

# **Neural Networks and Deep Learning**

**박달님 03.14.2021**

# Why need Non-Linear Activation Functions?

## 왜 비선형 함수가 필요한 것인가

- 신경망의 흥미로운 함수들을 산출하기 위해서는 비선형 activation function이 필요하다
- Forward propagation equations for the neural network ->
- 만약 선형 함수로 하게 되면 ->
  - 'Identity activation function' -> just outputting linear function of input
  - Activation 함수가 없는 경우에는 신경망이 몇개의 층으로 이루어졌다고 하더라도 단순히 선형 activation 함수를 산출하는 것
- linear function and sigmoid function - 표현적이지 않음. Hidden layer 쓸모 없음
  - 2개의 선형 함수의 구성요소 그 자체가 선형 함수이기 때문
  - 그러므로 비선형 특성을 기입하지 않는 이상은 더 흥미로운 함수를 산출하는 것이 아님
- 선형 함수를 쓰는 한가지의 경우:  $g(z) = z$  machine learning on regression problem
  - 예시) 집 값 예측.  $Y$ 는 0이나 1이 아니라 실수가 될 것
  - Hidden units가 activation function을 쓰면 안됨
    - ReLU나 tanh, Leaky ReLU etc. 집 값은 다 양수 이므로 여기에서는 ReLU 함수를 사용하면 됨
  - 주로 output layer에서 쓰인다. 하지만, 그 이외의 선형 activation 함수를 hidden layer에서 쓰는 경우는 아주 특별한 경우를 제외하고는 없다

Given  $x$ :

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

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

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

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

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

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

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

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

$$= w' x + b'$$

$$g(z) = z \quad \text{"linear activation function"}$$

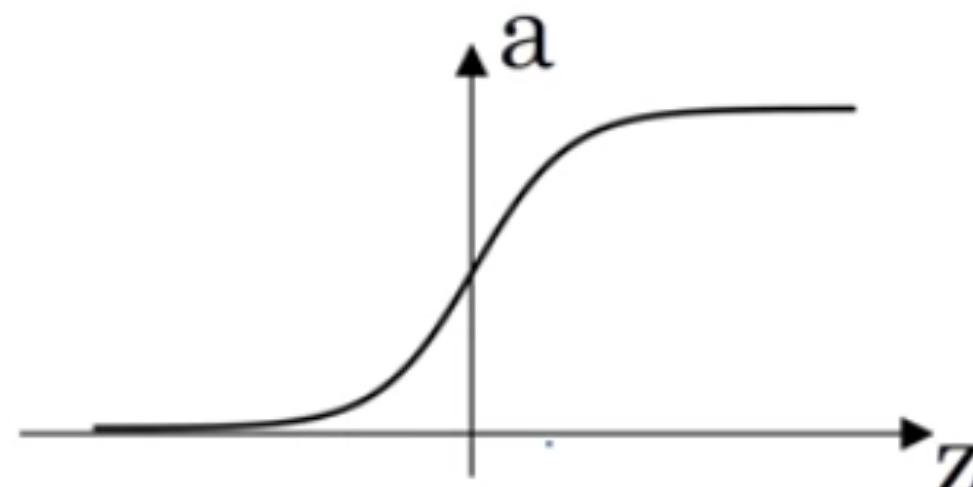
$$a^{[1]} = z^{[1]}$$

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

$$a^{[2]} = z^{[2]}$$

# Derivatives of Activation Functions - Sigmoid

## Back Propagation에서 필요한 activation function의 derivative



Sigmoid function derivative:

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

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

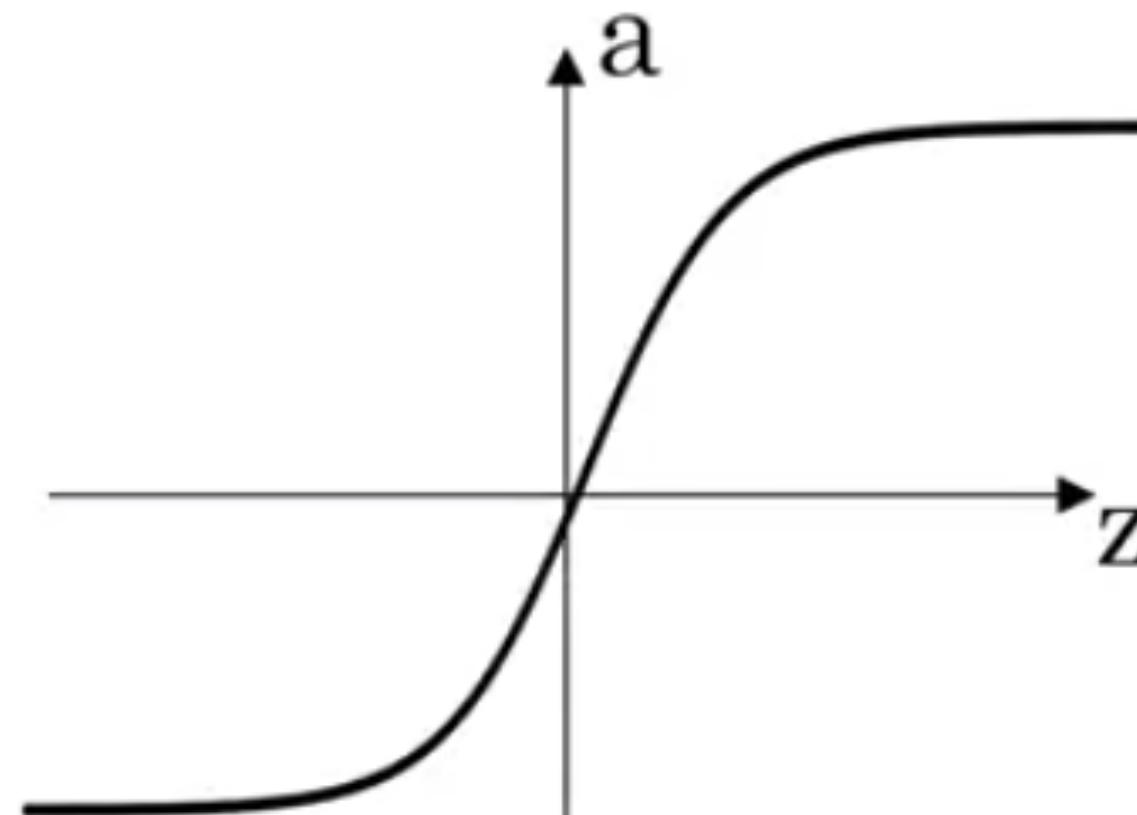
$$= a(1-a)$$

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

- If  $z$  is very large -> ex)  $z = 10$ ,  $g(z) = 1$
- $g(z)(1-g(z))$
- $d/dz g(z) = g'(z) = 1(1-1) = 0$
- If  $z$  is very small -> ex)  $z = -10$   $g(z) = 0$
- $g'(z) = 0(1-0) = 0$
- If  $z = 0$  ->  $g(z) = 1/2$
- $g'(z) = 1/2(1-1/2) = 1/4$
- In neural network:  $a = g(z) = \frac{1}{1 + e^{-z}}$

