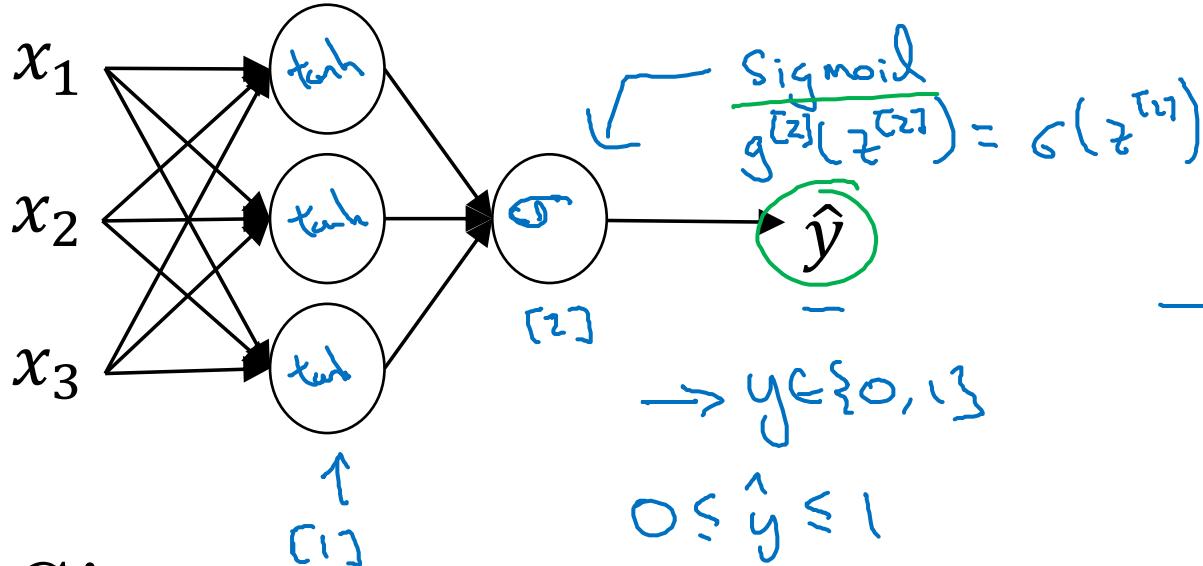
# One hidden layer Neural Network

---

## Activation functions

# Activation functions

$$\boxed{g^{(1)}(z^{(1)}) = \tanh(z^{(1)})}$$



Given x:

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

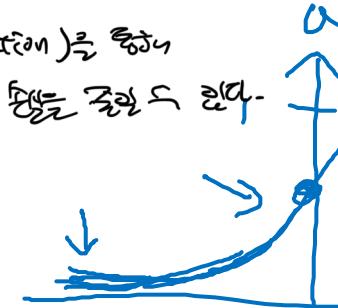
$$\rightarrow a^{[1]} = \sigma(z^{[1]})$$

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

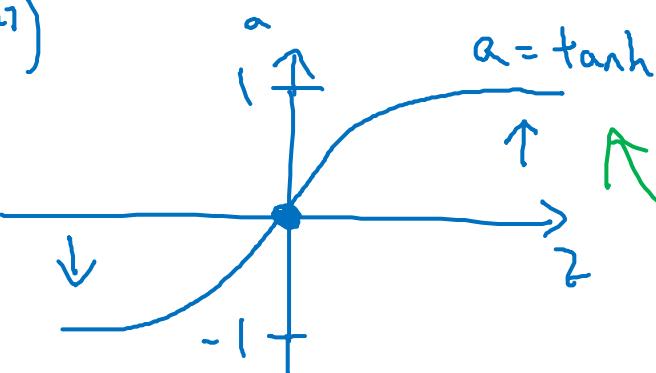
$$\rightarrow a^{[2]} = \sigma(z^{[2]}) g^{[2]}(z^{[2]})$$

• 활성화 함수 (Activation function)을 통해  
각 단위의  $w_i, b$ 를 계산하고 출력을 출력 한다.

