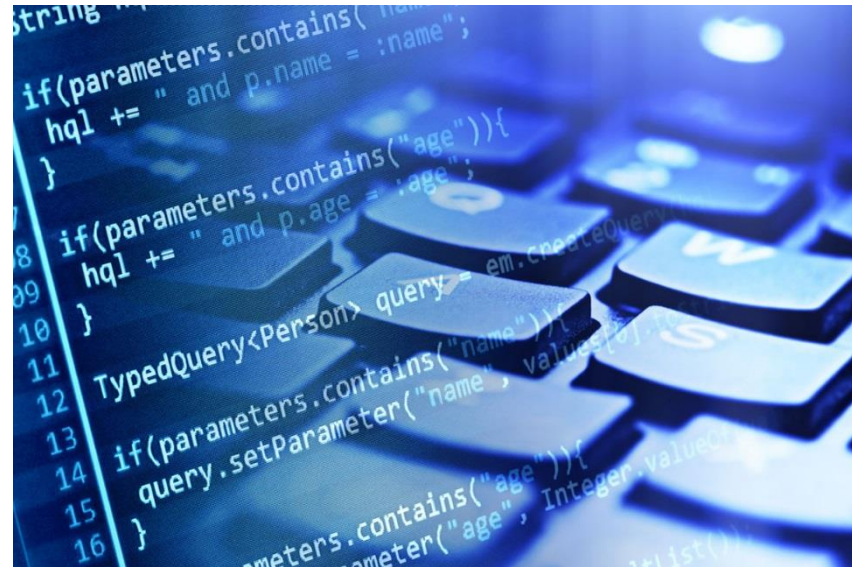


What is  
Machine Learning?

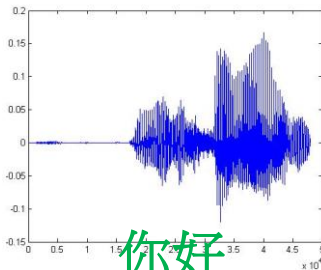
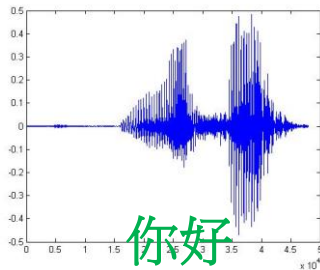
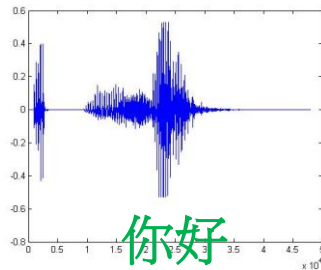
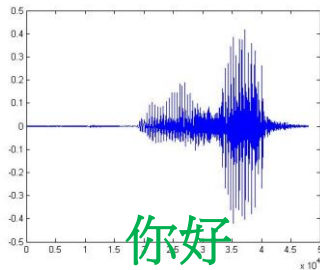
# You know how to program ...



- You can ask computers to do lots of things for you.
- However, computer can only do what you ask it to do.
- Computer can never solve the problem you can't solve.

# Some tasks are very complex

- One day, you are asked to write a program for speech recognition.

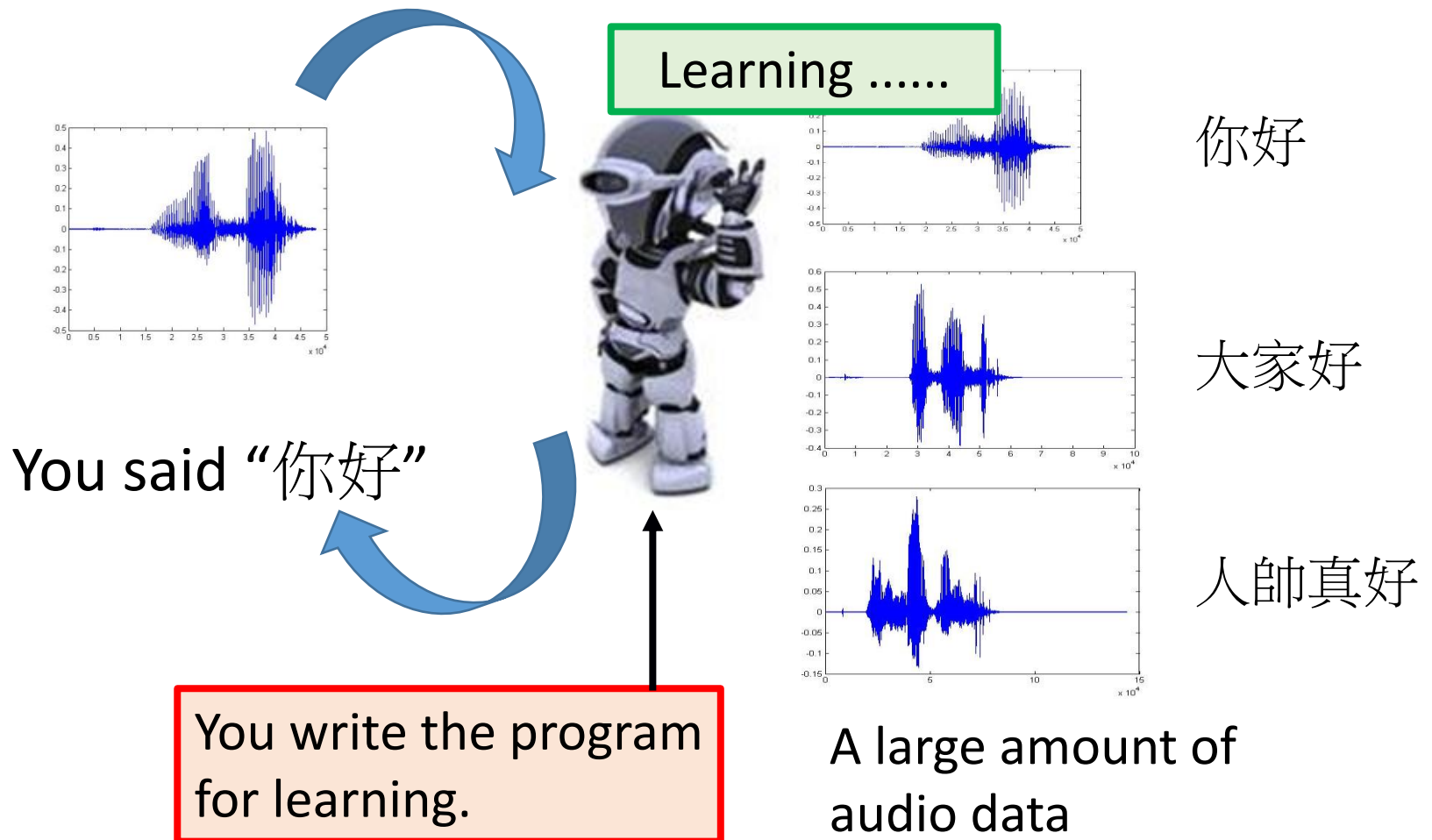


Find the common patterns from the left waveforms.

You quickly get lost in the exceptions and special cases.

It seems impossible to write a program for speech recognition.

# Let the machine learn by itself



# Learning $\approx$ Looking for a Function

- Speech Recognition

$$f\left(\text{[Waveform of '你好' (Hello)]}\right) = \text{“你好”}$$

- Handwritten Recognition

$$f\left(\text{[Handwritten digit '2']}\right) = \text{“2”}$$

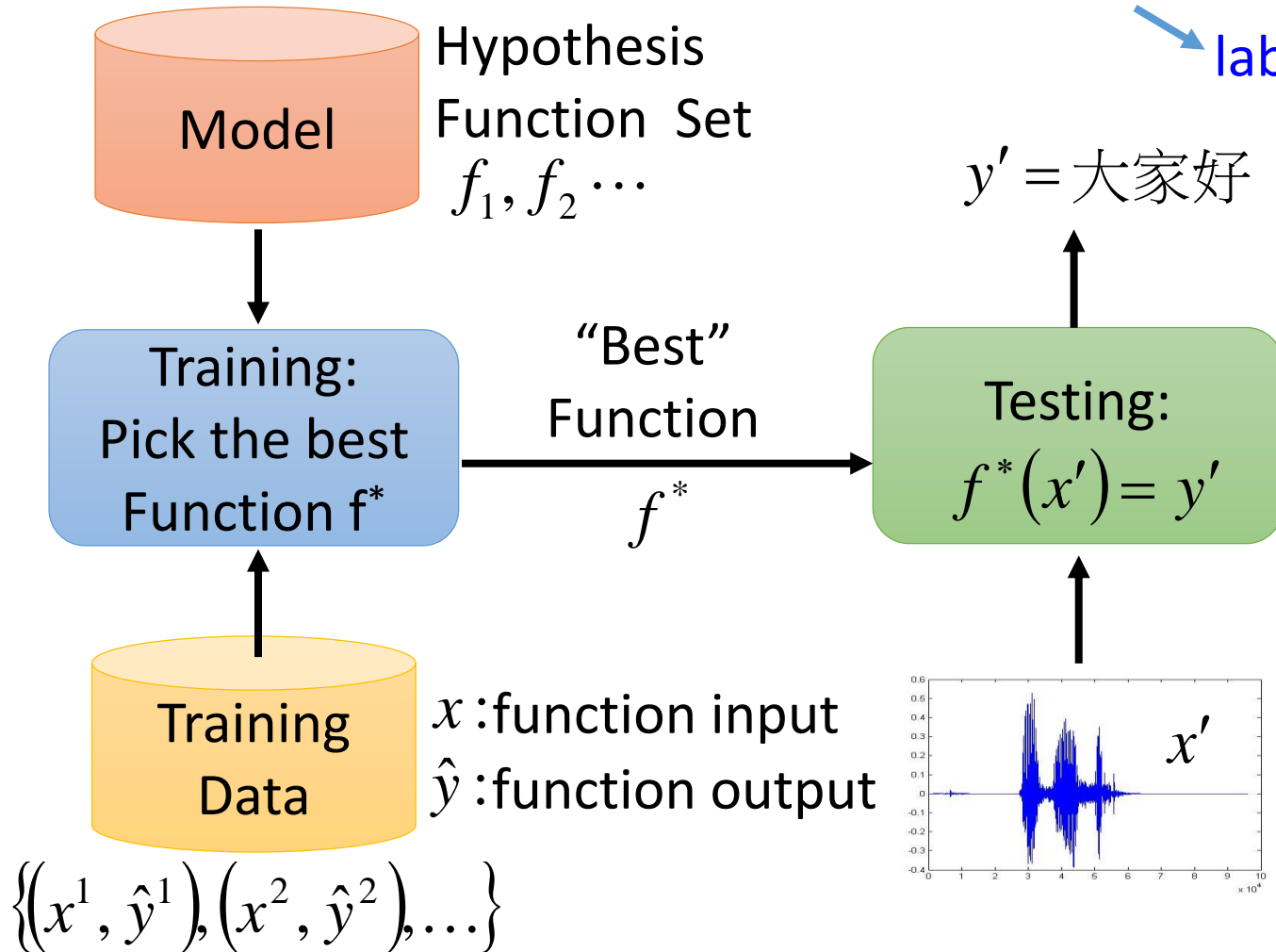
- Weather forecast

$$f\left(\text{weather today}\right) = \text{“sunny tomorrow”}$$

- Play video games

$$f\left(\begin{array}{l} \text{Positions and} \\ \text{number of enemies} \end{array}\right) = \text{“fire”}$$

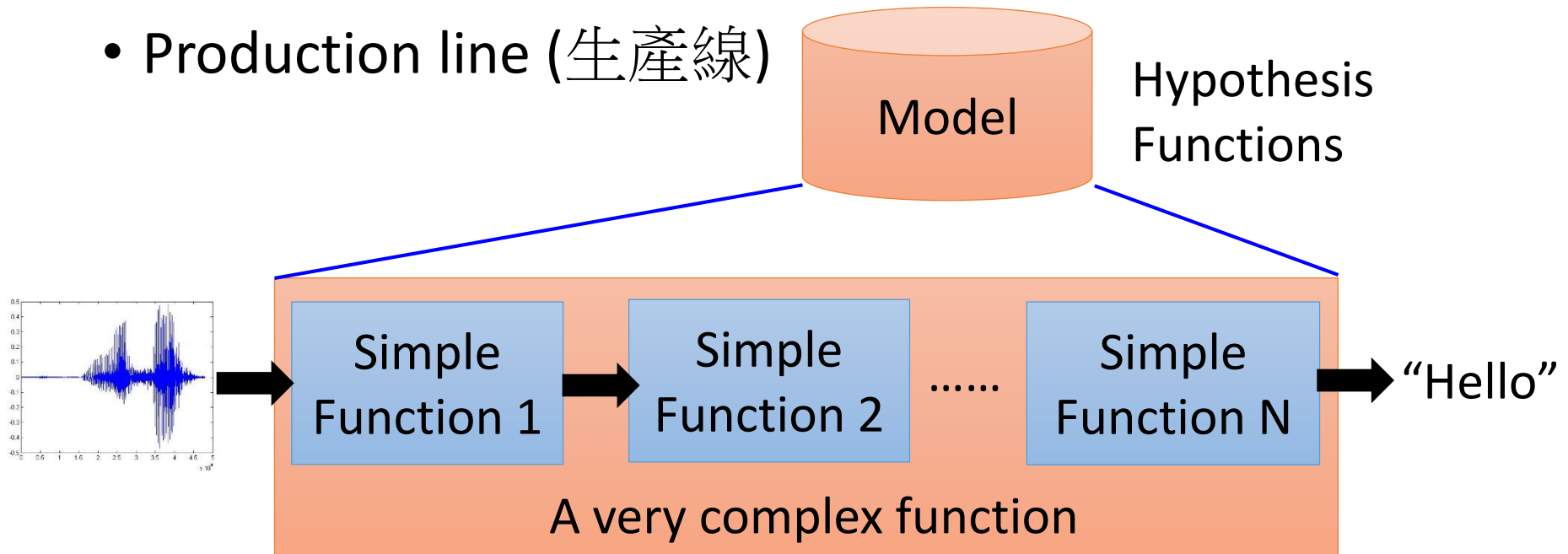
# Framework



# Deep Learning

# What is Deep Learning?

- Production line (生産線)



End-to-end training:

What each function should do is learned automatically

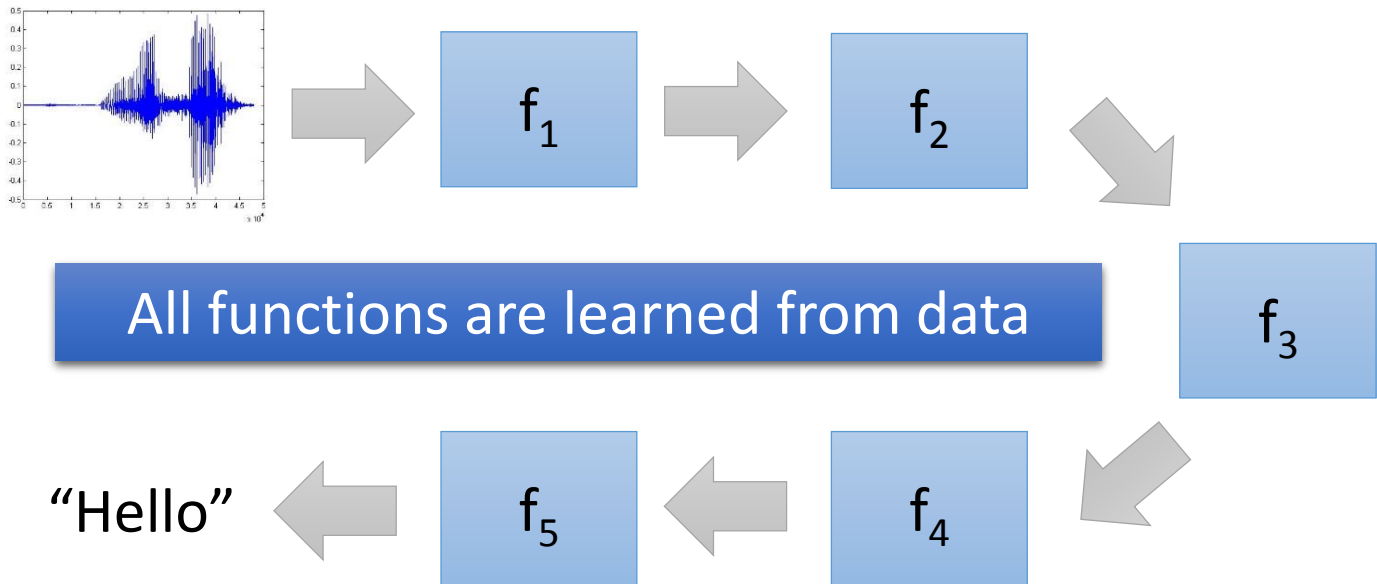


# Deep v.s. Shallow

## - Speech Recognition

- Deep Learning

“Bye bye, MFCC”  
- Deng Li in  
Interspeech 2014

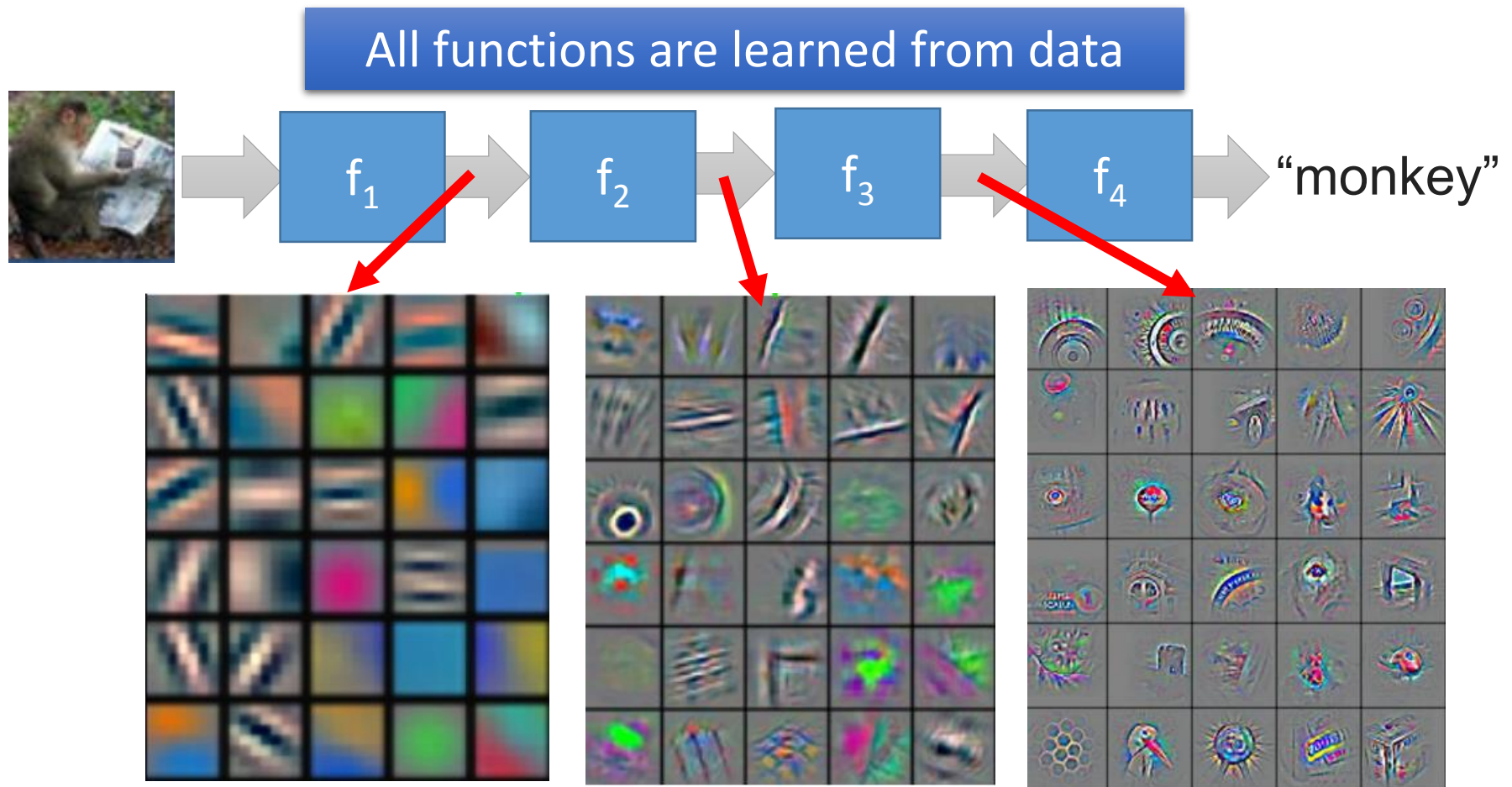


Less engineering labor, but machine learns more

# Deep v.s. Shallow - Image Recognition

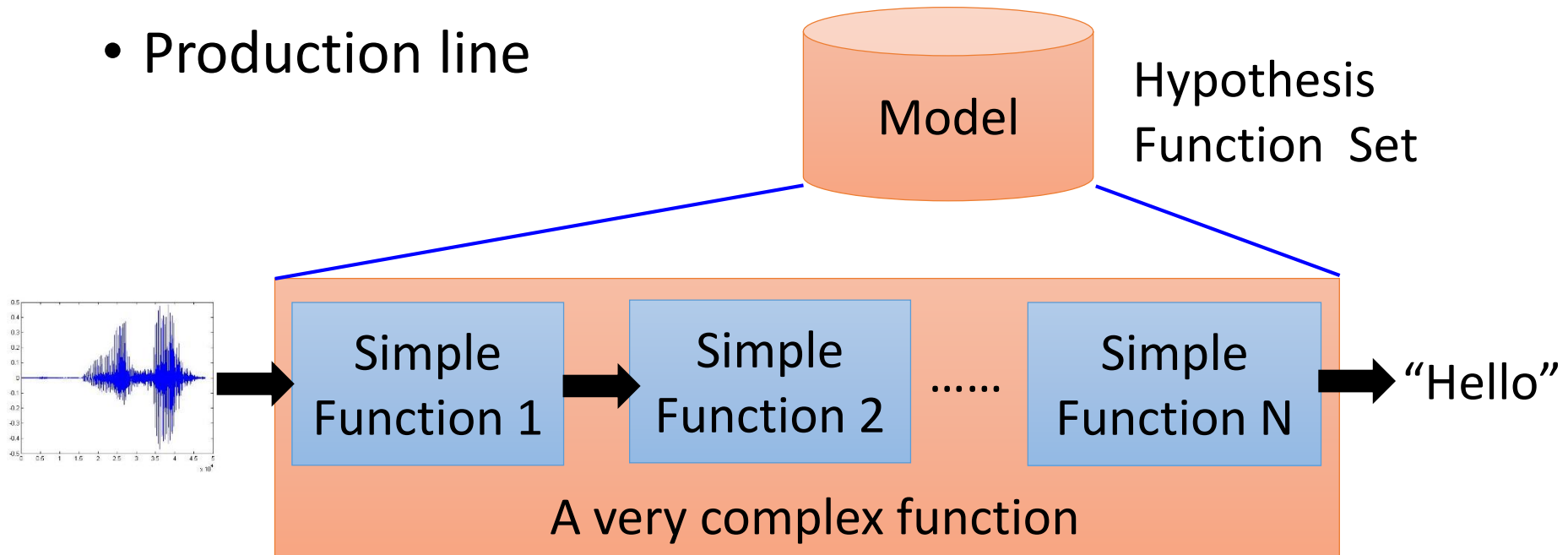
Reference: Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *Computer Vision—ECCV 2014* (pp. 818-833)

