



暑期培训



# CS229 Lecture 2

Machine Learning

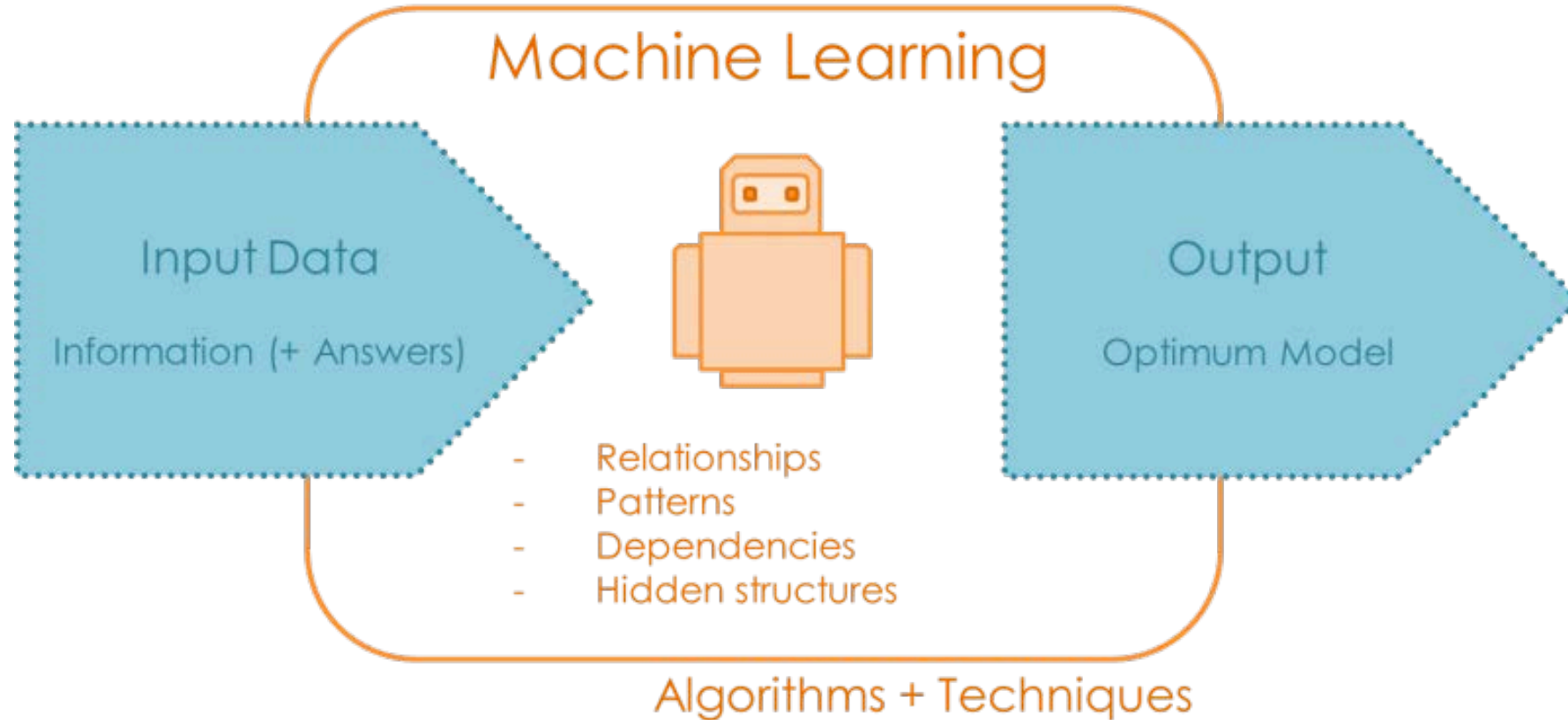
顾晓玲

[guxl@hdu.edu.cn](mailto:guxl@hdu.edu.cn)

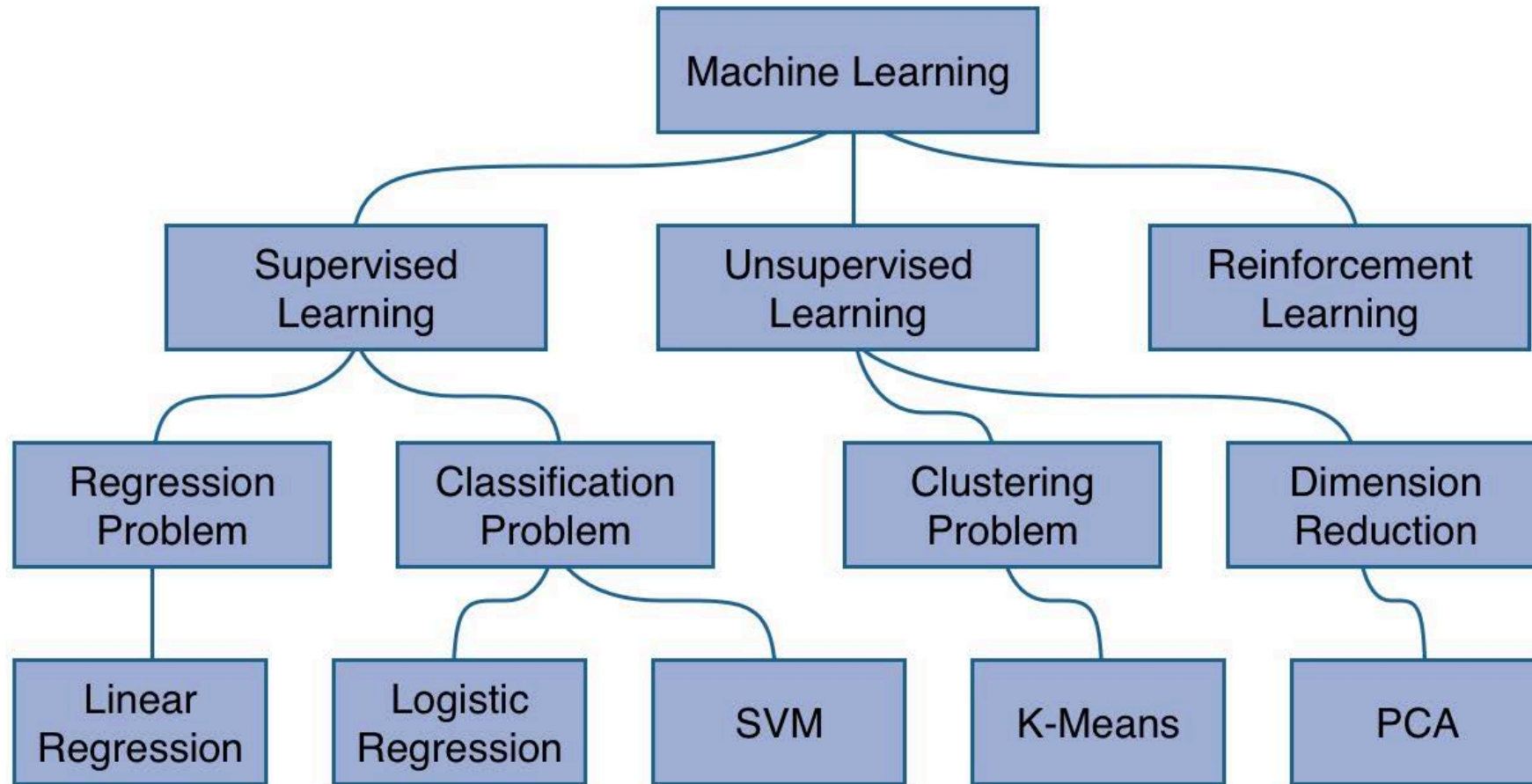
# Outlines

- Machine Learning: Overview
- Linear Regression
- Classification
- Logistic Regression
- Regularization
- Perceptron