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



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



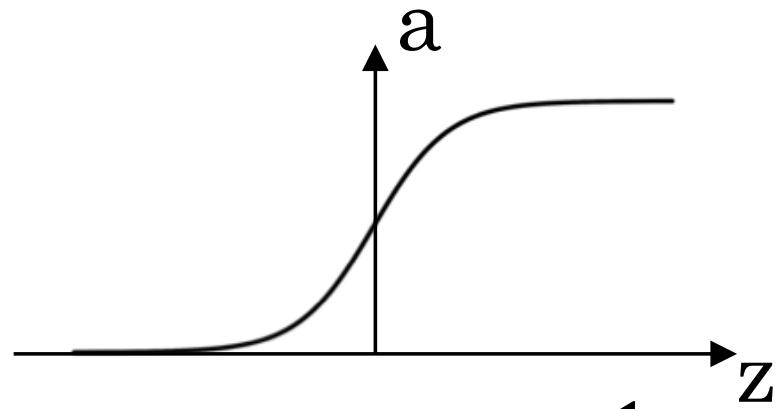
$$z > 0$$

Rely

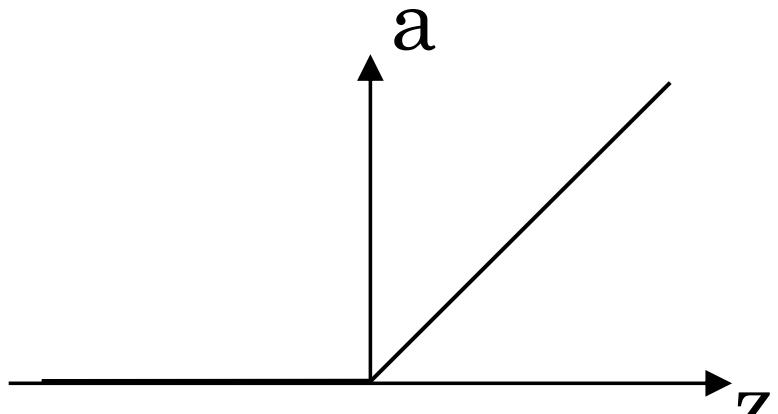
## Rectified Line Unit

## Leaky ReLU

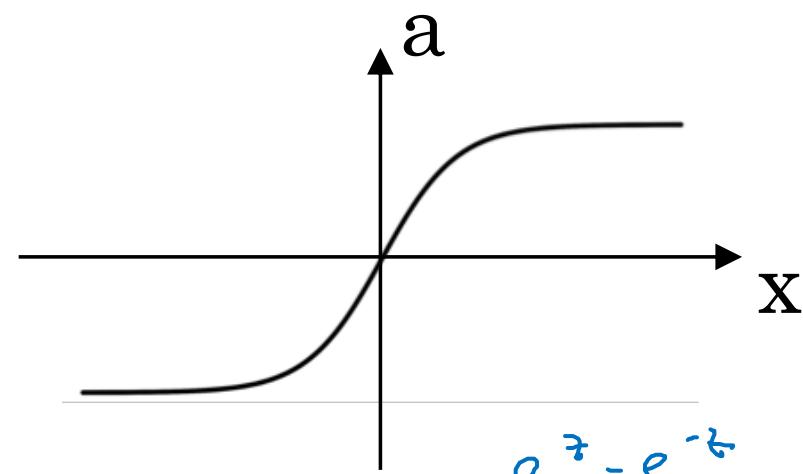
# Pros and cons of activation functions



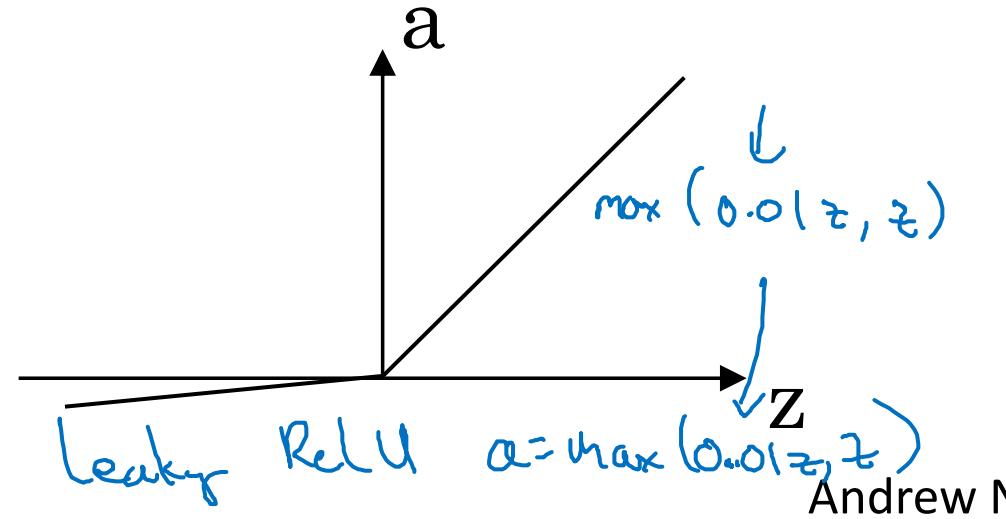
$$\text{sigmoid: } a = \frac{1}{1 + e^{-z}}$$



$$\text{ReLU} \quad a = \max(0, z)$$



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



$$\text{Leaky ReLU} \quad a = \max(0.01z, z) \quad \text{Andrew Ng}$$



deeplearning.ai

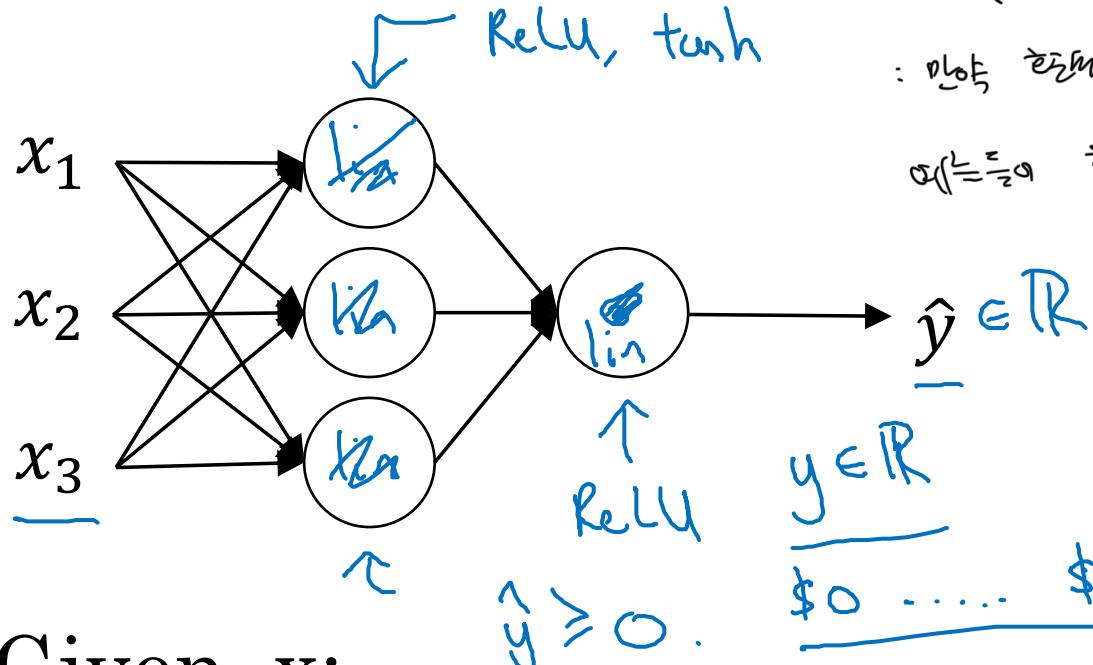
# One hidden layer Neural Network

---

Why do you  
need non-linear  
activation functions?

# Activation function

< 허용되는 활성화 함수가 모든 비선형 함수인 경우>



: 만약 허용되는 활성화 함수가 선형함수라면, 선형성을 가진 허용되는 모든 비선형 함수가 된다.

$$y = \alpha z + b \quad (y = mx + b)$$

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

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

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

Given  $x$ :

$$\rightarrow z^{[1]} = W^{[1]}x + b^{[1]} \quad \text{non-linear activation function of } g(z) \neq z$$

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

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

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

"linear activation function"

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

$$= \underline{w'x + b'} \quad \begin{array}{l} \text{선형성을} \\ \text{가진 허용되는} \\ \text{모든 비선형} \end{array}$$

$$g(z) = z \quad \begin{array}{l} \text{선형성을} \\ \text{가진 허용되는} \\ \text{모든 비선형} \end{array}$$

선형성을 가진 허용되는 모든 비선형 함수를 선형으로 만드는 것은 불가능하다.



