

深度神经网络的另一个视角

# 深度学习概述

Adept form Hung-yi Lee's tutorial

# Outline

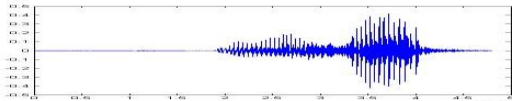
Introduction to Machine Learning

“Three Steps” for Deep Learning

# Machine Learning

## ≈ Looking for a Function


- Speech Recognition

$$f(\text{  ) = \text{“How are you”}$$

- Image Recognition

$$f(\text{  ) = \text{“Cat”}$$

- Playing Go

$$f(\text{  ) = \text{“5-5” (next move)}$$

- Dialogue System

$$f(\text{“Hi” (what the user said)} ) = \text{“Hello” (system response)}$$

# Framework

Image Recognition:

$$f\left(\text{img}\right) = \text{"cat"}$$



$$f_1\left(\text{img}\right) = \text{"cat"}$$

$$f_2\left(\text{img}\right) = \text{"money"}$$

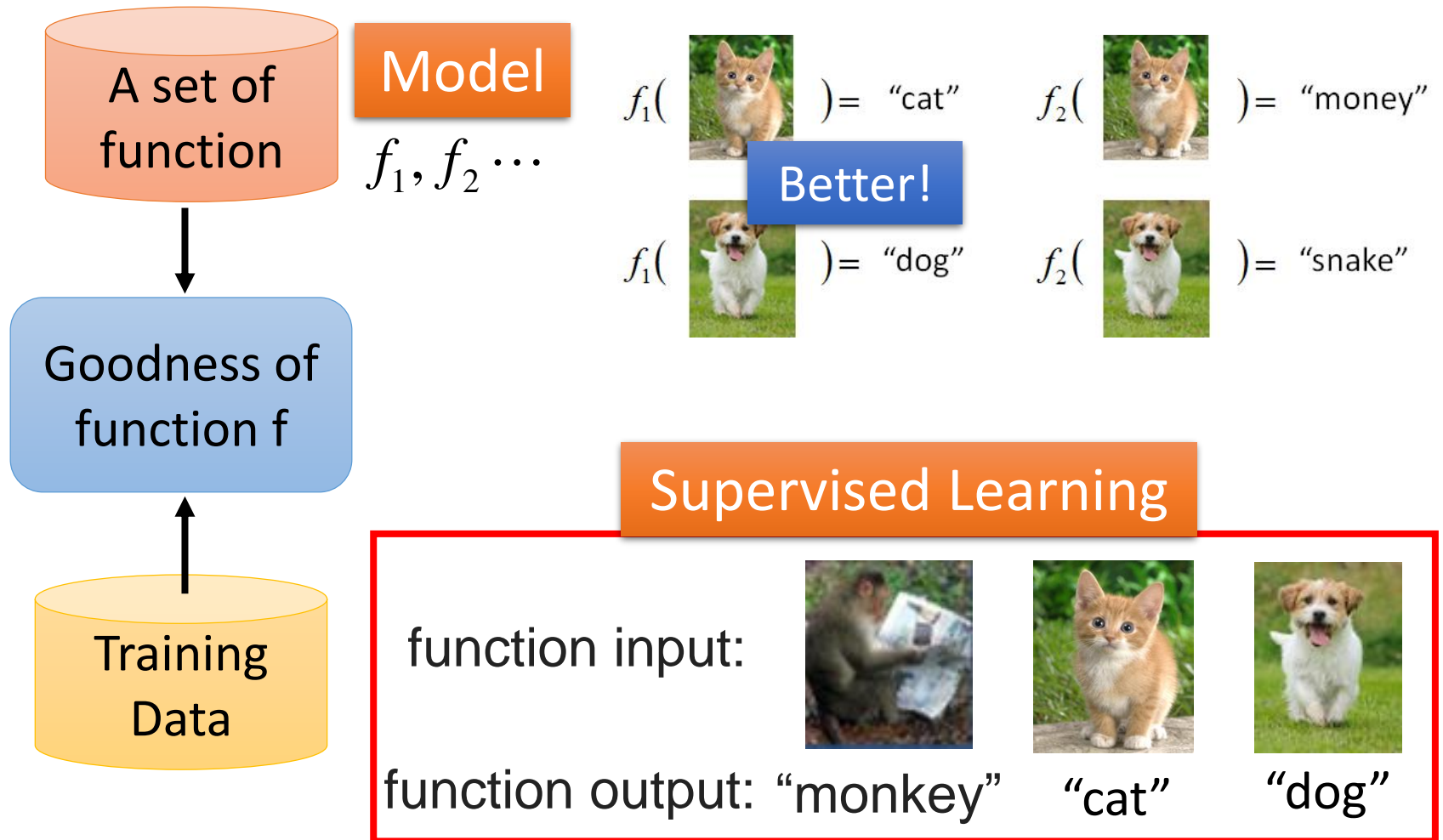
$$f_1\left(\text{img}\right) = \text{"dog"}$$

$$f_2\left(\text{img}\right) = \text{"snake"}$$

# Framework

Image Recognition:

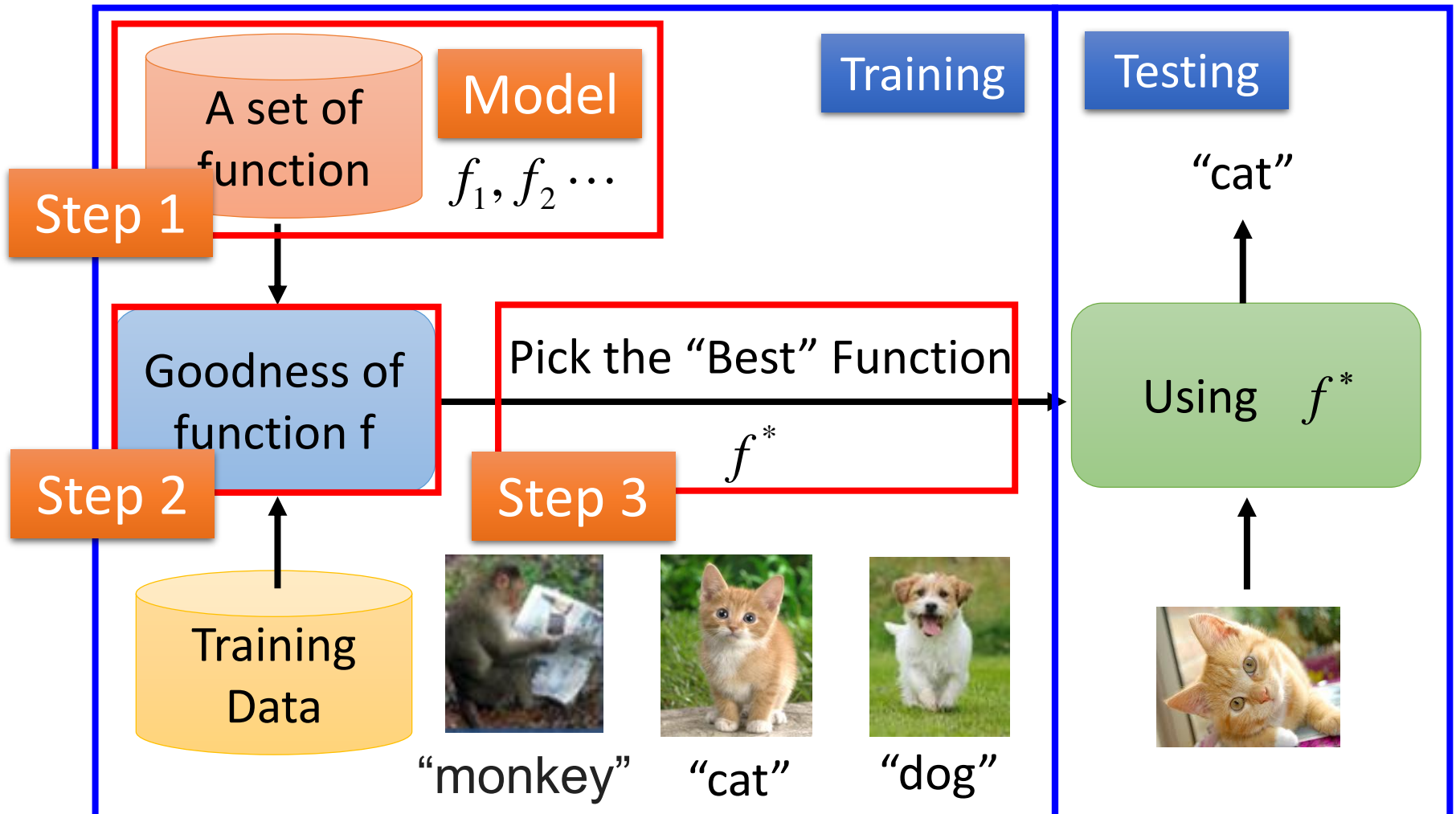
$$f(\text{img\_cat}) = \text{"cat"}$$



# Framework

Image Recognition:

$$f\left(\text{Image of a cat}\right) = \text{"cat"}$$



## II. Three Steps for Deep Learning

Step 1: define a set of function

Neural Network

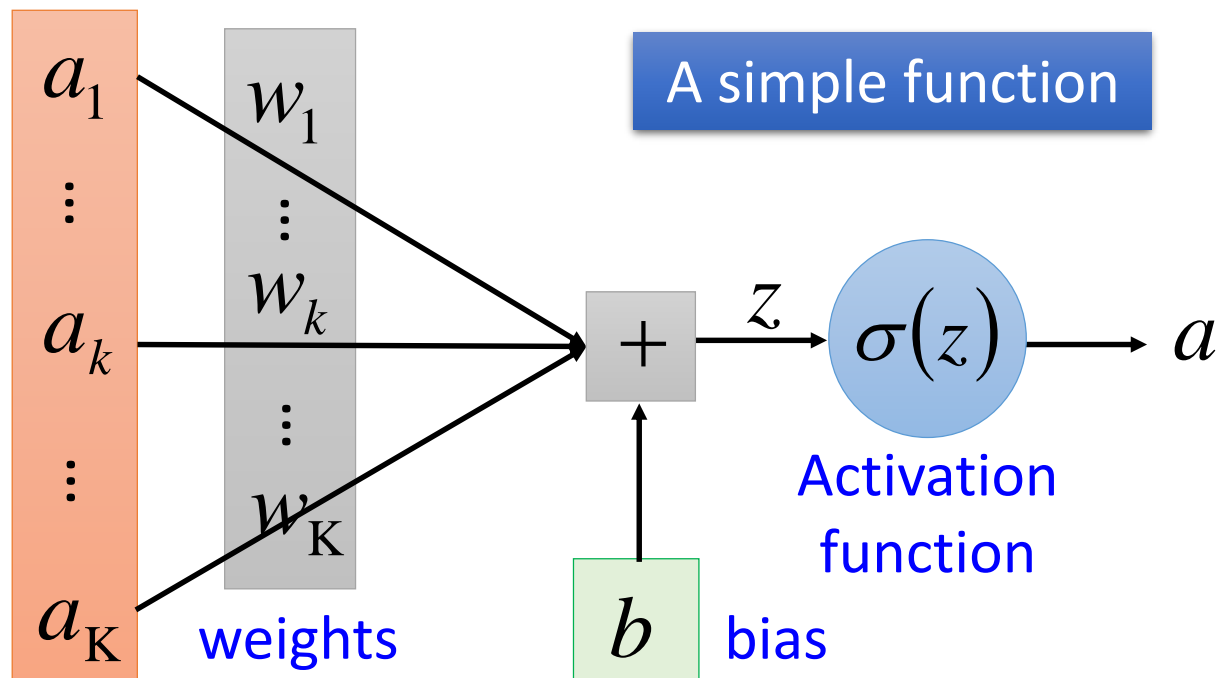
Step 2: goodness of function

Step 3: pick the best function

# Neural Network

## Neuron

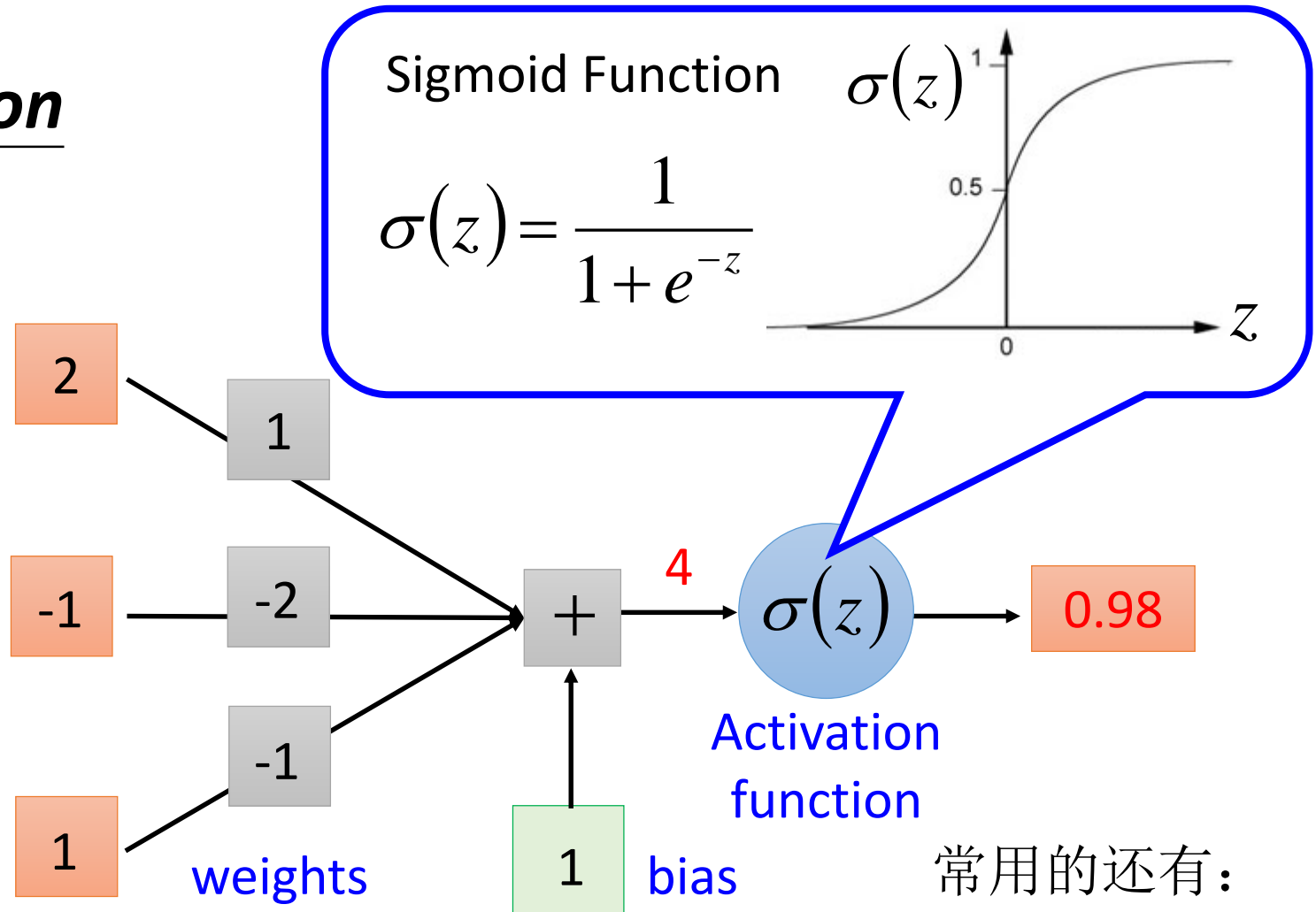
$$z = a_1 w_1 + \dots + a_k w_k + \dots + a_K w_K + b$$





# Neural Network

## Neuron



常用的还有:  
Relu, Tanh

# 神经元与Logistic Regression单元

- 线性回归 (Linear Regression)

- $y = \sum_i w_i x_i + b$

- Logistic回归 (Logistic Regression)

- Logistic函数:  $(-\infty, +\infty)$

$$y = \frac{L}{1 + e^{-k(z - z_0)}}$$

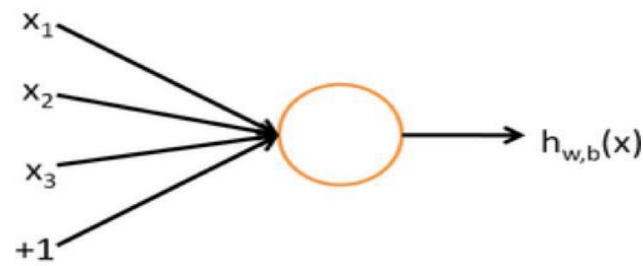
- 设  $z = \sum_i w_i x_i + b$

- Sigmoid函数 (Logistic函数的特例)

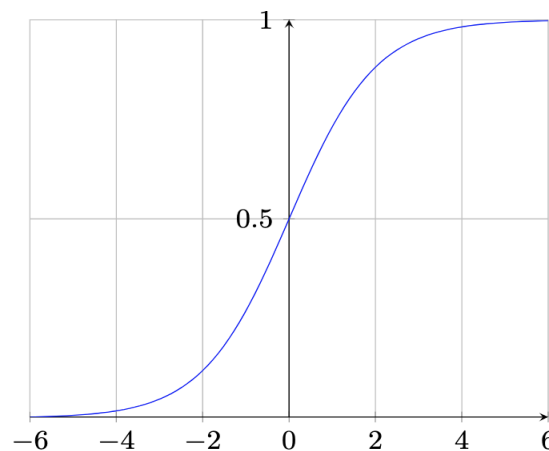
$$y = \frac{1}{1 + e^{-z}}$$

- 处理二元分类问题,  $y$  为输出的概率

- 垃圾邮件过滤、褒贬识别



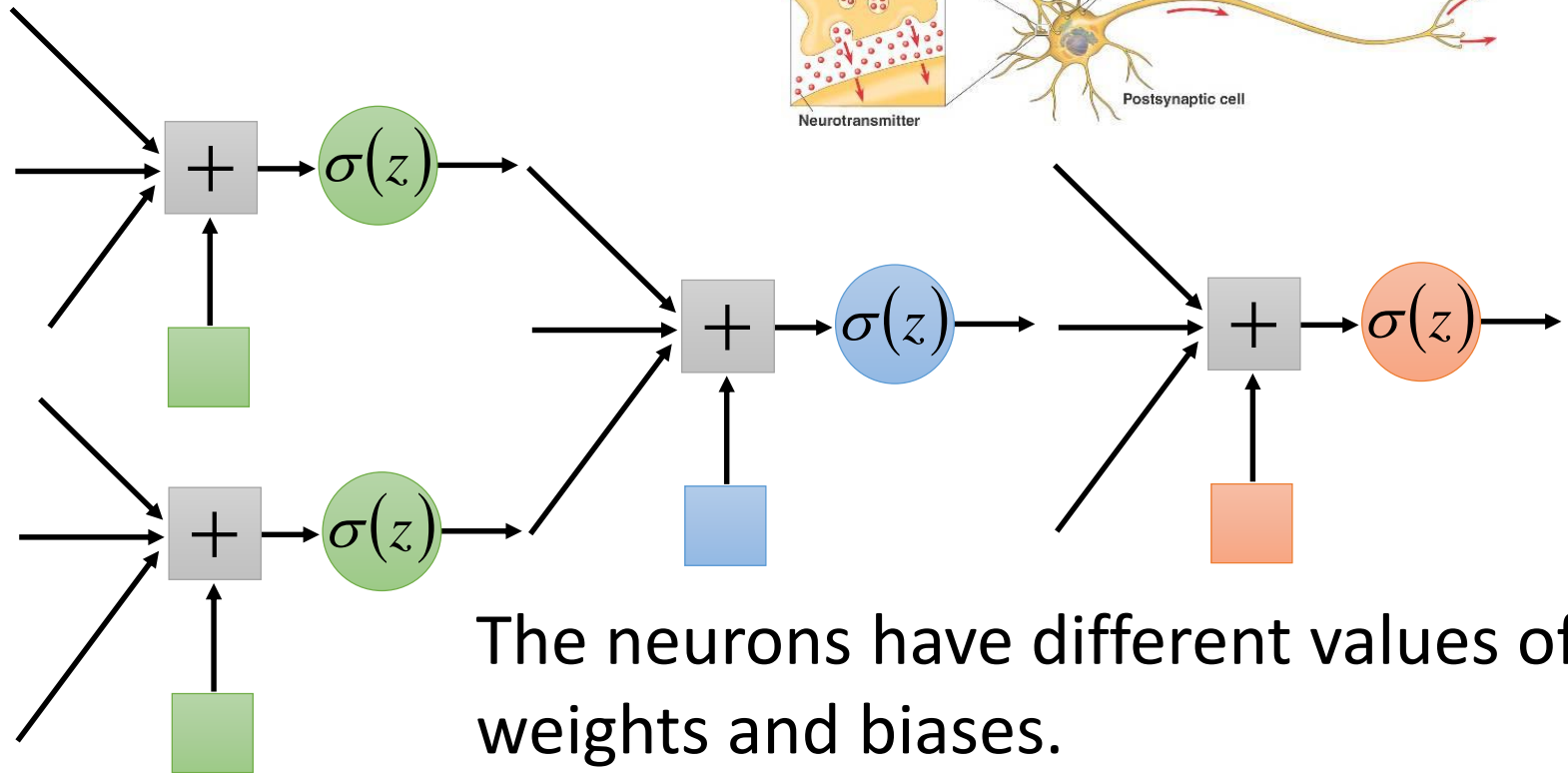
$w, b$  are the parameters of this neuron  
i.e., this logistic regression model



Sigmoid函数示意图

# Neural Network

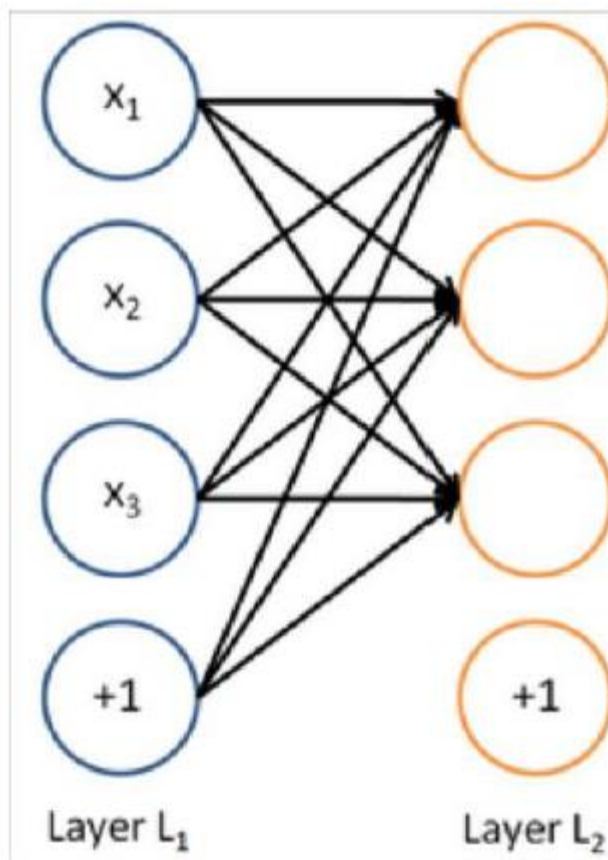
Different connections lead to different network structures



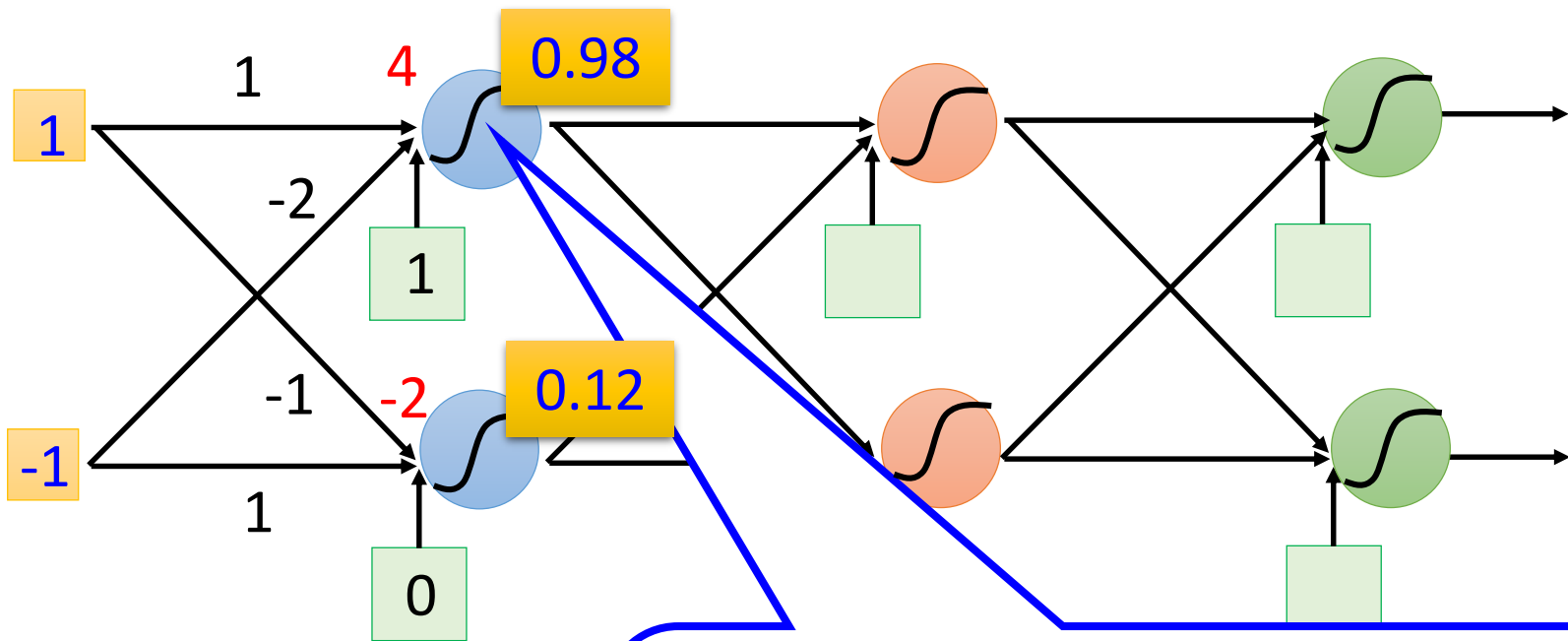
The neurons have different values of weights and biases.