# 1 Machine Learning: Overview



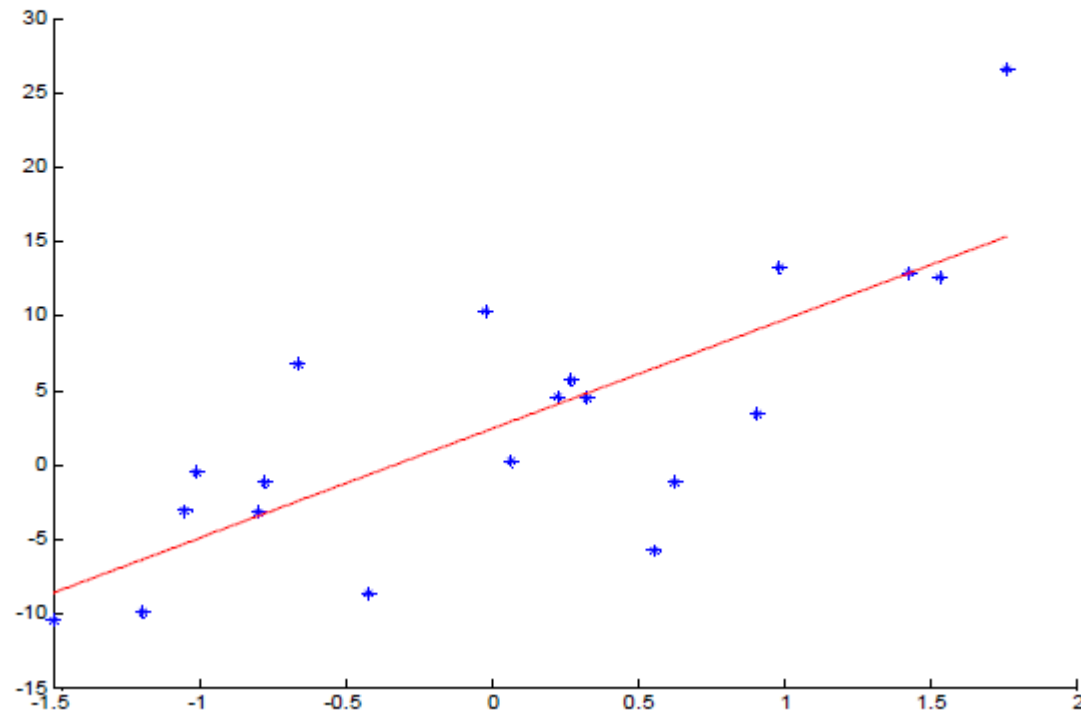
# 1 Machine Learning: Overview



# 2 Linear Regression: 1 Dimensional Input

## Linear regression. Example

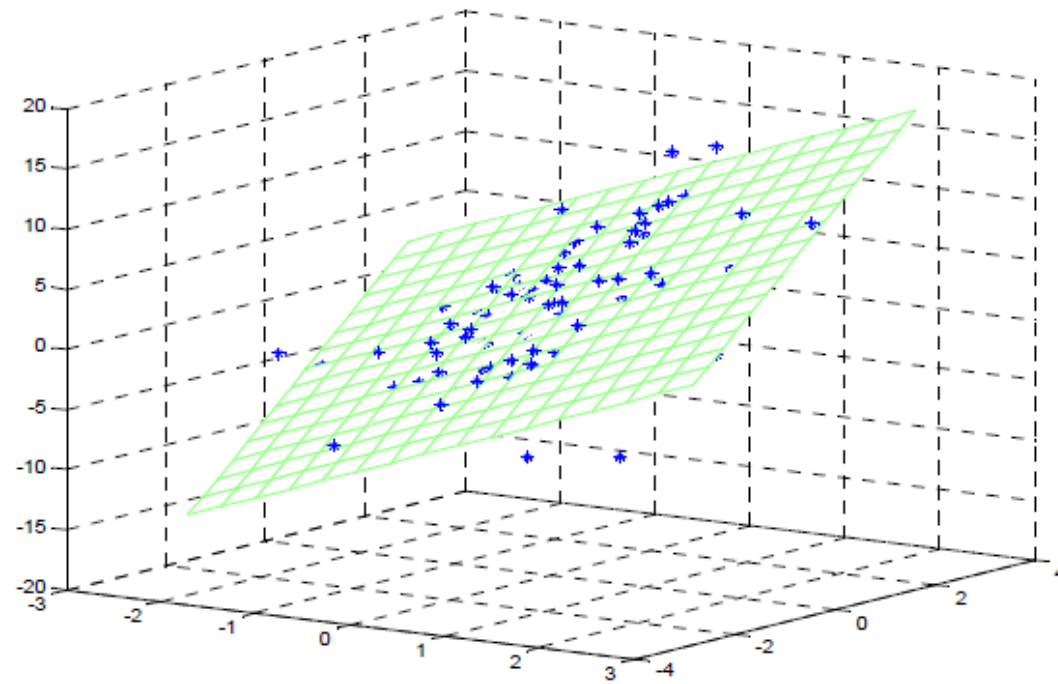
- 1 dimensional input  $\mathbf{x} = (x_1)$



# 2 Linear Regression: 2 Dimensional Input

## Linear regression. Example.

- 2 dimensional input  $\mathbf{x} = (x_1, x_2)$



## 2 Linear Regression: An Example

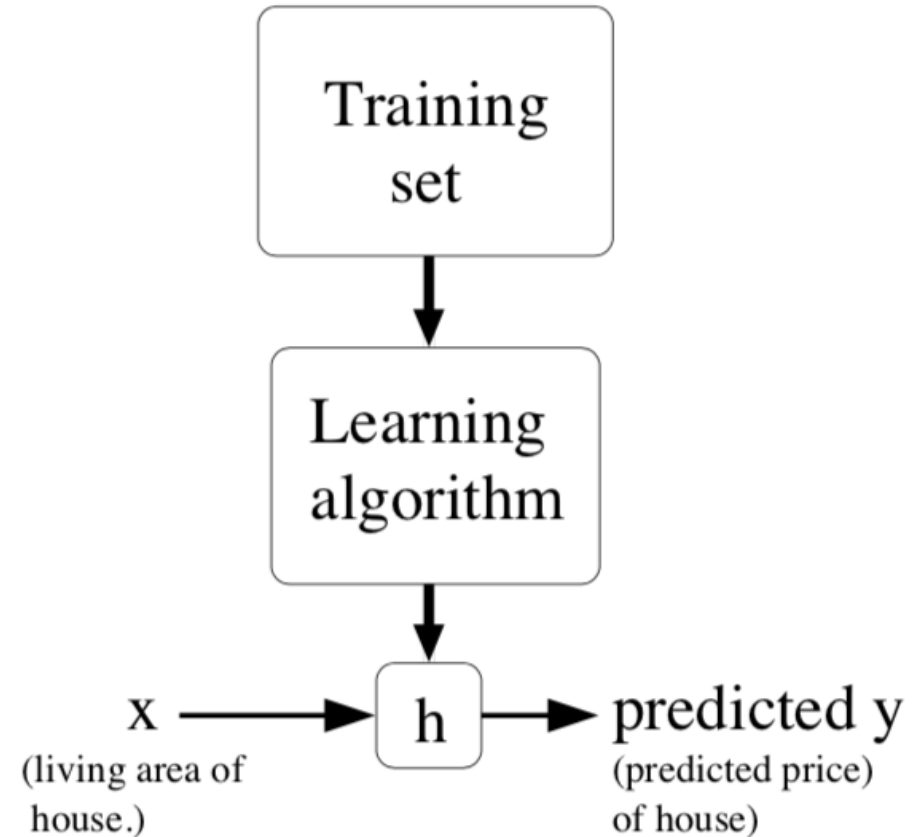
- Suppose we have a dataset giving the living areas, number of bedrooms and prices of 200 houses from a specific region:

Living area (feet <sup>2</sup> )	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮

- Given data like this, how can we learn to predict the prices of other houses, as a function of the size of their living areas and the number of bedrooms?