- Deep Learning



# What is Deep Learning?

- Production line



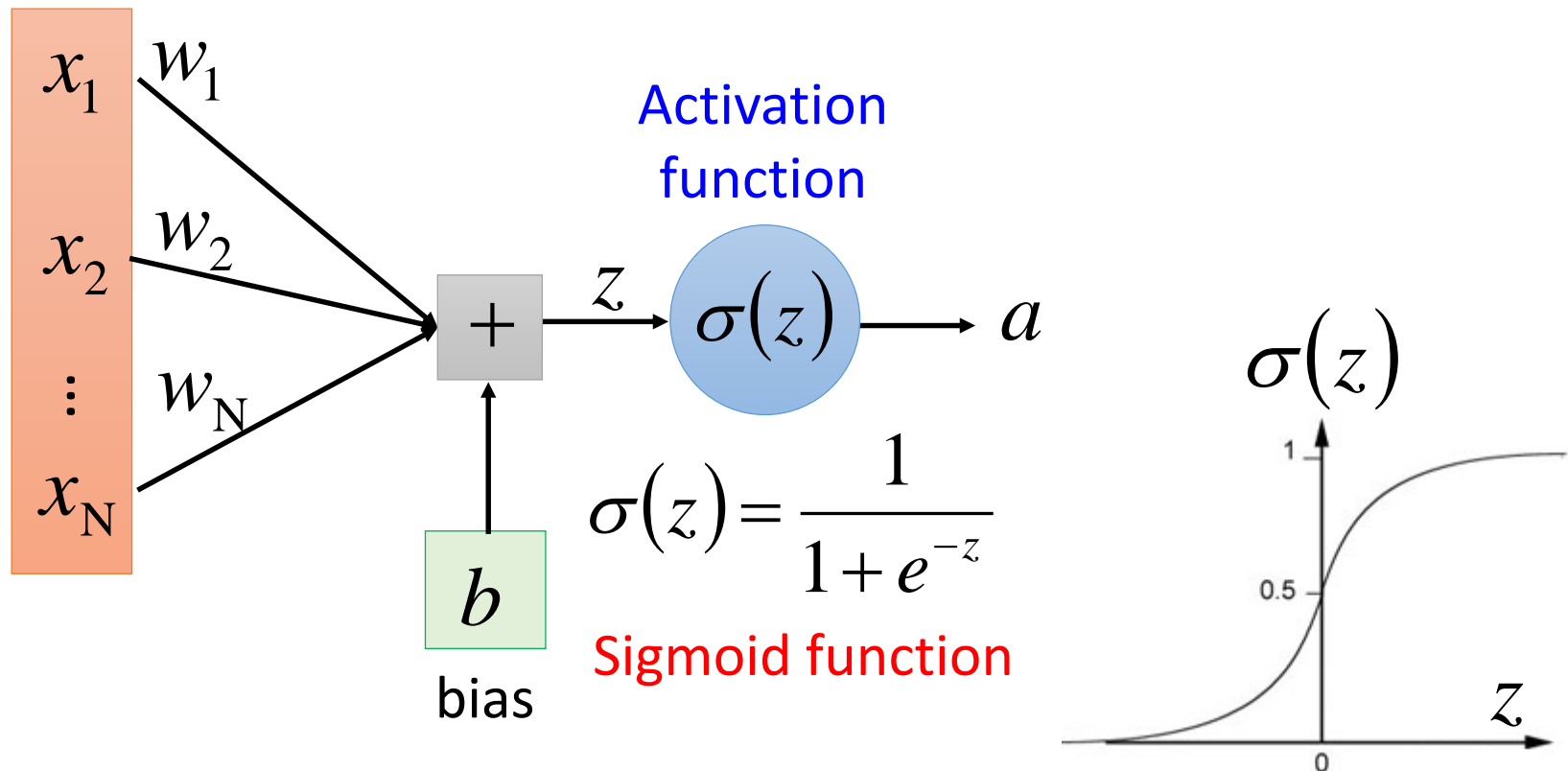
End-to-end training:

What each function should do is learned automatically

- Deep learning usually referred to neural network based approach

# A Neuron for Machine

Each neuron is a very simple function

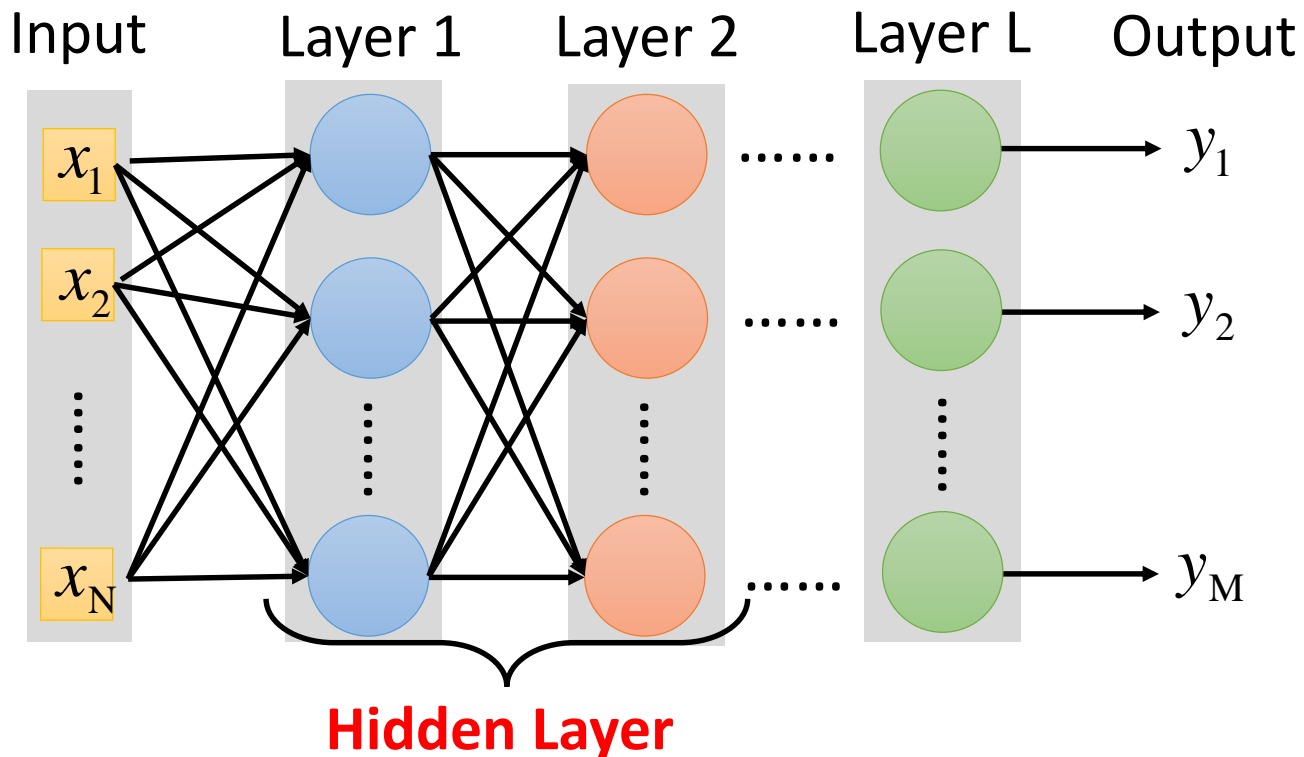


# Deep Learning

A neural network is a complex function:

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

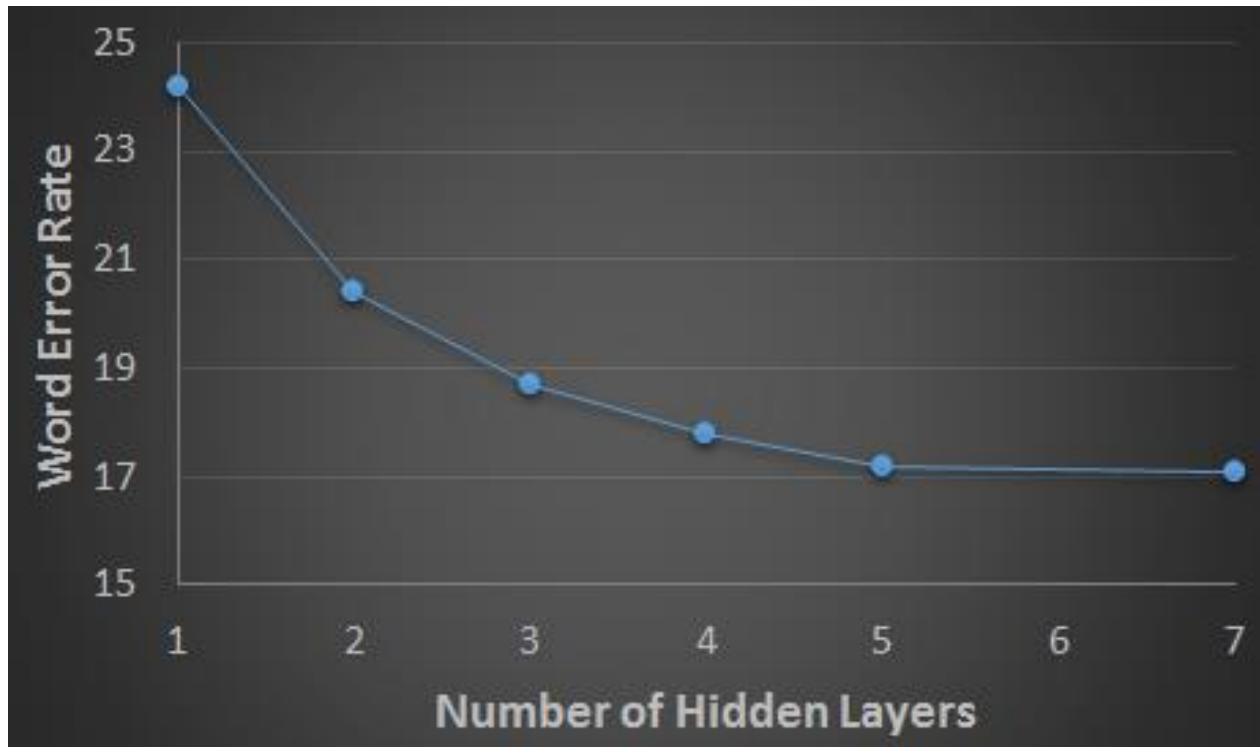
- Cascading the neurons to form a neural network.  
Each layer is a simple function in the production line.



# Why Deep Learning?

Deeper is Better.

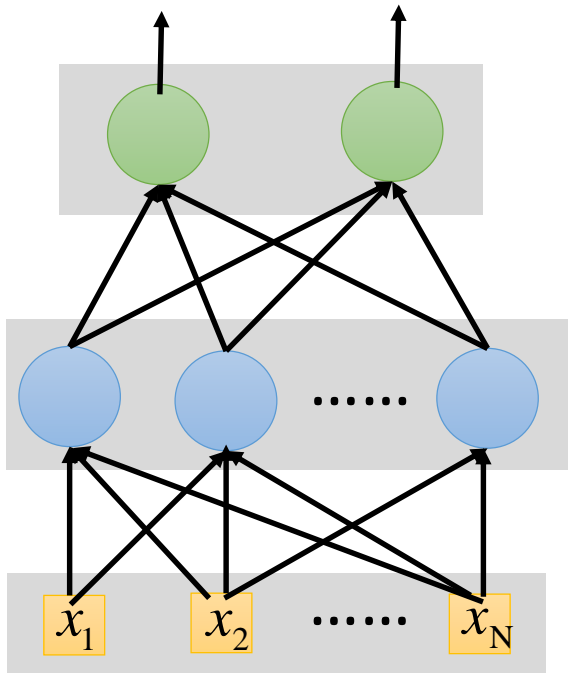
- Speech recognition



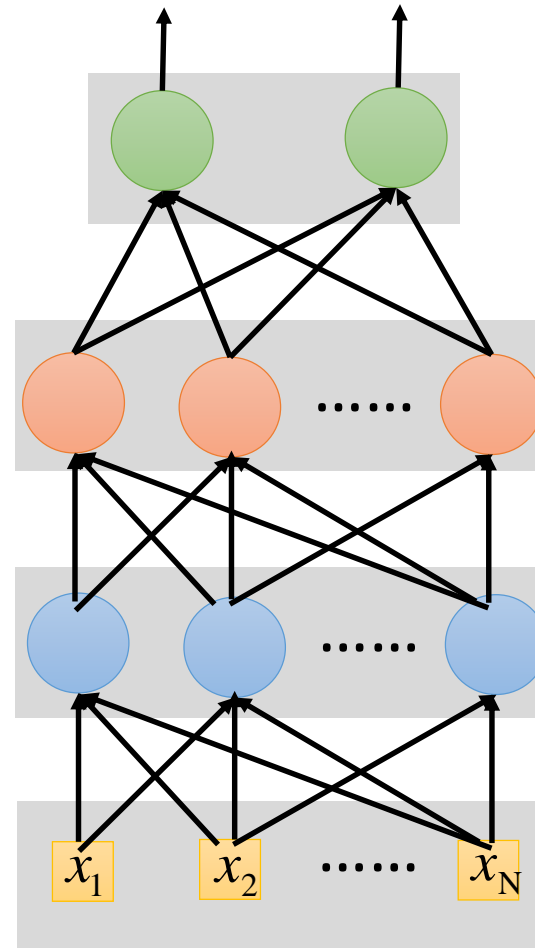
Seide, Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.

# Why Deeper is Better?

Deep works better simply because it uses more parameters.



Shallow

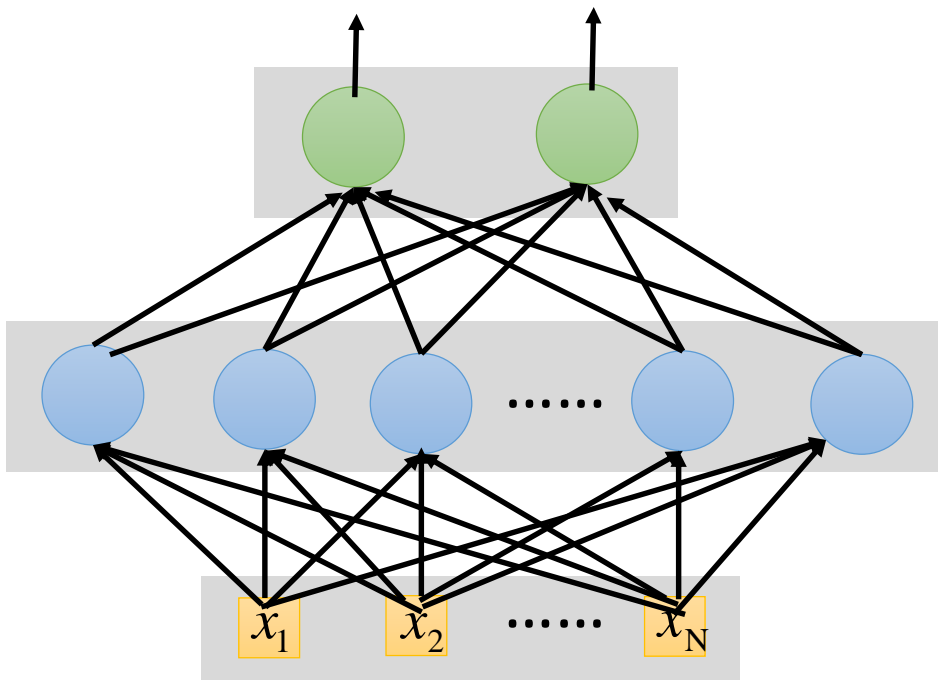


Deep

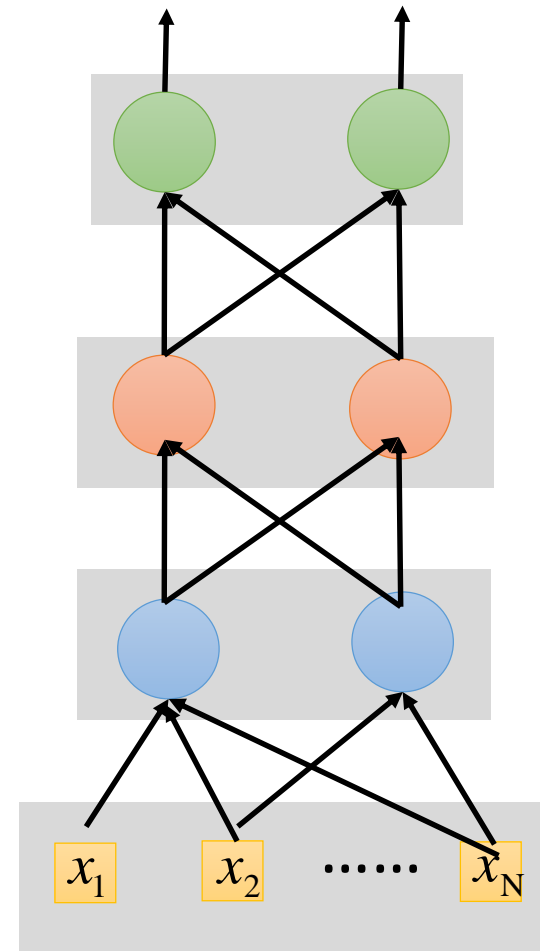
# Fat + Short v.s. Thin + Tall

If they have the same  
parameters,

Which one is better?