Weights and biases are network parameters  $\theta$

一个神经网络 i.e.  
多个logistic regression单元同时运行

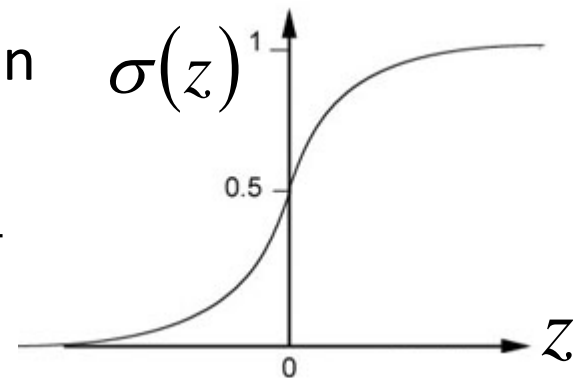


# Fully Connect Feedforward Network

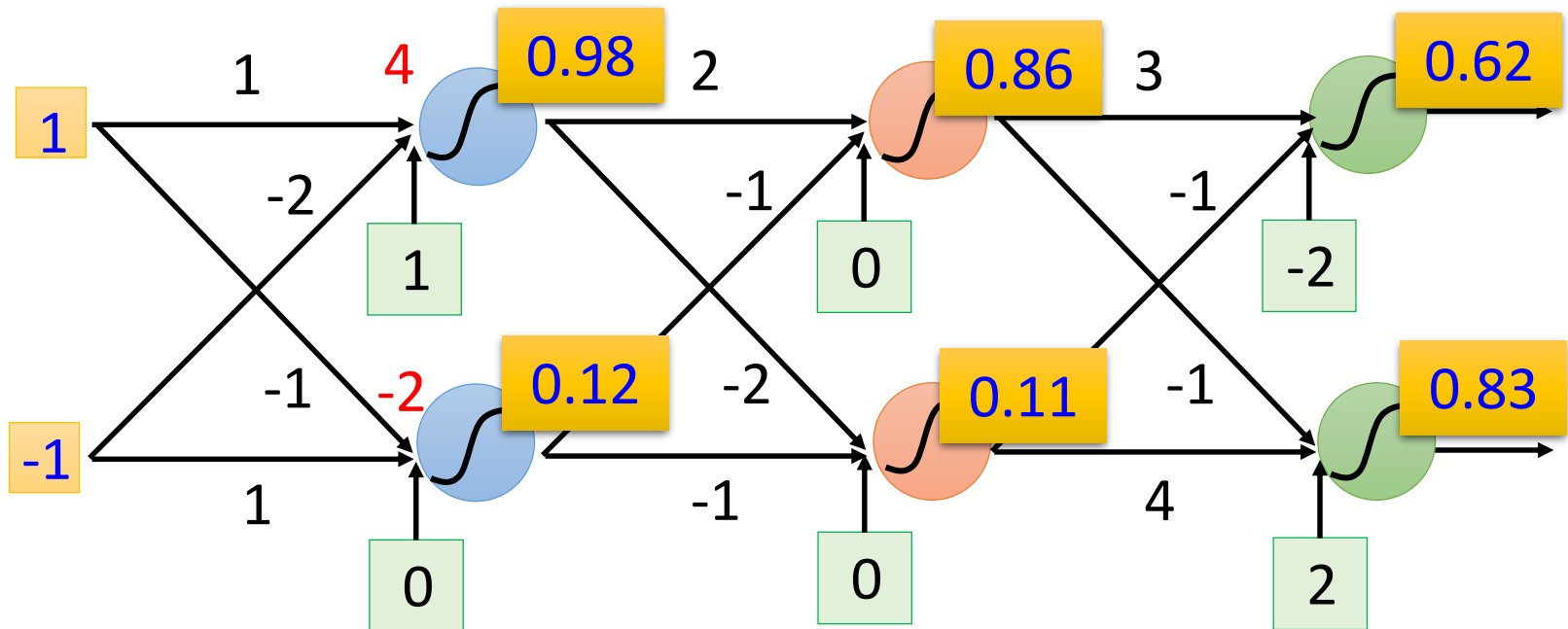


Sigmoid Function

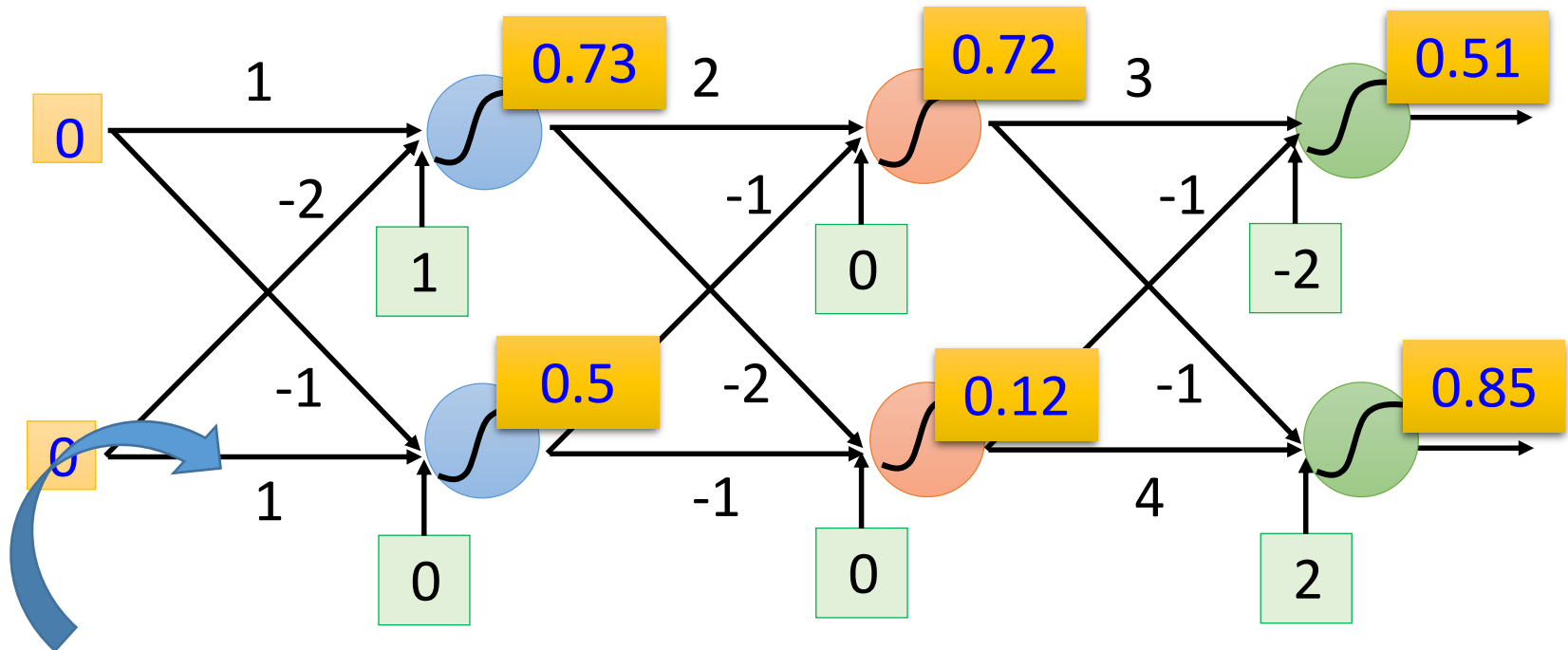
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



# Fully Connect Feedforward Network



# Fully Connect Feedforward Network



This is a function.

Input vector, output vector

$$f\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix} \quad f\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.51 \\ 0.85 \end{bmatrix}$$

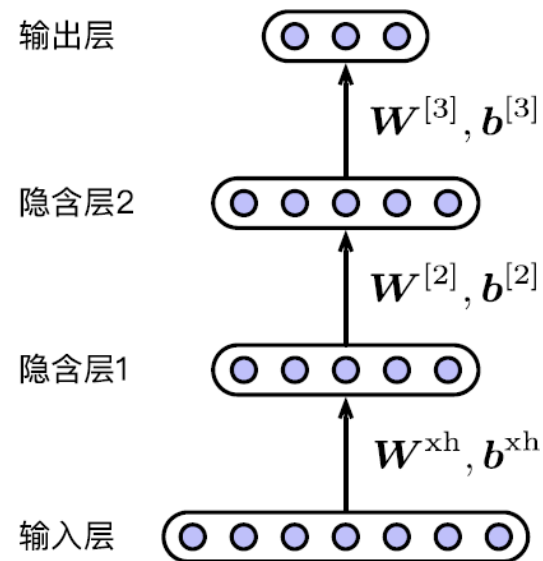
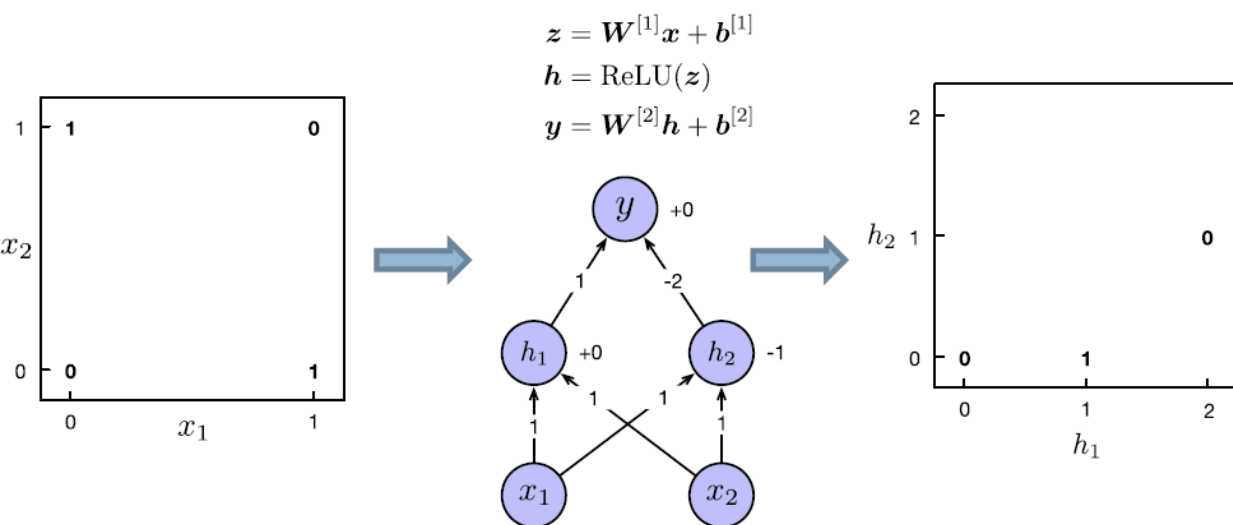
Given parameters  $\theta$ , define a function