# Derivatives of Activation Functions

## Tanh activation function



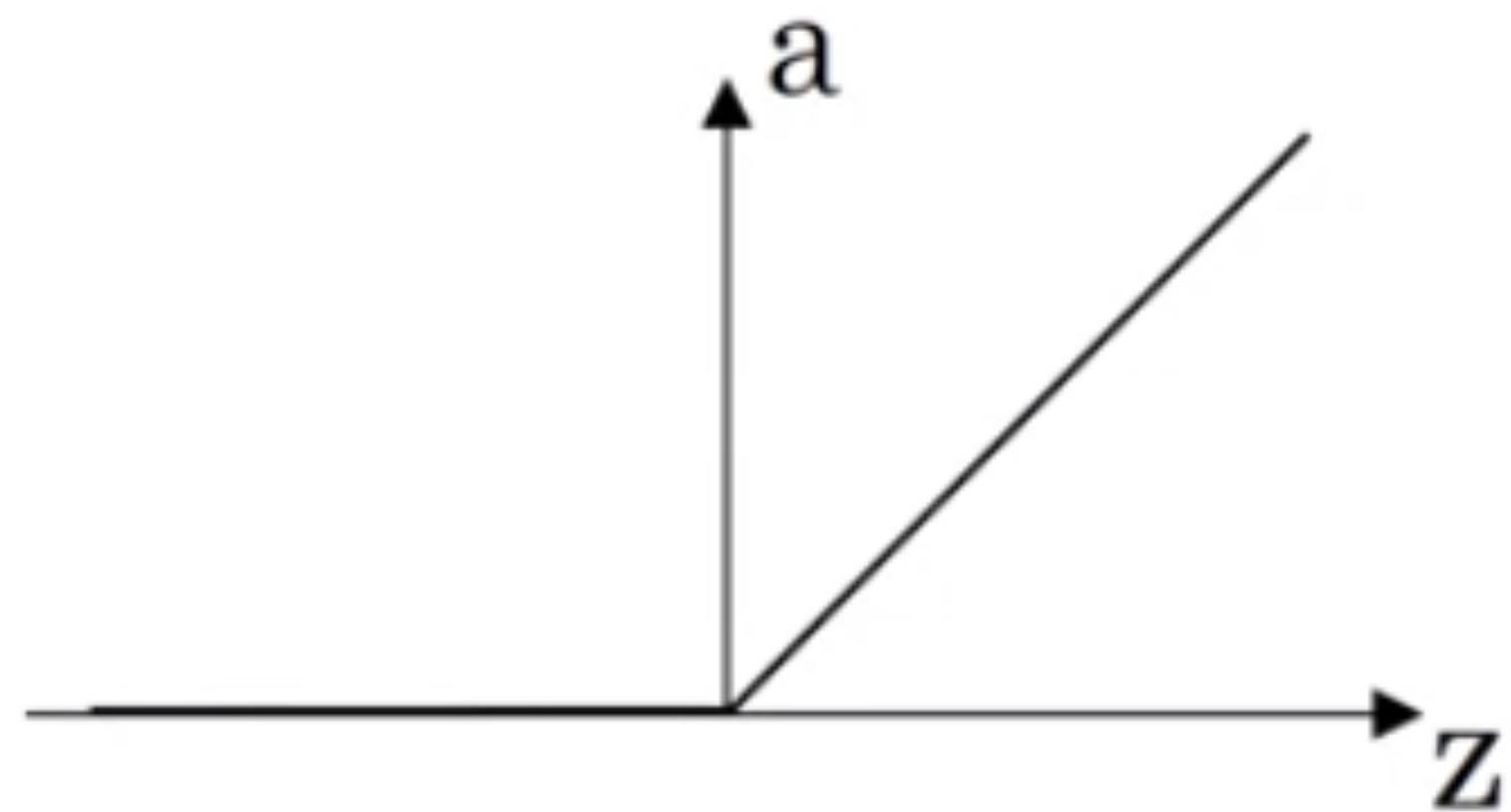
$$\begin{aligned}g(z) &= \tanh(z) \\&= \frac{(e^z - e^{-z})}{(e^z + e^{-z})}\end{aligned}$$

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

$$g'(z) = \frac{d}{dz} g(z) = 1 - (\tanh(z))^2$$

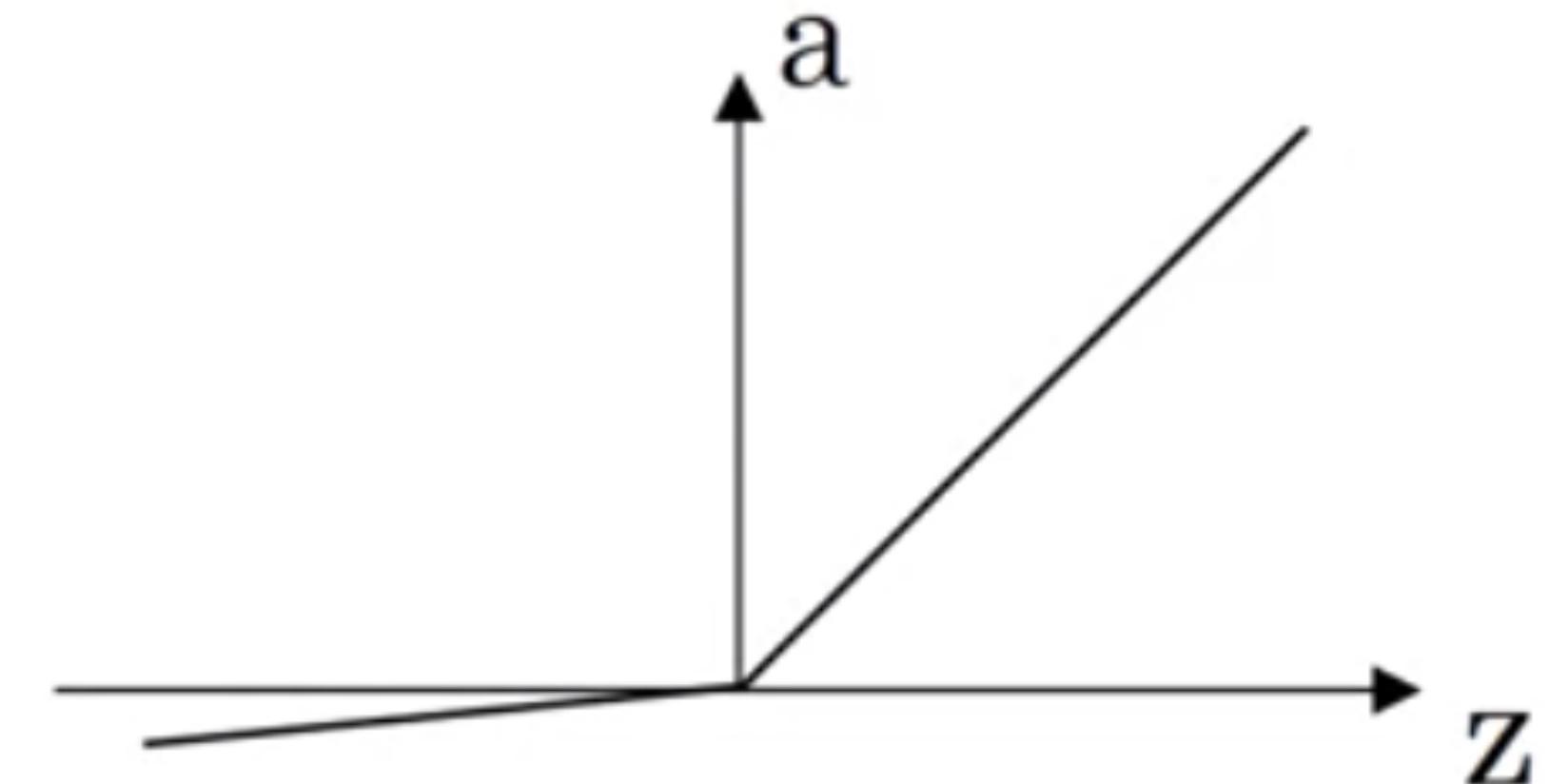
# Derivatives of Activation Functions

## ReLU and Leaky ReLU



ReLU

- $g(z) = \max(0, z)$
- $g'(z) = 0$  if  $z < 0$
- $g'(z) = 1$  if  $z \geq 0$
- ~~$g'(z) = \text{undefined}$  if  $z = 0$~~



Leaky ReLU

- $g(z) = \max(0.01z, z)$
- $g'(z) = 0.01$  if  $z < 0$
- $g'(z) = 1$  if  $z > 0$

# Gradient Descent for Neural Networks

## Equations needed to implement back-propagation/gradient descent

- Parameters:  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$
- $N_x = n^0], n^{[1]}, n^{[2]}=1$
- Cost Function:  $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_i^n L(\hat{y}, y)$
- Gradient descent:

```
Repeat {  
    Compute predictions ( $\hat{y}^{(i)}$ ,  $i=1,\dots,m$ )  
     $Dw^{[1]} = dJ/dw^{[1]}, db^{[1]} = dJ/db^{[1]}, \dots$   
     $w^{[1]} = w^{[1]} - \alpha dw^{[1]}$   
     $b^{[1]} := b^{[1]} - \alpha db^{[2]}$   
     $b^{[2]} := b^{[2]} - \alpha db^{[2]}$   
}
```

# Formulas for computing derivatives

Forward Propagation

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

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

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

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

Backward Propagation

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

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

$$db^{[2]} = \frac{1}{m} \underbrace{\text{np.sum}(dZ^{[2]}, axis=1, keepdims=True)}_{\substack{\text{summing horizontally} \\ \downarrow}}$$

$$dZ^{[0]} = \underbrace{W^{[2]T} dZ^{[2]}}_{\substack{\text{element-wise} \\ \text{product}}} \times \underbrace{g^{[1]}'(Z^{[1]})}_{\substack{\text{derivative of} \\ \text{activation function} \\ \downarrow}}$$