Shallow



Deep

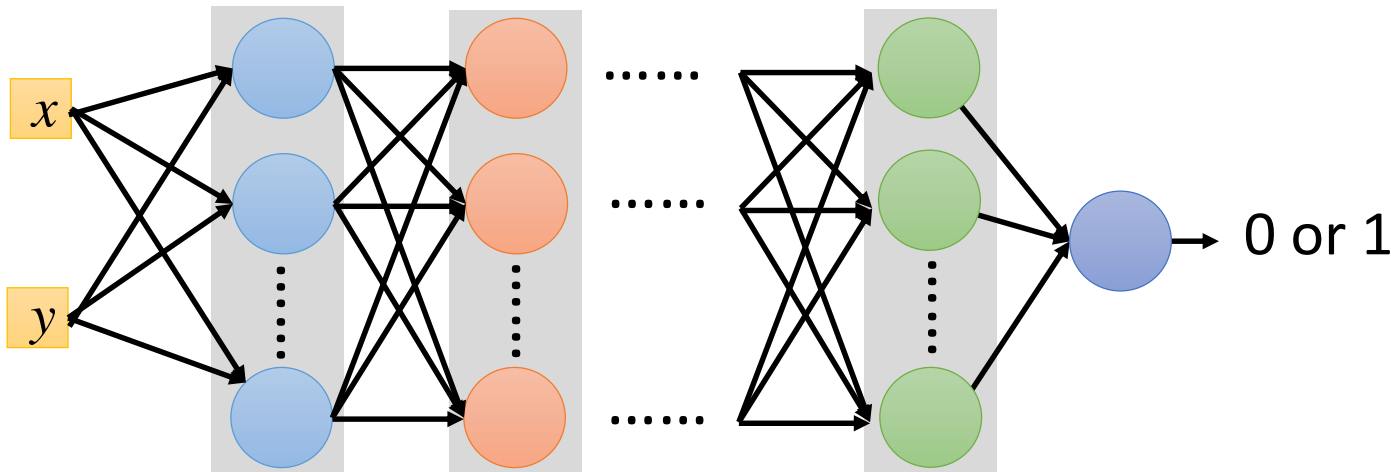
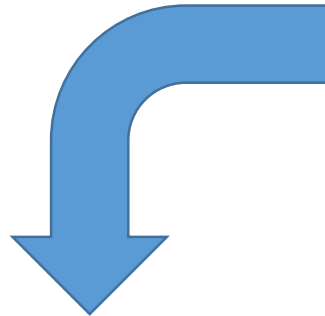
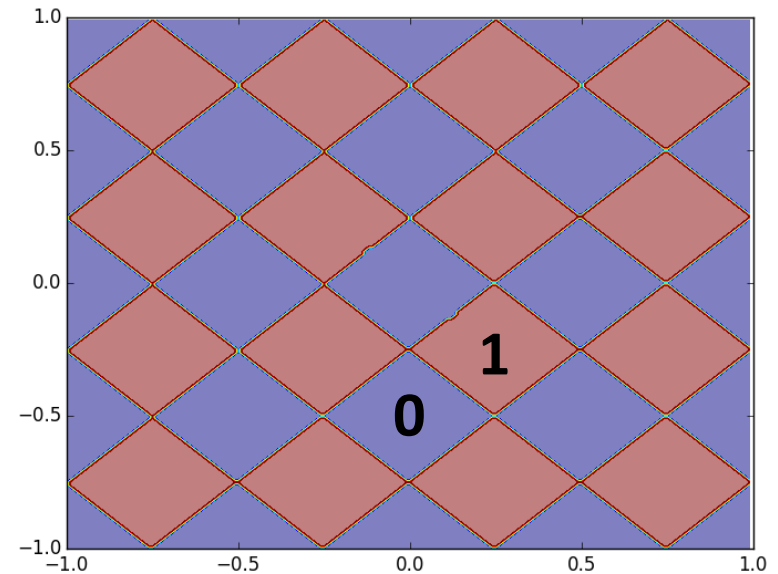


# Fat + Short v.s. Thin + Tall

## Toy Example

Sample 10,000  
points as training data

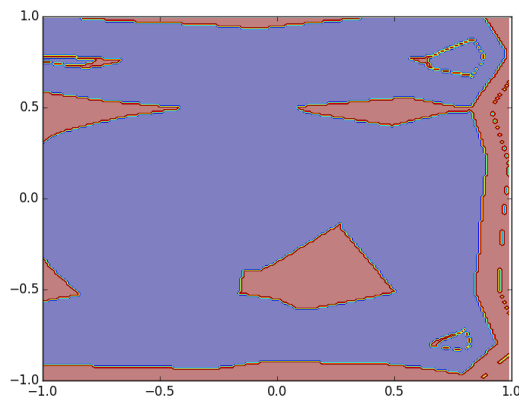
$$f : \mathbb{R}^2 \rightarrow \{0,1\}$$



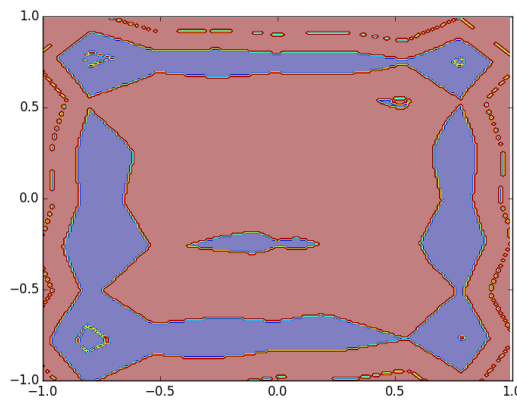
# Fat + Short v.s. Thin + Tall

## Toy Example

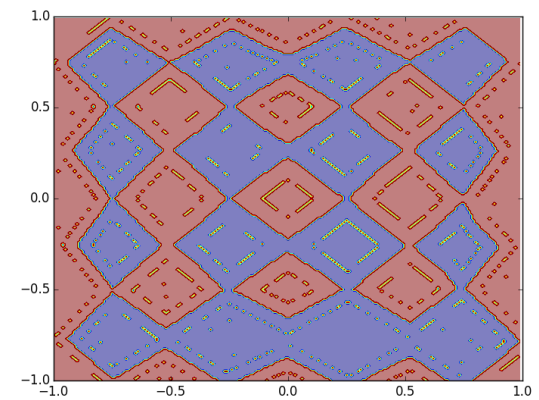
**1 hidden layer:**



(A) 125 neurons

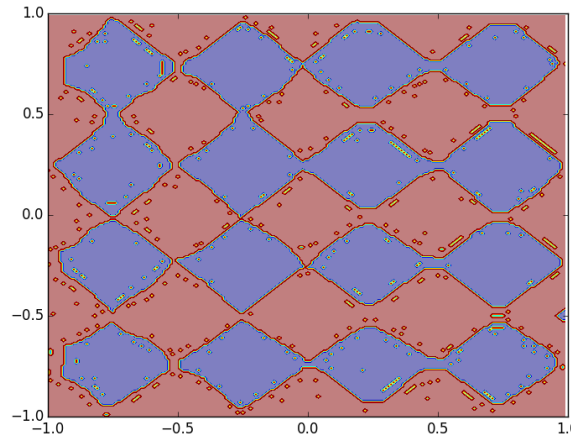


(B) 500 neurons



(C) 2500 neurons

**3 hidden layers:**

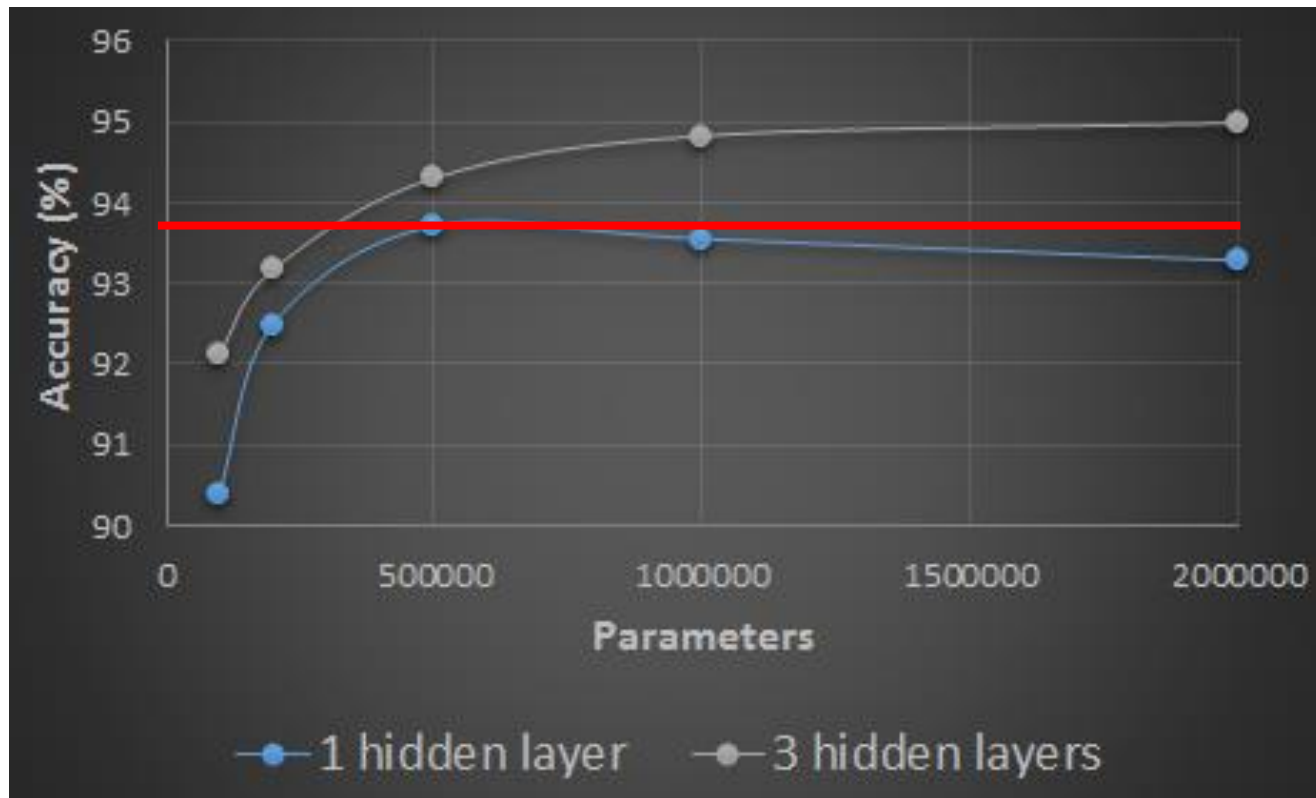


Q: the number of  
parameters close to  
(A), (B) or (C)?

# Fat + Short v.s. Thin + Tall

## Hand-writing digit classification

- Same parameters



Deeper: Using less parameters to achieve the same performance

# Fat + Short v.s. Thin + Tall Speech Recognition

- Word error rate (WER)

Multiple layers		1 hidden layer	
LxN	DBN-PT (%)	1xN	DBN-PT (%)
1×2k	24.2		
2×2k	20.4		
3×2k	18.4		
4×2k	17.8		
5×2k	17.2	1×3,772	22.5
7×2k	17.1	1×4,634	22.6
		1×16K	22.1

Seide, Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.

# Size of Training Data

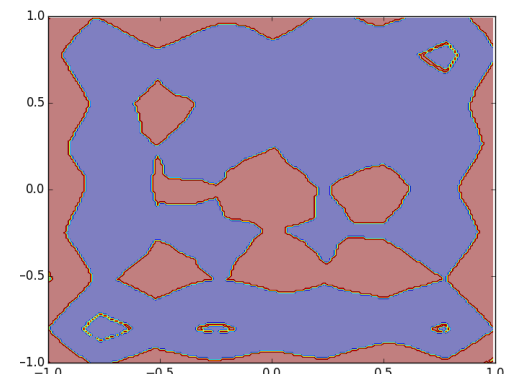
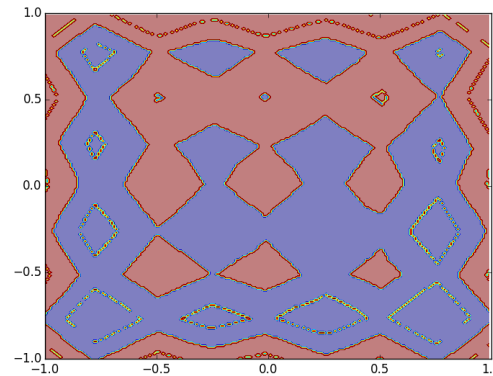
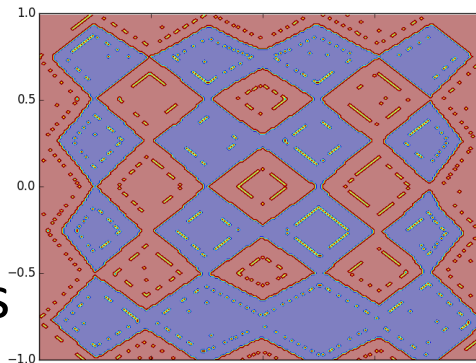
- Different numbers of training examples

10,000

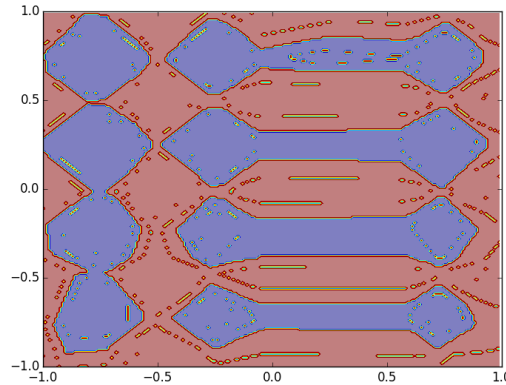
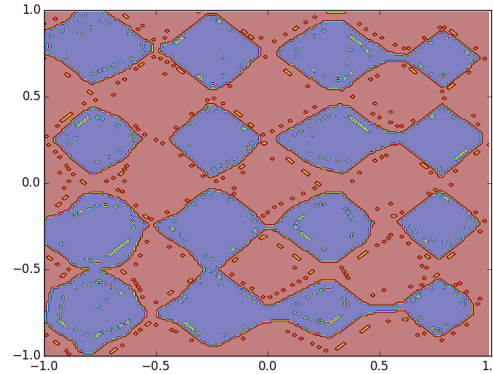
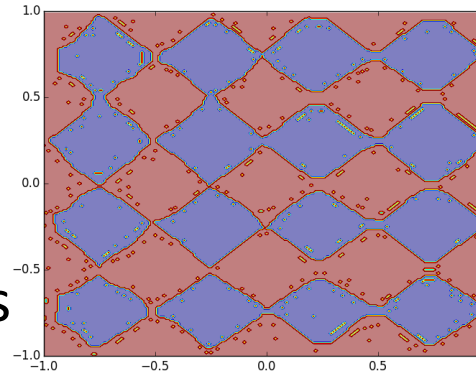
5,000

2,000

**1 hidden  
layer**  
More  
parameters

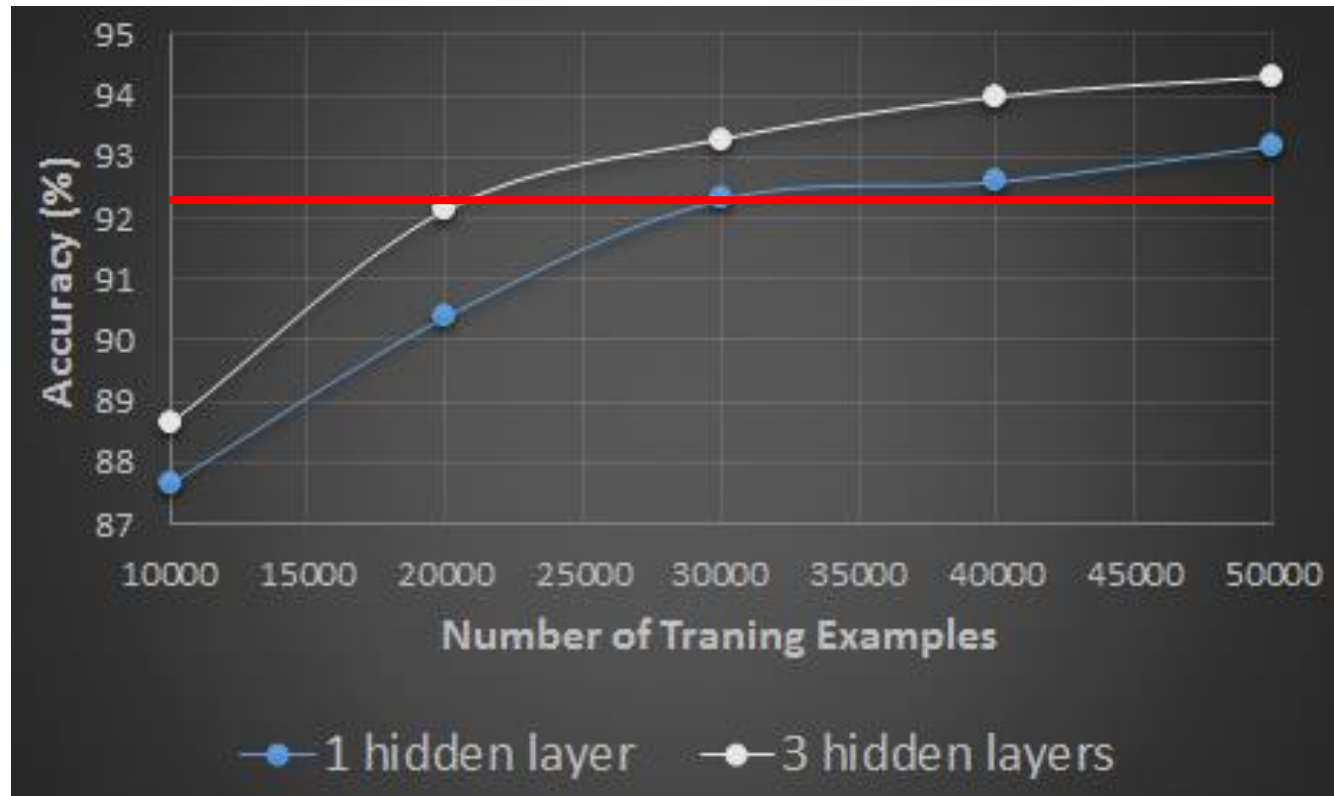


**3 hidden  
layers**  
Less  
parameters



# Size of Training Data

- Hand-writing digit classification



Deeper: Using less training data to achieve the same performance

# Learning $\approx$ Looking for a Function

- Speech Recognition

$$f\left(\text{[Waveform of '你好' (Hello)]}\right) = \text{“你好”}$$

- Handwritten Recognition

$$f\left(\text{[Handwritten digit '2']}\right) = \text{“2”}$$

- Weather forecast

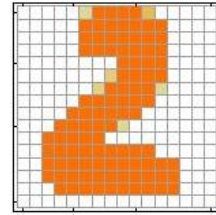
$$f\left(\text{weather today}\right) = \text{“sunny tomorrow”}$$

- Play video games

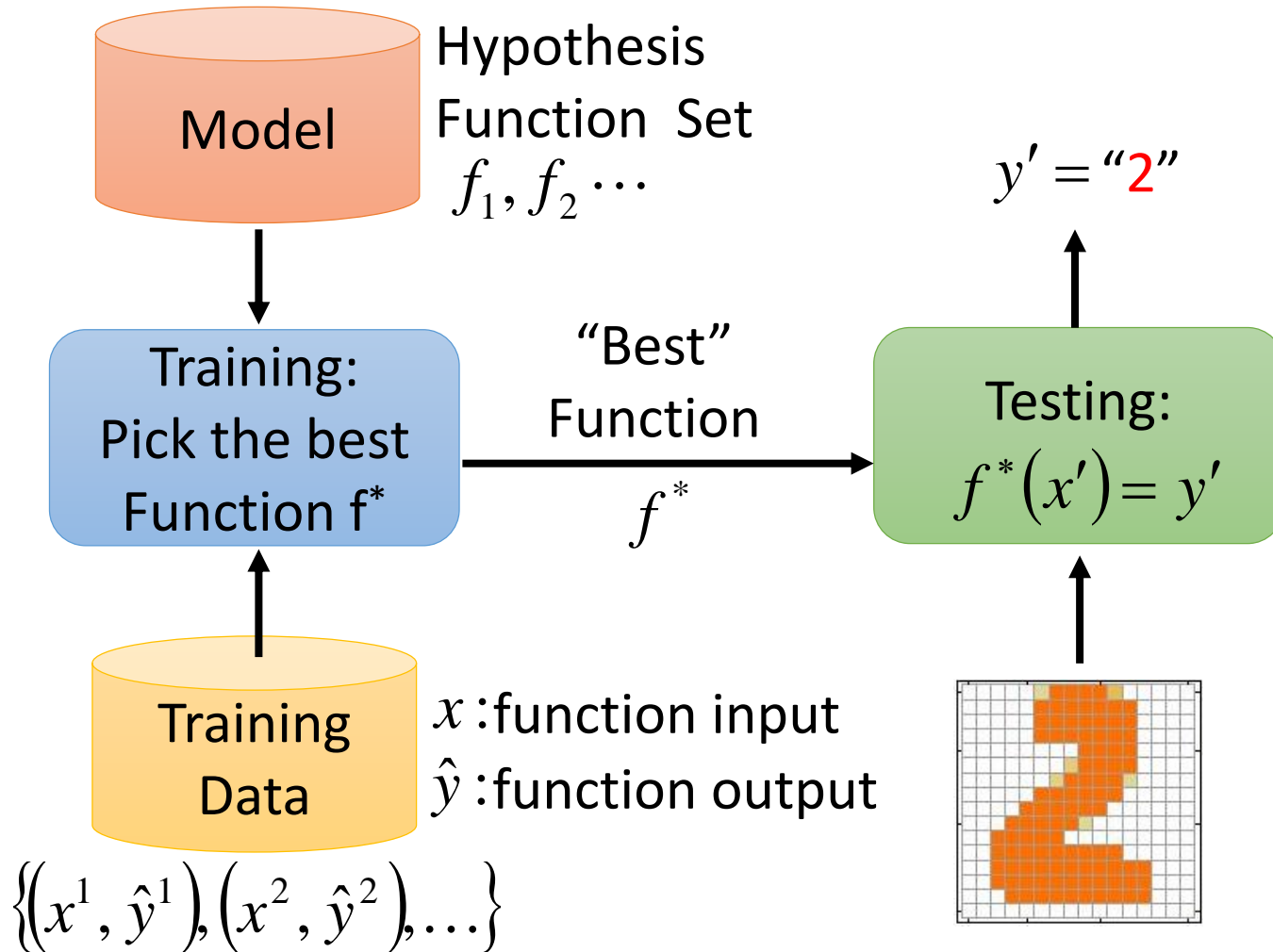
$$f\left(\begin{array}{l} \text{Positions and} \\ \text{number of enemies} \end{array}\right) = \text{“fire”}$$

# Framework

$x$ :



$\hat{y}$ : "2"  
(label)





# Outline

1. What is the model (function hypothesis set)?



2. What is the “best” function?



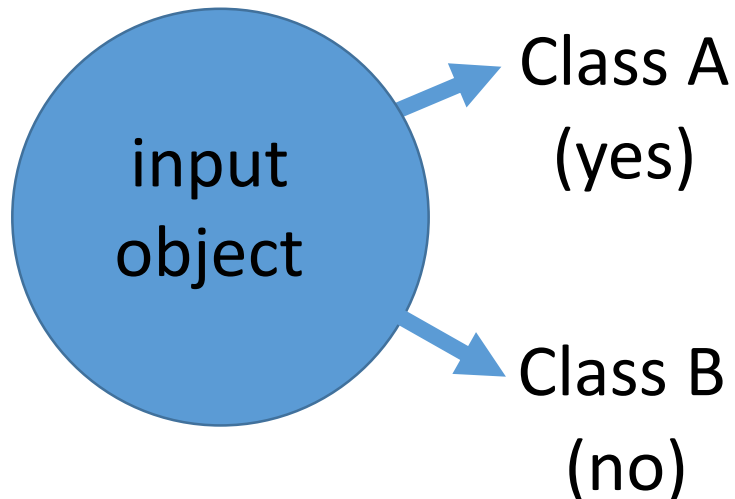
3. How to pick the “best” function?

# Task Considered Today

- **Classification**

- **Binary Classification**

Only two classes



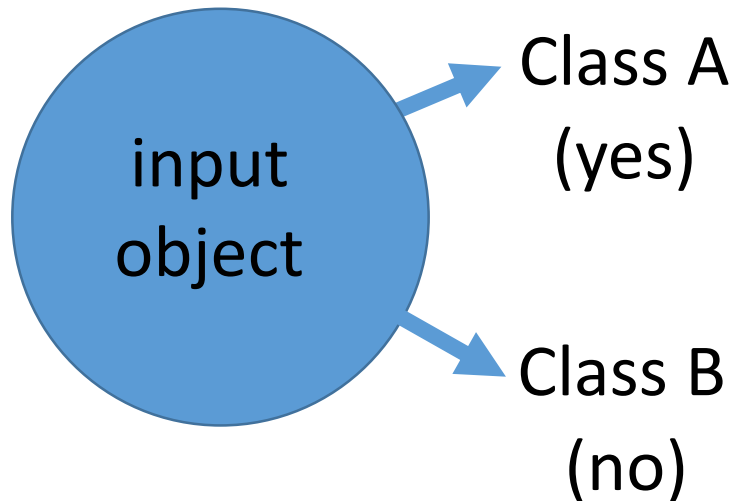
- **Spam filtering**
  - Is an e-mail spam or not?
- **Recommendation systems**
  - recommend the product to the customer or not?
- **Malware detection**
  - Is the software malicious or not?
- **Stock prediction**
  - Will the future value of a stock increase or not?