# 2 Linear Regression: Terms and Concepts

- Sample, Example
- Feature
- Target
- Hypothesis
- Training Data (Training Set)
- Test data (Test Set)
- Training Error & Test Error





# 2 Linear Regression: Hypothesis

Living area (feet <sup>2</sup> )	#bedrooms	Price (1000\$)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
$\vdots$	$\vdots$	$\vdots$

$x_1^{(i)}$  : the living area of the i-th house  
in the training set

$x_2^{(i)}$  : the number of bedrooms of the i-th house  
in the training set

$y^{(i)}$  : the price of the i-th house in the training  
set

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$\theta_i$ 's : **parameters (weights )**

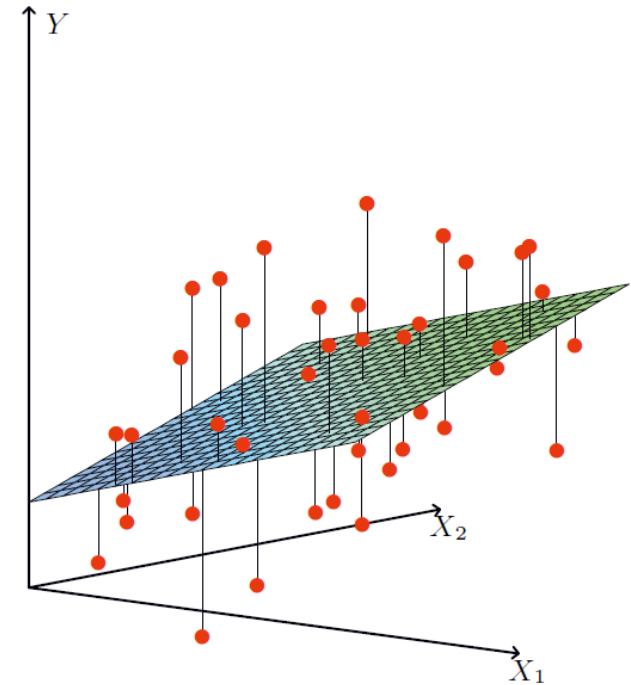
$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

# 2 Linear Regression: Cost Function

- Now, given a training set, how do we learn the parameters  $\theta$ ? One reasonable method seems to be to make  $h(x)$  close to  $y$ .

- **Cost Function (Loss Function):**

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



# 2 Linear Regression: Cost Function

Hypothesis:  $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$

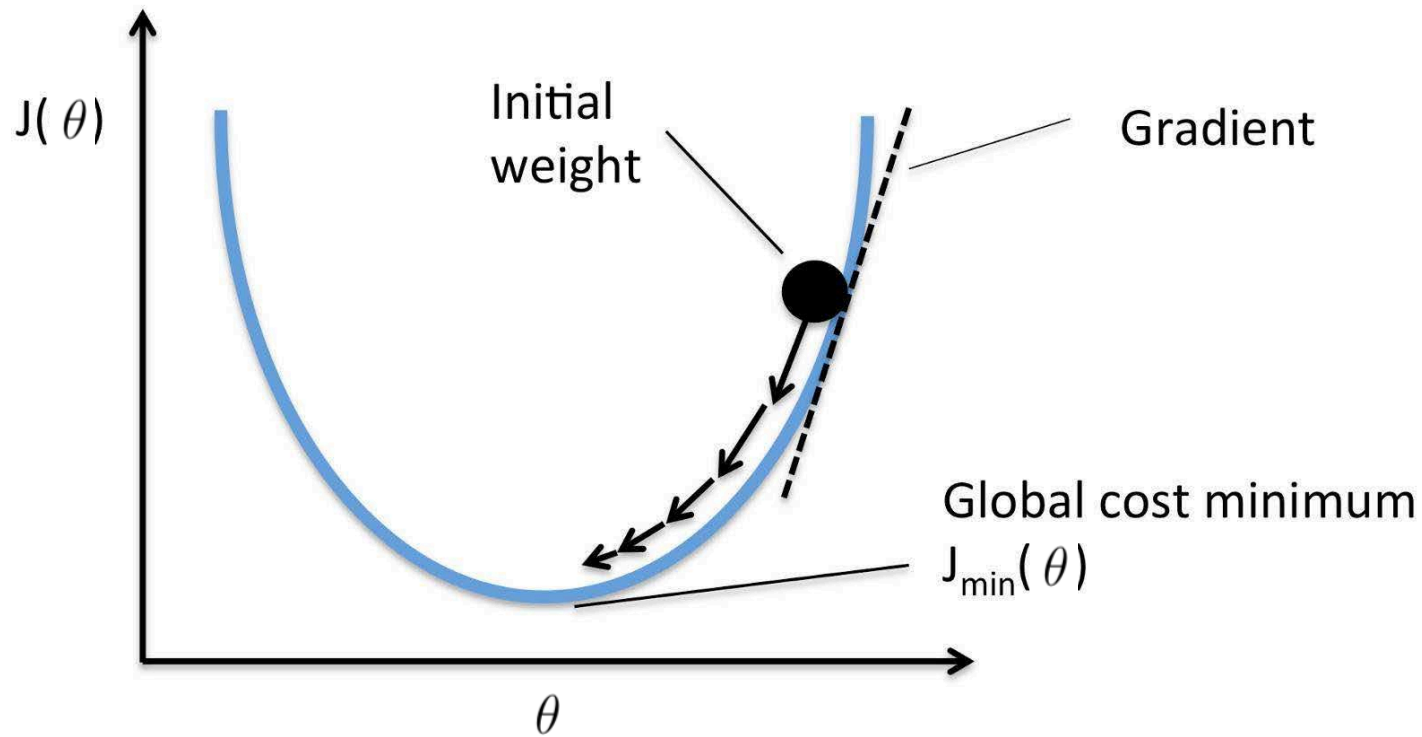
Parameters:  $\theta_0$  ,  $\theta_1$  ,  $\theta_2$

Cost Function:  $J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Goal: Minimize  $J(\theta)$