Given network structure, define a function set

# 神经网络特长1

□ 可解决线性不可分问题

▣ 如：异或 (XOR) 问题





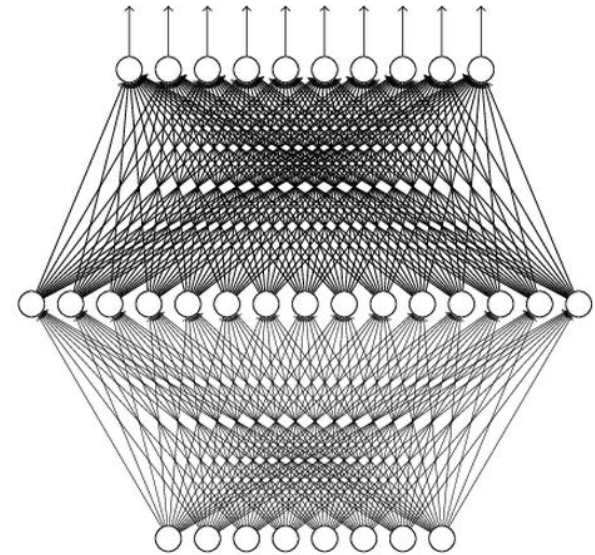
# 神经网络特长2: Universality Theorem

Any continuous function  $f$

$$f : R^N \rightarrow R^M$$

Can be realized by a network  
with one hidden layer

(given **enough** hidden  
neurons)

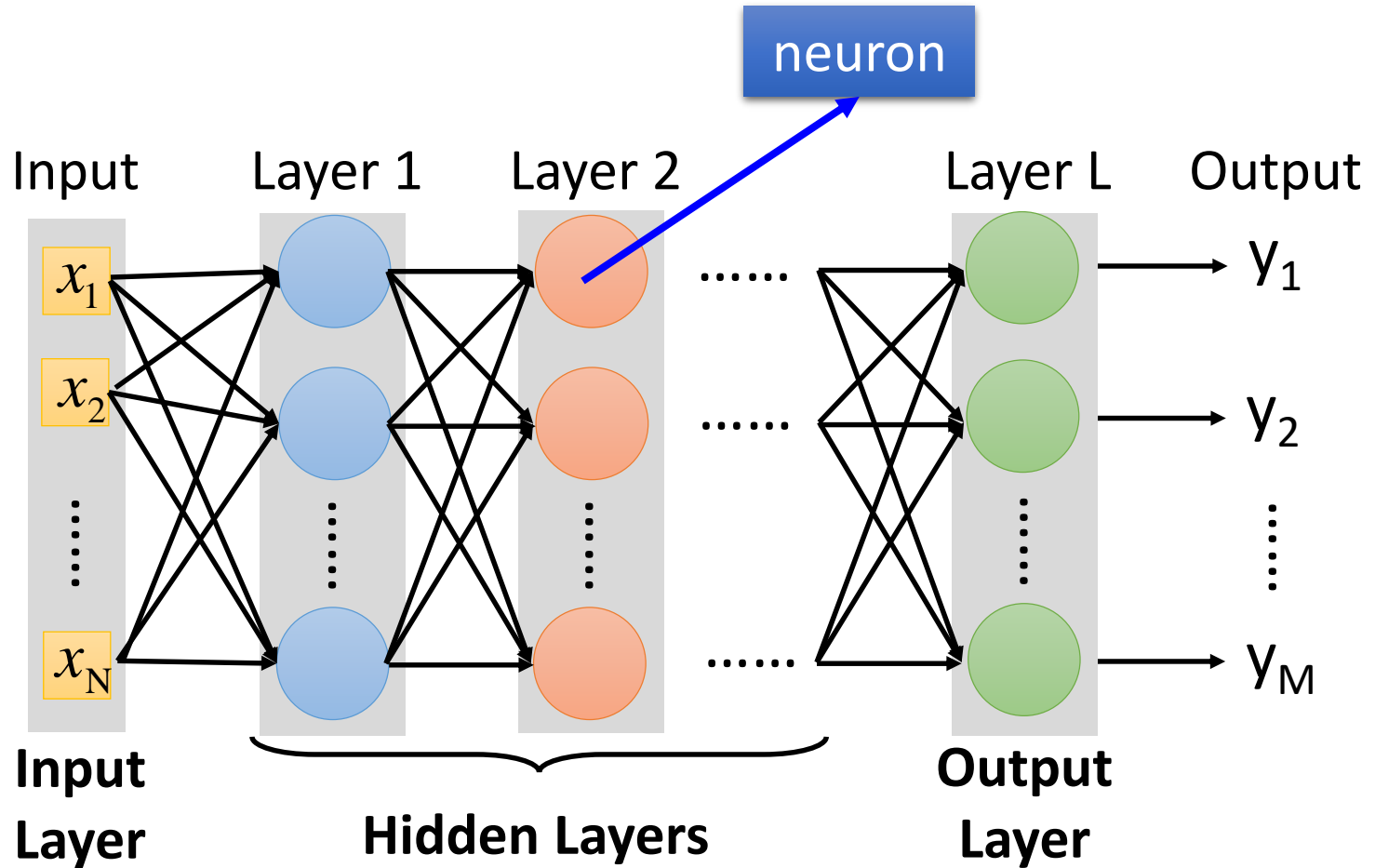


Reference for the reason:

<http://neuralnetworksanddeeplearning.com/chap4.html>

通用近似定理：即使是二层神经网络，都可以以任意精度近似任意一个连续函数。

# Deep means many hidden layers



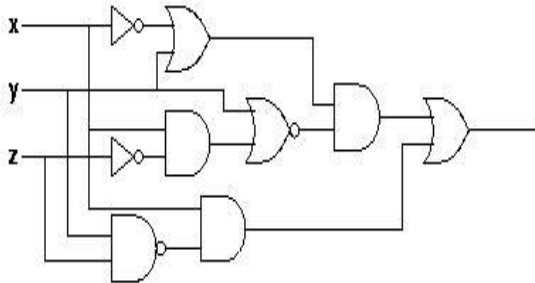
# Why Deep? Analogy

## Logic circuits

- Logic circuits consists of **gates**
- **A two layers of logic gates** can represent **any Boolean function**.
- Using multiple layers of logic gates to build some functions are much simpler



less gates needed



## Neural network

- Neural network consists of **neurons**
- **A hidden layer network** can represent **any continuous function**.
- Using multiple layers of neurons to represent some functions are much simpler



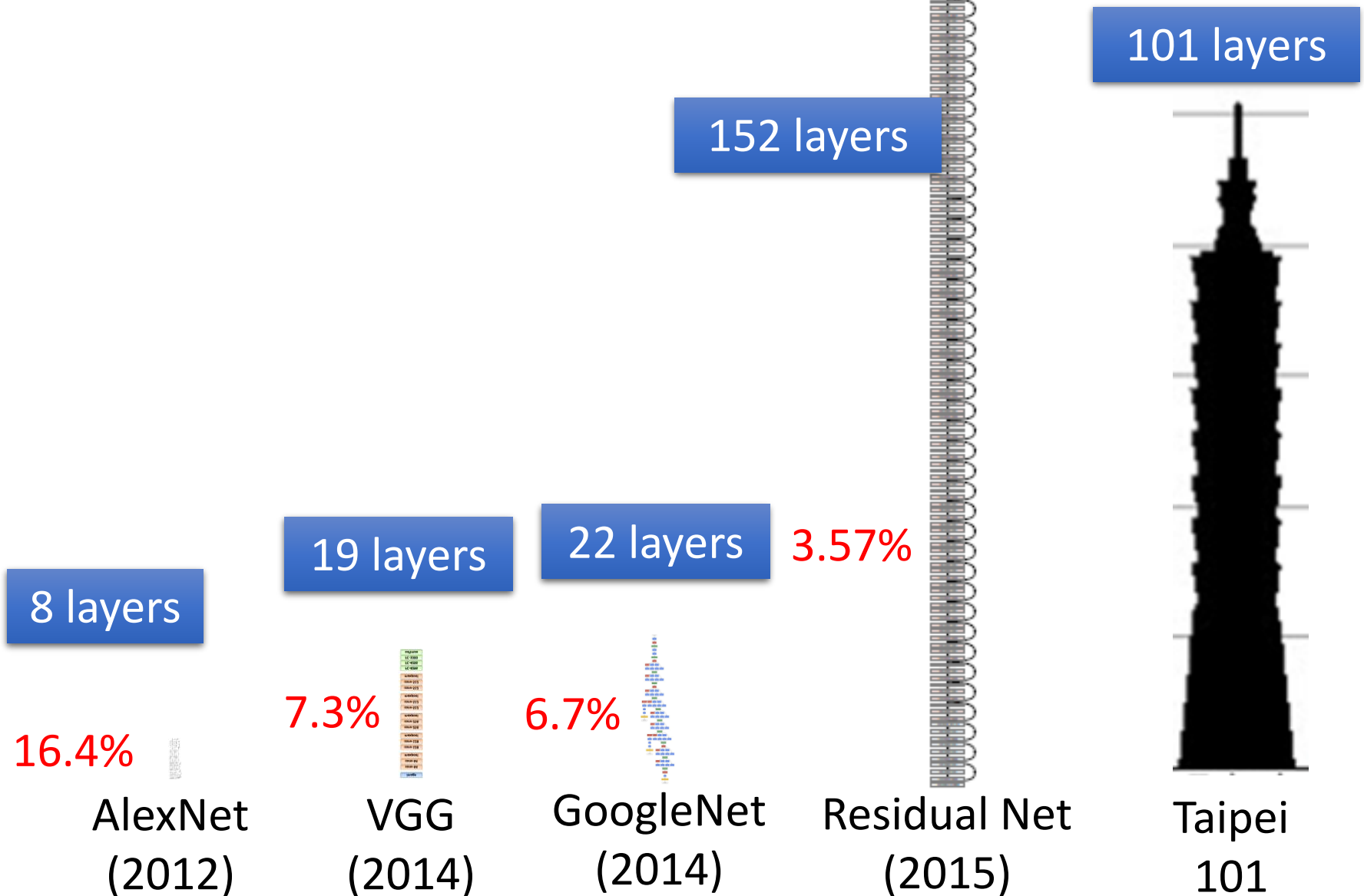
less parameters



less data?

定理 2. 存在这样的布尔函数，可以由深度为  $k$  的多项式复杂度的逻辑门电路表示，等价的深度为  $k-1$  的逻辑门电路变为指数复杂度。这里深度指从输入到输出的最长路径的长度，复杂度指逻辑门的个数。

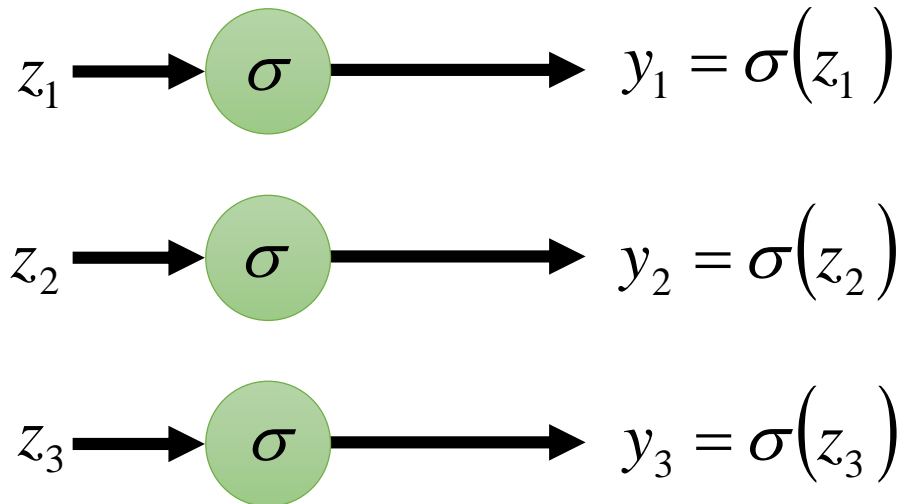
# Deep = Many hidden layers



# Output Layer

- Softmax layer as the output layer

## Ordinary Layer



In general, the output of network can be any value.