↓  
prevents producing  
(n, )  
insures producing  
(n<sup>[2]</sup>, 1)

↓  
both (n<sup>[1]</sup>, m) matrix

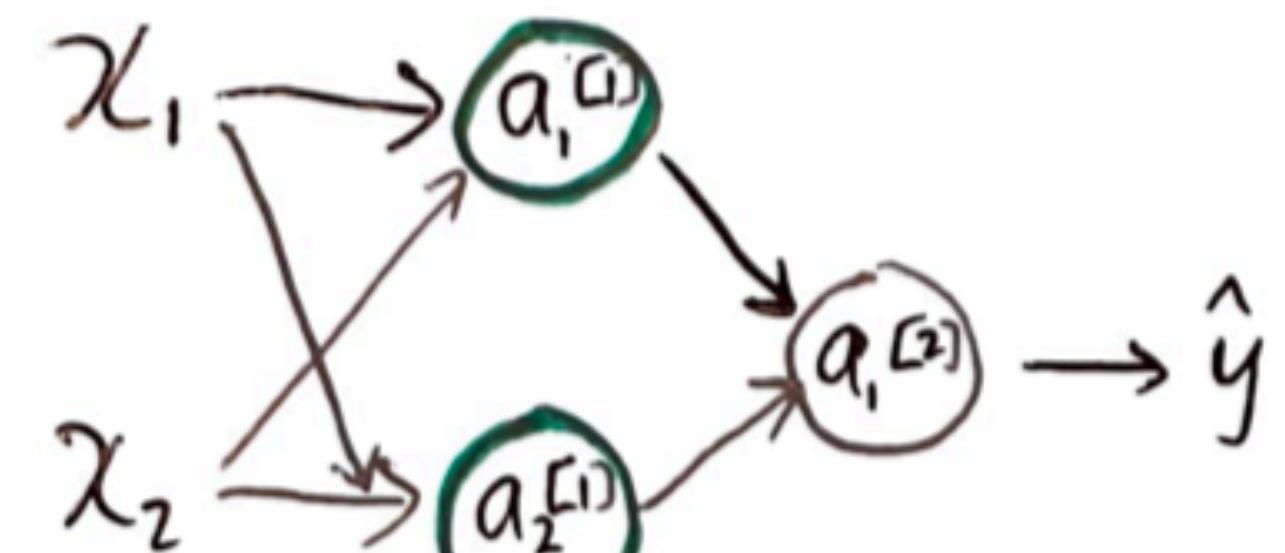
$$dw^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \underbrace{\text{np.sum}(dZ^{[1]}, axis=1, keepdims=True)}_{\substack{\downarrow \\ (n^{[0]}, 1)}}$$

# Random Initialization

- Neural network를 훈련시킬때 weights를 랜덤하게 initialize 해주는게 중요함
- Logistic regression은 weights를 0으로 initialize 해줌
- Neural network는 parameter weight를 0으로 initialize 해주고 gradient descent를 적용하면 작동하지 않을 것
  - 그렇게 하면 어떻게 되는가?

# Random Initialization



$$n^{[0]} = 2 \quad n^{[1]} = 2$$

$$W^{[1]} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]}$$

어떤 example 미던 둘의  
값이 같다.

둘의 같은 function을  
compute 하기 때문에.

$$dZ_1^{[1]} = dZ_2^{[1]}$$

Hidden units initialize the same way.  
outgoing weights are also identical.

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

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

weight update :  $W^{[1]} = W^{[1]} - \alpha dw$

몇번을 iterate 하던 상관 없이 두개의 hidden units는 같은 함수를 계산한다.

그러므로 두개의 hidden units는 output unit에 똑같은 영향을 준다.

다른 hidden units들을 이용해서 다른 함수들을 계산하려고 필요한건데 이 경우에는 Hidden unit이 필요가 없어진다.

Solution: parameter를 random하게 initialize한다.

**Symmetric** - 0 으로 initialize 하면 둘의 hidden units  $a_1^{[1]}, a_2^{[1]}$ 가 완전히 같다.

# Random Initialization

$\rightarrow W^{(1)} = \text{np.random.randn}(2, 2) * \underline{0.01}$

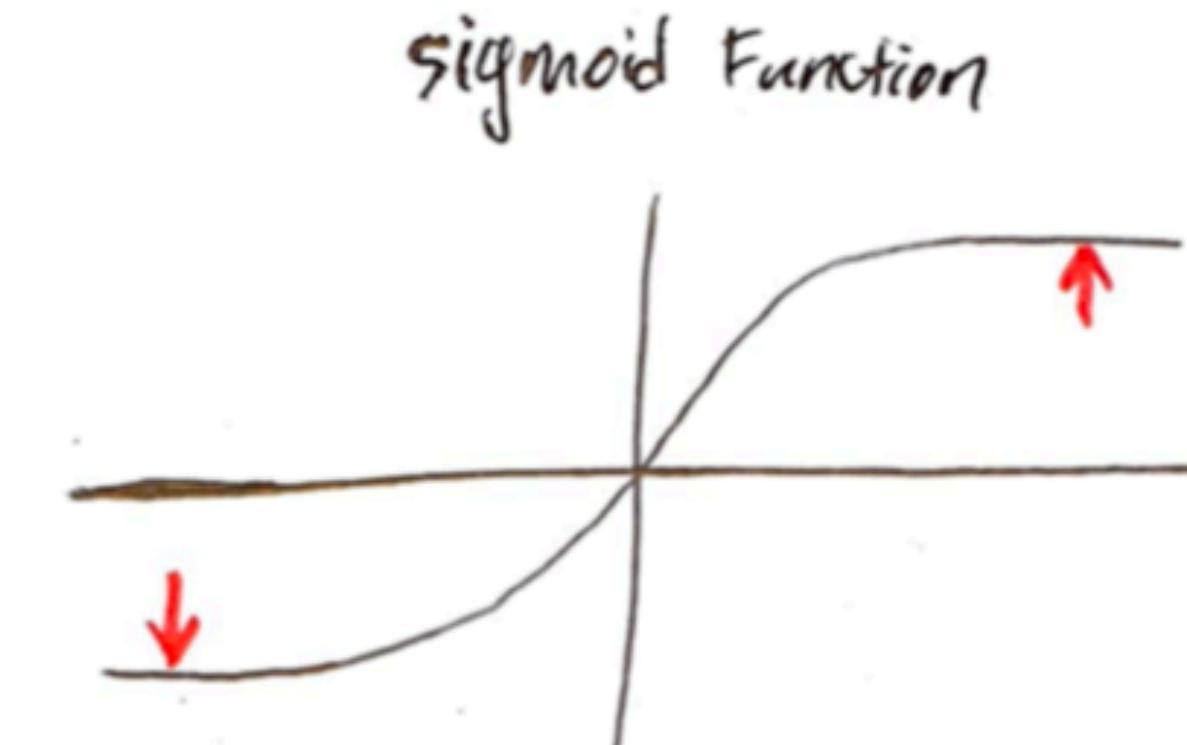
$b^{(1)} = \text{np.zeros}(2, 1)$

$\uparrow W^{(2)} = \text{np.random.randn}(2, 2)$

$b^{(2)} = 0$   
b does not have symmetry  
breaking problem

As long as w is initialized randomly,  
start off with different hidden units  
computing different things  
 $\therefore$  No longer have symmetry breaking problem

prefer to initialize with  
small random values.



if  $w$  is big  
 $z^{(1)} = W^{(1)}x + b^{(1)}$   
 $a^{(1)} = g^{(1)}(z^{(1)})$   
 $\uparrow$   
 $z$  will be big