# Task Considered Today

- ***Classification***

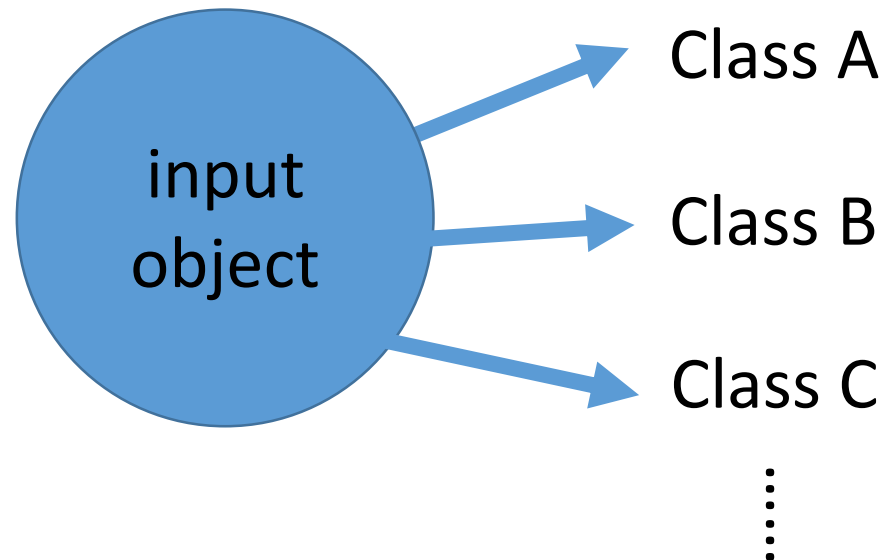
## ***Binary Classification***

Only two classes



## ***Multi-class Classification***

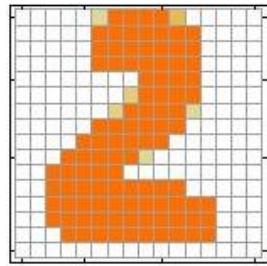
More than two classes



# Multi-class Classification

- Handwriting Digit Classification

Input:



Class: "1", "2", ..., "9", "0"  
10 classes

- Image Recognition

Input:



Class: "dog", "cat", "book", ....  
Thousands of classes

1. What is the model?

# What is the function we are looking for?

- **classification**

$$y = f(x) \quad \longrightarrow \quad f: R^N \rightarrow R^M$$

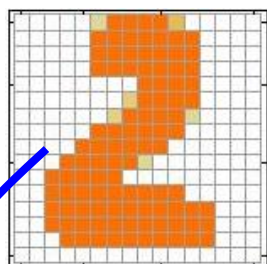
- x: input object to be classified
- y: class
- ***Assume both  $x$  and  $y$  can be represented as fixed-size vector***
  - x is a vector with N dimensions, and y is a vector with M dimensions

# What is the function we are looking for?

- Handwriting Digit Classification**

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

**x: image**



16 x 16

Each pixel corresponds to an element in the vector

$\begin{bmatrix} 0 \\ 1 \\ \vdots \end{bmatrix}$

1: for ink,  
0: otherwise  
16 x 16 = 256  
dimensions

**y: class**

10 dimensions for digit recognition

"1"

$\begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$

"1"  
"2"  
"3"

"2"

$\begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix}$

"1" → "1" or not  
"2" → "2" or not  
"3" → "3" or not

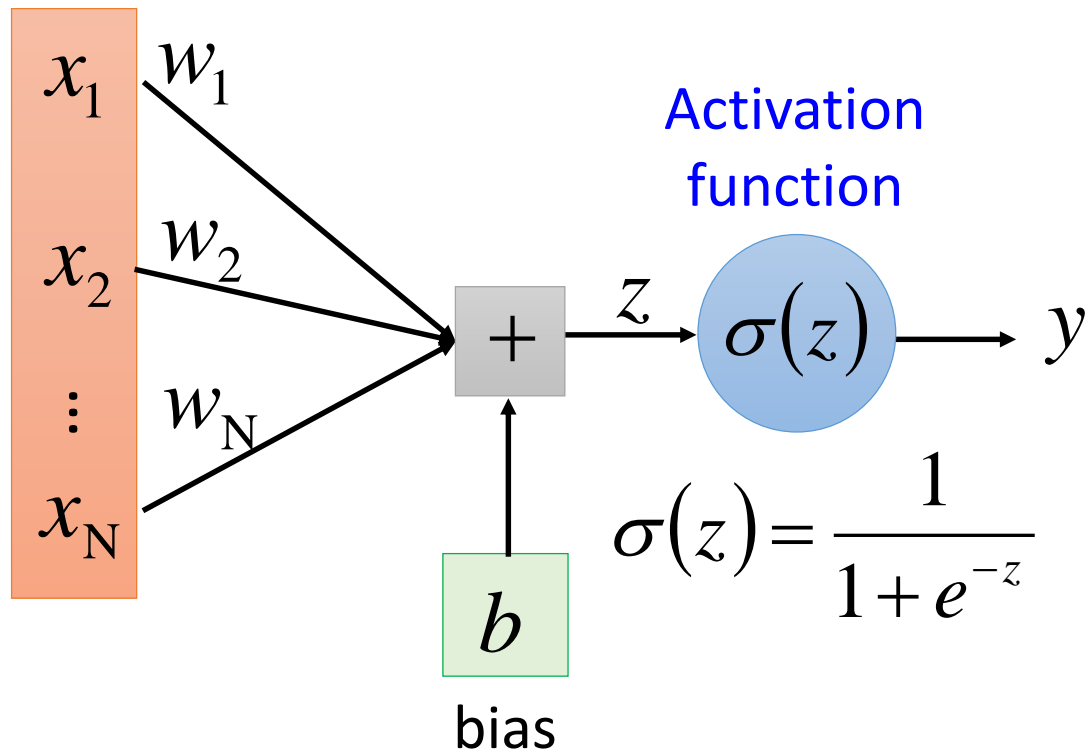
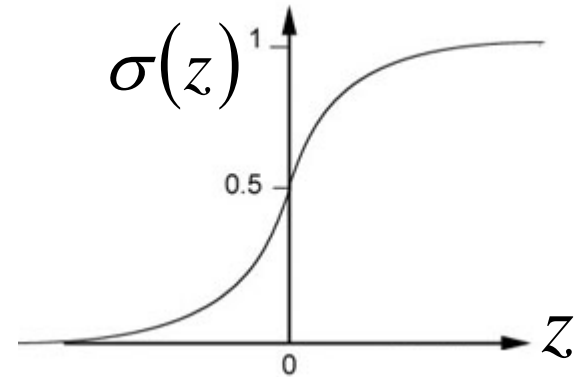
1. What is the model?

A Layer of Neuron



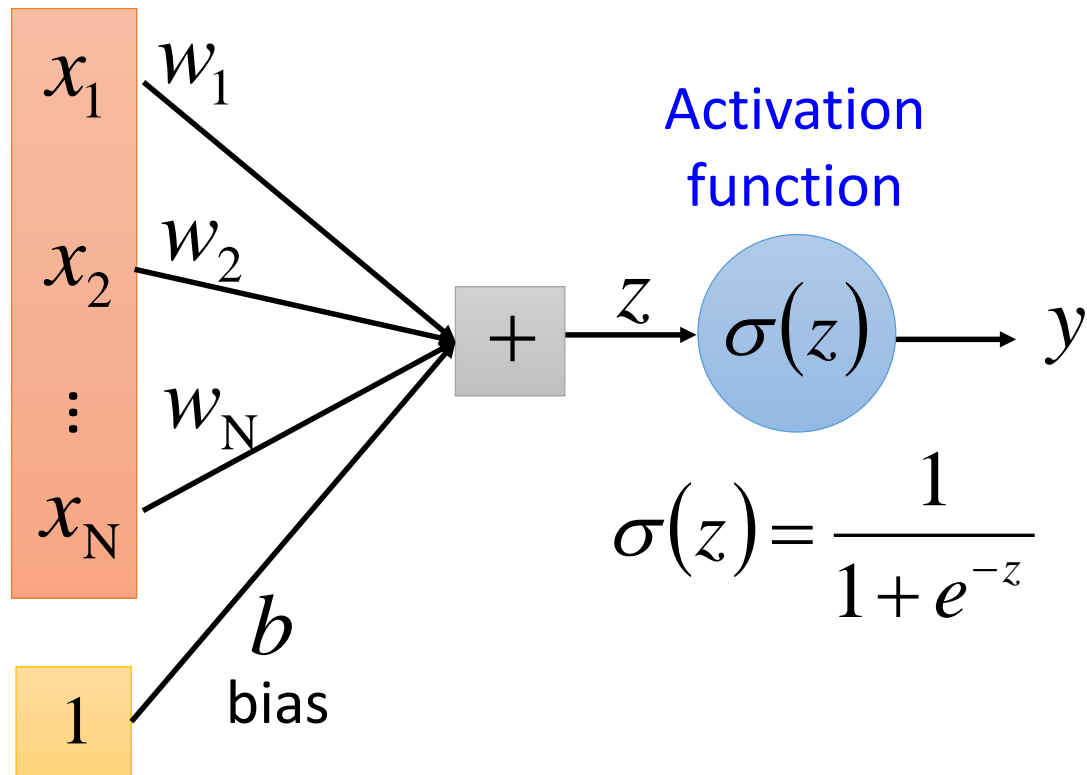
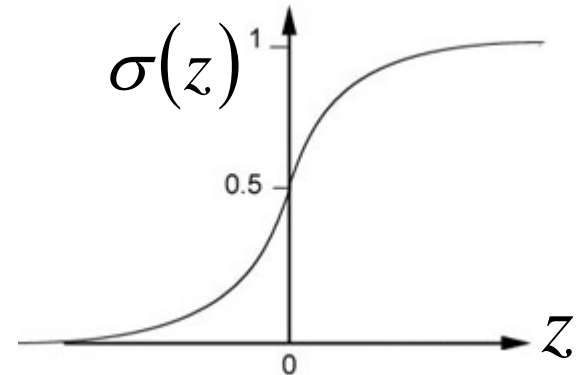
# Single Neuron

$$f: \mathbb{R}^N \rightarrow \mathbb{R}$$



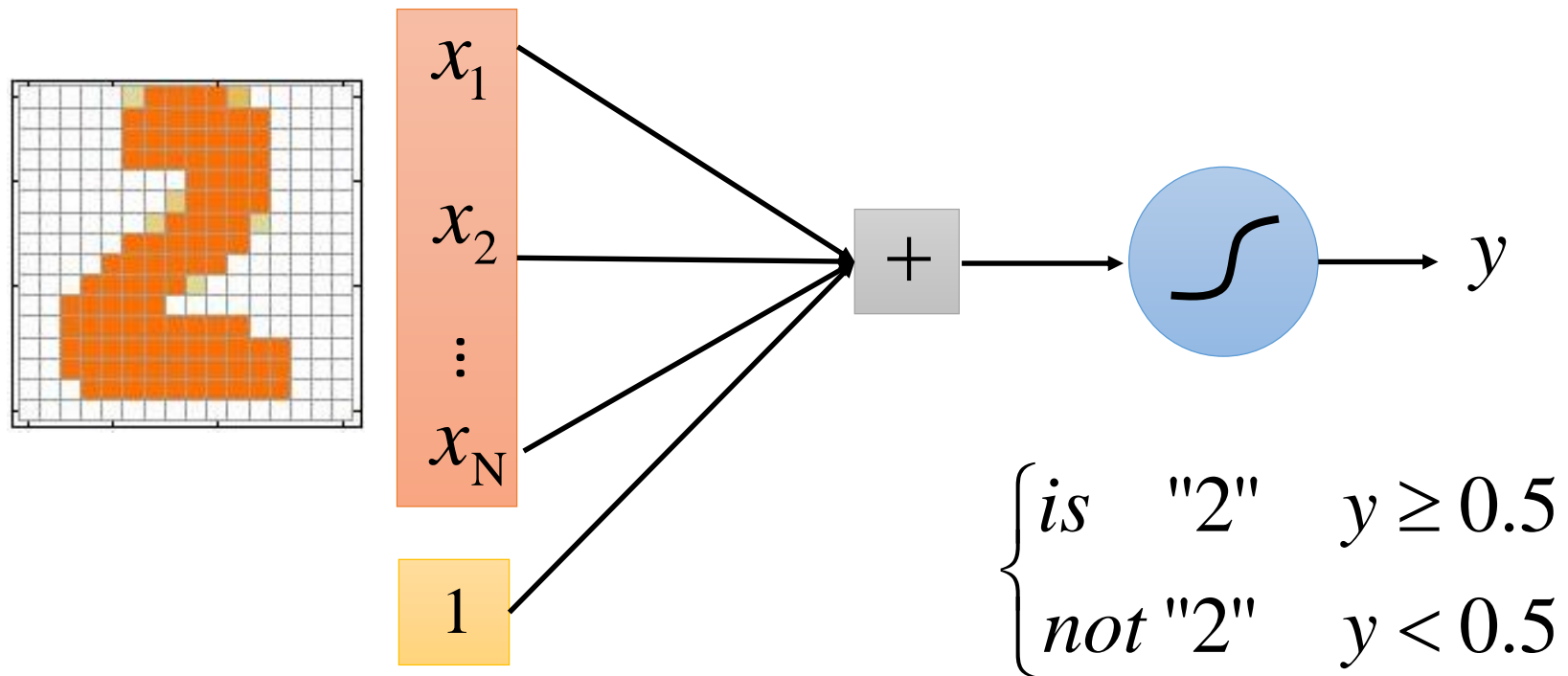
# Single Neuron

$$f: \mathbb{R}^N \rightarrow \mathbb{R}$$



# Single Neuron $f: R^N \rightarrow R$

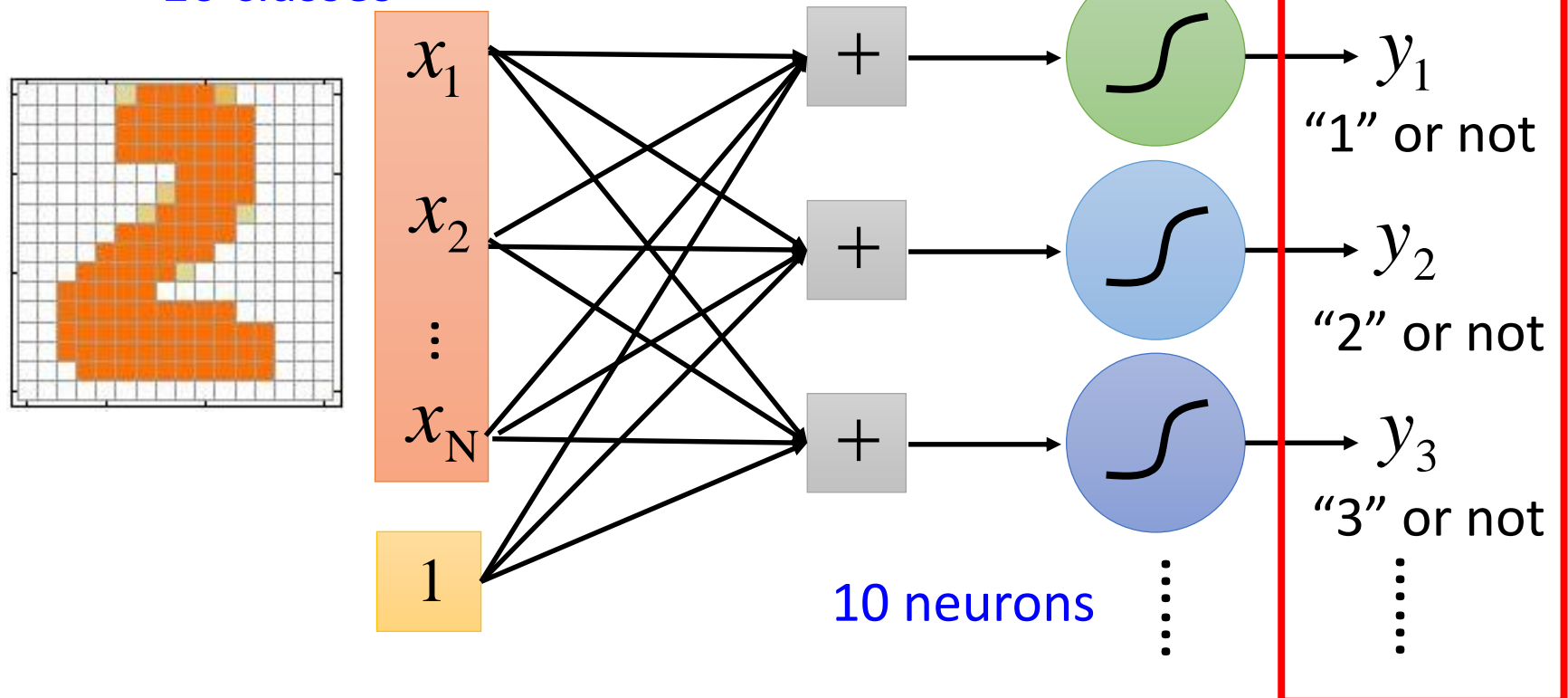
- Single neuron can only do binary classification, cannot handle multi-class classification



# A Layer of Neuron $f: R^N \rightarrow R^M$

- Handwriting digit classification
  - Classes: "1", "2", ..., "9", "0"
  - 10 classes

If  $y_2$  is the max, then the image is "2".