# 2 Linear Regression: Gradient Descent

- **Gradient Descent Algorithm**



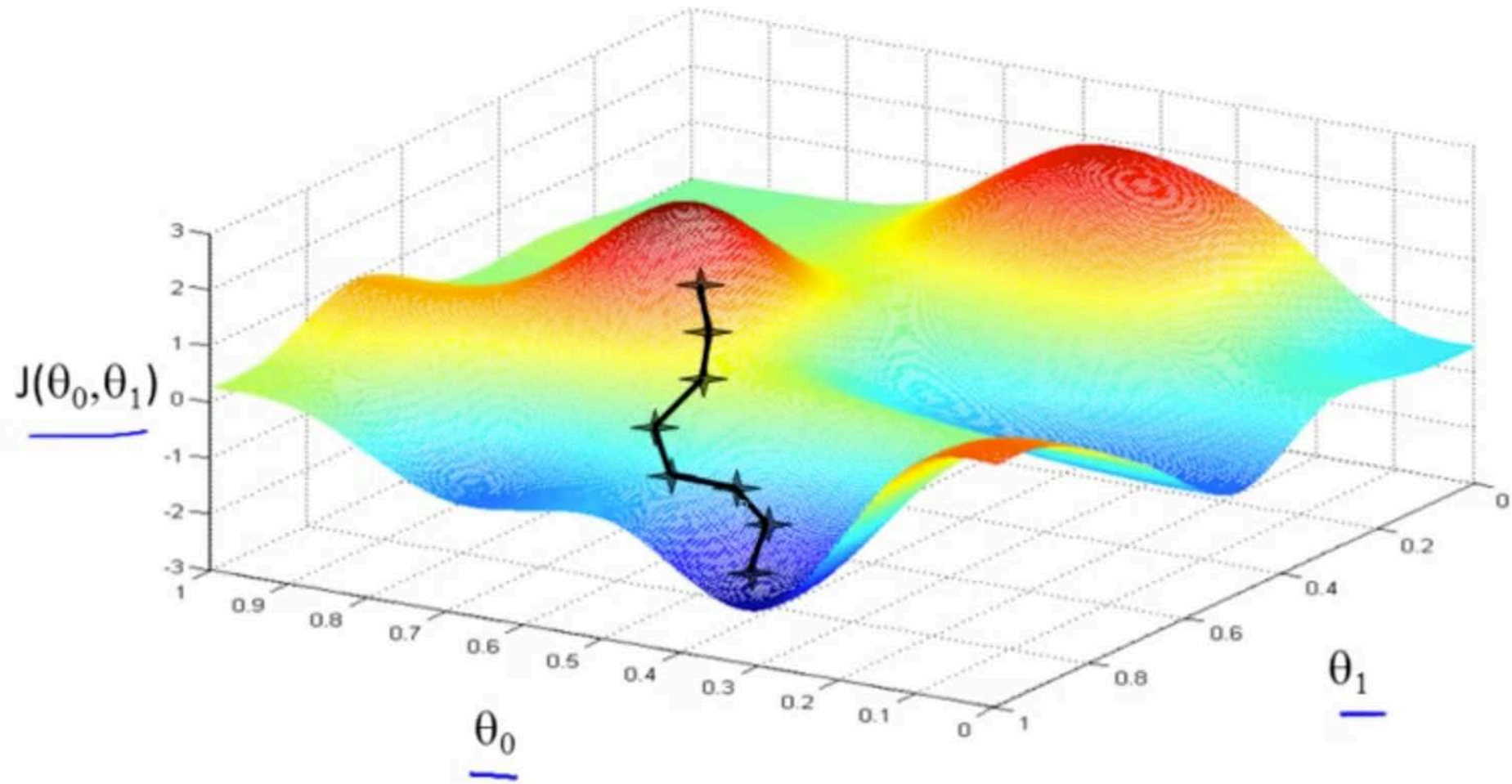
**Update Rule:**

$$\theta := \theta - \alpha \nabla_{\theta} J(\theta)$$

**In Our Case:**

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

# 2 Linear Regression: Gradient Descent



# 2 Linear Regression : Gradient Descent

- Let's assume we have only one training example  $(x, y)$ :

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y) \\ &= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^n \theta_i x_i - y \right) \\ &= (h_{\theta}(x) - y) x_j\end{aligned}$$

- For a single training example, this gives the update rule:

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

# 2 Linear Regression: Gradient Descent

- **Batch Gradient Descent:** looks at every example in the entire training set on every step.

**Old Version**

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \quad (\text{for every } j).$$

}

**New Version**

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

# 2 Linear Regression: Gradient Descent

- **Stochastic Gradient Descent:** update the parameters according to the gradient of the error with respect to that single training example only.

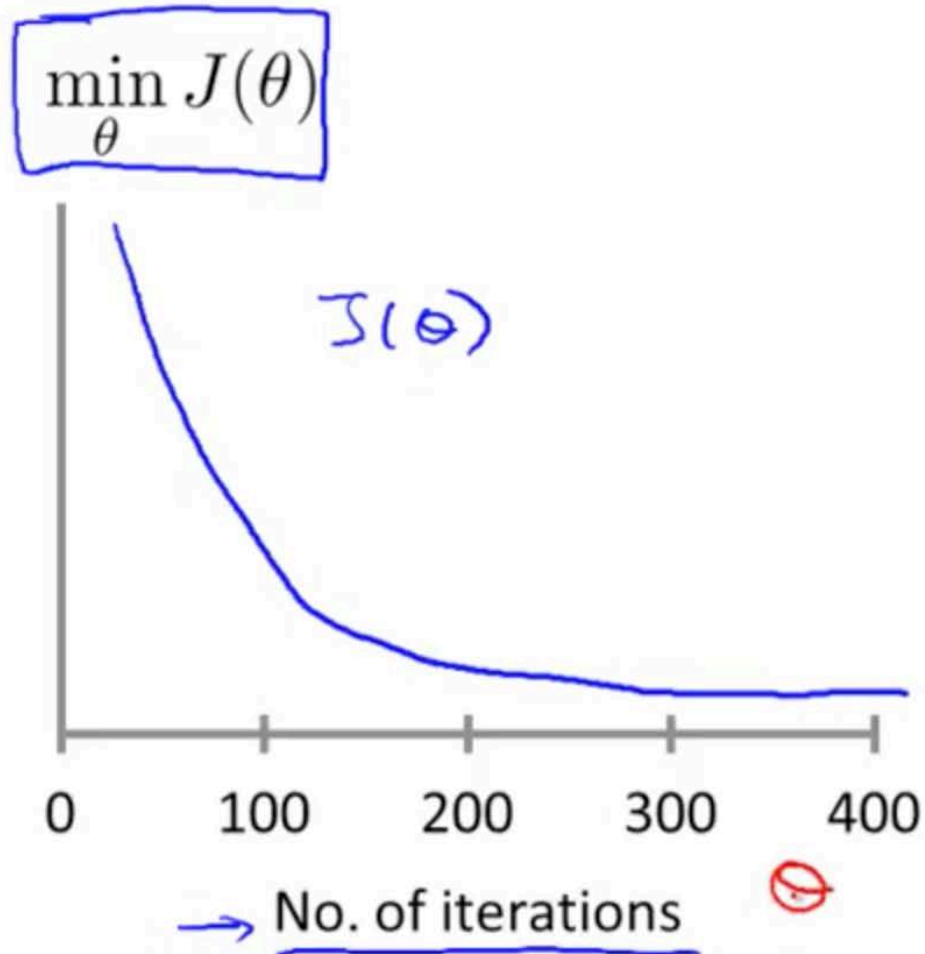
```
Loop {  
    for i=1 to m, {  
         $\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$     (for every  $j$ ).  
    }  
}
```

- Often, stochastic gradient descent gets  $\theta$  “close” to the minimum much faster than batch gradient descent.