May not be easy to interpret

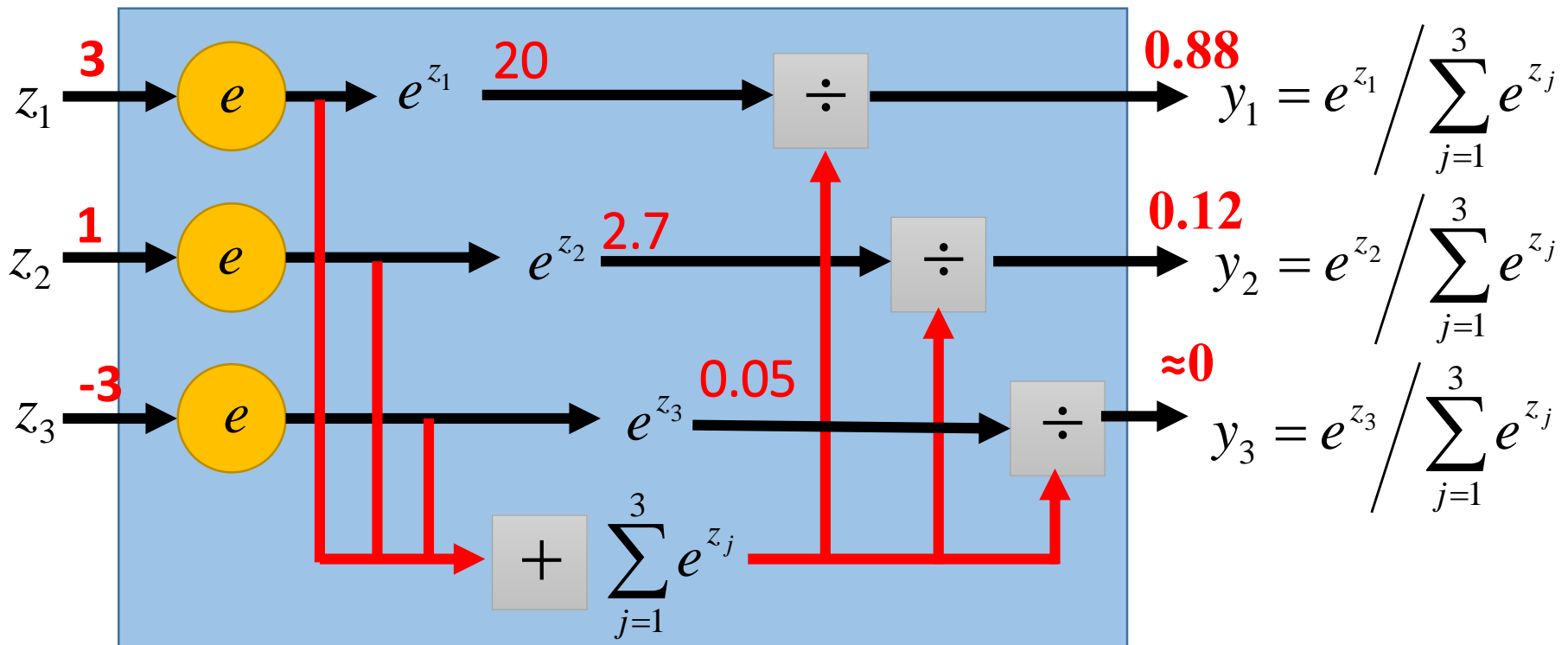
# Output Layer

- Softmax layer as the output layer

**Probability:**

- $1 > y_i > 0$
- $\sum_i y_i = 1$

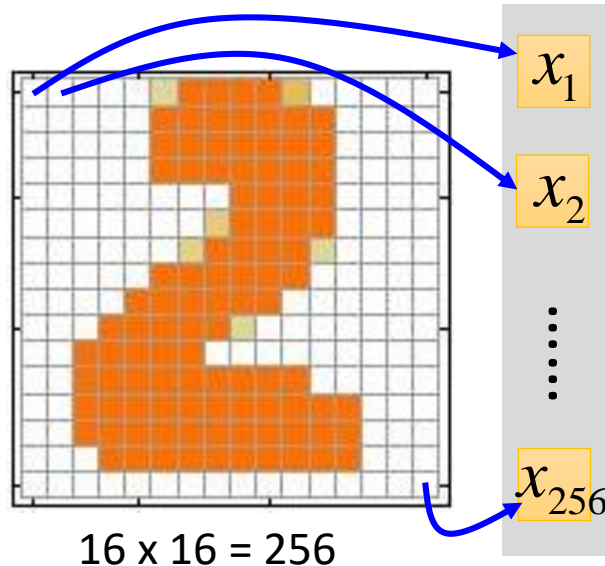
## Softmax Layer



# Example Application

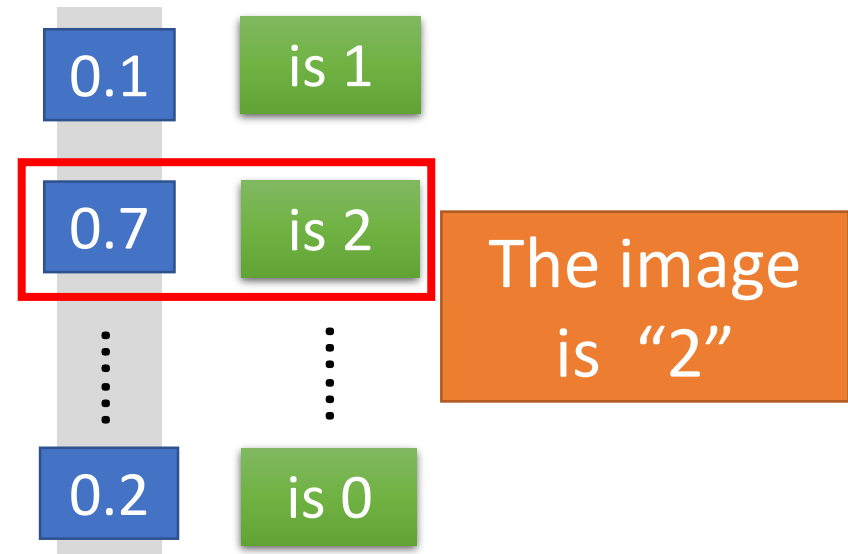


## Input



with Ink  $\rightarrow 1$   
w/o ink  $\rightarrow 0$

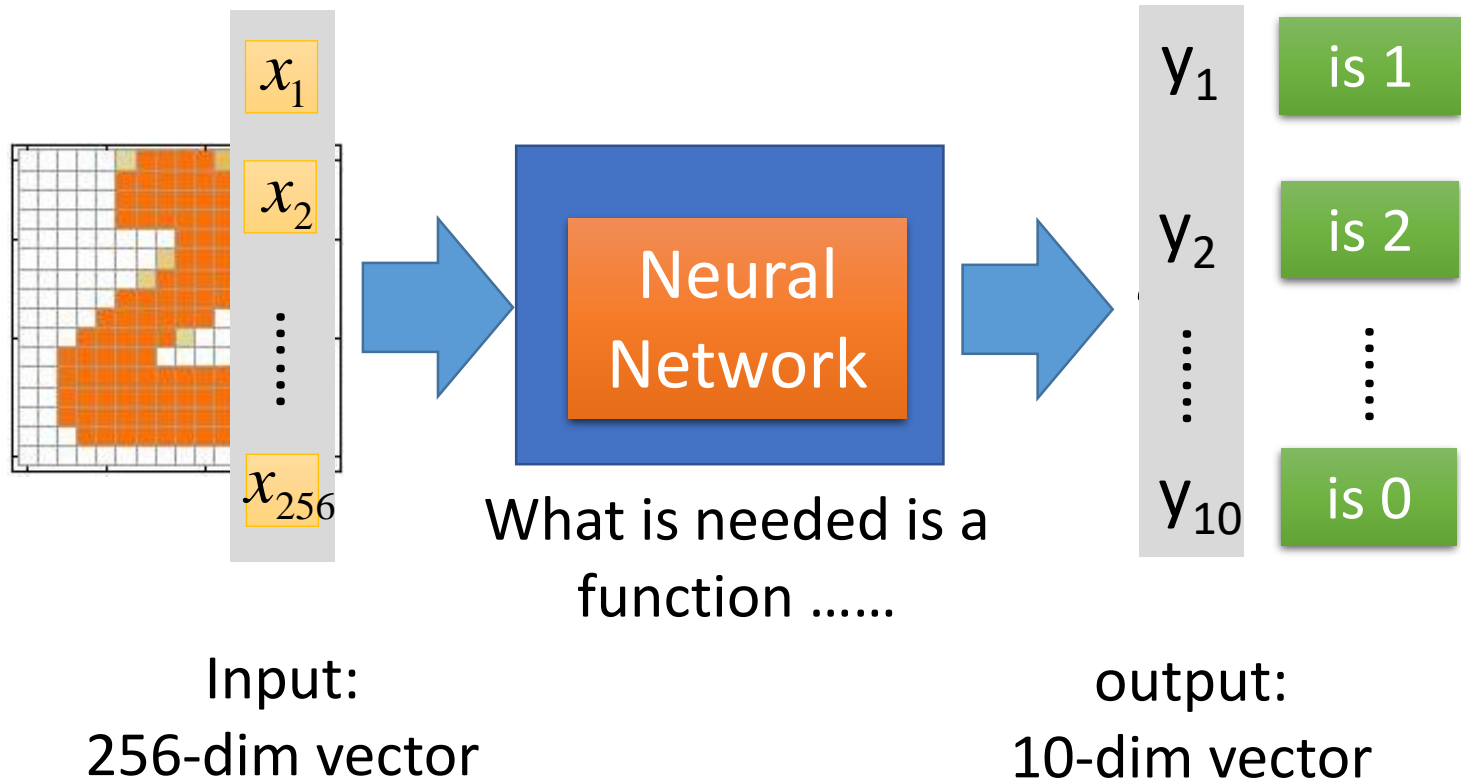
## Output



Each dimension represents the confidence of a digit.

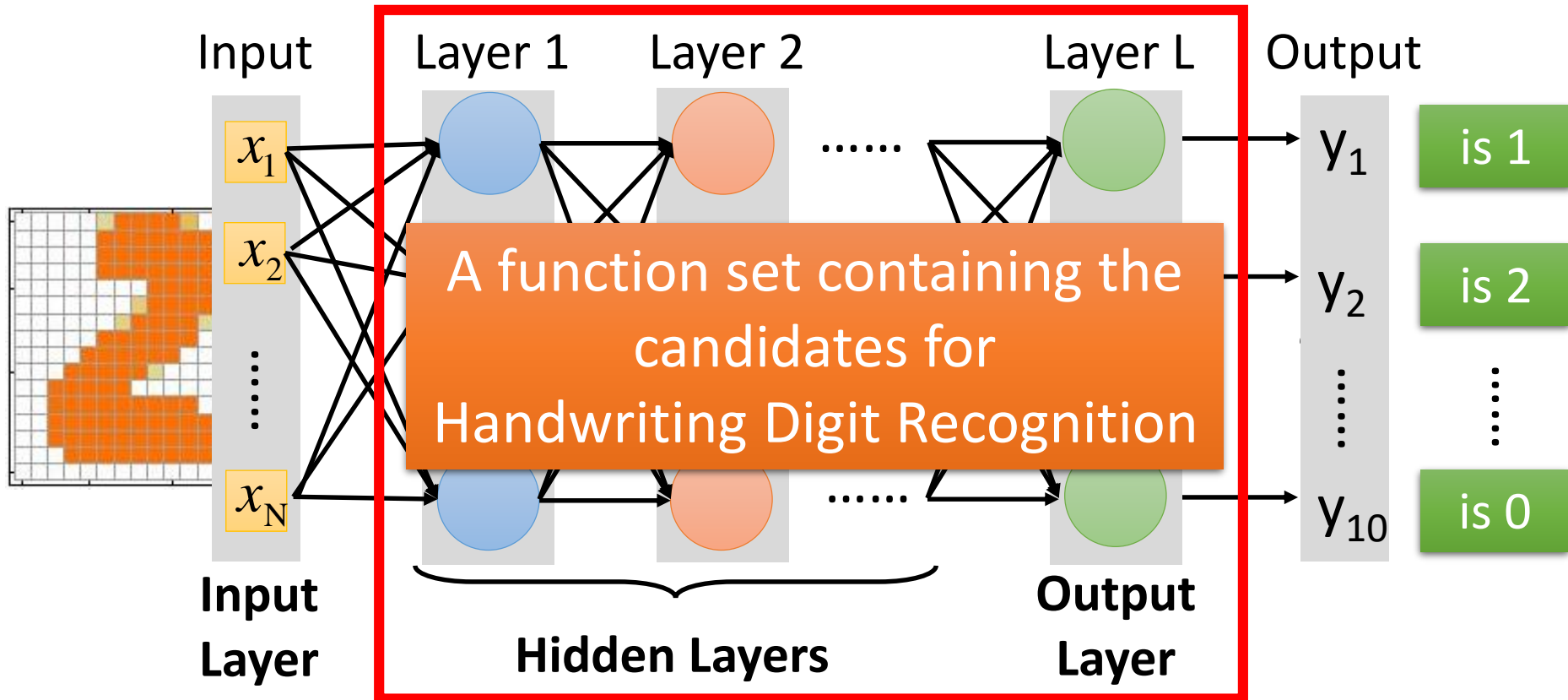
# Example Application

- Handwriting Digit Recognition





# Example Application



You need to decide the network structure to let a good function in your function set.