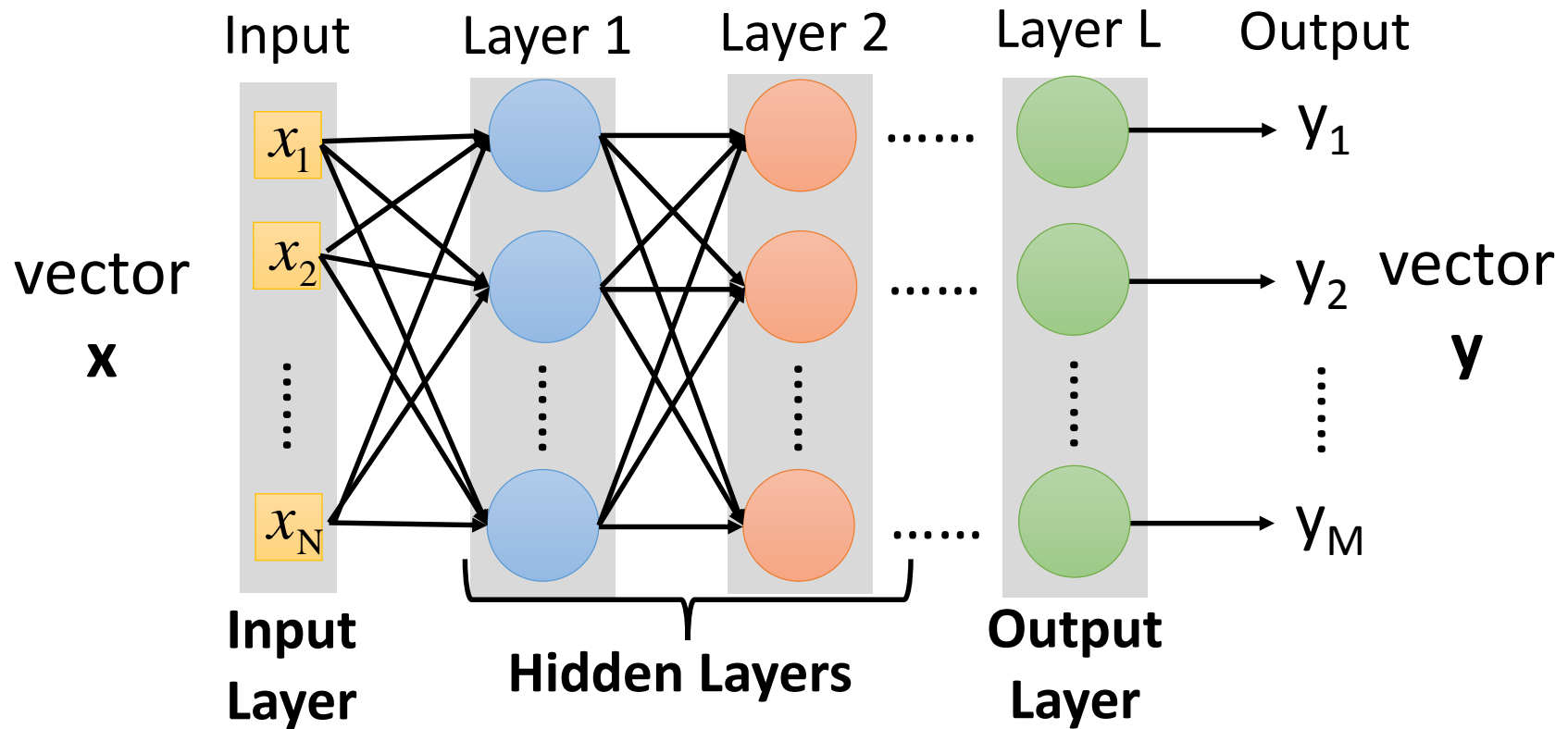


1. What is the model?

Neural Network

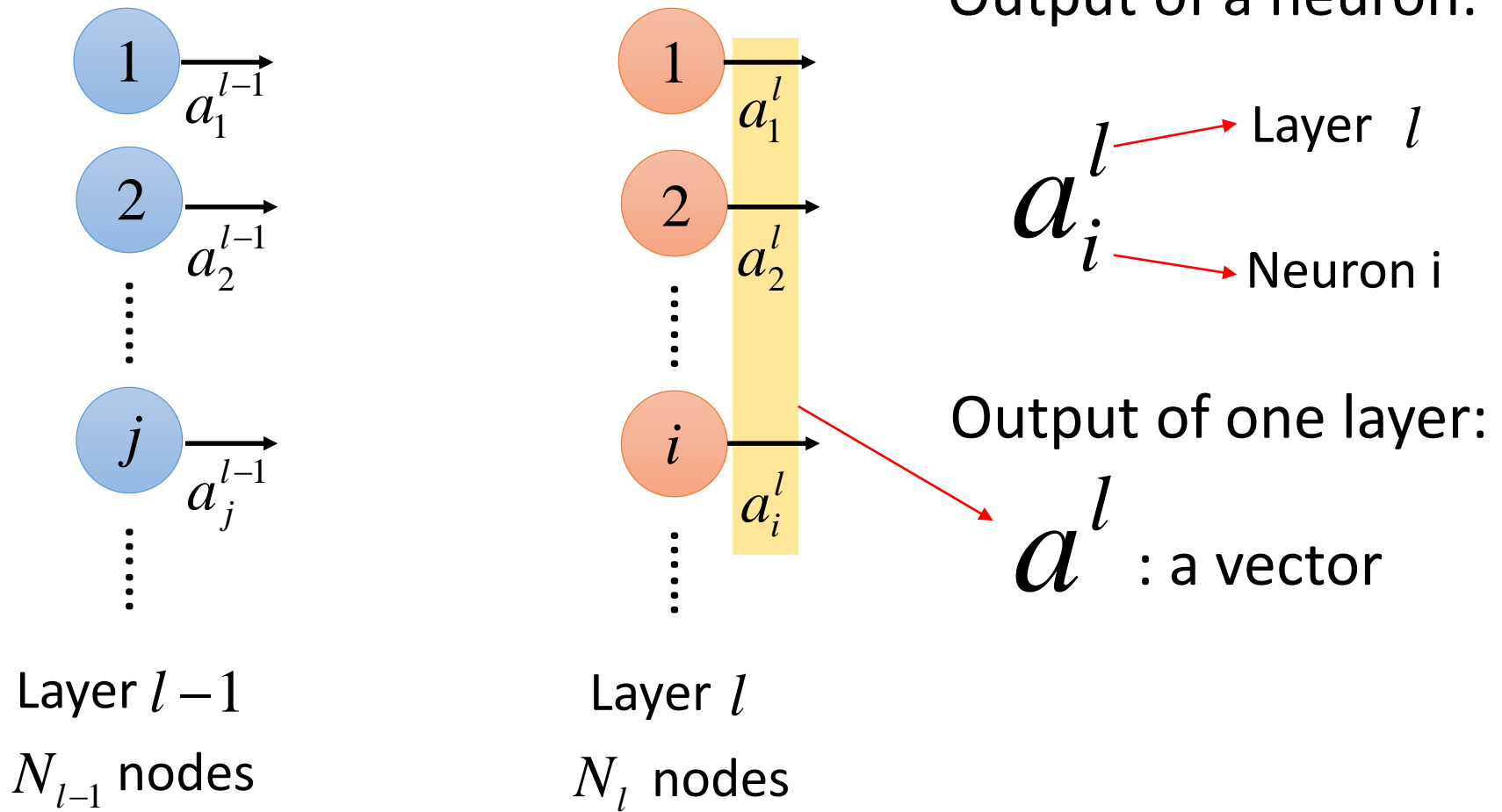
# Neural Network as Model

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

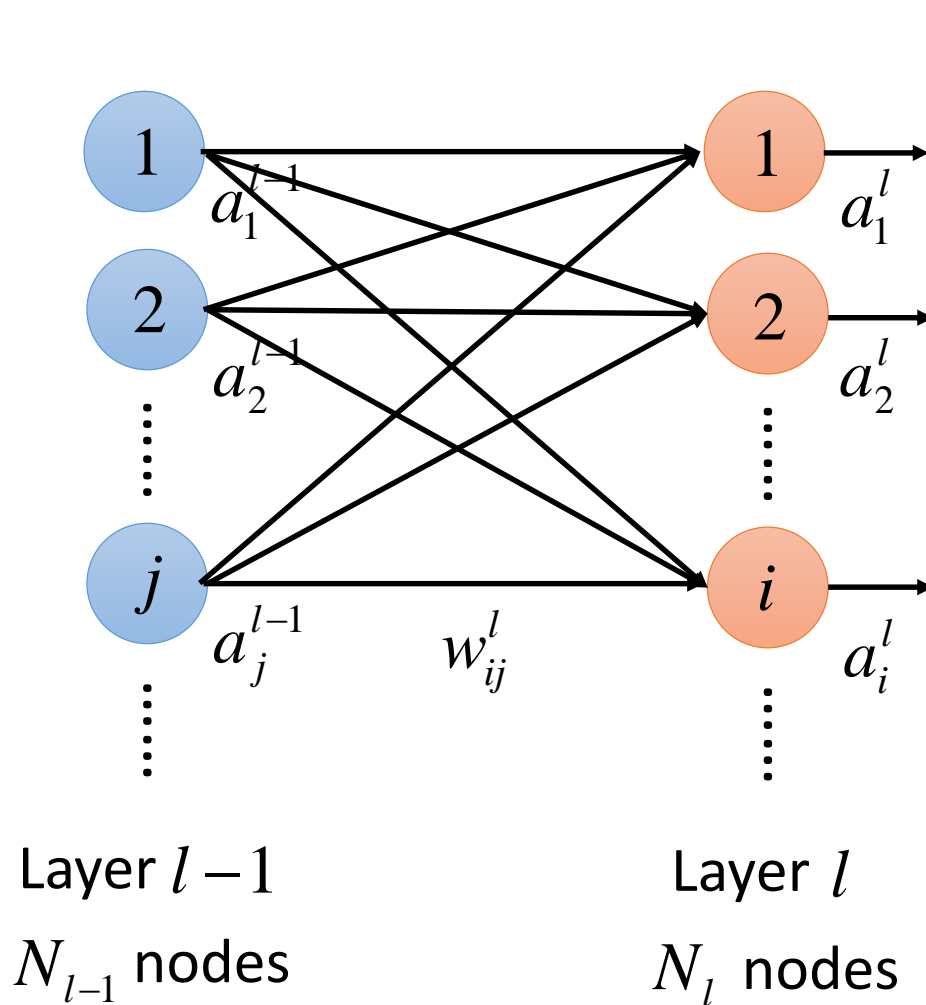


- Fully connected feedforward network
- Deep Neural Network: many hidden layers

# Notation



# Notation

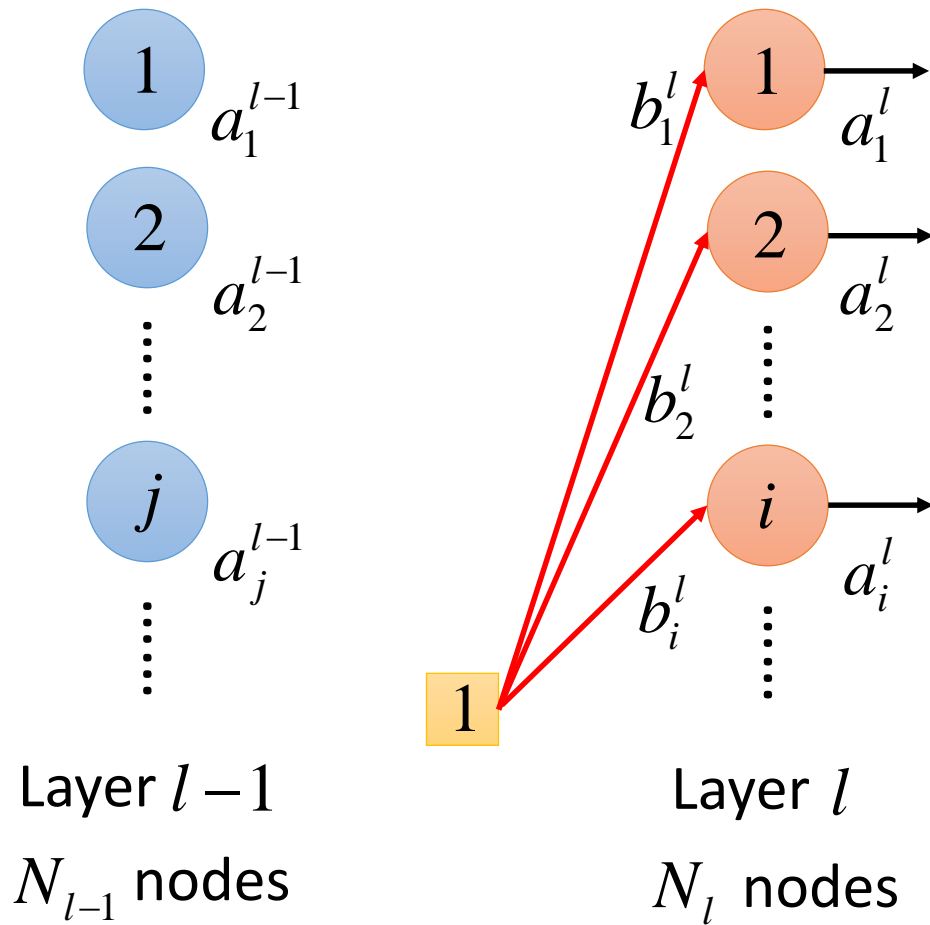


$w_{ij}^l$   $\xrightarrow{\text{red arrow}}$  Layer  $l-1$   
to Layer  $l$   
 $\downarrow$   
from neuron  $j$  (Layer  $l-1$ )  
to neuron  $i$  (Layer  $l$ )

$$W^l = \underbrace{\begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \\ \vdots & & \ddots \end{bmatrix}}_{N_{l-1} \times N_l} \underbrace{\quad}_{N_l}$$



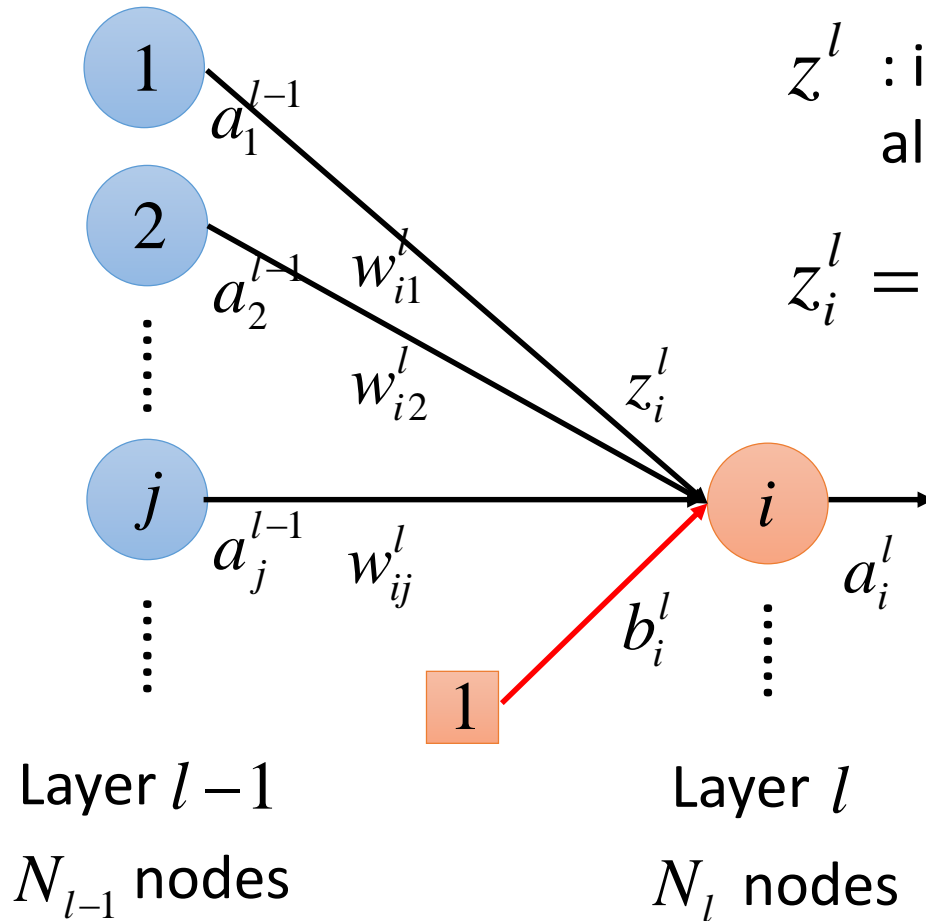
# Notation



$b_i^l$  : bias for neuron  $i$  at layer  $l$

$$b^l = \begin{bmatrix} b_1^l \\ b_2^l \\ \vdots \\ b_i^l \\ \vdots \end{bmatrix} \quad \text{bias for all neurons in layer } l$$

# Notation



$z_i^l$  : input of the activation function for neuron  $i$  at layer  $l$

$z^l$  : input of the activation function all the neurons in layer  $l$

$$z_i^l = w_{i1}^l a_1^{l-1} + w_{i2}^l a_2^{l-1} \dots + b_i^l$$

$$z_i^l = \sum_{j=1}^{N_{l-1}} w_{ij}^l a_j^{l-1} + b_i^l$$

# Notation - Summary

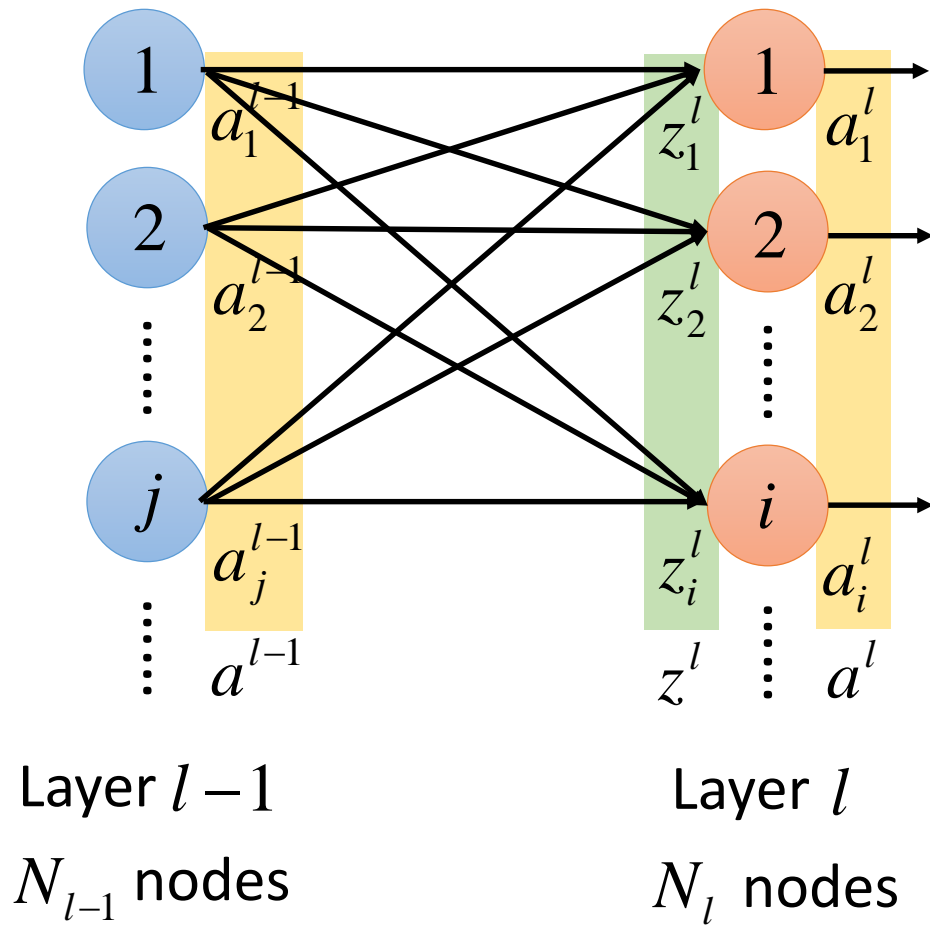
$a_i^l$  : output of a neuron       $w_{ij}^l$  : a weight

$a^l$  : output of a layer       $\mathbf{W}^l$  : a weight matrix

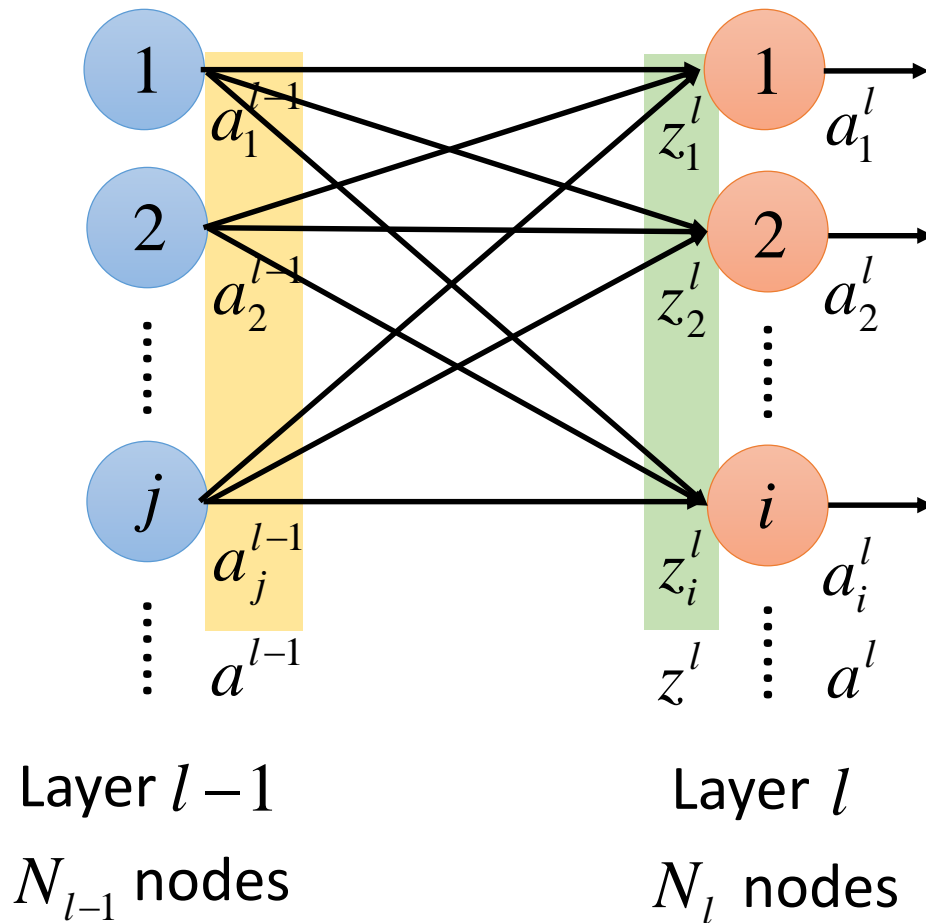
$z_i^l$  : input of activation  
function       $b_i^l$  : a bias