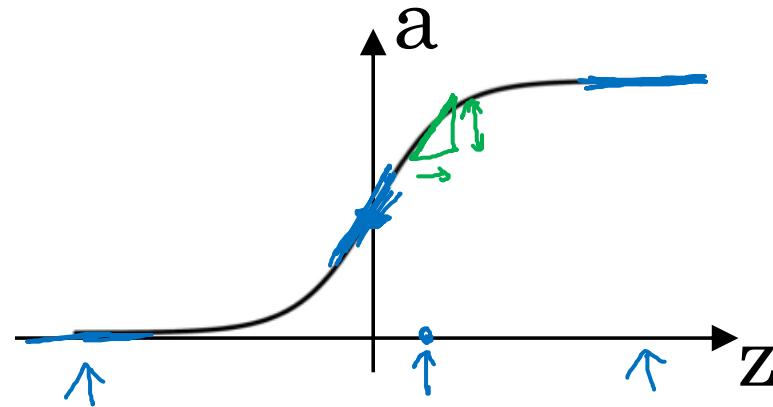
deeplearning.ai

# One hidden layer Neural Network

---

## Derivatives of activation functions

# Sigmoid activation function



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

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

$$\begin{aligned} g'(z) &= \boxed{\frac{d}{dz} g(z)} \\ &= \frac{1}{1 + e^{-z}} \left( 1 - \frac{1}{1 + e^{-z}} \right) \\ &= g(z) \left( 1 - g(z) \right) \quad \leftarrow \\ &= \boxed{a(1-a)} \quad \left| \begin{array}{l} g'(z) = a(1-a) \\ \uparrow \end{array} \right. \end{aligned}$$

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

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

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

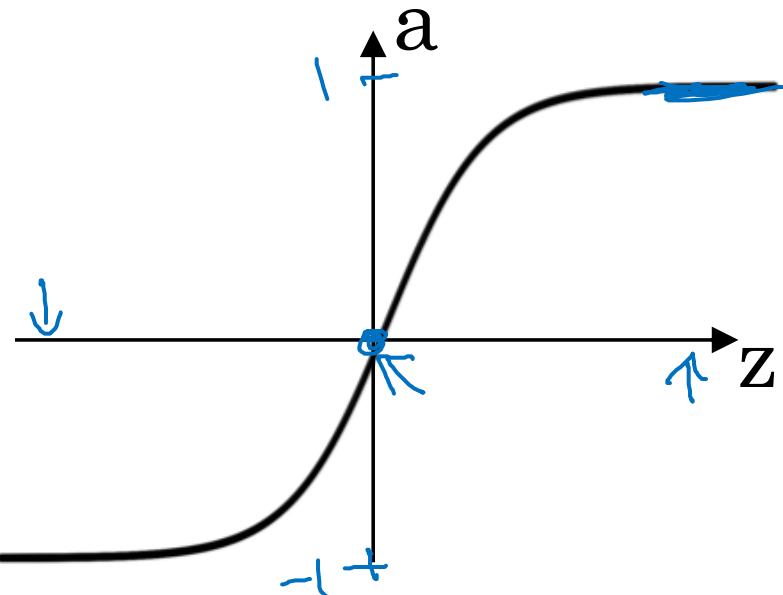
$$\frac{d}{dz} g(z) \approx 0 \cdot (1-0) \approx 0$$