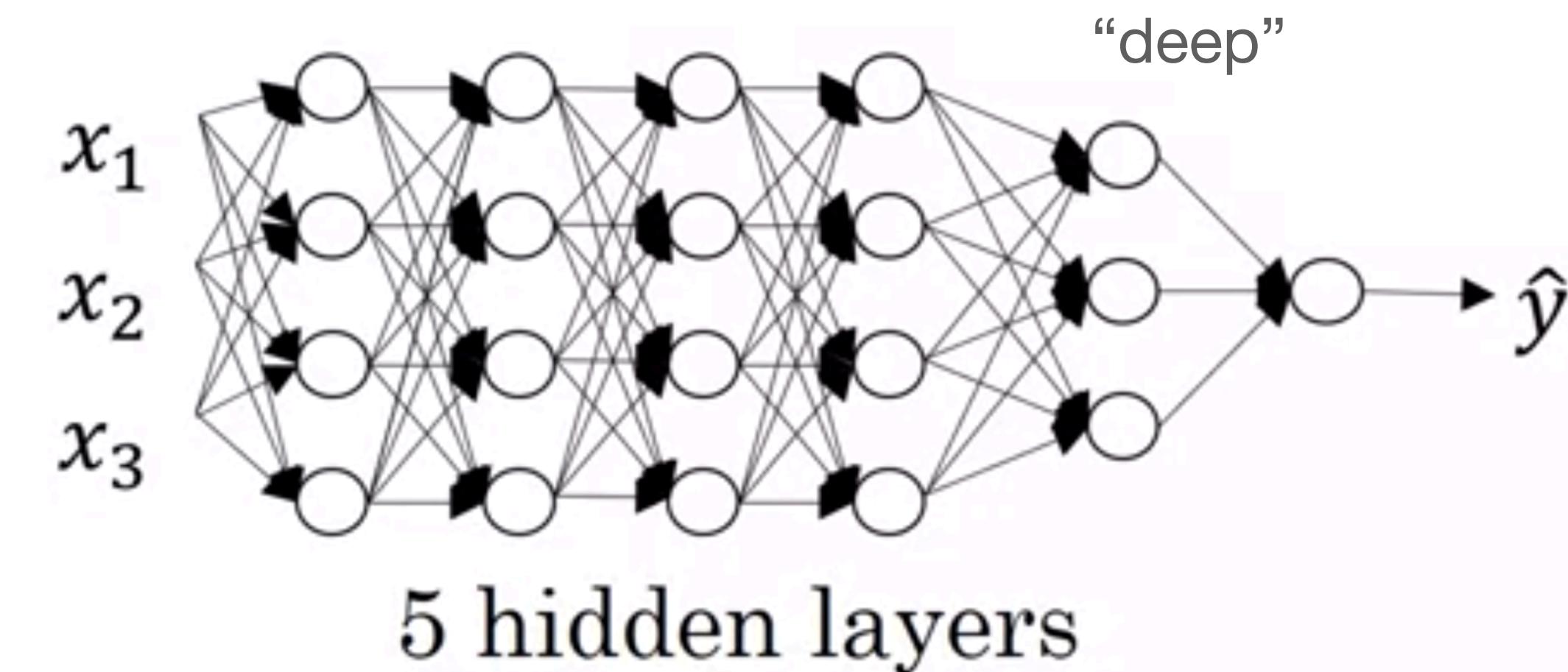
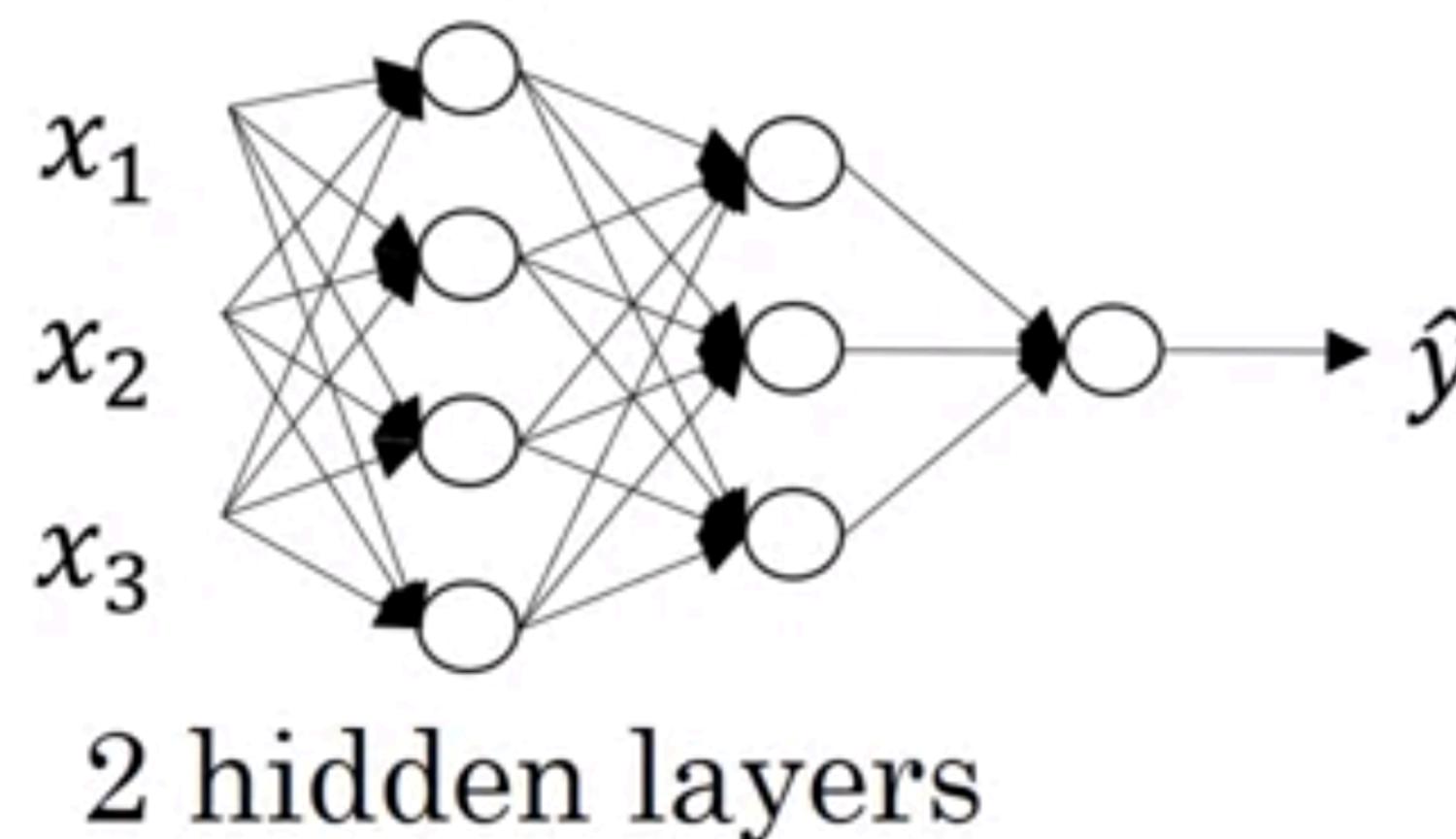
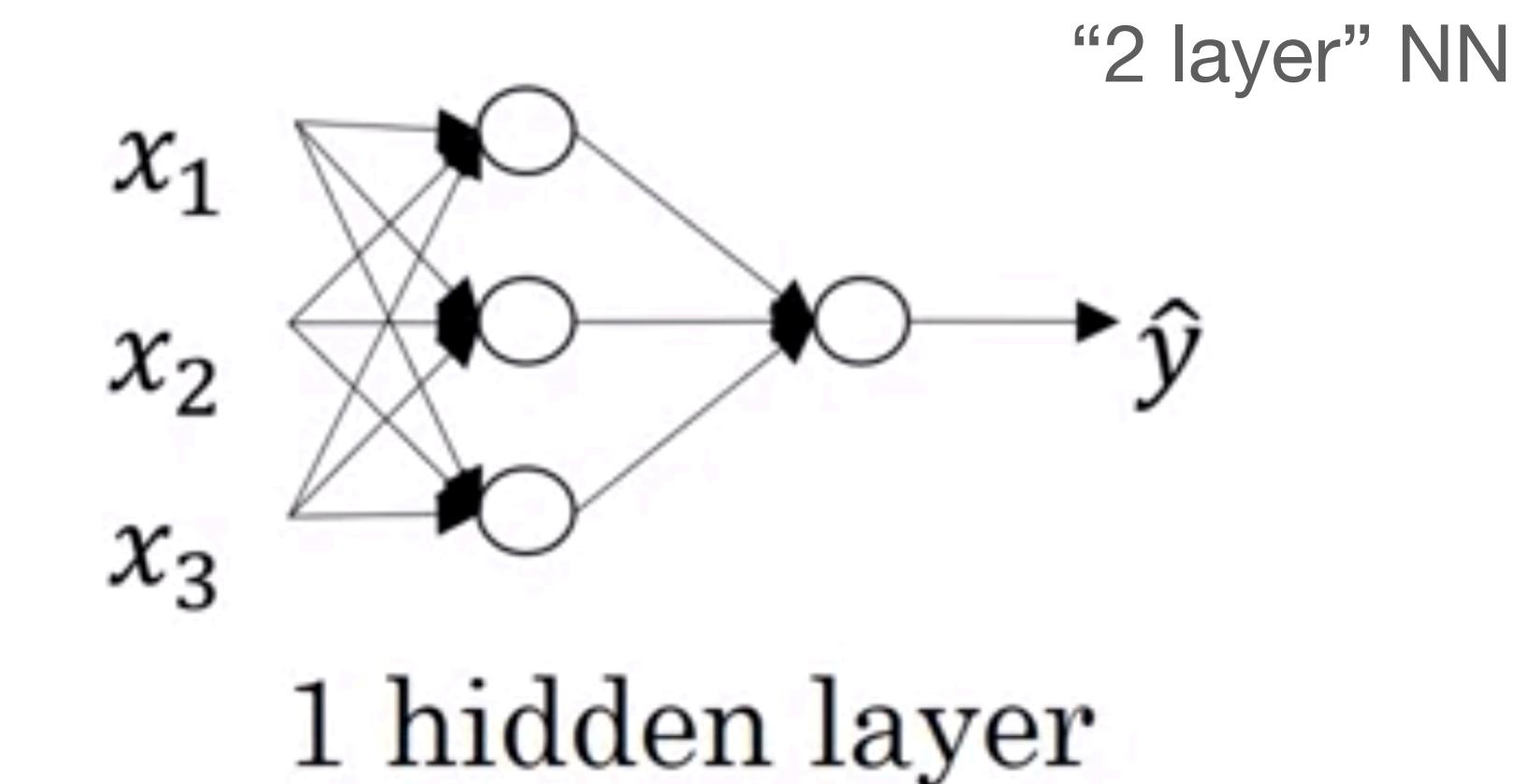
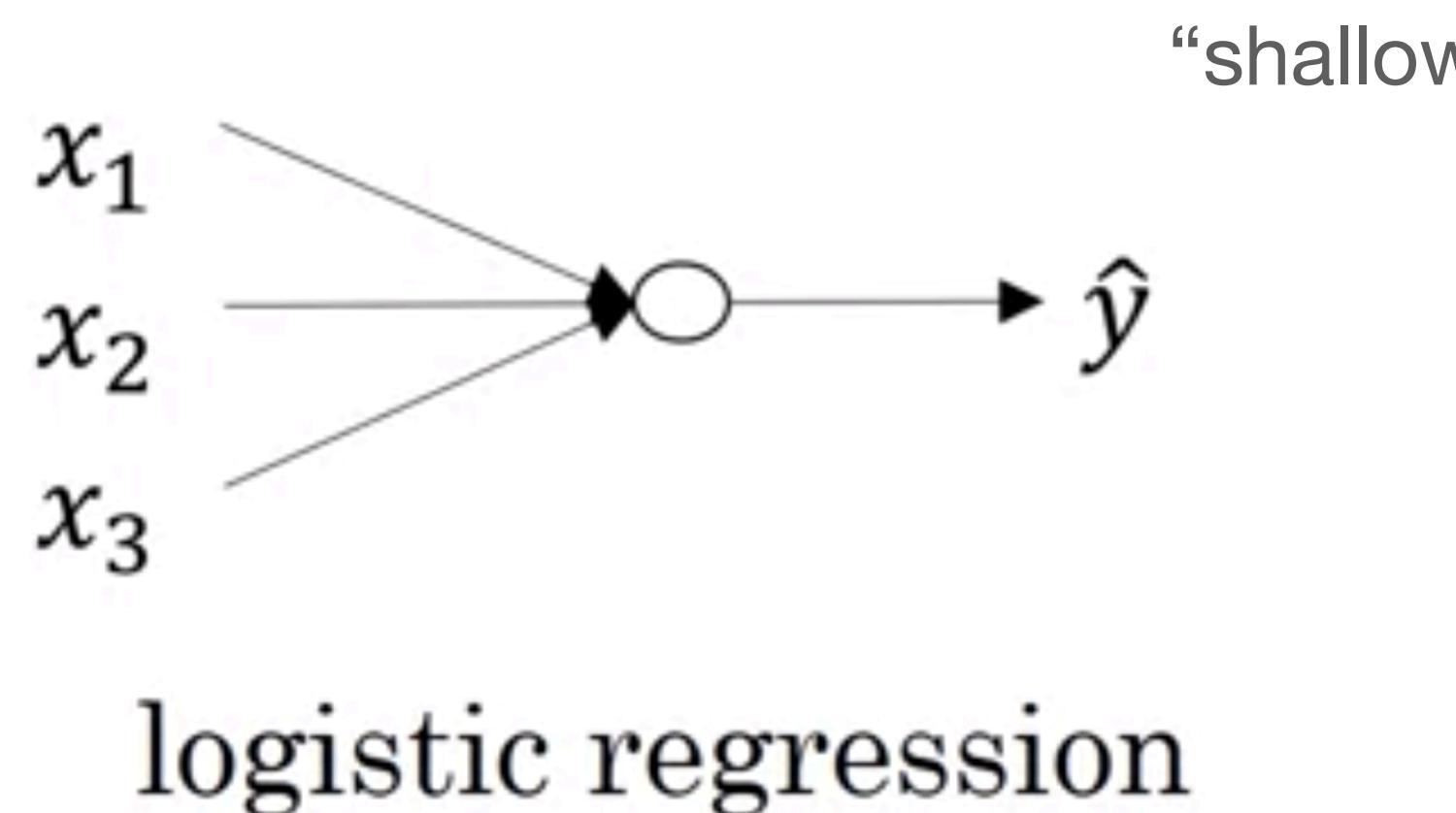
slope of gradient function is very small.  
 $\therefore$  learning rate will be very small.  
ex) binary classification.tanh/sigmoid.