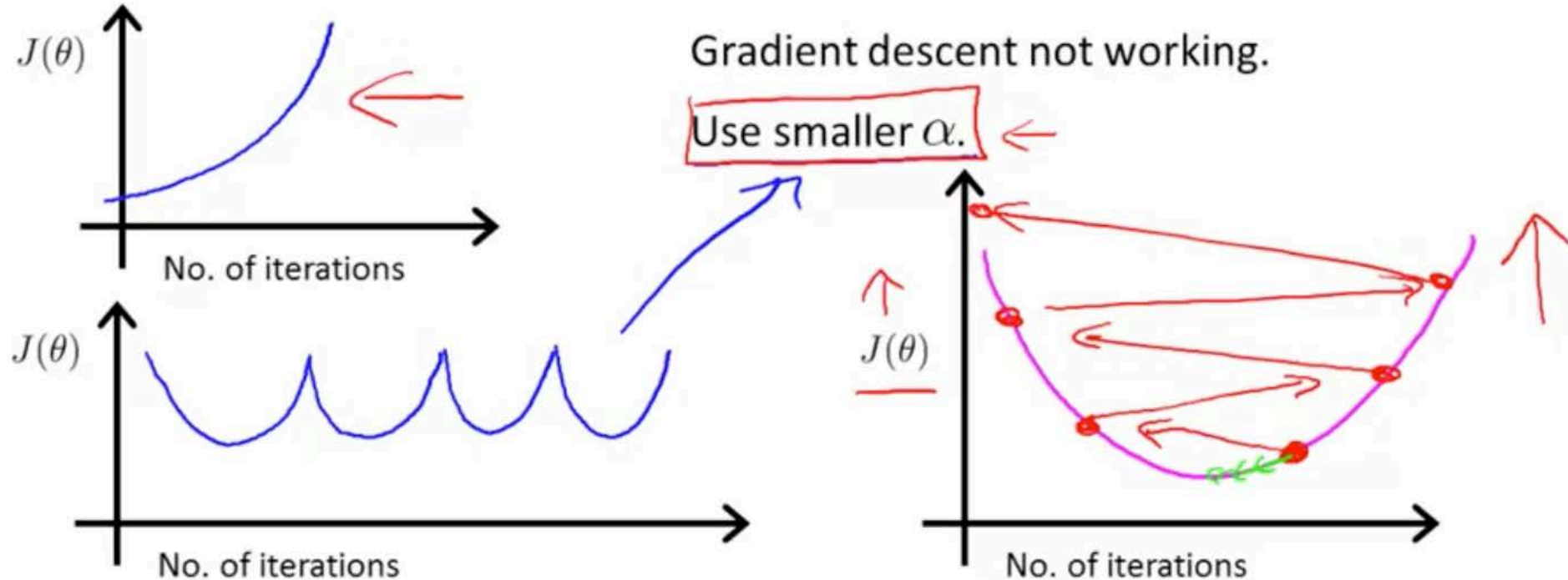
# 2 Linear Regression: Gradient Descent

Making sure gradient descent is working correctly.



# 2 Linear Regression: Gradient Descent

Making sure gradient descent is working correctly.



- For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration. ←
- But if  $\alpha$  is too small, gradient descent can be slow to converge.

# 2 Linear Regression: Normal Equations

- Represent  $J$  as matrix-vectorial notation

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



$$J(\theta) = \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y})$$

# 2 Linear Regression: Normal Equations

- Represent  $J$  as matrix-vectorial notation

$$X = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \vdots \\ \text{---} (x^{(m)})^T \text{---} \end{bmatrix}$$

**m-by-n matrix, actually m-by-n + 1**

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

**m-dimensional vector**

# 2 Linear Regression: Normal Equations

- For example

Living area (feet <sup>2</sup> )	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
$\vdots$	$\vdots$	$\vdots$

$$X = \begin{bmatrix} 1 & 2104 & 3 \\ 1 & 1600 & 3 \\ 1 & 2400 & 3 \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} 400 \\ 330 \\ 369 \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

**Note: more precisely, the features should be scaled.**

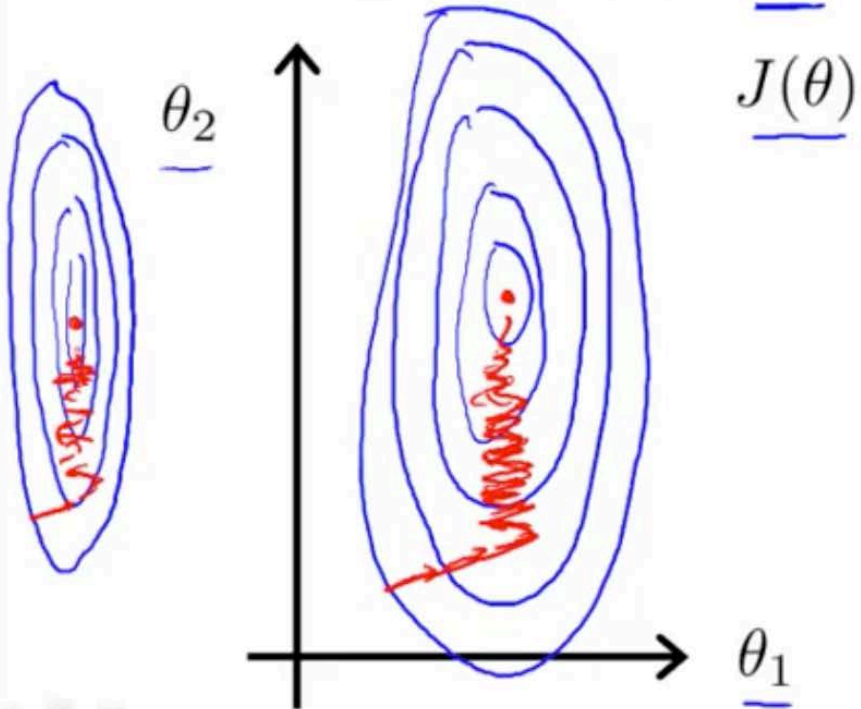
# 2 Linear Regression: Feature Scaling

## Feature Scaling

Idea: Make sure features are on a similar scale.