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

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

Andrew Ng

# Tanh activation function



$$g(z) = \tanh(z)$$

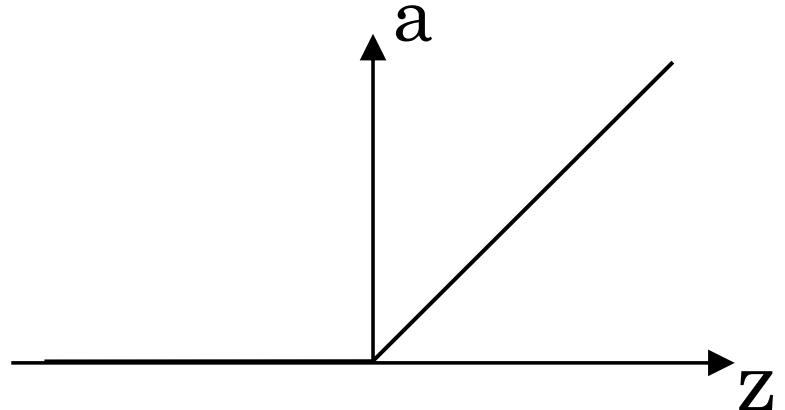
$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\begin{aligned} g'(z) &= \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z \\ &= \underline{\underline{1 - (\tanh(z))^2}} \leftarrow \end{aligned}$$

$$a = g(z), \quad g'(z) = 1 - a^2$$

$$\left| \begin{array}{ll} z = 10 & \tanh(z) \approx 1 \\ & g'(z) \approx 0 \\ z = -10 & \tanh(z) \approx -1 \\ & g'(z) \approx 0 \\ z = 0 & \tanh(z) = 0 \\ & g'(z) = 1 \end{array} \right.$$

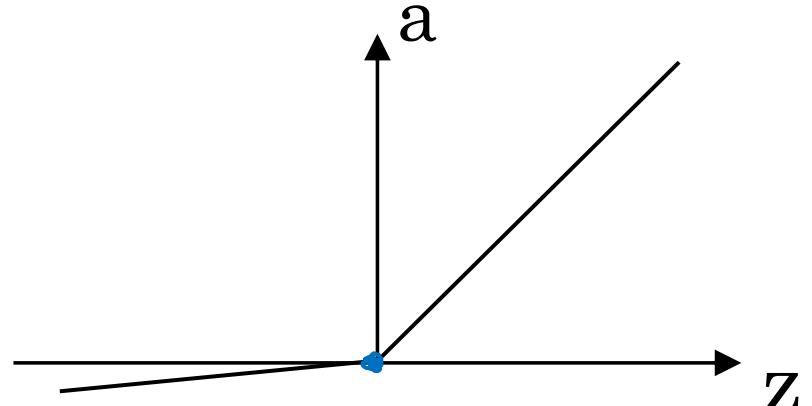
# ReLU and Leaky ReLU



ReLU

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

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



## Leaky ReLU

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

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



deeplearning.ai

One hidden layer  
Neural Network

---

Gradient descent for  
neural networks

# Gradient descent for neural networks

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

Cost function:  $J(\underline{\omega}^{[1]}, \underline{b}^{[1]}, \underline{\omega}^{[2]}, \underline{b}^{[2]}) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}_i, y_i)$

$n_x = n^{[0]}$ ,  $n^{[1]}$ ,  $n^{[2]} = 1$   
 ↓  
 input feature  
 + hidden feature  
 ↓  
 output feature

Gradient Descent:

→ Repeat {

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

$\frac{\partial J}{\partial \omega^{[1]}} = \frac{\partial J}{\partial \omega^{[1]}} , \quad \frac{\partial J}{\partial b^{[1]}} = \frac{\partial J}{\partial b^{[1]}} , \dots$

$\omega^{[1]} := \omega^{[1]} - \alpha \frac{\partial J}{\partial \omega^{[1]}}$

$b^{[1]} := b^{[1]} - \alpha \frac{\partial J}{\partial b^{[1]}}$

$\omega^{[2]} := \dots$

$b^{[2]} := \dots$

# Formulas for computing derivatives

Forward propagation:

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

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

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

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

Back propagation:

$$dz^{[2]} = A^{[2]} - Y \leftarrow$$

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

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

$$dz^{[1]} = \underbrace{w^{[2]T} dz^{[2]}}_{(n^{[2]}, m)} \times \underbrace{g^{[2]'}(z^{[2]})}_{\text{element-wise product}} \quad (n^{[1]}, m)$$

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

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

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

$$(n^{[1]}) \leftarrow$$

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



deeplearning.ai

One hidden layer  
Neural Network

---

Backpropagation  
intuition (Optional)