When training a neural network with just one  
hidden layer, 0.01 fine.

When training a deep neural network,  
pick a different #.  
Relatively small #.

# Deep L-layer Neural Network

What is a deep neural network?



# Deep neural network notation

$L = 4$  (# of layers)

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

$a^{[l]}$  = activations in layer  $l$

$a^{[l]} = g^{[l]}(z^{[l]})$

$w^{[l]}$  = weights for  $z^{[l]}$

$b^{[l]}$  = used to compute  $z^{[l]}$

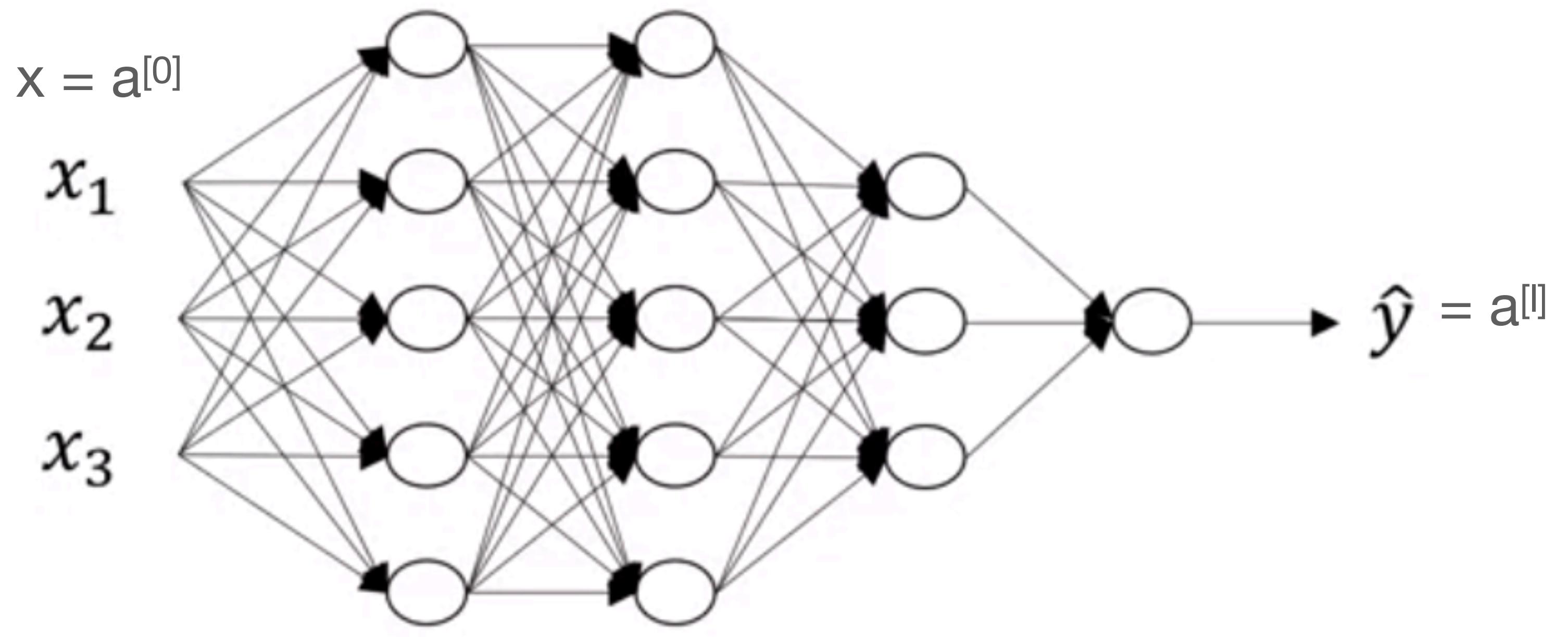
$n^{[1]} = 5$

$n^{[2]} = 5$

$n^{[3]} = 3$

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

$n^{[0]} = n_x = 3$



# Forward Propagation in a Deep Network

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

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

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

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

$$z^{[4]} = w^{[4]} a^{[3]} + b^{[4]}$$

$$a^{[4]} = g^{[4]}(z^{[4]})$$

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

$$a^{[l]} = g^{[l]}(z^{[l]})$$

Vectorized:

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

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

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

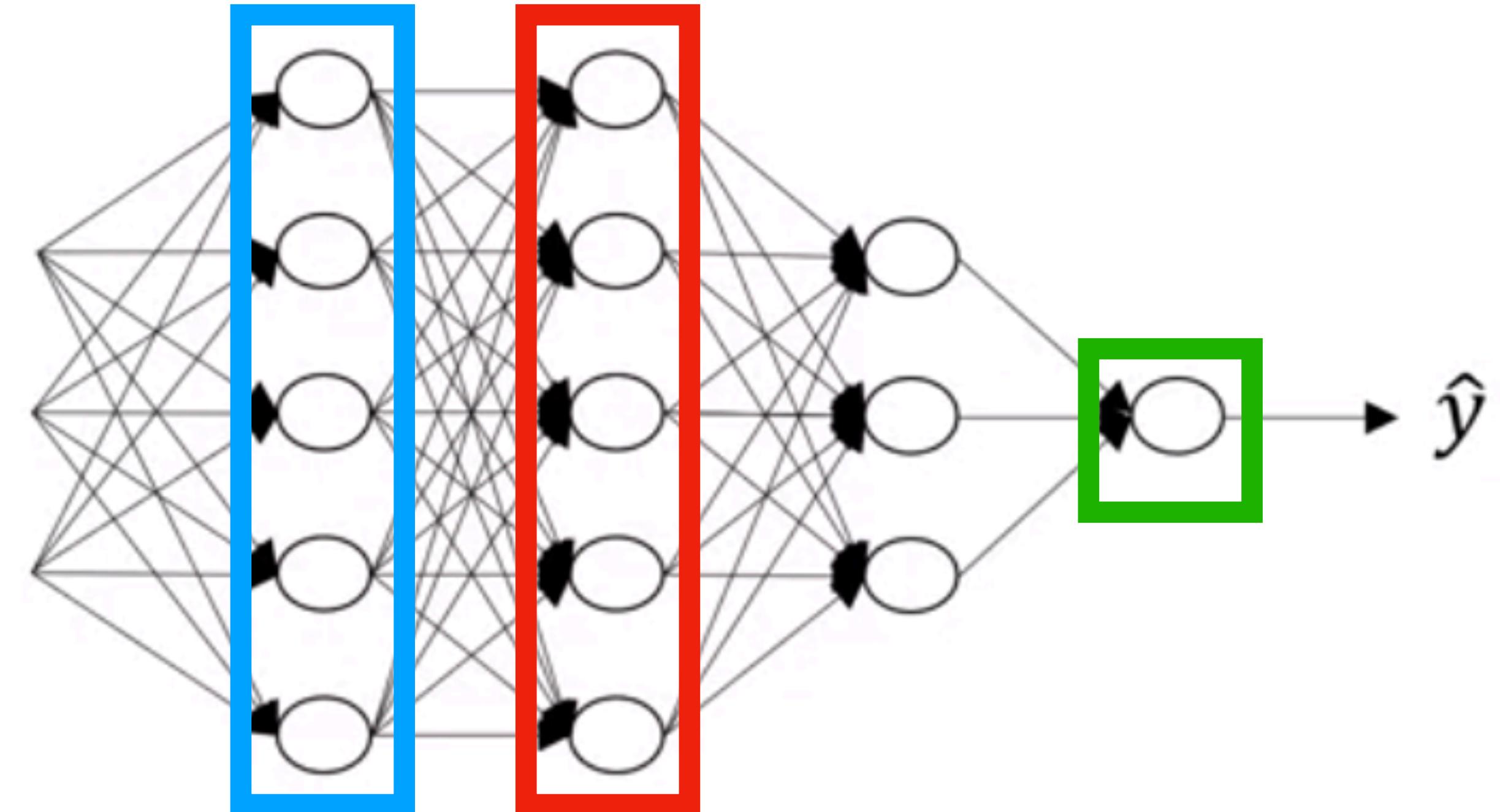
$$A^{[2]} = g^{[2]}(Z^{[2]})$$

$$\hat{y} = g(Z^{[4]}) = A^{[4]}$$

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

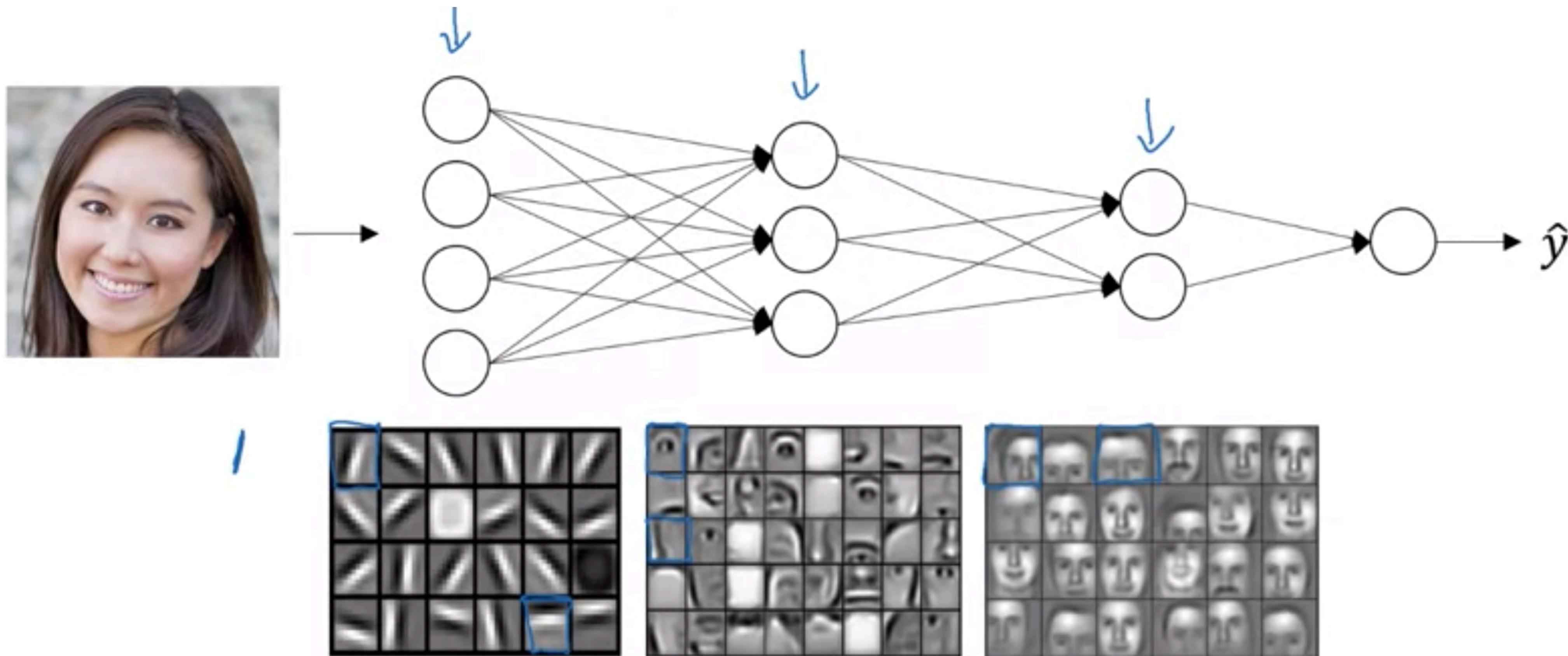
$$A^{[l]} = g^{[l]}(Z^{[l]})$$

$x_1$   
 $x_2$   
 $x_3$



# Why Deep Representations?

What is the deep network computing?