$z^l$  : input of activation  
function for a layer       $b^l$  : a bias vector

# Relations between Layer Outputs



# Relations between Layer Outputs

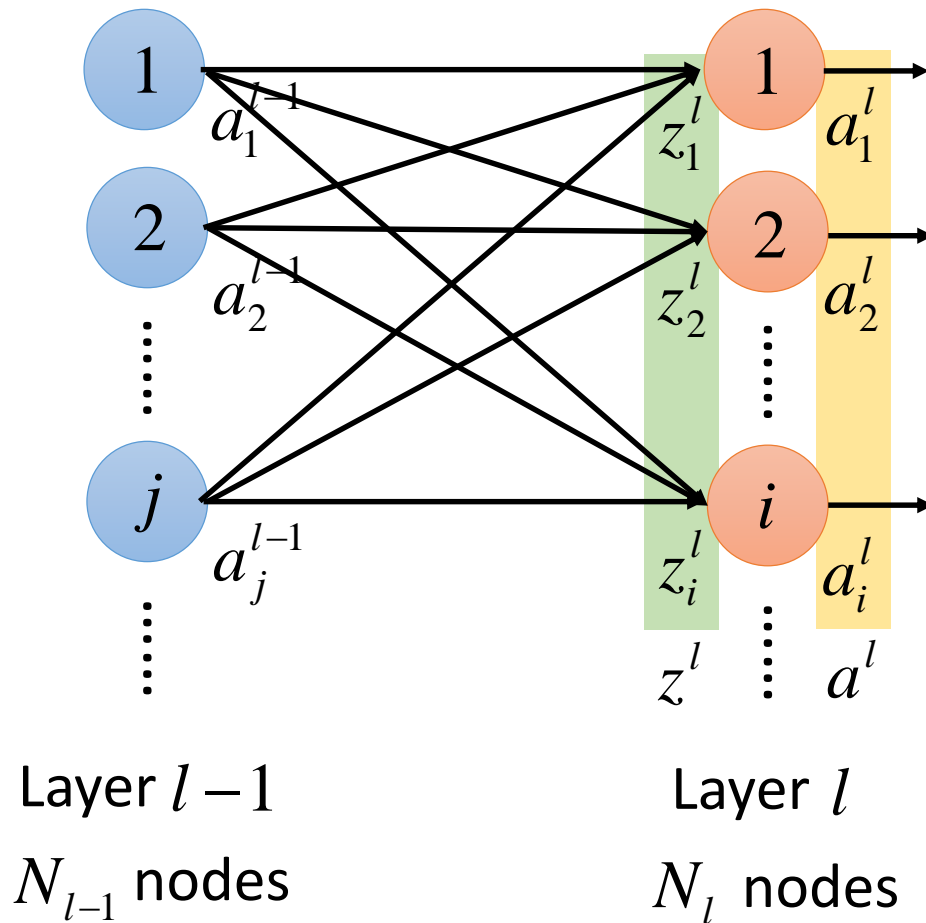


$$\begin{aligned} z_1^l &= w_{11}^l a_1^{l-1} + w_{12}^l a_2^{l-1} + \cdots + b_1^l \\ z_2^l &= w_{21}^l a_1^{l-1} + w_{22}^l a_2^{l-1} + \cdots + b_2^l \\ &\vdots \\ z_i^l &= w_{i1}^l a_1^{l-1} + w_{i2}^l a_2^{l-1} + \cdots + b_i^l \\ &\vdots \end{aligned}$$

$$\begin{bmatrix} z_1^l \\ z_2^l \\ \vdots \\ z_i^l \\ \vdots \end{bmatrix} = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \\ \vdots & & \ddots \end{bmatrix} \begin{bmatrix} a_1^{l-1} \\ a_2^{l-1} \\ \vdots \\ a_i^{l-1} \\ \vdots \end{bmatrix} + \begin{bmatrix} b_1^l \\ b_2^l \\ \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

$$z^l = W^l a^{l-1} + b^l$$

# Relations between Layer Outputs

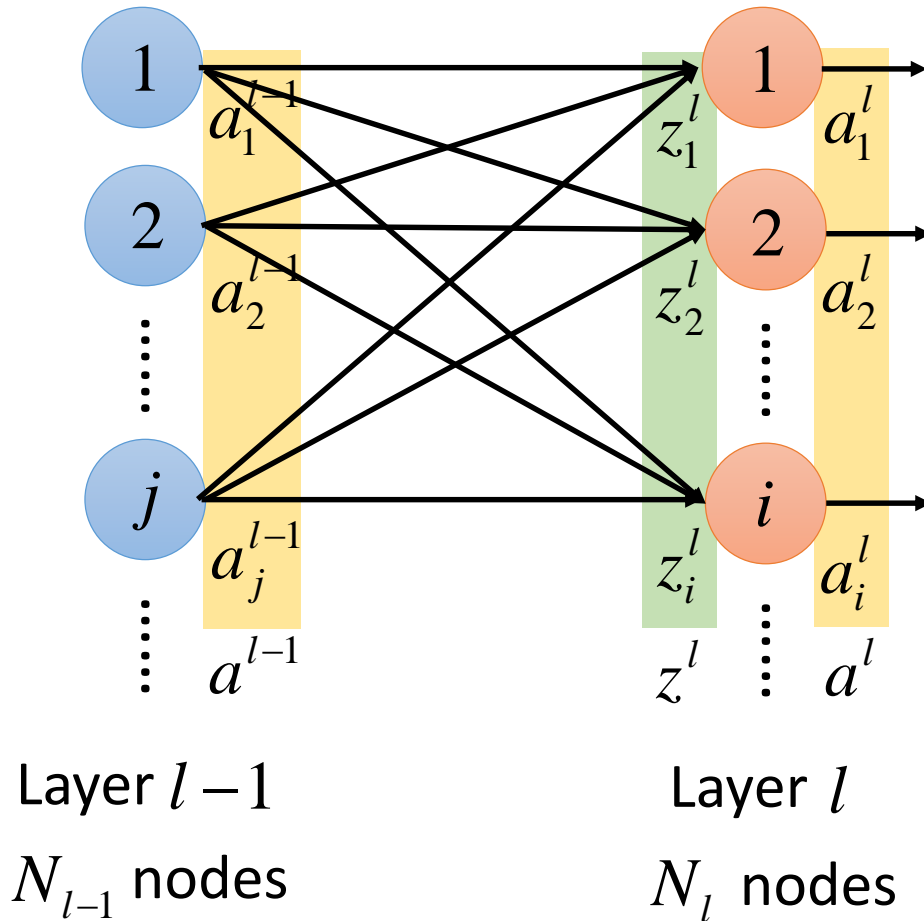


$$a_i^l = \sigma(z_i^l)$$

$$\begin{bmatrix} a_1^l \\ a_2^l \\ \vdots \\ a_i^l \\ \vdots \end{bmatrix} = \begin{bmatrix} \sigma(z_1^l) \\ \sigma(z_2^l) \\ \vdots \\ \sigma(z_i^l) \\ \vdots \end{bmatrix}$$

$$a^l = \sigma(z^l)$$

# Relations between Layer Outputs

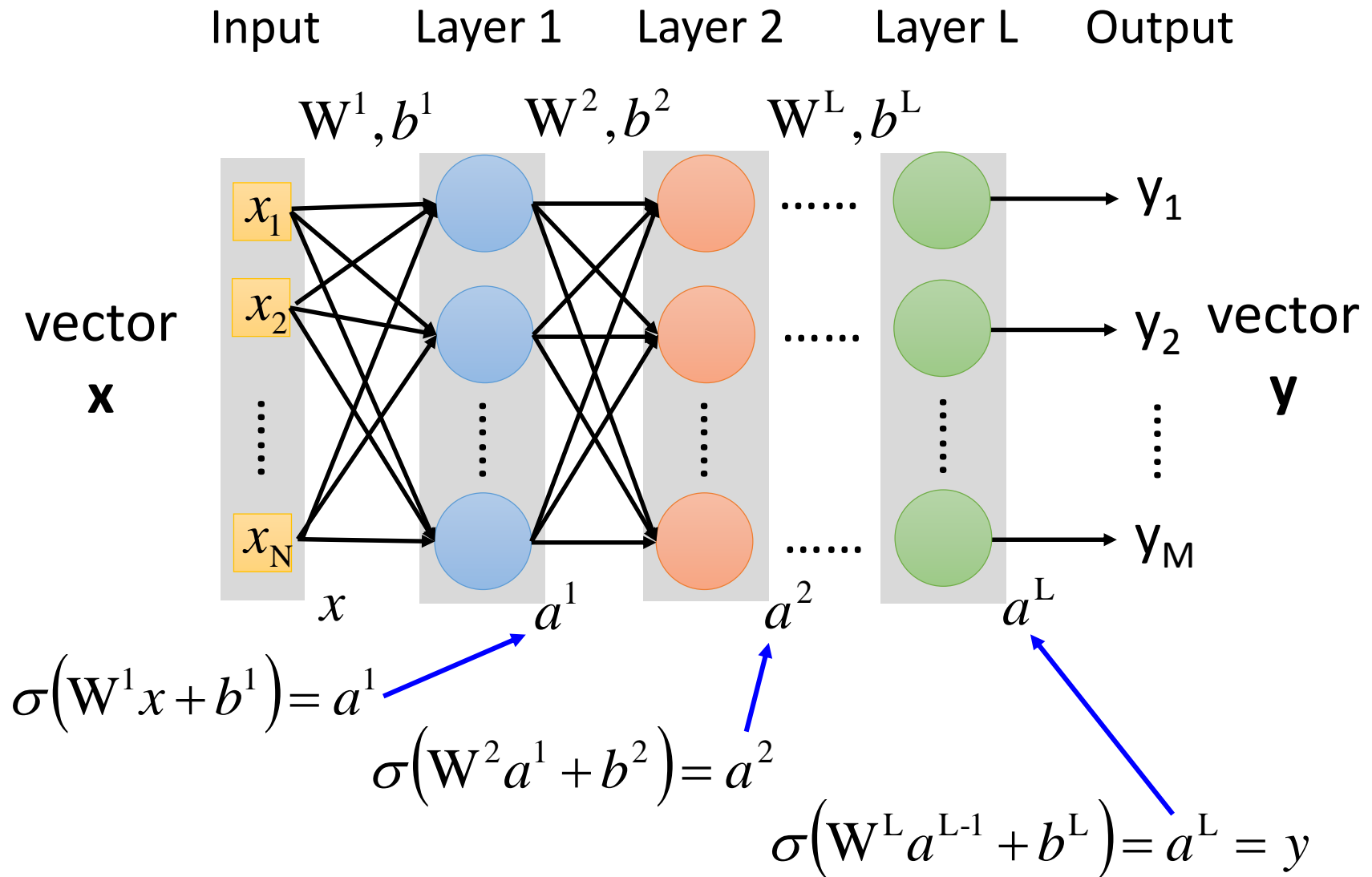


$$z^l = W^l a^{l-1} + b^l$$

$$a^l = \sigma(z^l)$$

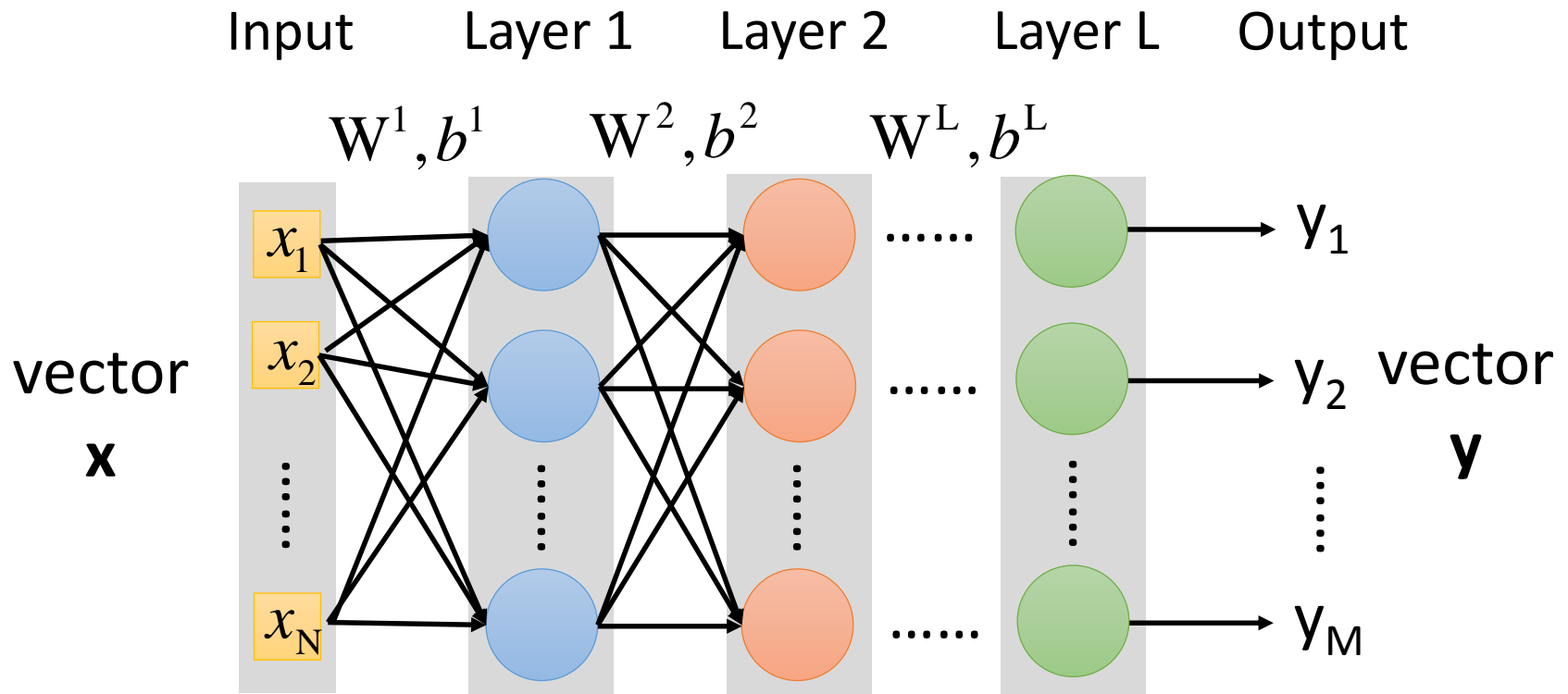
$$a^l = \sigma(W^l a^{l-1} + b^l)$$

# Function of Neural Network





# Function of Neural Network



$$y = f(x)$$

$$= \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

2. What is the “best”  
function?

# Best Function = Best Parameters

$$y = f(x) = \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

function set

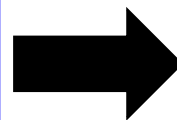
because different parameters  $W$  and  $b$  lead to different function

Formal way to define a function set:

$f(x; \theta) \rightarrow$  parameter set

$$\theta = \{W^1, b^1, W^2, b^2 \dots W^L, b^L\}$$

Pick the “best”  
function  $f^*$



Pick the “best”  
parameter set  $\theta^*$

# Cost Function

- Define a function for parameter set  $C(\theta)$ 
  - $C(\theta)$  evaluate how bad a parameter set is
  - The best parameter set  $\theta^*$  is the one that minimizes  $C(\theta)$

$$\theta^* = \arg \min_{\theta} C(\theta)$$

- $C(\theta)$  is called ***cost/loss/error function***
  - If you define the goodness of the parameter set by another function  $O(\theta)$
  - $O(\theta)$  is called objective function

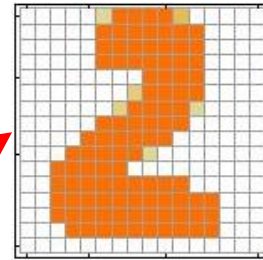
# Cost Function

Given training data:

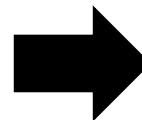
$$\{(x^1, \hat{y}^1) \dots (x^r, \hat{y}^r) \dots (x^R, \hat{y}^R)\}$$

- Handwriting Digit Classification

sum over all  
training examples



$$C(\theta) = \frac{1}{R} \sum_r \| f(x^r; \theta) - \hat{y}^r \|$$



Minimize distance

$$\begin{bmatrix} 0.1 \\ 0.4 \\ 0.2 \\ \vdots \end{bmatrix} \begin{matrix} \text{"1"} \\ \text{"2"} \\ \text{"3"} \end{matrix}$$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \begin{matrix} \text{"1"} \\ \text{"2"} \\ \text{"3"} \end{matrix}$$

3. How to pick  
the “best” function?

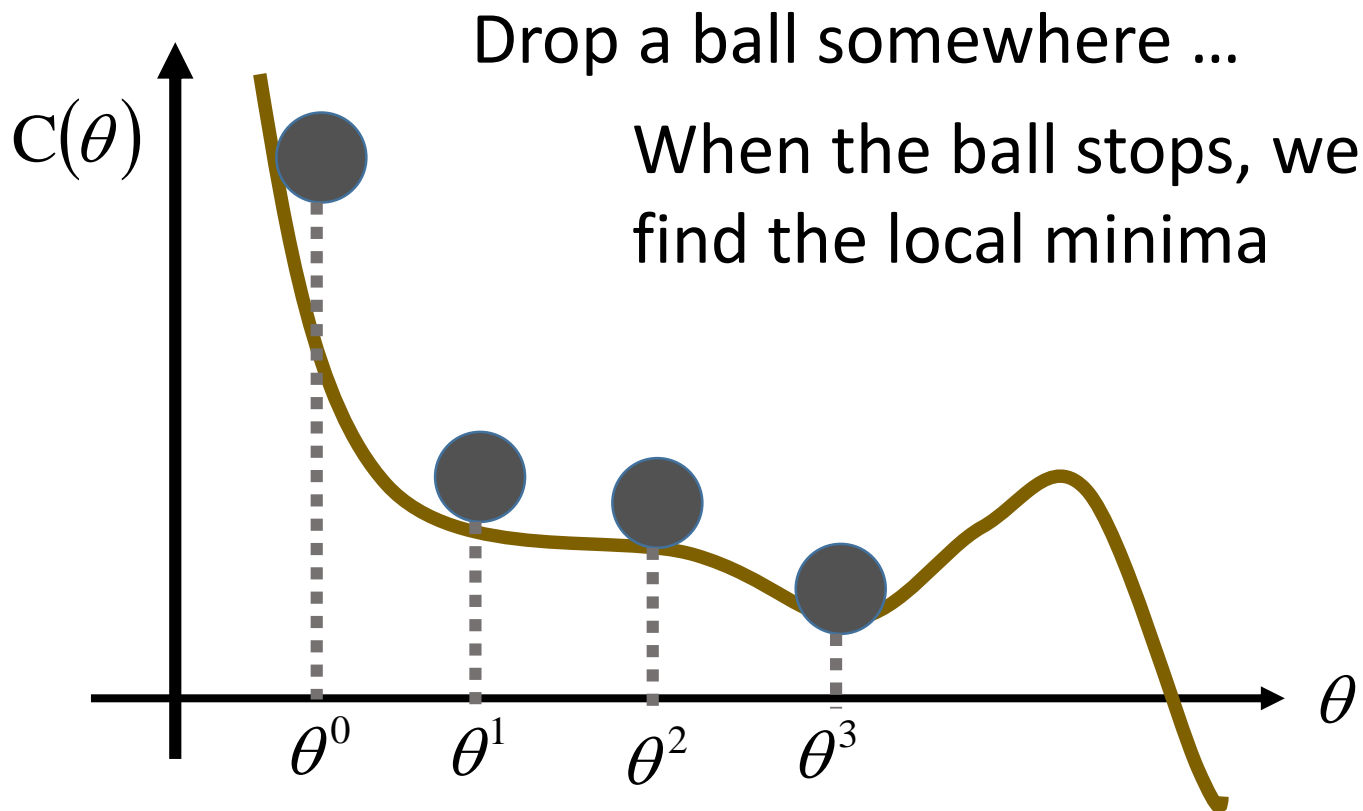
Gradient Descent

# Statement of Problems

- Statement of problems:
  - There is a function  $C(\theta)$ 
    - $\theta$  represents parameter set
    - $\theta = \{\theta_1, \theta_2, \theta_3, \dots\}$
  - Find  $\theta^*$  that minimizes  $C(\theta)$
- Brute force?
  - Enumerate all possible  $\theta$
- Calculus?
  - Find  $\theta^*$  such that  $\left. \frac{\partial C(\theta)}{\partial \theta_1} \right|_{\theta=\theta^*} = 0, \left. \frac{\partial C(\theta)}{\partial \theta_2} \right|_{\theta=\theta^*} = 0, \dots$

# Gradient Descent – Idea

- For simplification, first consider that  $\theta$  has only one variable

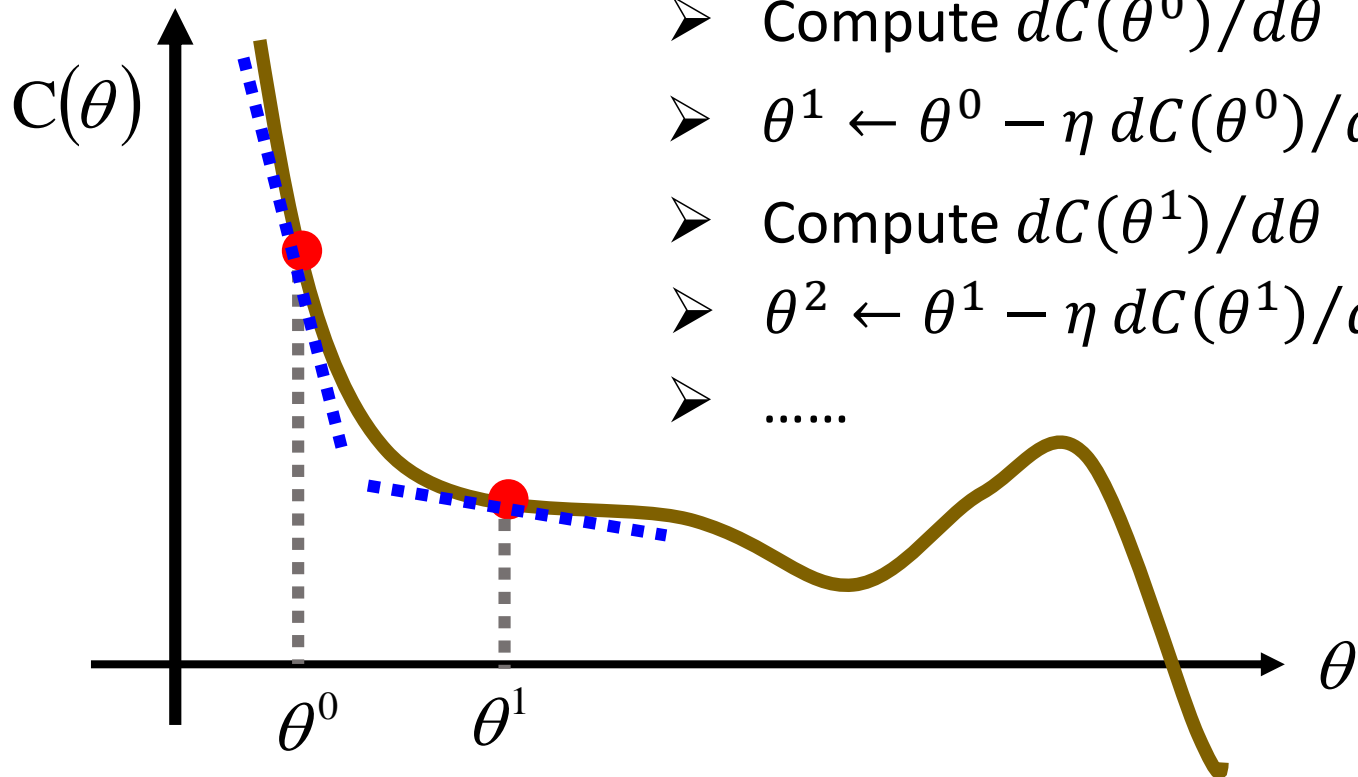




# Gradient Descent – Idea

$\eta$  is called  
*“learning rate”*

- For simplification, first consider that  $\theta$  has only one variable



- Randomly start at  $\theta^0$
- Compute  $dC(\theta^0)/d\theta$
- $\theta^1 \leftarrow \theta^0 - \eta dC(\theta^0)/d\theta$
- Compute  $dC(\theta^1)/d\theta$
- $\theta^2 \leftarrow \theta^1 - \eta dC(\theta^1)/d\theta$
- .....