E.g.  $x_1 = \text{size (0-2000 feet}^2\text{)}$  ←

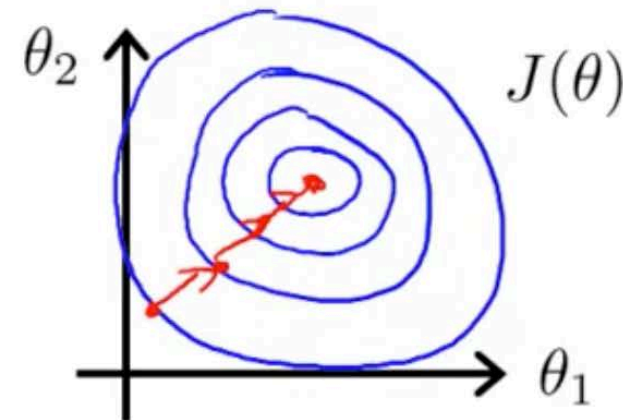
$x_2 = \text{number of bedrooms (1-5)}$  ←



$$\rightarrow x_1 = \frac{\text{size (feet}^2\text{)}}{2000} \quad \checkmark$$

$$\rightarrow x_2 = \frac{\text{number of bedrooms}}{5} \quad \checkmark$$

$$0 \leq x_1 \leq 1 \quad 0 \leq x_2 \leq 1$$



# 2 Linear Regression: Normal Equations

$$J(\theta) = \frac{1}{2}(X\theta - \vec{y})^T(X\theta - \vec{y})$$

$$\mathbf{X}: \mathbf{m} \times (\mathbf{n} + 1)$$

$$\vec{y}: \mathbf{m} \times 1$$

$$\theta: (\mathbf{n} + 1) \times 1$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2}(X\theta - \vec{y})^T(X\theta - \vec{y})$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T X^T X \theta - \theta^T X^T \vec{y} - \vec{y}^T X \theta + \vec{y}^T \vec{y})$$

a real number

$$= \frac{1}{2} \nabla_{\theta} \text{tr} (\theta^T X^T X \theta - \theta^T X^T \vec{y} - \vec{y}^T X \theta + \vec{y}^T \vec{y})$$

$$= \frac{1}{2} \nabla_{\theta} (\text{tr} \theta^T X^T X \theta - 2 \text{tr} \vec{y}^T X \theta)$$

$$\text{tr} A = \text{tr} A^T$$

$$= \frac{1}{2} (X^T X \theta + X^T X \theta - 2 X^T \vec{y})$$

$$= X^T X \theta - X^T \vec{y}$$

$$\nabla_{A^T} \text{tr} A B A^T C = B^T A^T C^T + B A^T C$$

$$\nabla_A \text{tr} A B = B^T$$

# 2 Linear Regression: Normal Equations

- Normal Equations

$$X^T X \theta = X^T \vec{y}$$

$$\theta = (X^T X)^{-1} X^T \vec{y}.$$

What if  $X^T X$  is *non-invertible*?



# 2 Linear Regression: Normal Equations

What if  $X^T X$  is non-invertible?

- Redundant features (linearly dependent).

E.g.  $x_1$  = size in feet<sup>2</sup>

$x_2$  = size in m<sup>2</sup>

$$1\text{m} = 3.28 \text{ feet}$$

$$\underline{x_1 = (3.28)^2 x_2}$$

$$m = 10$$

$$n = 100$$

- Too many features (e.g.  $m \leq n$ ).

$$\Theta \in \underline{\mathbb{R}^{101}}$$

- Delete some features, or use regularization.

# 2 Gradient Descent Vs. Normal Equations

$m$  training examples,  $n$  features.

## Gradient Descent

- • Need to choose  $\alpha$ .
- • Needs many iterations.
- Works well even when  $n$  is large.



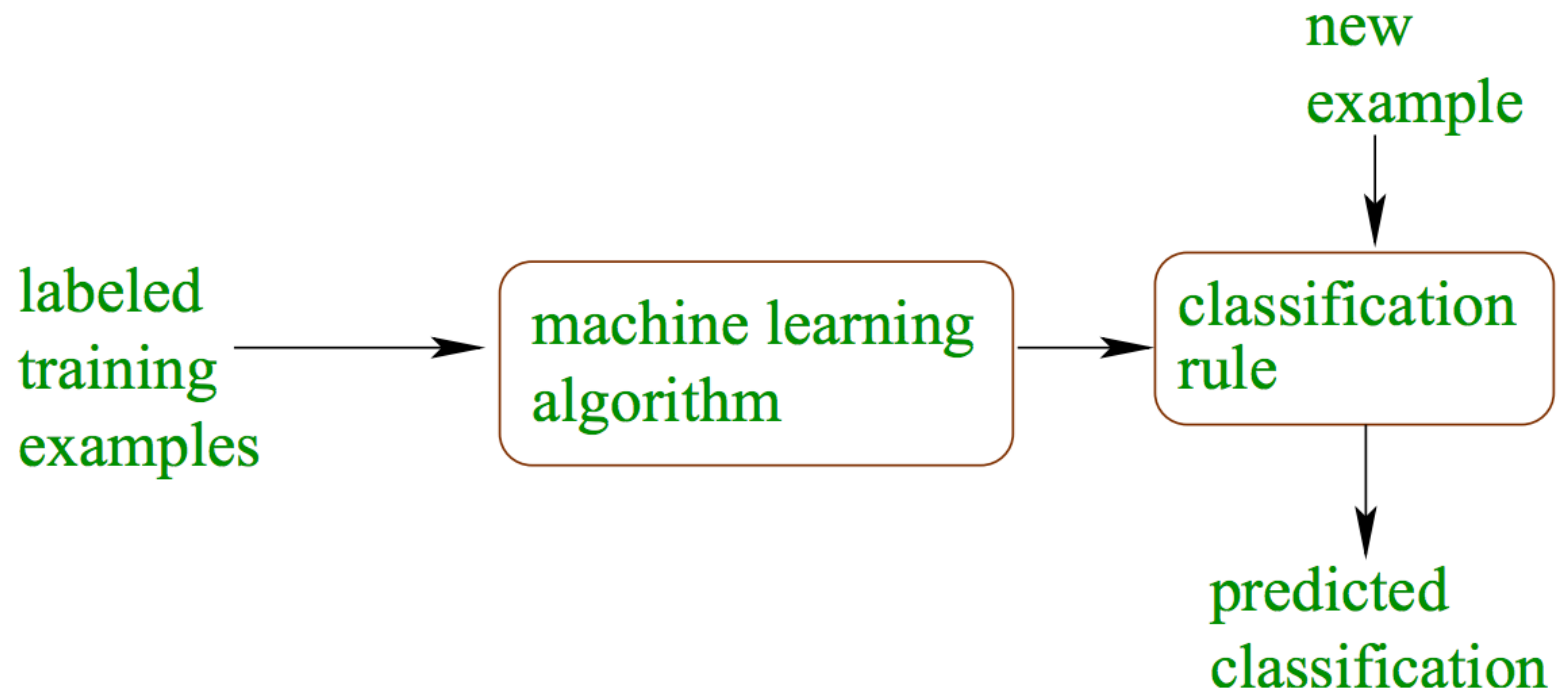
## Normal Equation

- • No need to choose  $\alpha$ .
- • Don't need to iterate.
- Need to compute  $(X^T X)^{-1}$   $\frac{n \times n}{\quad}$   $\underline{O(n^3)}$
- Slow if  $n$  is very large.