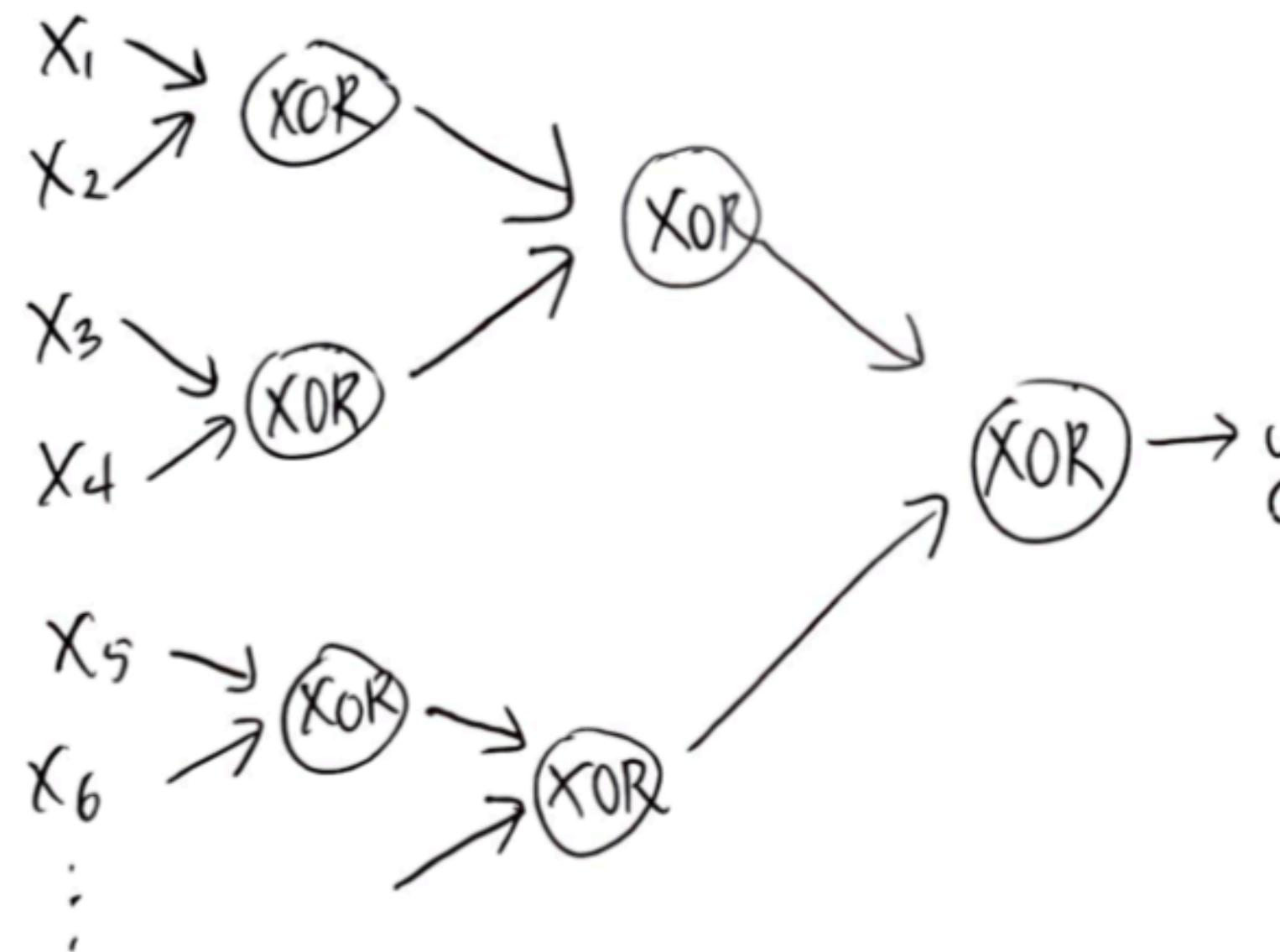


# Circuit theory and deep learning

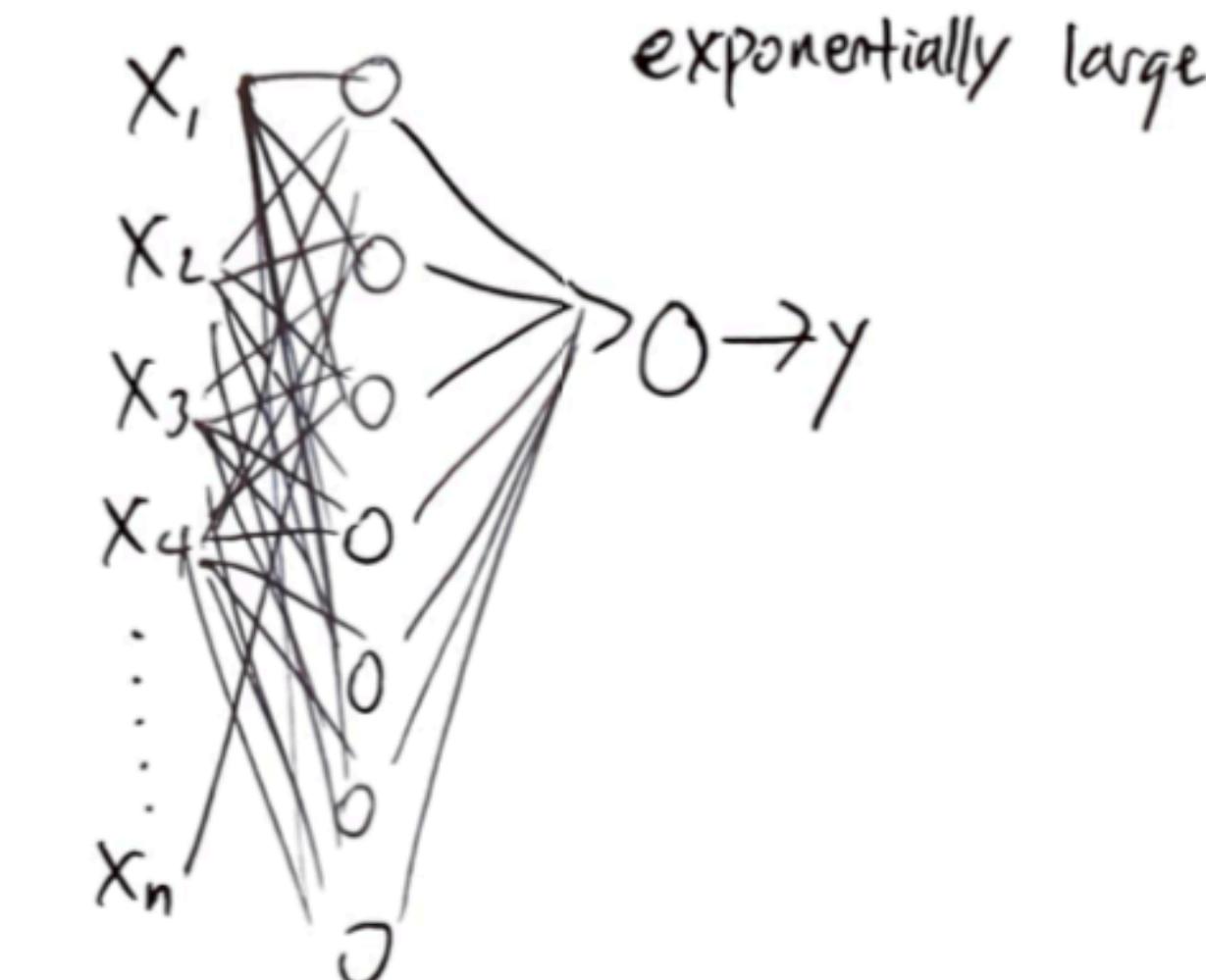
Informally: There are functions you can compute with a “small” L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