# Computing gradients

## Logistic regression

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

$$\delta w = \delta z \cdot x$$

$$\delta b = \delta z$$

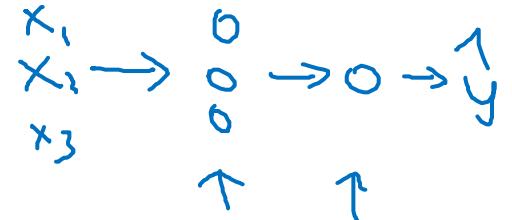
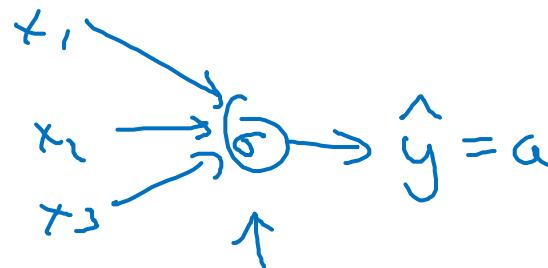
$$\underline{\delta z} = a - y$$

$$\delta z = \delta a \cdot g'(z)$$

$$g(z) = \sigma(z)$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z}$$

$$\text{"d}z" = \text{"d}a$$



$$a = \sigma(z)$$

$$\mathcal{L}(a, y)$$

$$\delta a$$

$$\frac{\partial \mathcal{L}(a, y)}{\partial a} = -y \log a - (1-y) \log(1-a)$$

$$\frac{\partial \mathcal{L}(a, y)}{\partial a} = -y + 1 - a$$

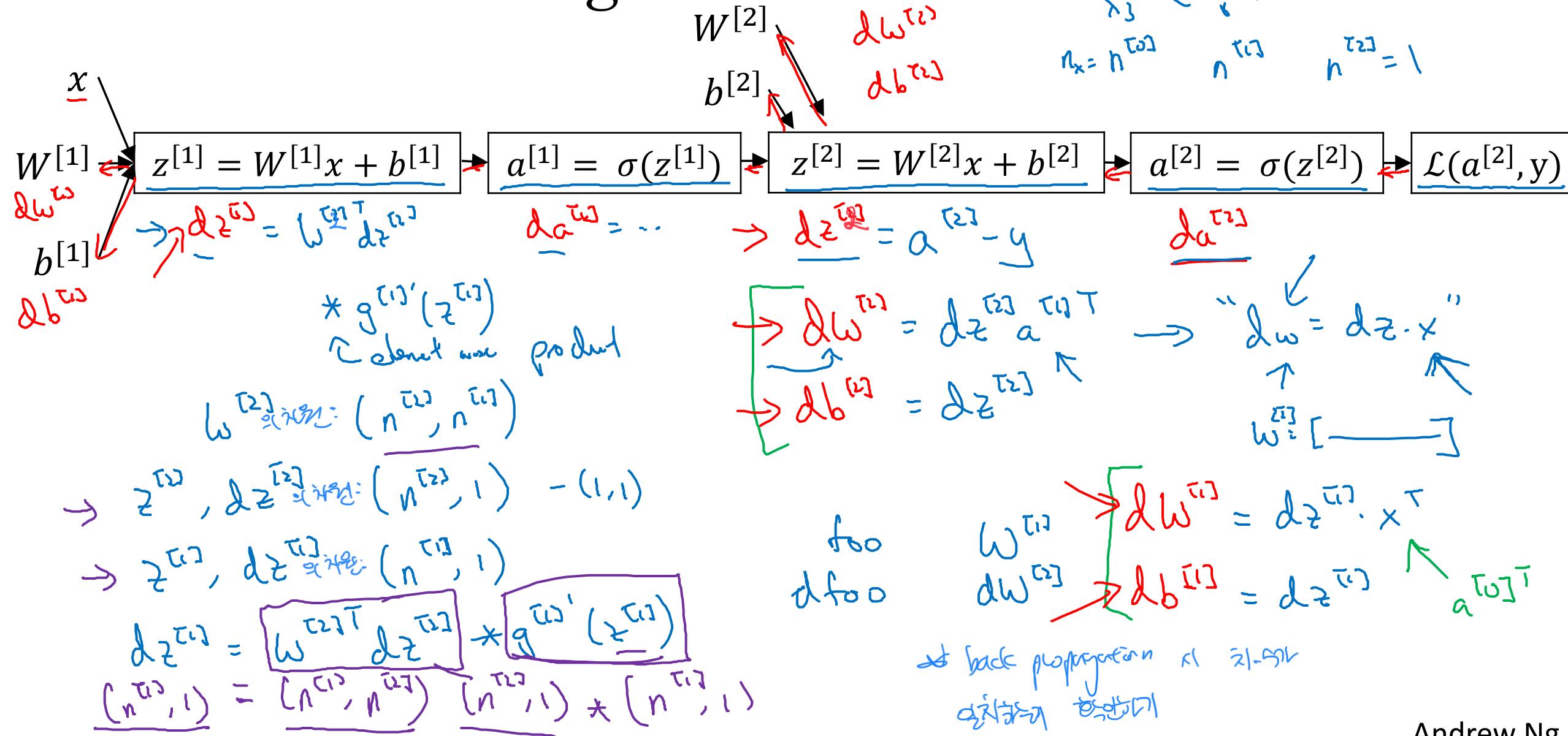
$$\mathcal{L}(a, y) = -y \log a$$

$$- (1-y) \log(1-a)$$

$$\frac{\partial \mathcal{L}(a, y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{\partial}{\partial z} g(z) = g'(z)$$

# Neural network gradients



# Summary of gradient descent

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

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

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

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

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

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

Vectorized implementation:

$$\begin{aligned} z^{[1]} &= \underbrace{w^{[1]} x + b^{[1]}}_{\text{Implementation}} \\ a^{[1]} &= g^{[1]}(z^{[1]}) \end{aligned}$$

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

$$\begin{aligned} z^{[1]} &= w^{[1]} x + b^{[1]} \\ A^{[1]} &= g^{[1]}(z^{[1]}) \end{aligned}$$

# ✓ Summary of gradient descent

< Vectorized Implementation >

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

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

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

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

$(n^{[1]}, 1)$

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

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

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

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

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

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]'}(Z^{[1]})}_{\text{elementwise product}} \quad (n^{[1]}, m)$$