$$n = 100$$

# 3 Classification: Overview

- Classification: classifies examples into given set of categories



# 3 Classification: Application

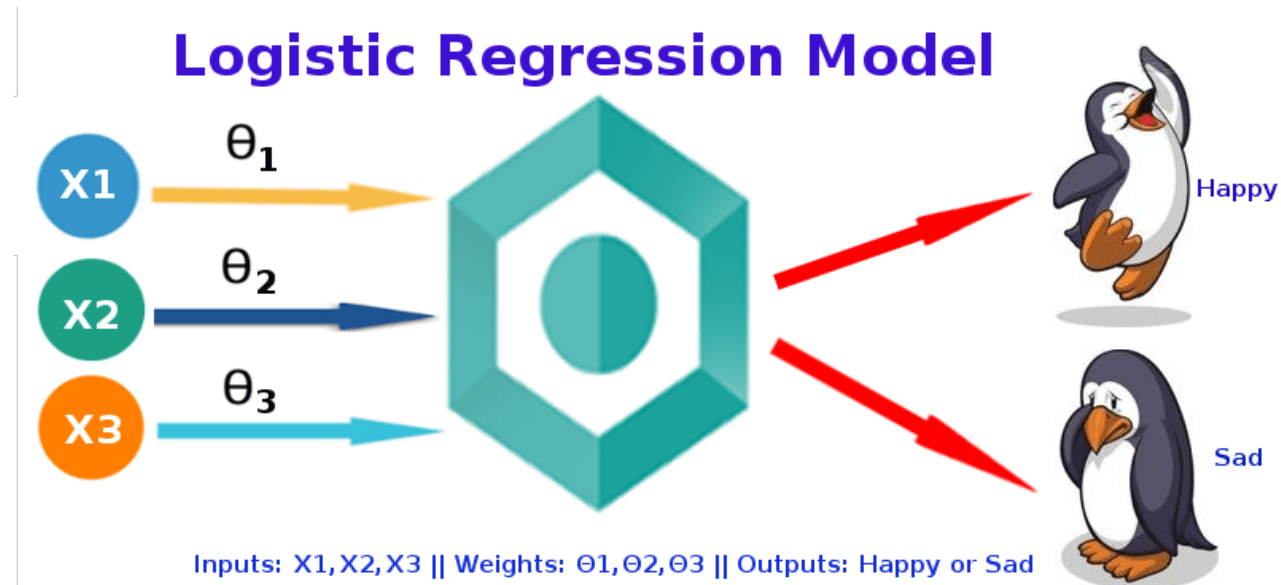
- Image/Document Classification
- Spam Email Filtering
- Optical Character Recognition
- Machine Vision (e.g., face detection)
- Natural Language Processing (e.g., spoken language understanding)
- Fraud Detection (e.g., if a transaction is fraudulent)
- Market Segmentation (e.g.: predict if customer will respond to promotion)
- Bioinformatics (e.g., classify proteins according to their function)

# 3 Classification: Algorithms

- Logistic Regression
- Support Vector Machine (SVM)
- Naïve Bayes Classifier
- K-Nearest Neighbor Classifier
- Decision Tree, Random Forest
- (Deep) Neural Network

# 4 Logistic Regression: Overview

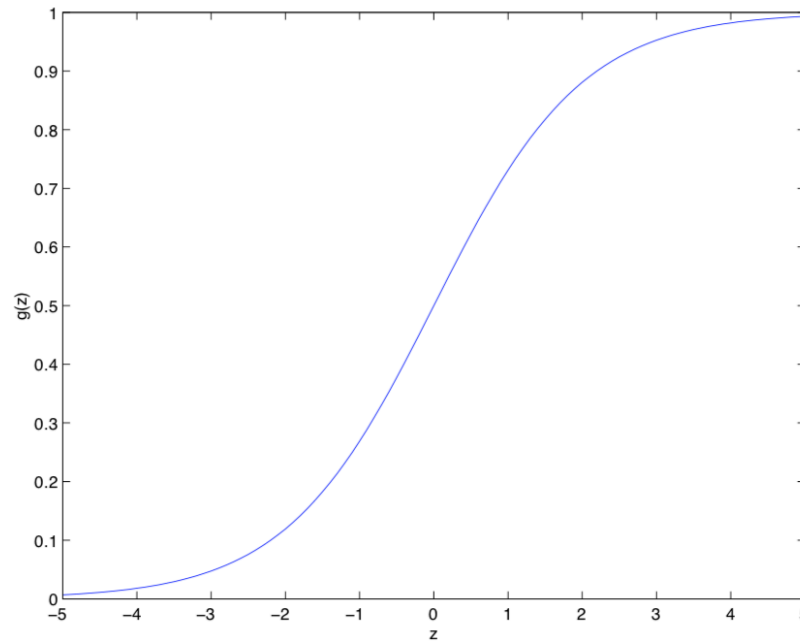
- Logistic Regression was used in the biological sciences in early twentieth century. Logistic Regression is used when the dependent variable(target) is categorical. For example:
  - To predict whether an email is spam (1) or (0)
  - Whether the tumor is malignant (1) or not (0)



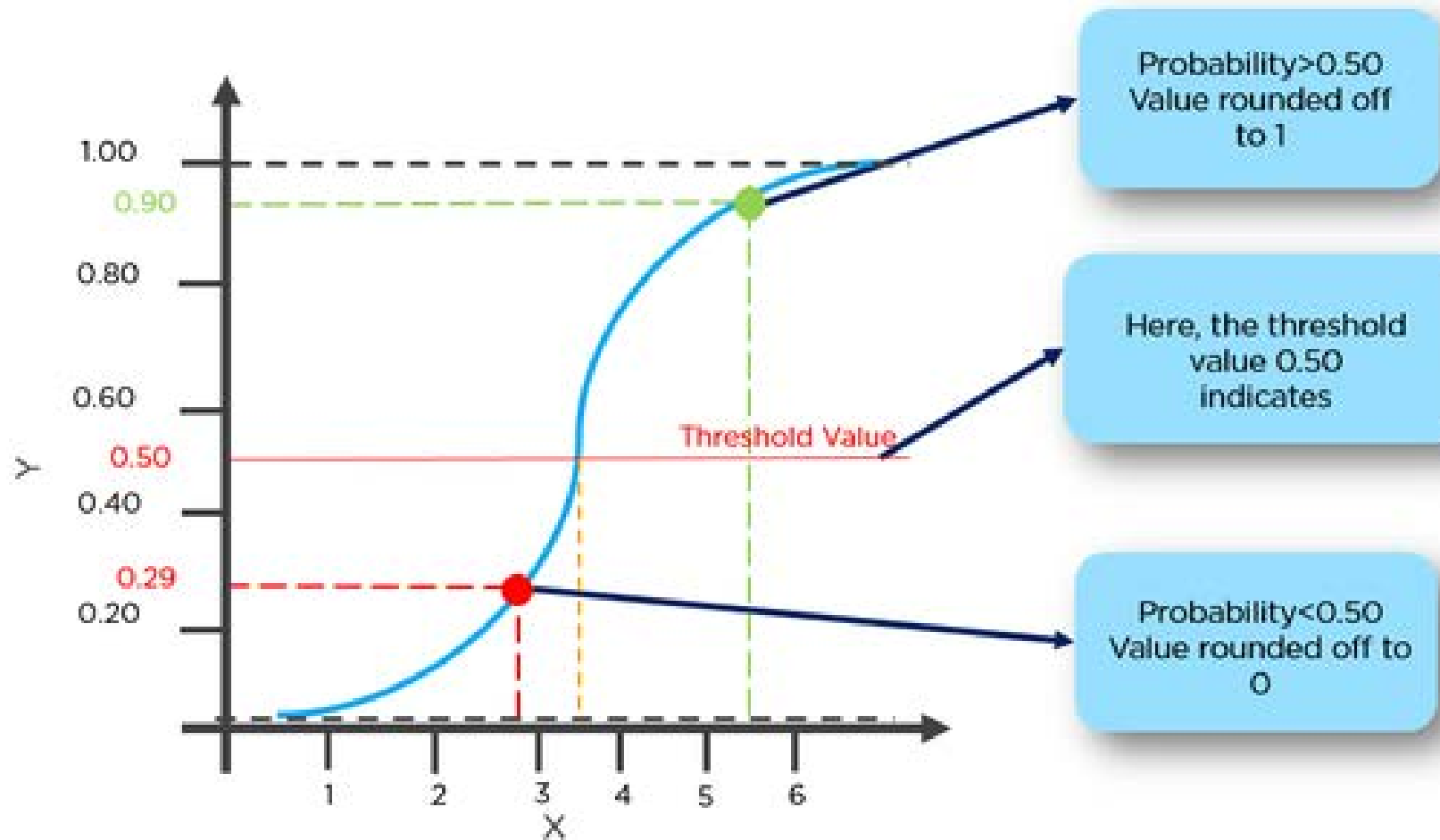
# 4 Logistic Regression: : Hypothesis

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

$g(z) = \frac{1}{1 + e^{-z}}$  is called the logistic function or sigmoid function