# 小结：损失函数

- 损失函数（Loss Function）
  - 评估参数好坏的**准则**
  - 为什么不直接使用**准确率**等指标进行评估？
- 两种常用的损失函数
  - **均方误差**（Mean Squared Error, MSE）

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

- **交叉熵**（Cross-Entropy, CE）

$$\text{CE} = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^c y_j^{(i)} \log \hat{y}_j^{(i)}$$



$$\text{CE} = -\frac{1}{m} \sum_{i=1}^m \log \hat{y}_t^{(i)}$$

**负对数似然损失**  
Negative Log  
Likelihood, NLL

# Three Steps for Deep Learning

Step 1: define a set of function



Step 2: goodness of function

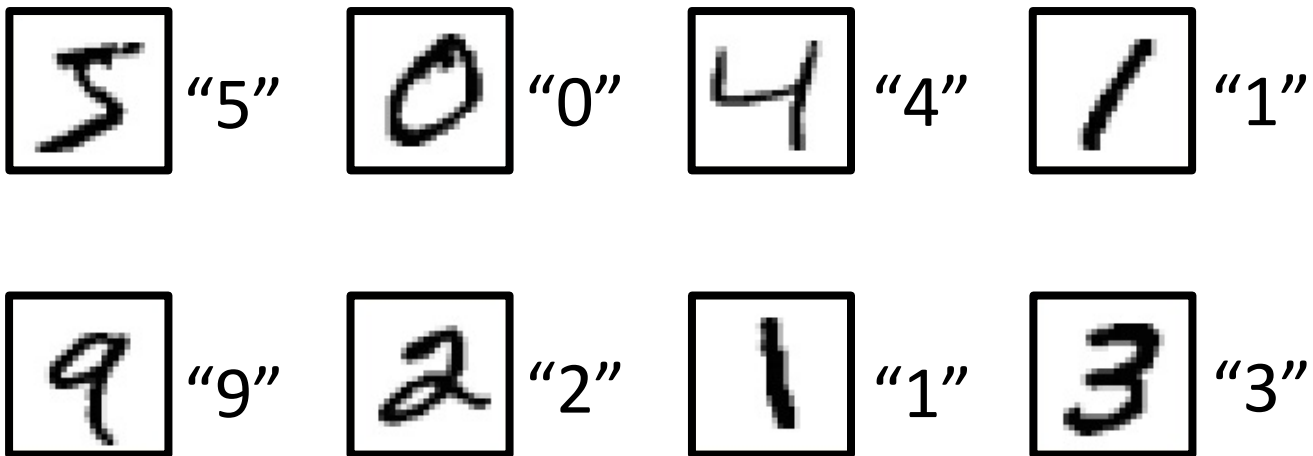


Step 3: pick the best function



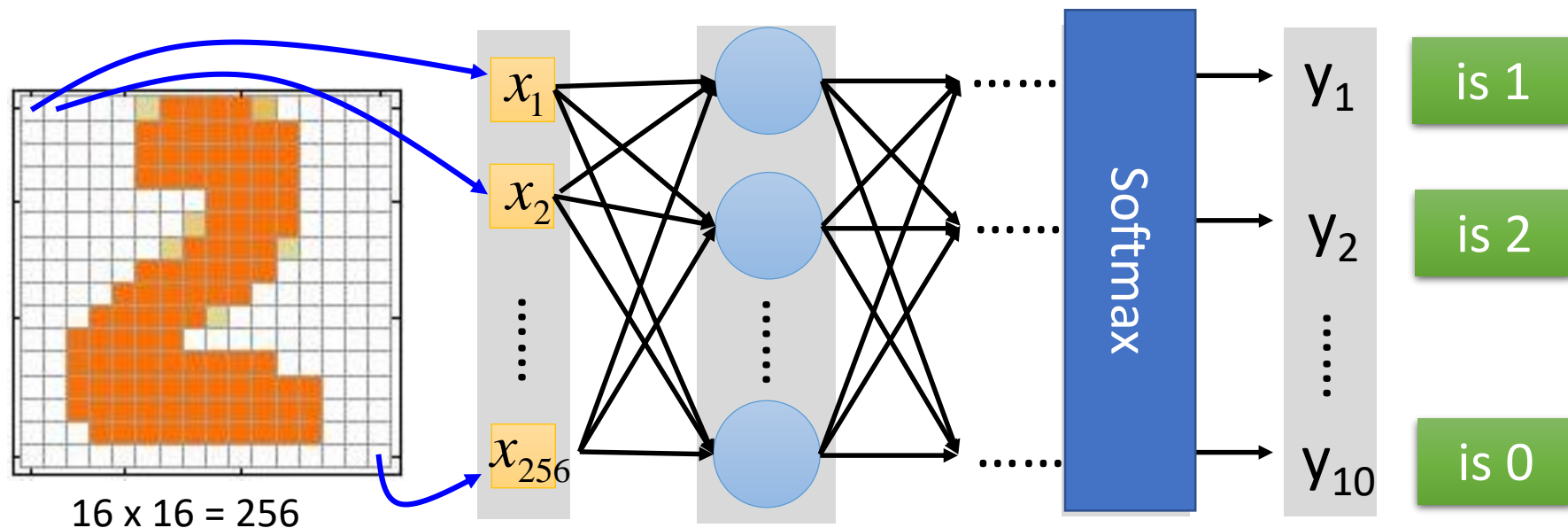
# Training Data

- Preparing training data: images and their labels



The learning target is defined on the training data.


# Learning Target




Ink  $\rightarrow$  1

No ink  $\rightarrow$  0

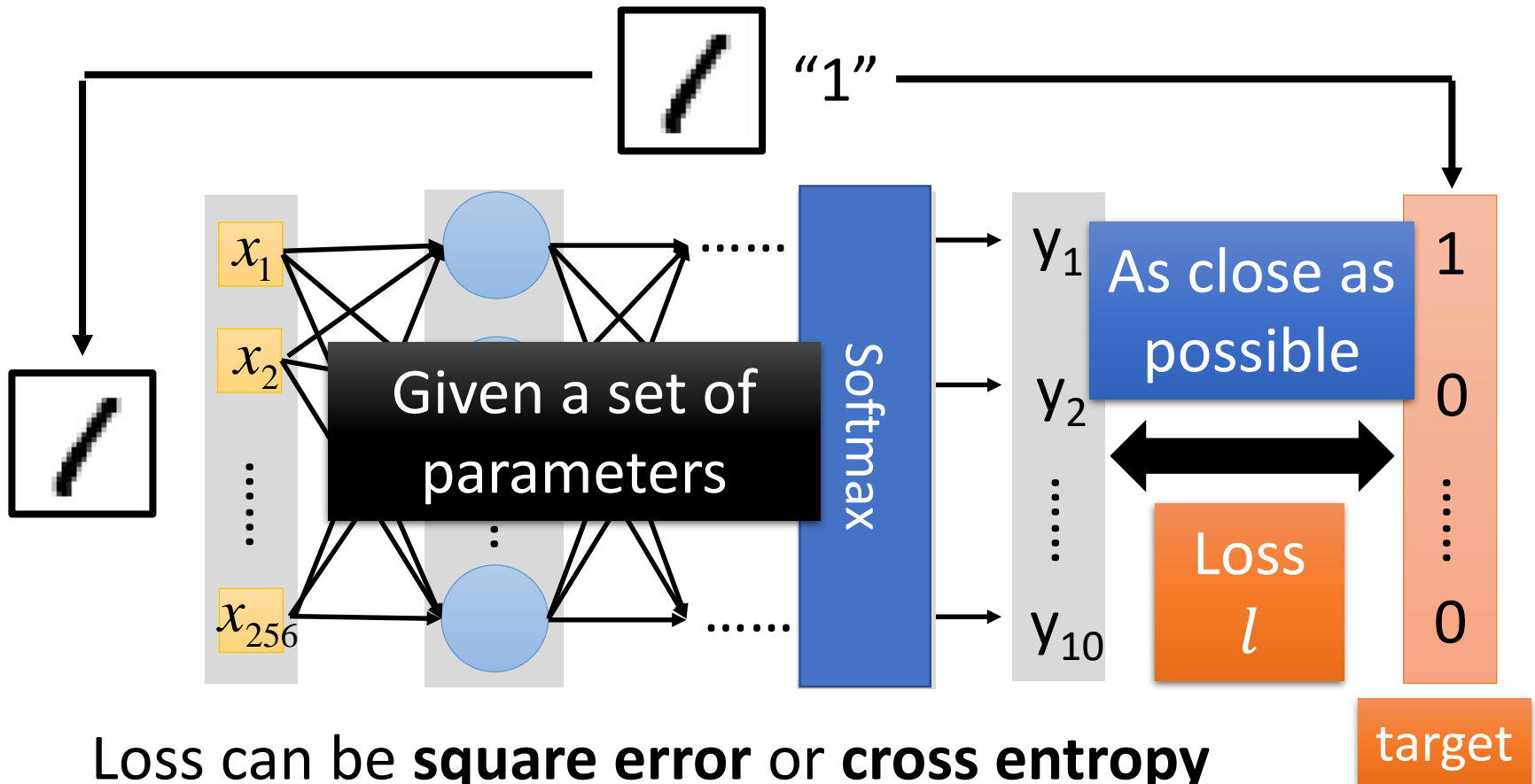
The learning target is .....

Input:   $\rightarrow$   $y_1$  has the maximum value

Input:   $\rightarrow$   $y_2$  has the maximum value

# Loss

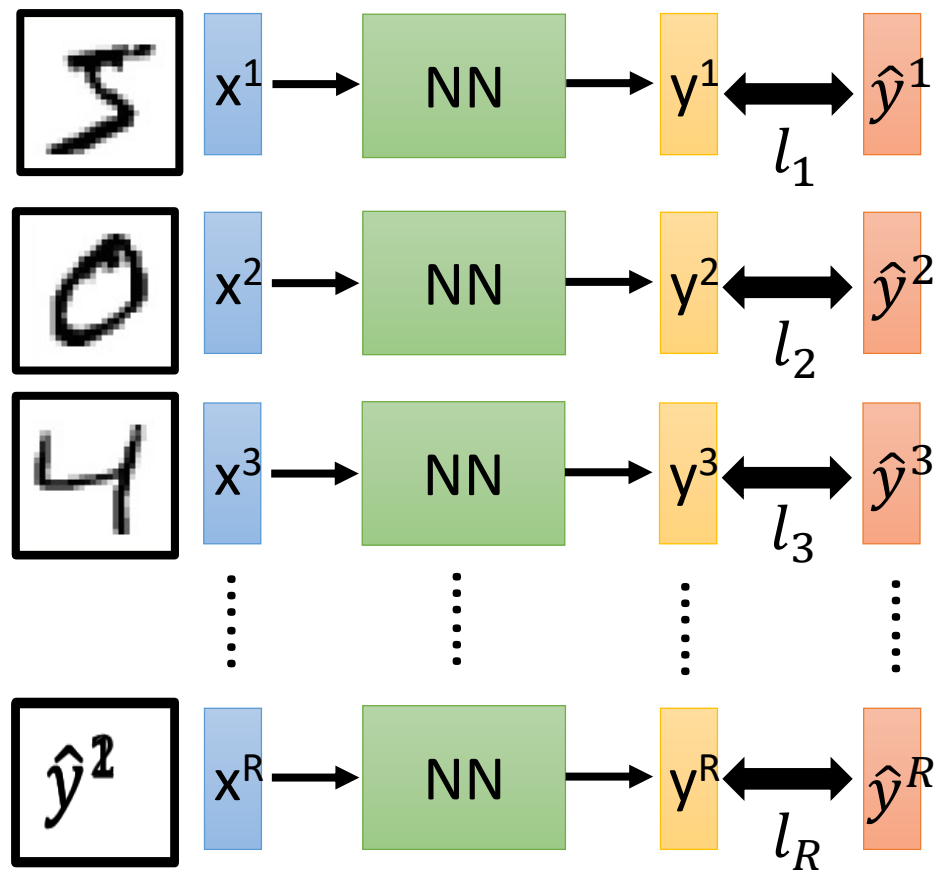
A good function should make the loss of all examples as small as possible.



Loss can be **square error** or **cross entropy** between the network output and target

# Total Loss

For all training data ...