# Gradient Descent

- Suppose that  $\theta$  has two variables  $\{\theta_1, \theta_2\}$

➤ Randomly start at  $\theta^0 = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix}$

➤ Compute the gradients of  $C(\theta)$  at  $\theta^0$ :  $\nabla C(\theta^0) = \begin{bmatrix} \partial C(\theta_1^0)/\partial \theta_1 \\ \partial C(\theta_2^0)/\partial \theta_2 \end{bmatrix}$

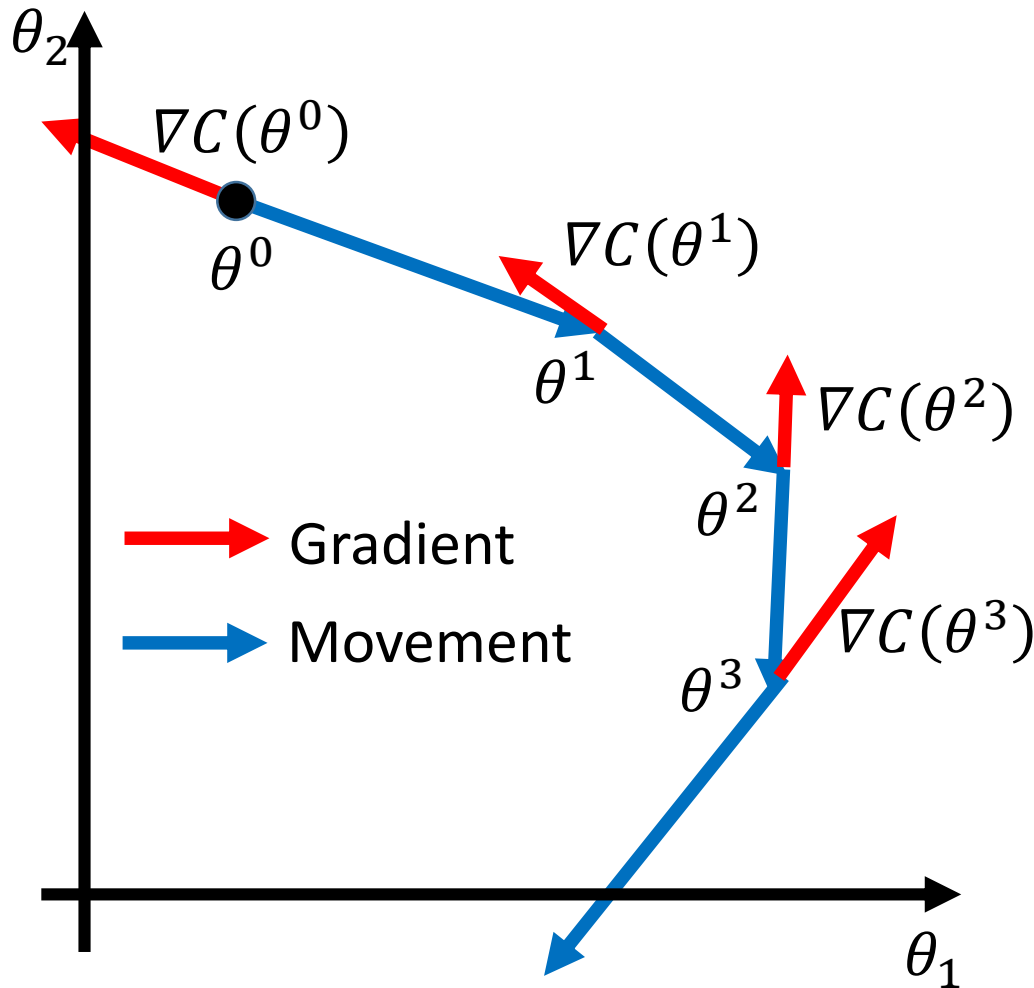
➤ Update parameters

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \end{bmatrix} = \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \end{bmatrix} - \eta \begin{bmatrix} \partial C(\theta_1^0)/\partial \theta_1 \\ \partial C(\theta_2^0)/\partial \theta_2 \end{bmatrix} \Rightarrow \theta^1 = \theta^0 - \eta \nabla C(\theta^0)$$

➤ Compute the gradients of  $C(\theta)$  at  $\theta^1$ :  $\nabla C(\theta^1) = \begin{bmatrix} \partial C(\theta_1^1)/\partial \theta_1 \\ \partial C(\theta_2^1)/\partial \theta_2 \end{bmatrix}$

➤ .....

# Gradient Descent



Start at position  $\theta^0$

Compute gradient at  $\theta^0$

Move to  $\theta^1 = \theta^0 - \eta \nabla C(\theta^0)$

Compute gradient at  $\theta^1$

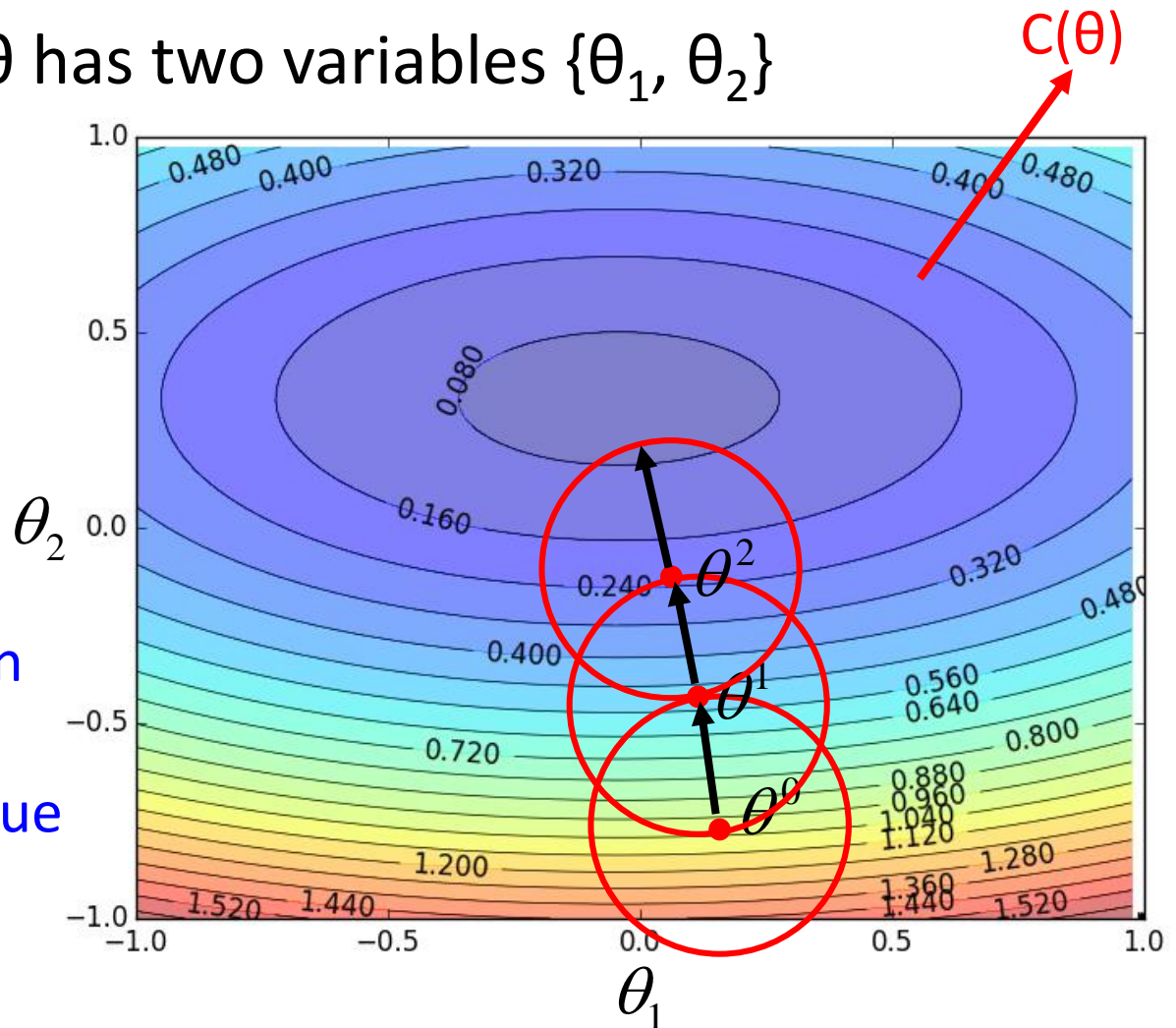
Move to  $\theta^2 = \theta^1 - \eta \nabla C(\theta^1)$

$\vdots$

# Formal Derivation of Gradient Descent

- Suppose that  $\theta$  has two variables  $\{\theta_1, \theta_2\}$

Given a point, we can easily find the point with the smallest value nearby.



# Gradient Descent for Neural Network

$$\begin{aligned} & \text{Compute } \nabla C(\theta^0) \\ & \theta^1 = \theta^0 - \eta \nabla C(\theta^0) \\ & \text{Compute } \nabla C(\theta^1) \\ & \theta^2 = \theta^1 - \eta \nabla C(\theta^1) \end{aligned}$$

Starting  
Parameters

$$\theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2 \longrightarrow \dots$$

$$\nabla C(\theta)$$

$$\theta = \{W^1, b^1, W^2, b^2, \dots, W^l, b^l, \dots, W^L, b^L\}$$

$$= \begin{bmatrix} \vdots \\ \frac{\partial C(\theta)}{\partial w_{ij}^l} \\ \vdots \\ \frac{\partial C(\theta)}{\partial b_i^l} \\ \vdots \end{bmatrix}$$

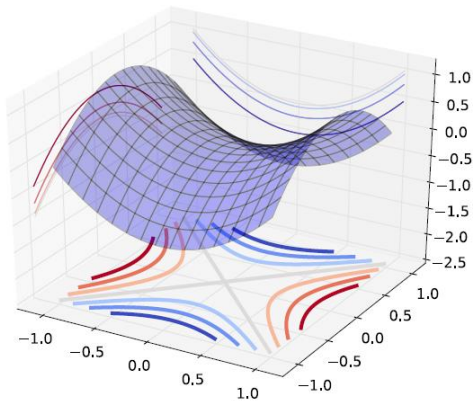
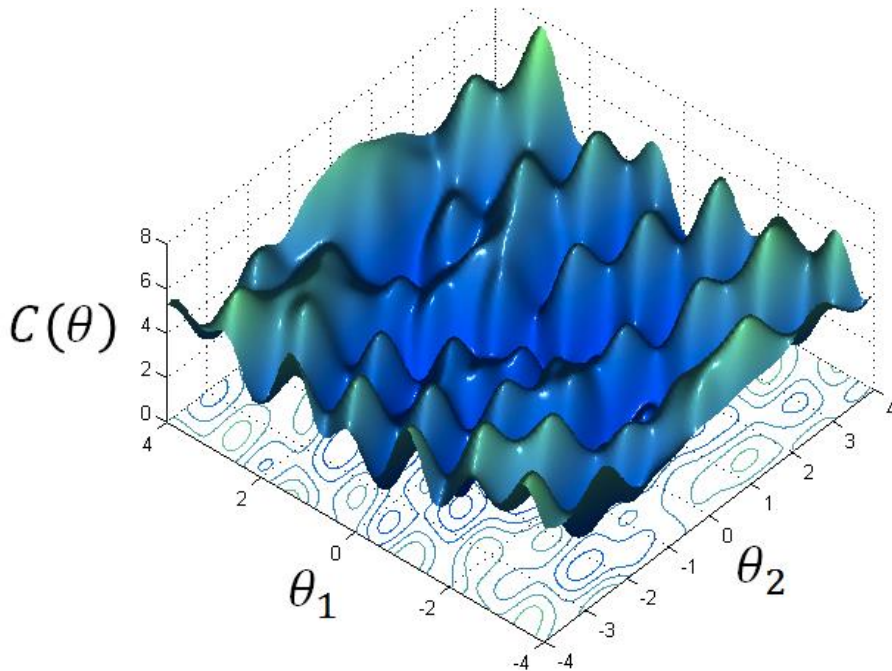
$$\begin{bmatrix} w_{11}^l & w_{12}^l & \dots \\ w_{21}^l & w_{22}^l & \\ \vdots & & \ddots \end{bmatrix}$$

$$\begin{bmatrix} b_1^l \\ b_2^l \\ \vdots \\ b_i^l \\ \vdots \end{bmatrix}$$

Millions of parameters .....

To compute the gradients efficiently, we  
use backpropagation.

# Stuck at local minima?



Saddle  
point

- Who is Afraid of Non-Convex Loss Functions?
- [http://videolectures.net/eml07\\_lecun\\_wia/](http://videolectures.net/eml07_lecun_wia/)
- Deep Learning: Theoretical Motivations
- [http://videolectures.net/deeplearning2015\\_bengio\\_theoretical\\_motivations/](http://videolectures.net/deeplearning2015_bengio_theoretical_motivations/)

3. How to pick  
the “best” function?

Practical Issues  
for neural network

# Practical Issues for neural network

- Parameter Initialization
- Learning Rate
- Stochastic gradient descent and Mini-batch
- Recipe for Learning



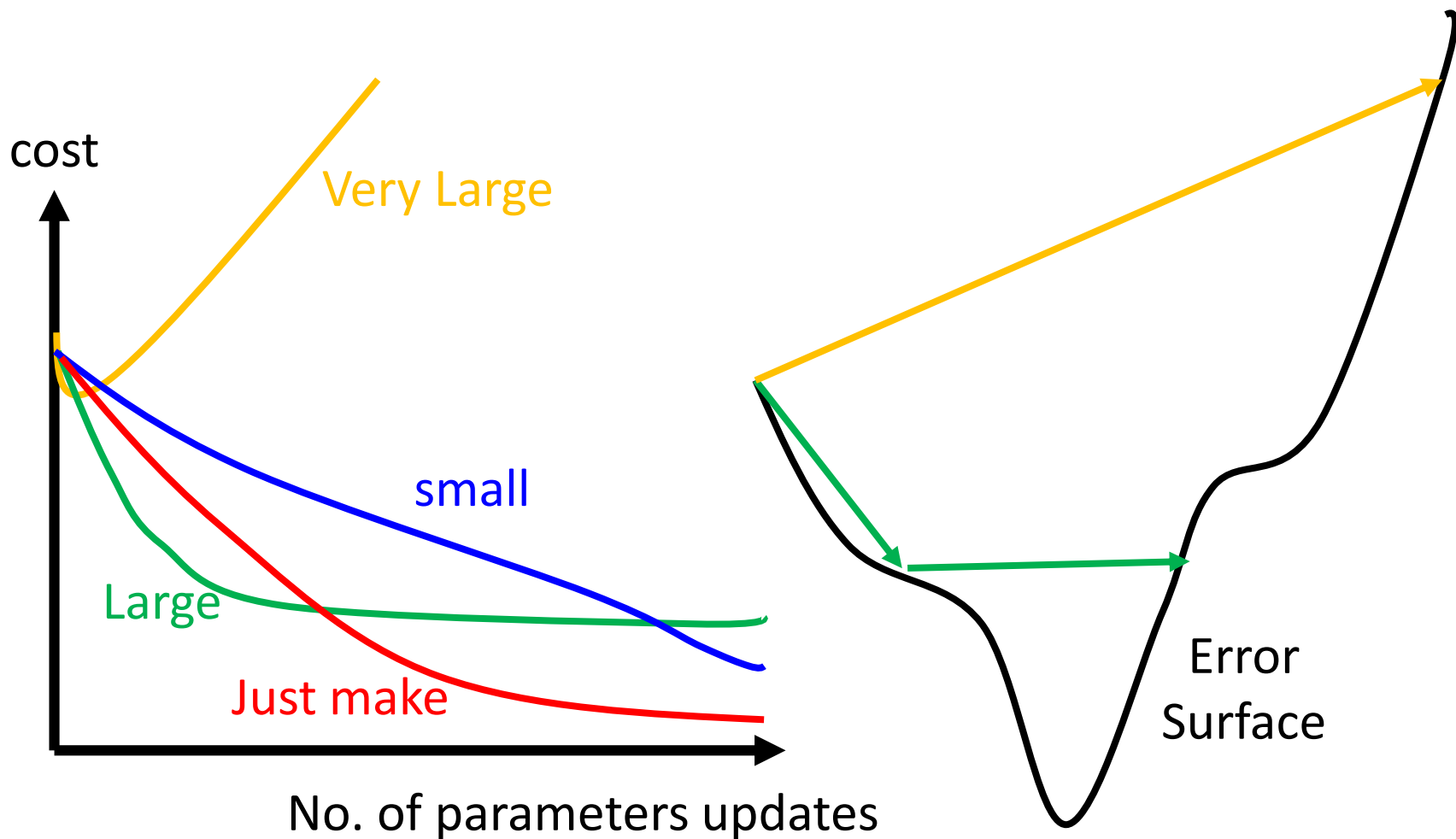
# Parameter Initialization

- For gradient Descent, we need to pick an initialization parameter  $\theta^0$ .
- The initialization parameters have some influence to the training.
  - We will go back to this issue in the future.
- Suggestion today:
  - Do not set all the parameters  $\theta^0$  equal
  - Set the parameters in  $\theta^0$  randomly

# Learning Rate

$$\theta^i = \theta^{i-1} - \eta \nabla C(\theta^{i-1})$$

- Set the learning rate  $\eta$  carefully

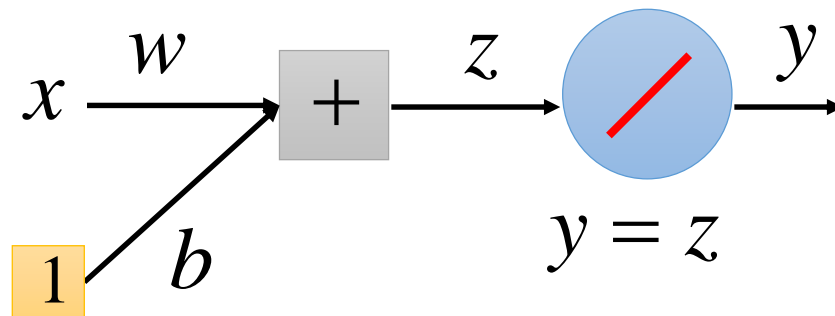


# Learning Rate

$$\theta^i = \theta^{i-1} - \eta \nabla C(\theta^{i-1})$$

- Set the learning rate  $\eta$  carefully

- **Toy Example**



$$\theta^* = \begin{bmatrix} w = 1 \\ b = 0 \end{bmatrix}$$

Training Data (20 examples)

$x = [0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5, 4.0, 4.5, 5.0, 5.5, 6.0, 6.5, 7.0, 7.5, 8.0, 8.5, 9.0, 9.5]$

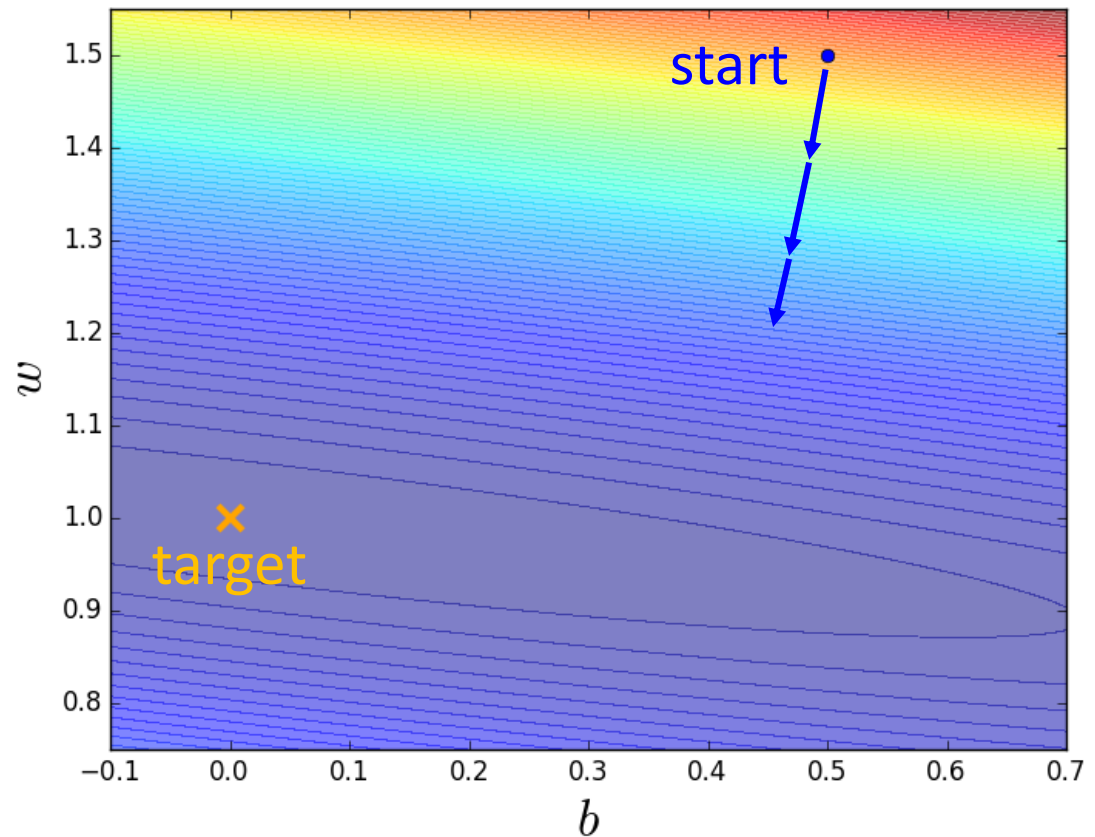
$y = [0.1, 0.4, 0.9, 1.6, 2.2, 2.5, 2.8, 3.5, 3.9, 4.7, 5.1, 5.3, 6.3, 6.5, 6.7, 7.5, 8.1, 8.5, 8.9, 9.5]$