DeepNN- XOR binary tree

$$y = x_1 \text{XOR} (x_2 \text{XOR} (x_3 \text{XOR} \dots (x_n \text{XOR}$$



Shallow NN - one single layer  
enumerate all  $2^n$  configuration



# Building Blocks of Deep Neural Networks

## One step forward and backward propagation

Layer L:  $w^{[L]}, b^{[L]}$

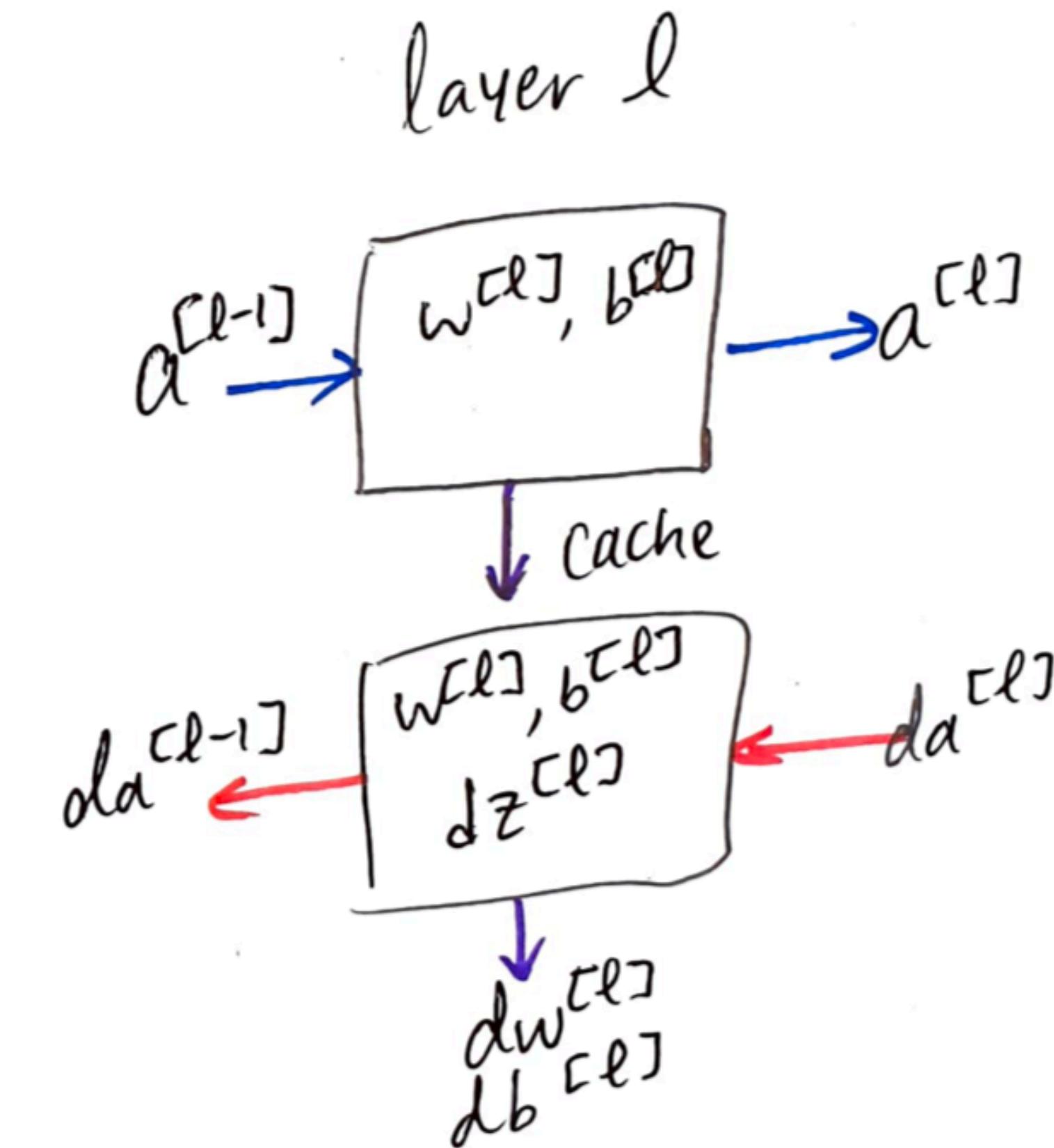
Forward : Input  $a^{[l-1]}$ , output  $a^{[l]}$

$$z^{[l]}: w^{[l]}a^{[l-1]} + b^{[l]}$$

$$a^{[l]}: g^{[l]}(z^{[l]})$$

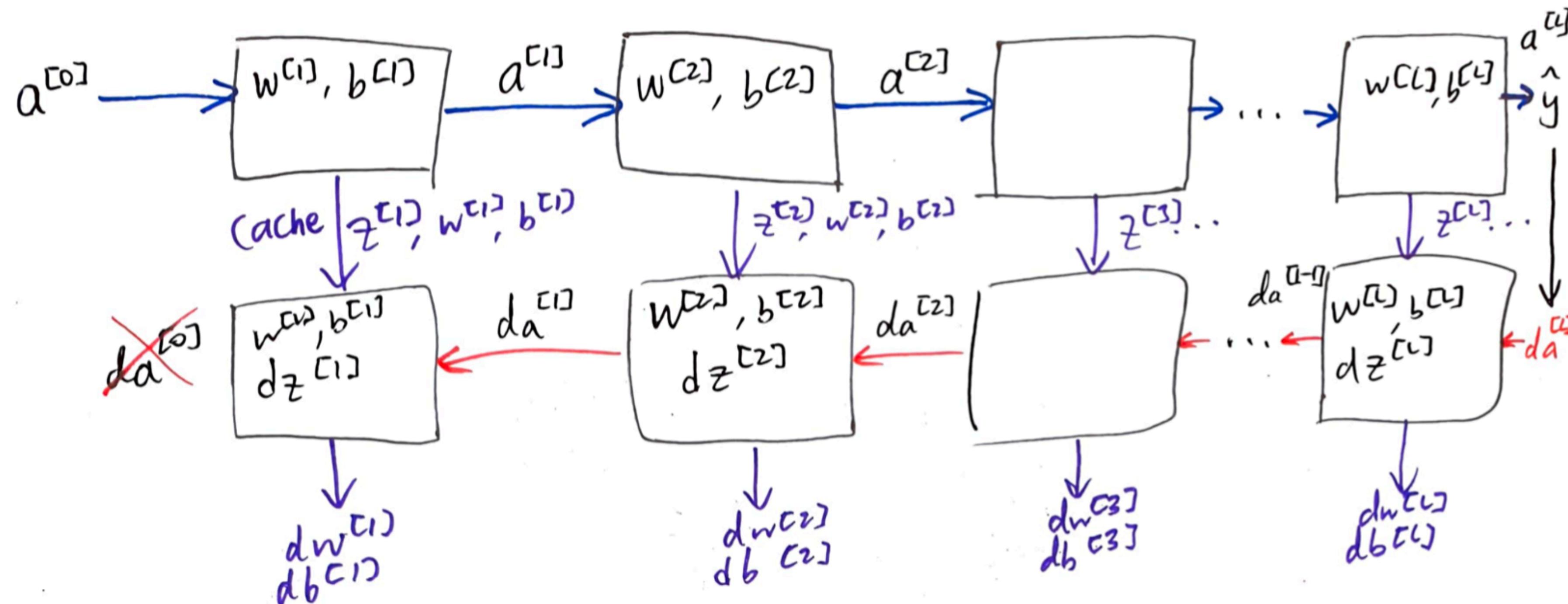
Backward: Input  $da^{[l]}$ , output  $da^{[l-1]}$   
cache( $z^{[l]}$ )

$$\frac{dw^{[l]}}{db^{[l]}}$$



# Building Blocks of Deep Neural Networks

## One Iteration forward and backward propagation



$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}$$
$$b^{[l]} := b^{[l]} - \alpha db^{[l]}$$

# Parameters vs. Hyperparameters

- parameters:  $W^{[l]}$  and  $b^{[l]}$  → trained from data
  - **hyperparameters:**
    - **alpha (learning\_rate), number of iterations, L,  $n^{[l]}$  size of each layer,  $g^{[l]}$  at each layer...**
    - **momentum, mini-batch, regularization...**
- finally decides what params will be.

empirical: try out different hyperparameters