Total Loss:

$$L = \sum_{r=1}^R l_r$$

As small as possible

Find a function in function set that minimizes total loss  $L$

Find the network parameters  $\theta^*$  that minimize total loss  $L$

# Three Steps for Deep Learning

Step 1: define a set of function



Step 2: goodness of function



Step 3: pick the best function





# How to pick the best function

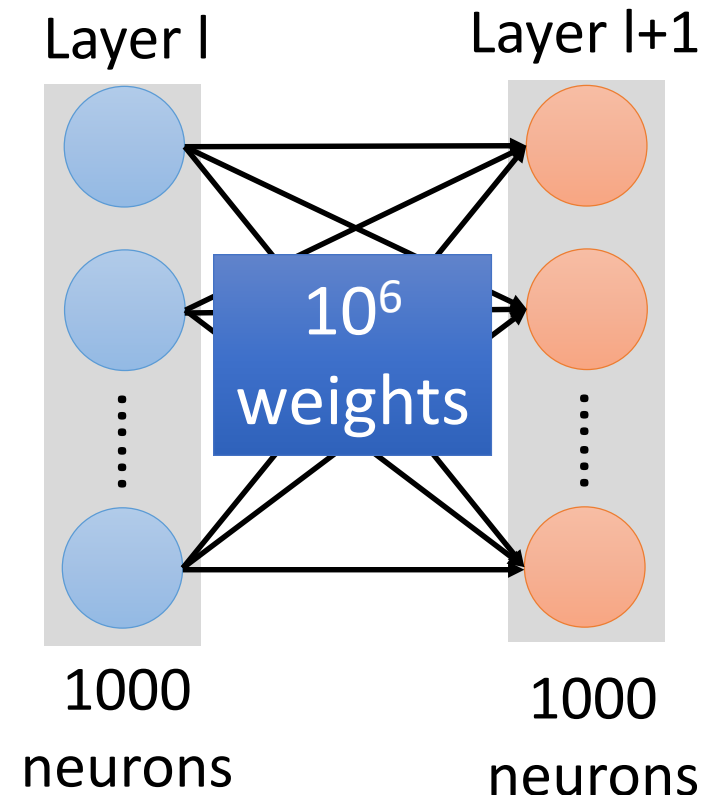
Find network parameters  $\theta^*$  that minimize total loss  $L$

Enumerate all possible values

Network parameters  $\theta =$   
 $\{w_1, w_2, w_3, \dots, b_1, b_2, b_3, \dots\}$

Millions of parameters

E.g. speech recognition: 8 layers and  
1000 neurons each layer



# Gradient Descent

Network parameters  $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

Find network parameters  $\theta^*$  that minimize total loss  $L$

- Pick an initial value for  $w$
- Compute  $\partial L / \partial w$

Total  
Loss  $L$

Negative

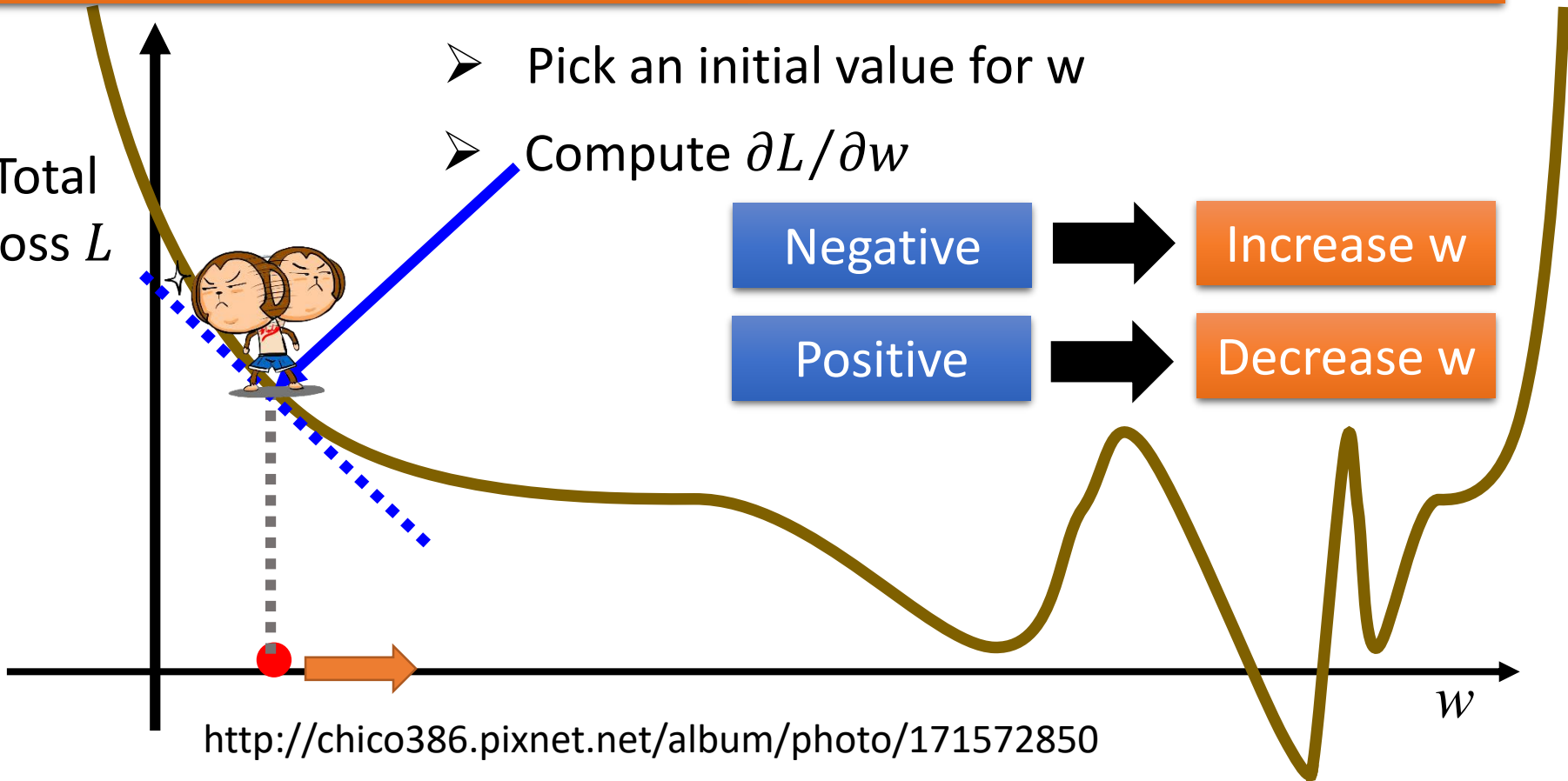


Increase  $w$

Positive



Decrease  $w$



# Gradient Descent

Network parameters  $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

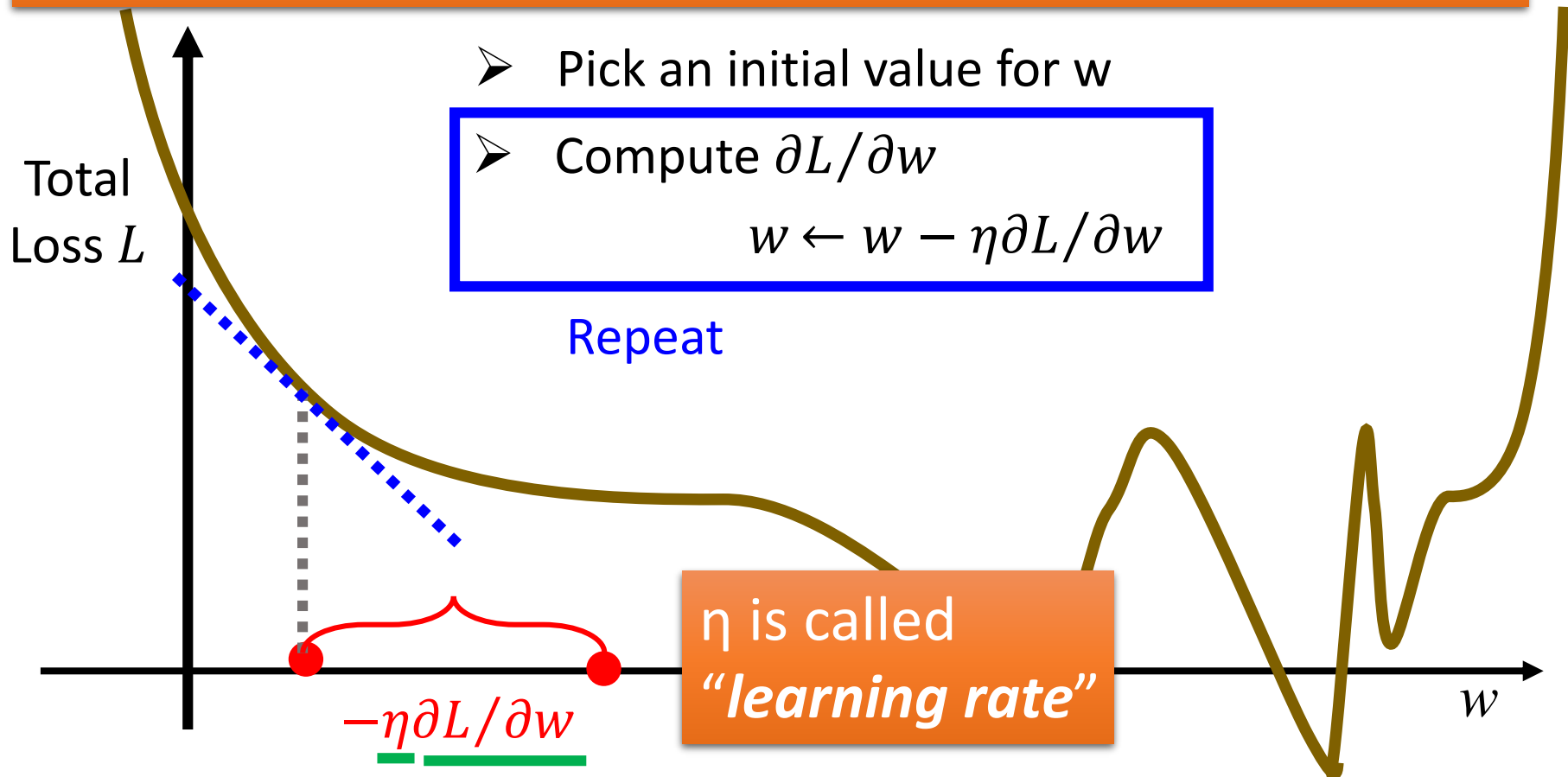
Find network parameters  $\theta^*$  that minimize total loss  $L$