# Learning Rate

$$\theta^i = \theta^{i-1} - \eta \nabla C(\theta^{i-1})$$

- ***Toy Example***

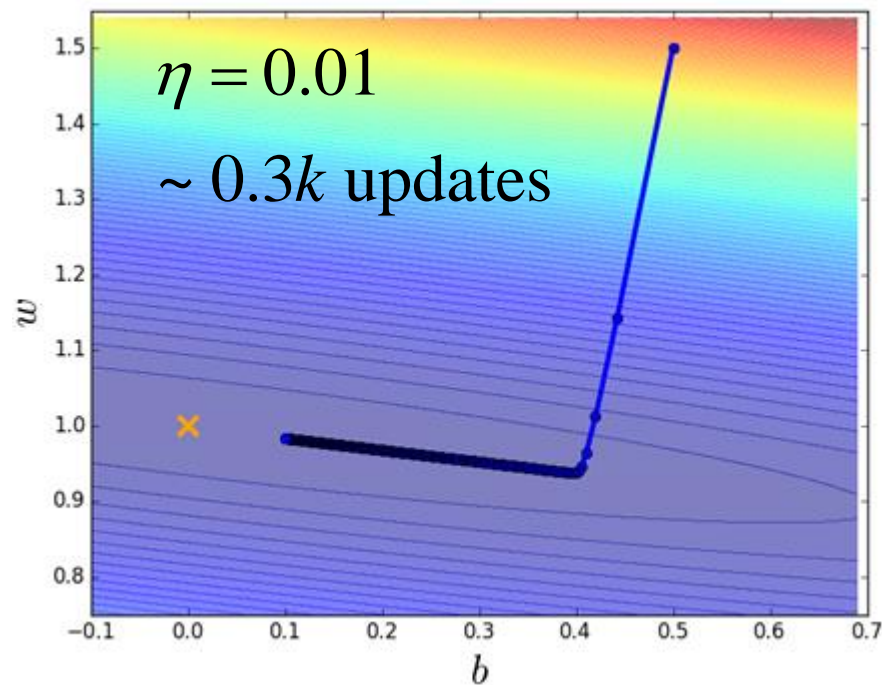
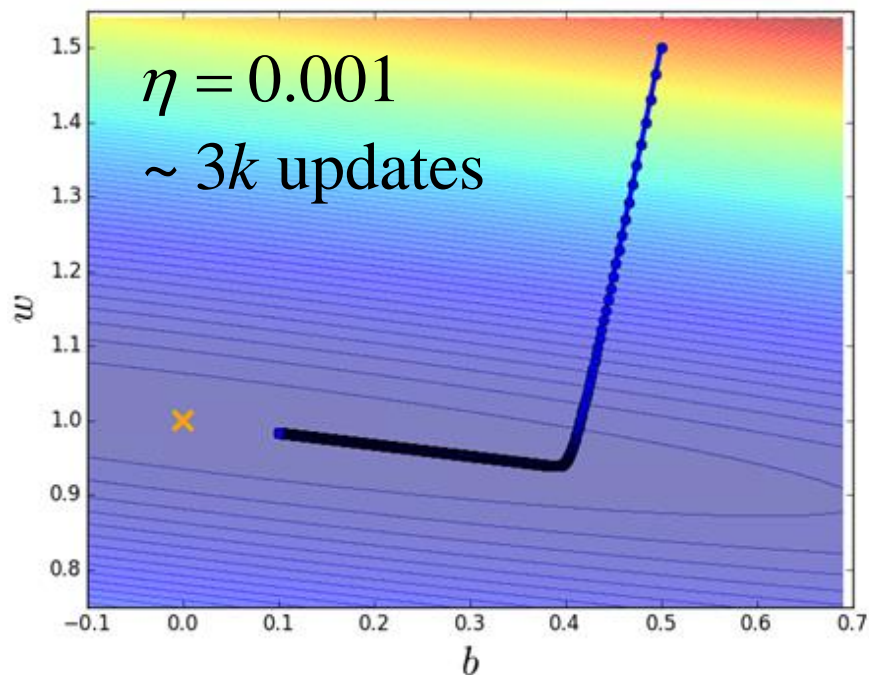
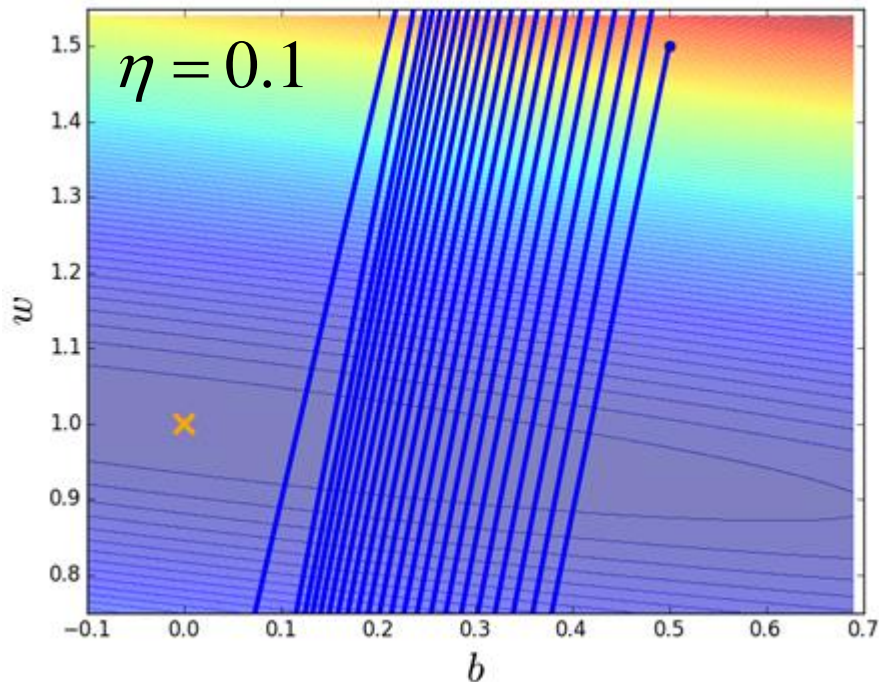
Error Surface:  $C(w,b)$



# Learning Rate

- **Toy Example**

Different learning rate  $\eta$



# Stochastic Gradient Descent and Mini-batch

$$C(\theta) = \frac{1}{R} \sum_r \|f(x^r; \theta) - \hat{y}^r\|$$
$$= \frac{1}{R} \sum_r C^r(\theta)$$

## ◆ Gradient Descent

$$\theta^i = \theta^{i-1} - \eta \nabla C(\theta^{i-1})$$
$$\nabla C(\theta^{i-1}) = \frac{1}{R} \sum_r \nabla C^r(\theta^{i-1})$$

## ◆ Stochastic Gradient Descent

Faster!

Better!

Pick an example  $x^r$

$$\theta^i = \theta^{i-1} - \eta \nabla C^r(\theta^{i-1})$$

If all example  $x^r$  have  
equal probabilities to  
be picked

$$E[\nabla C^r(\theta^{i-1})] = \frac{1}{R} \sum_r \nabla C^r(\theta^{i-1})$$

# Stochastic Gradient Descent and Mini-batch

What is epoch?

Training Data:  $\{(x^1, \hat{y}^1), (x^2, \hat{y}^2), \dots (x^r, \hat{y}^r), \dots (x^R, \hat{y}^R)\}$

When using stochastic gradient descent

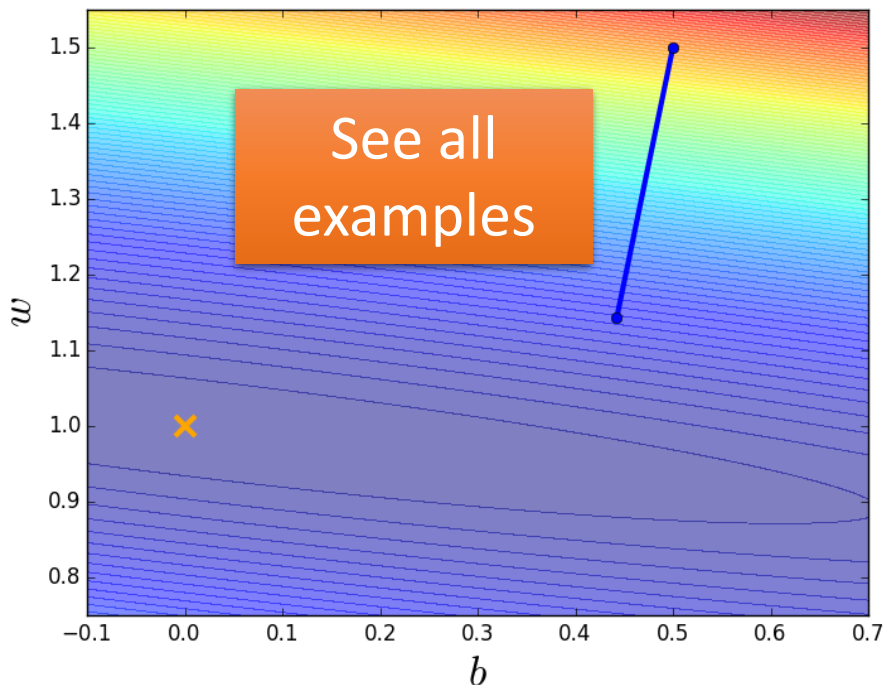
Starting at $\theta_0$	pick $x^1$	$\theta^1 = \theta^0 - \eta \nabla C^1(\theta^0)$	Seen all the examples once <b>One epoch</b>
	pick $x^2$	$\theta^2 = \theta^1 - \eta \nabla C^2(\theta^1)$	
	$\vdots$	$\vdots$	
	pick $x^r$	$\theta^r = \theta^{r-1} - \eta \nabla C^r(\theta^{r-1})$	
	$\vdots$	$\vdots$	
	pick $x^R$	$\theta^R = \theta^{R-1} - \eta \nabla C^R(\theta^{R-1})$	
<hr/>			
	pick $x^1$	$\theta^{R+1} = \theta^R - \eta \nabla C^1(\theta^R)$	

# Stochastic Gradient Descent and Mini-batch

- ***Toy Example***

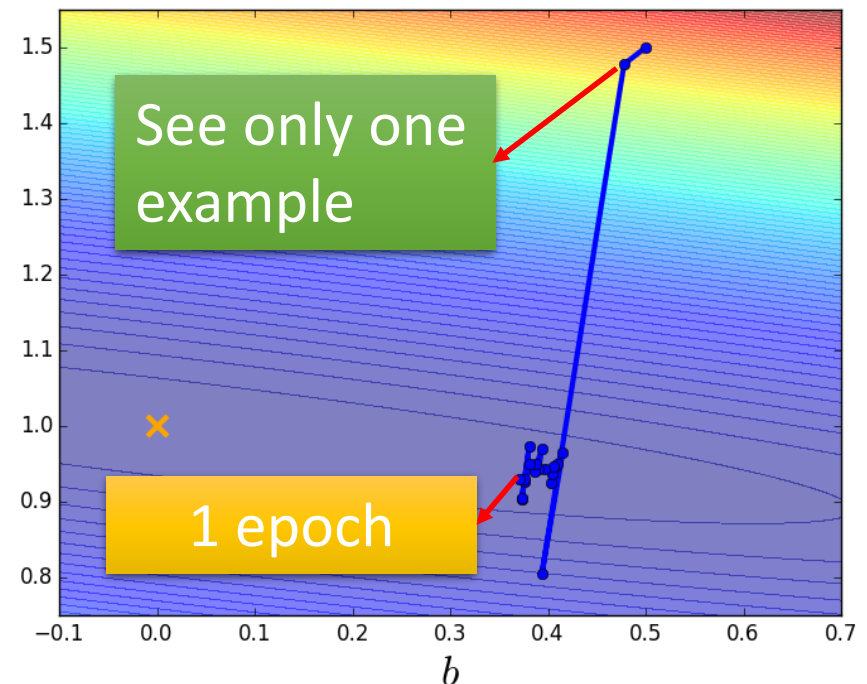
## ***Gradient Descent***

Update after seeing all examples



## ***Stochastic Gradient Descent***

If there are 20 examples, update 20 times in one epoch.





# Stochastic Gradient Descent and Mini-batch

## ◆ Gradient Descent

$$\theta^i = \theta^{i-1} - \eta \nabla C(\theta^{i-1}) \quad \nabla C(\theta^{i-1}) = \frac{1}{R} \sum_r \nabla C^r(\theta^{i-1})$$

## ◆ Stochastic Gradient Descent

Pick an example  $x_r$

$$\theta^i = \theta^{i-1} - \eta \nabla C^r(\theta^{i-1})$$

## ◆ Mini Batch Gradient Descent

Pick B examples as  
a batch b

B is batch size

Shuffle your data

$$\theta^i = \theta^{i-1} - \eta \frac{1}{B} \sum_{x_r \in b} \nabla C^r(\theta^{i-1})$$

Average the gradient of the  
examples in the batch b

# Stochastic Gradient Descent and Mini-batch

- **Handwriting Digit Classification**

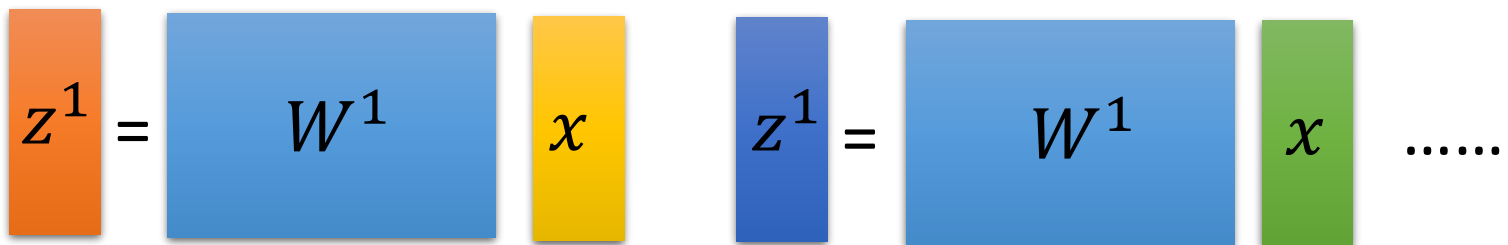


Gradient Descent

# Stochastic Gradient Descent and Mini-batch

- Why mini-batch is faster than stochastic gradient descent?

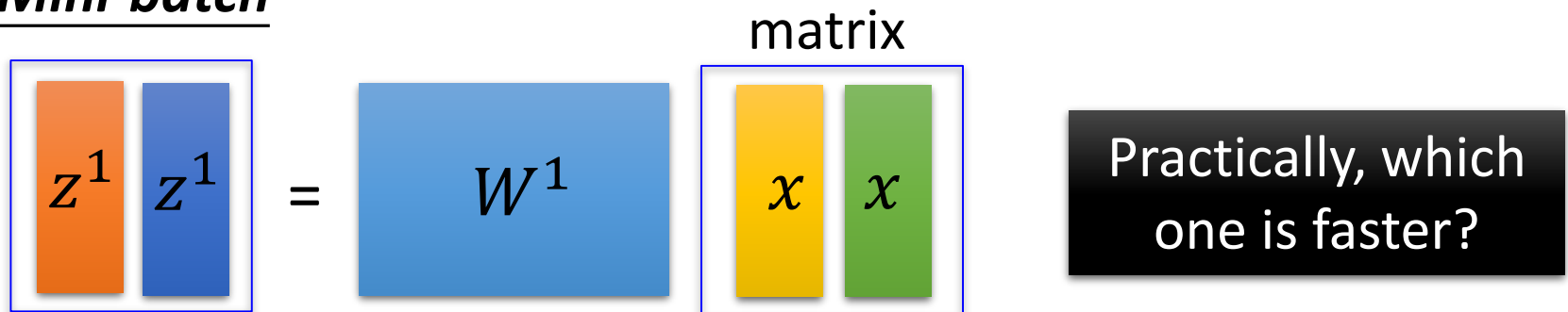
## Stochastic Gradient Descent



The diagram illustrates the Stochastic Gradient Descent process. It shows two separate calculations. In the first, an orange vertical rectangle labeled  $z^1$  is equal to a blue square labeled  $W^1$  multiplied by a yellow vertical rectangle labeled  $x$ . In the second, a blue vertical rectangle labeled  $z^1$  is equal to the same blue square  $W^1$  multiplied by a green vertical rectangle labeled  $x$ . This is followed by an ellipsis, indicating that these calculations are repeated for many different data points.

$$z^1 = W^1 x \quad z^1 = W^1 x \quad \dots$$

## Mini-batch



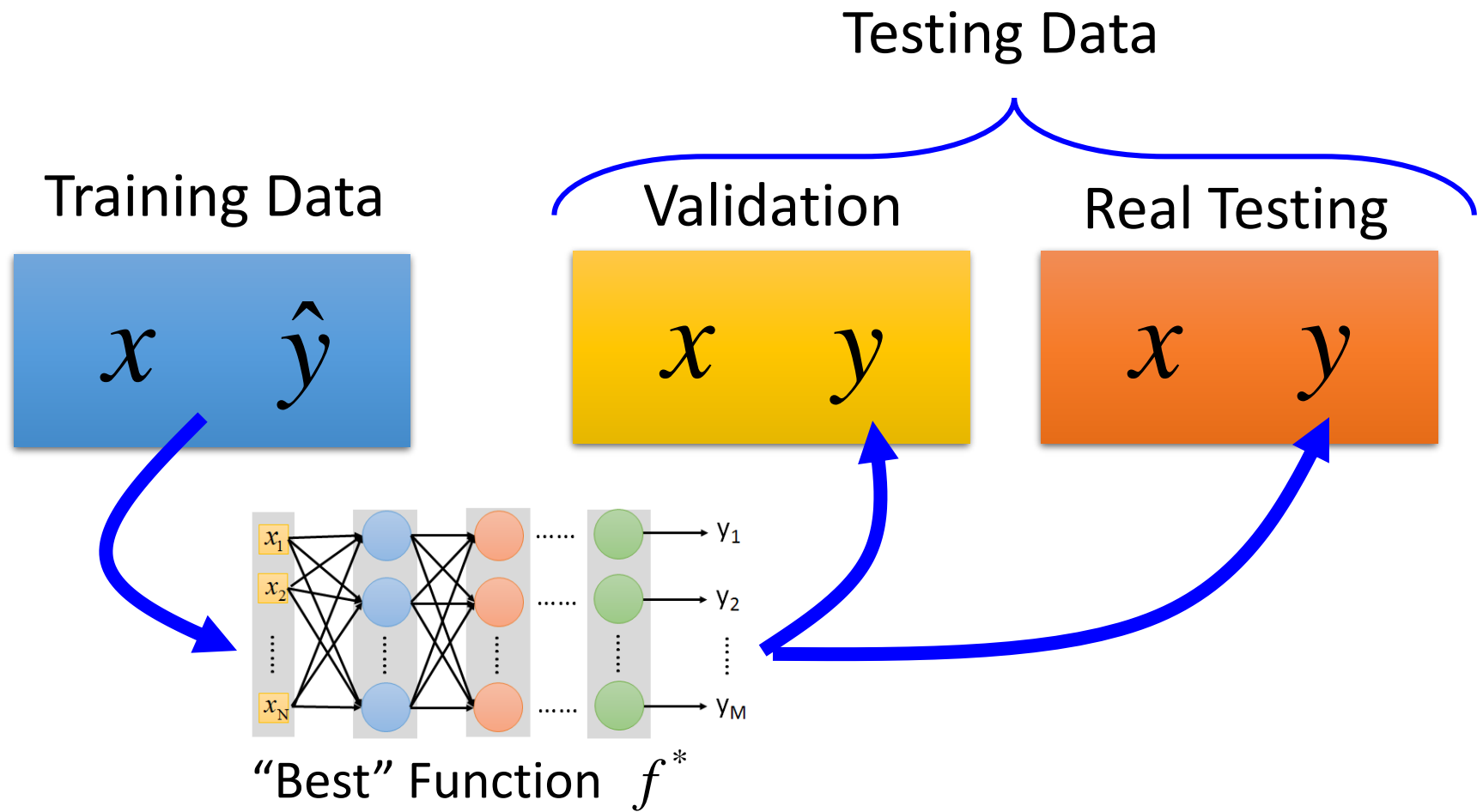
The diagram illustrates the Mini-batch process. On the left, two vertical rectangles (one orange labeled  $z^1$  and one blue labeled  $z^1$ ) are grouped together in a blue-bordered box. This is followed by an equals sign and a blue square labeled  $W^1$ . To the right of the matrix is another blue-bordered box containing two vertical rectangles (one yellow labeled  $x$  and one green labeled  $x$ ). The word "matrix" is written above this second box. To the right of the entire equation is a black box with white text that asks, "Practically, which one is faster?"

$$\begin{bmatrix} z^1 & z^1 \end{bmatrix} = W^1 \begin{bmatrix} x & x \end{bmatrix}$$

matrix

Practically, which one is faster?

# Recipe for Learning



# Recipe for Learning - Overfitting

- You pick a “best” parameter set  $\theta^*$

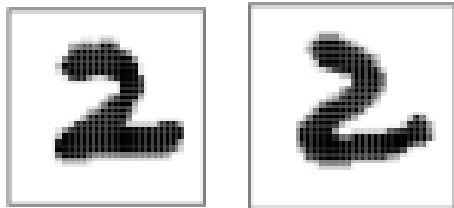
Training Data:  $\{\dots(x^r, \hat{y}^r)\dots\} \Rightarrow \forall r: f(x^r; \theta^*) = \hat{y}^r$

However,

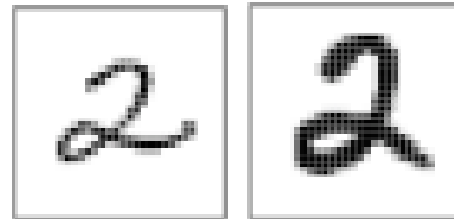
Testing Data:  $\{\dots x^u \dots\} \Rightarrow f(x^u; \theta^*) \neq \hat{y}^u$

Training data and testing data have different distribution.

Training Data:



Testing Data:



# Recipe for Learning - Overfitting

- Panacea: Have more training data
- We will go back to this issue in the future.

# Concluding Remarks

1. What is the model (function hypothesis set)?

Neural Network

2. What is the “best” function?

Cost Function

3. How to pick the “best” function?

Gradient Descent

- Parameter Initialization
- Learning Rate
- Stochastic gradient descent, Mini-batch
- Recipe for Learning