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

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



deeplearning.ai

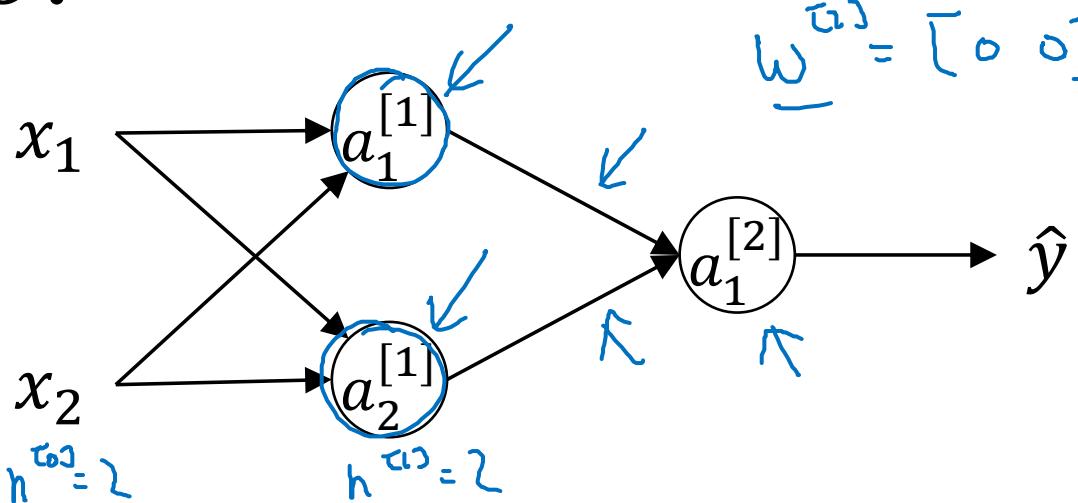
One hidden layer  
Neural Network

---

Random Initialization

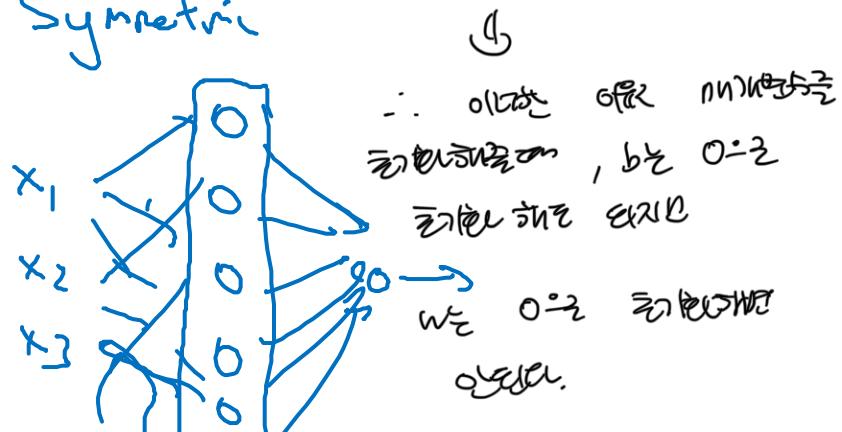
# What happens if you initialize weights to zero?

< 어떤경우  $W$ 를 0으로 초기화하는게 안되는가? >



이전경우 초기화후,  $w, b$ 의 값은 0이면  
차이점은, (전체경우 전부 0으로 초기화되는경우)  
신경망 가시성을 갖지 못해 이해가 어렵게 된다.

Symmetric

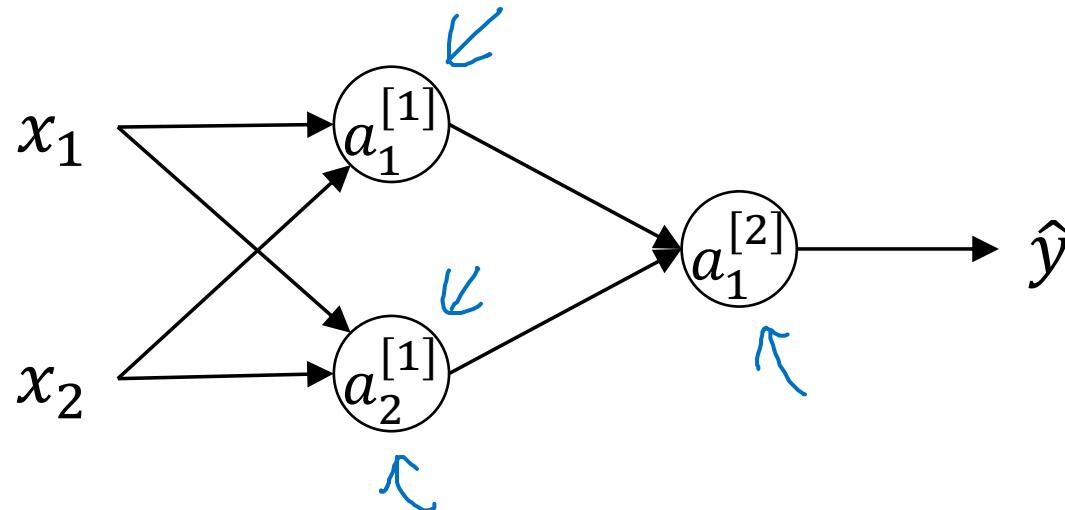


$$W^{[1]} = \begin{bmatrix} \dots & \dots & \dots \end{bmatrix}$$

$$\Delta W = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

$$W^{[1]} = W^{[1]} - \lambda \Delta W$$

# Random initialization



$$\rightarrow w^{[1]} = \text{np.random.randn}(2, 2)$$

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

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

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

<  $w^{[1]}$  를 초기화할 때 0으로 하는 경우, 학습률이 0.01일 때 흐름이 멈춰버린다.  
 →  $w^{[1]}$  를 초기화하는 경우  
 초기화 값이 0인 경우 학습률이 0.01일 때 흐름이 멈춰버린다.  
 학습률이 0.01일 때 흐름이 멈춰버린다.  
 학습률이 0.01일 때 흐름이 멈춰버린다.  
 학습률이 0.01일 때 흐름이 멈춰버린다.