➤ Pick an initial value for  $w$

➤ Compute  $\partial L / \partial w$

$$w \leftarrow w - \eta \partial L / \partial w$$

Repeat



# Gradient Descent

Network parameters  $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

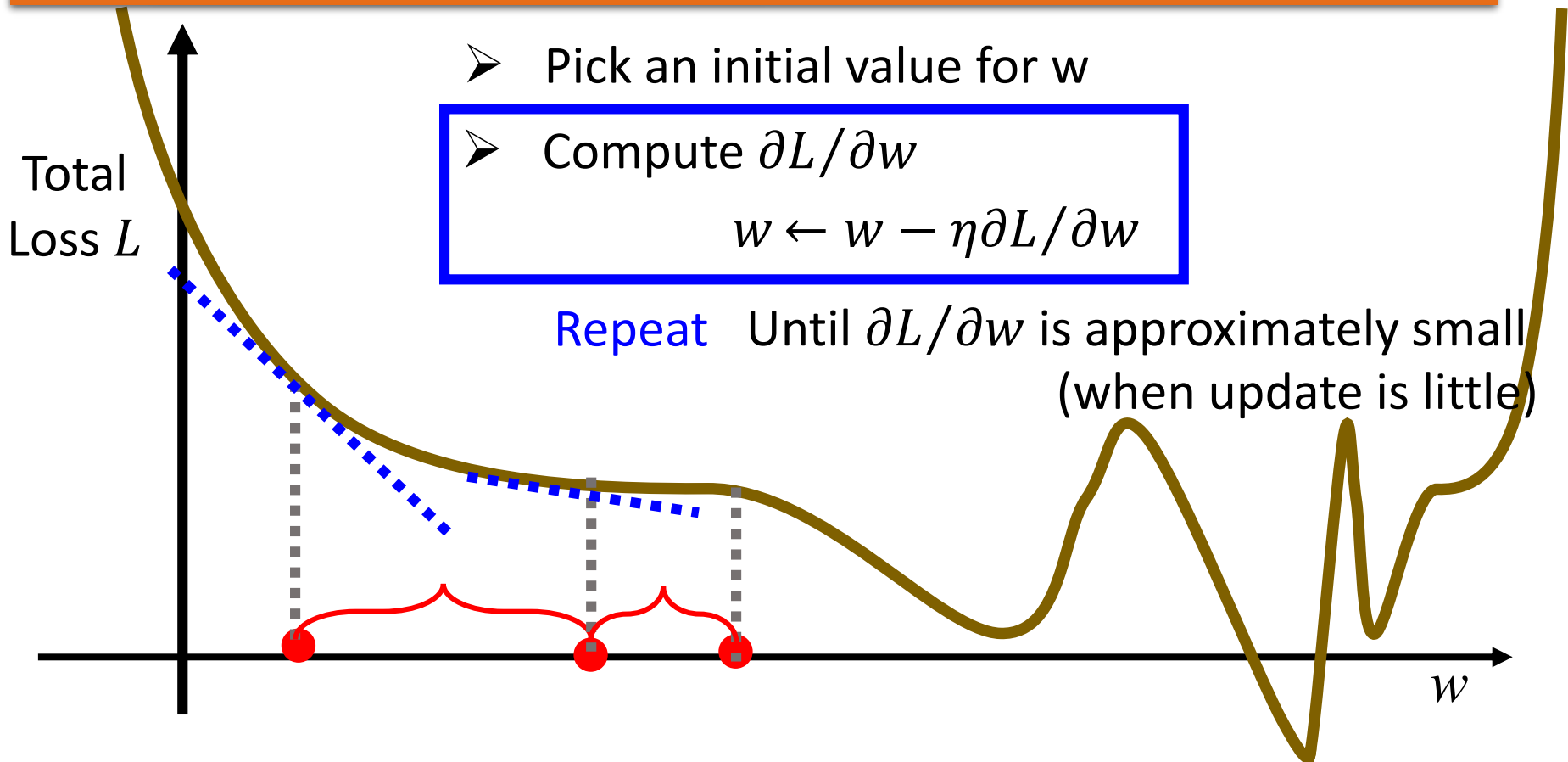
Find network parameters  $\theta^*$  that minimize total loss  $L$

➤ Pick an initial value for  $w$

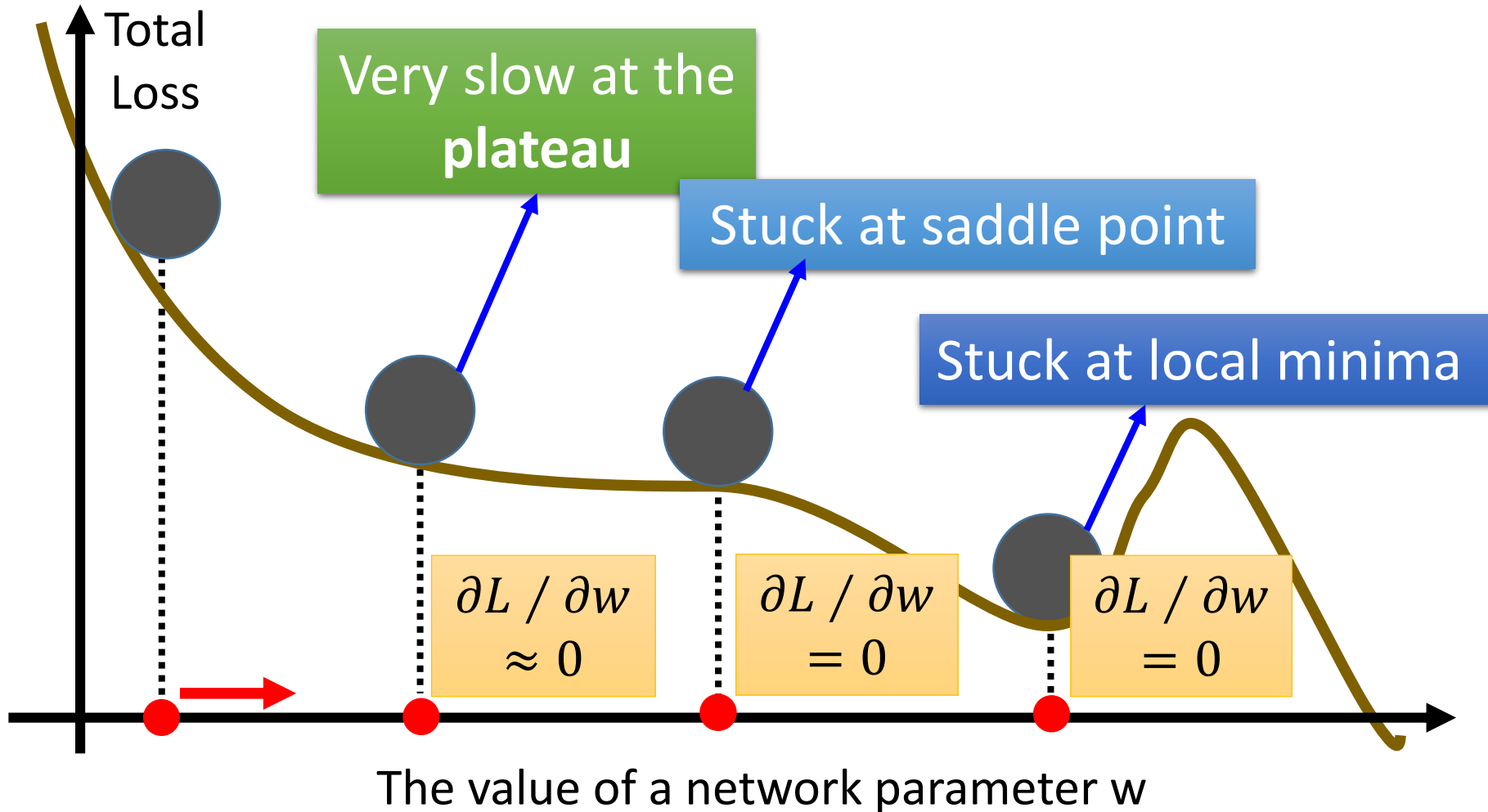
➤ Compute  $\partial L / \partial w$

$$w \leftarrow w - \eta \partial L / \partial w$$

Repeat Until  $\partial L / \partial w$  is approximately small  
(when update is little)

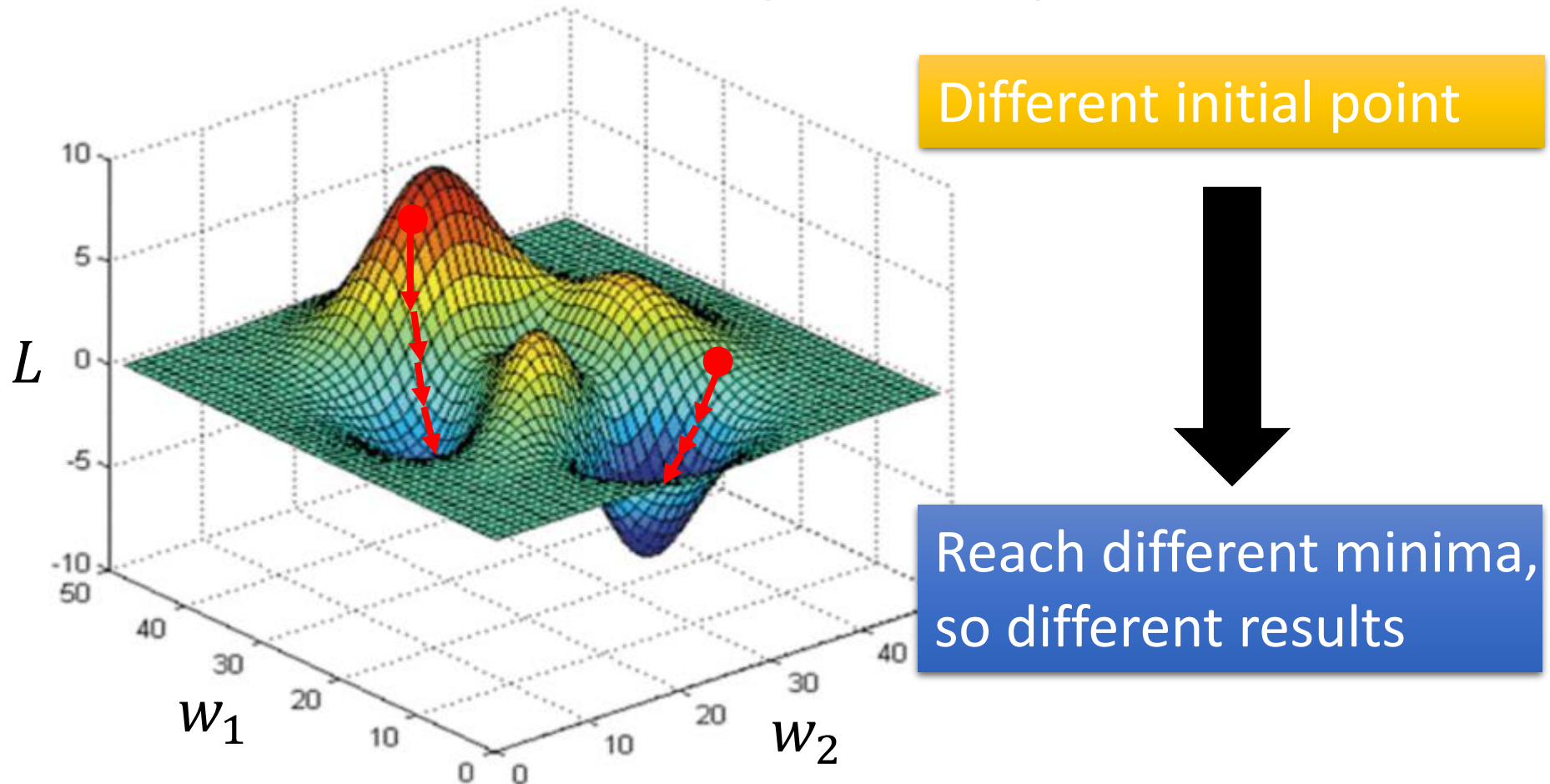


# Local Minima



# Local Minima

- Gradient descent never guarantee global minima



# 小结：梯度下降

- 梯度（Gradient）

- 以向量的形式写出的对多元函数各个参数求得的偏导数
- 是函数值增加最快的方向
- 沿着梯度相反的方向，更加容易找到函数的极小值

- 梯度下降算法

（Gradient Descent, GD）

---

**算法 4.1** 梯度下降算法

---

**Input:** 学习率  $\alpha$ ; 含有  $m$  个样本的训练数据

**Output:** 优化参数  $\theta$

1. 设置损失函数为  $L(f(x; \theta), y)$ ;
  2. 初始化参数  $\theta$ 。
  3. **while** 未达到终止条件 **do**
  4.     计算梯度  $g = \frac{1}{m} \nabla_{\theta} \sum_i^m L(f(x^{(i)}; \theta), y^{(i)})$ ;
  5.      $\theta = \theta - \alpha g$ 。
  6. **end**
- 

- 小批次梯度下降法（Mini-batch Gradient Descent）

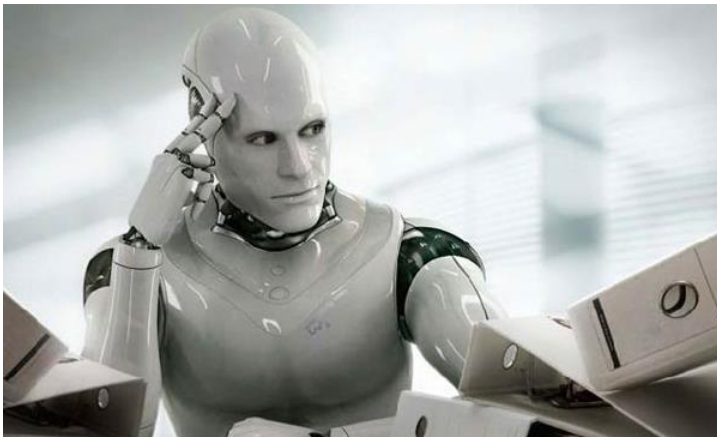
- 每次随机采样小规模的训练数据来估计梯度
- 提高算法的运行速度

# Gradient Descent

This is the “learning” of machines in deep learning .....

➡ Even alpha go using this approach.

People imagine .....



Actually .....



I hope you are not too disappointed :p



# Three Steps for Deep Learning



Deep Learning is so simple .....

Now If you want to find a function

If you have lots of function input/output (?) as training data

 You can use deep learning

Next

Network structures!