# 4 Logistic regression: Hypothesis





# 4 Logistic Regression: Sigmoid Function

A useful property of the derivative of the sigmoid function:

$$\begin{aligned}g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\&= \frac{1}{(1 + e^{-z})^2} (e^{-z}) \\&= \frac{1}{(1 + e^{-z})} \cdot \left(1 - \frac{1}{(1 + e^{-z})}\right) \\&= g(z)(1 - g(z)).\end{aligned}$$

# 4 Logistic Regression: How to fit $\theta$ ?

Let's assume that:

$$P(y = 1 \mid x; \theta) = h_{\theta}(x)$$

$$P(y = 0 \mid x; \theta) = 1 - h_{\theta}(x)$$

This can be written as:

$$p(y \mid x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

# 4 Logistic Regression: Maximize the Likelihood

- Assuming that the  $m$  training examples were generated independently, we can then write down the likelihood of the parameters as:

$$\begin{aligned} L(\theta) &= p(\vec{y} \mid X; \theta) \\ &= \prod_{i=1}^m p(y^{(i)} \mid x^{(i)}; \theta) \\ &= \prod_{i=1}^m (h_{\theta}(x^{(i)}))^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}} \end{aligned}$$

- It will be easier to maximize the log likelihood:

$$\begin{aligned} \ell(\theta) &= \log L(\theta) \\ &= \sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)})) \end{aligned}$$

# 4 Logistic Regression: Maximize the likelihood

- Maximize the likelihood

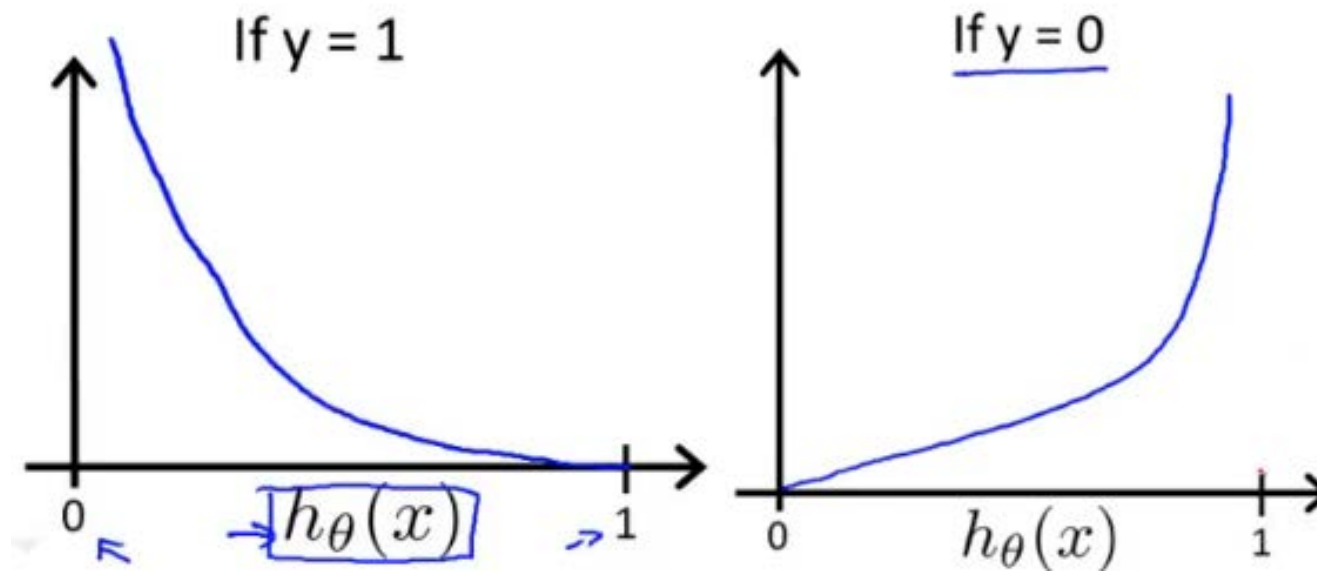
$$\begin{aligned}\frac{\partial}{\partial \theta_j} \ell(\theta) &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x)(1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= (y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x)) x_j \\ &= (y - h_\theta(x)) x_j\end{aligned}$$

- The stochastic gradient ascent rule

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

## 4 Logistic Regression: Minimize Cross-Entropy Loss

$$Cost(h_{\theta}(x^{(i)}), y^{(i)}) = \begin{cases} -\log(h_{\theta}(x^{(i)})), & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x^{(i)})), & \text{if } y = 0 \end{cases}$$



## 4 Logistic Regression: Minimize Cross-Entropy Loss

$$Cost(h_{\theta}(x^{(i)}), y^{(i)}) = \begin{cases} -\log(h_{\theta}(x^{(i)})), & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x^{(i)})), & \text{if } y = 0 \end{cases}$$

- **Cross-Entropy Loss**

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m Cost(h_{\theta}(x^{(i)}), y^{(i)}) = -\frac{1}{m} [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$$

# 4 Logistic regression: Decision Boundary

## Logistic regression

$$\rightarrow h_{\theta}(x) = g(\theta^T x) = p(y=1|x;\theta)$$

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

Suppose predict " $y = 1$ " if  $h_{\theta}(x) \geq 0.5$

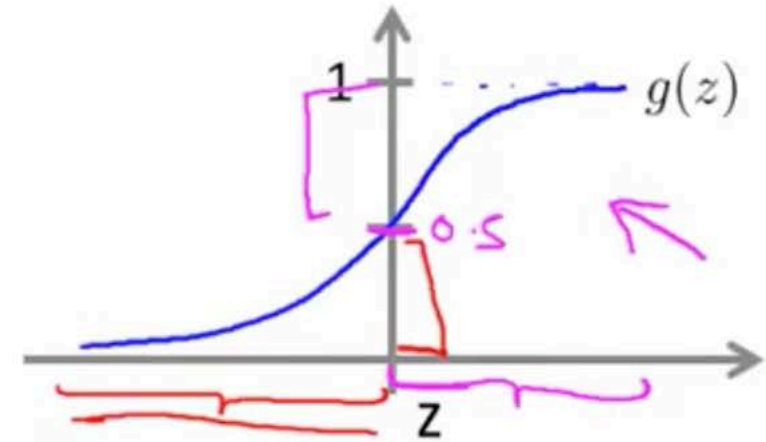
$$\theta^T x \geq 0$$

predict " $y = 0$ " if  $h_{\theta}(x) < 0.5$

$$h_{\theta}(x) = g(\theta^T x)$$

$$\theta^T x < 0$$

$$\underline{g(z) < 0.5}$$



$$g(z) \geq 0.5$$

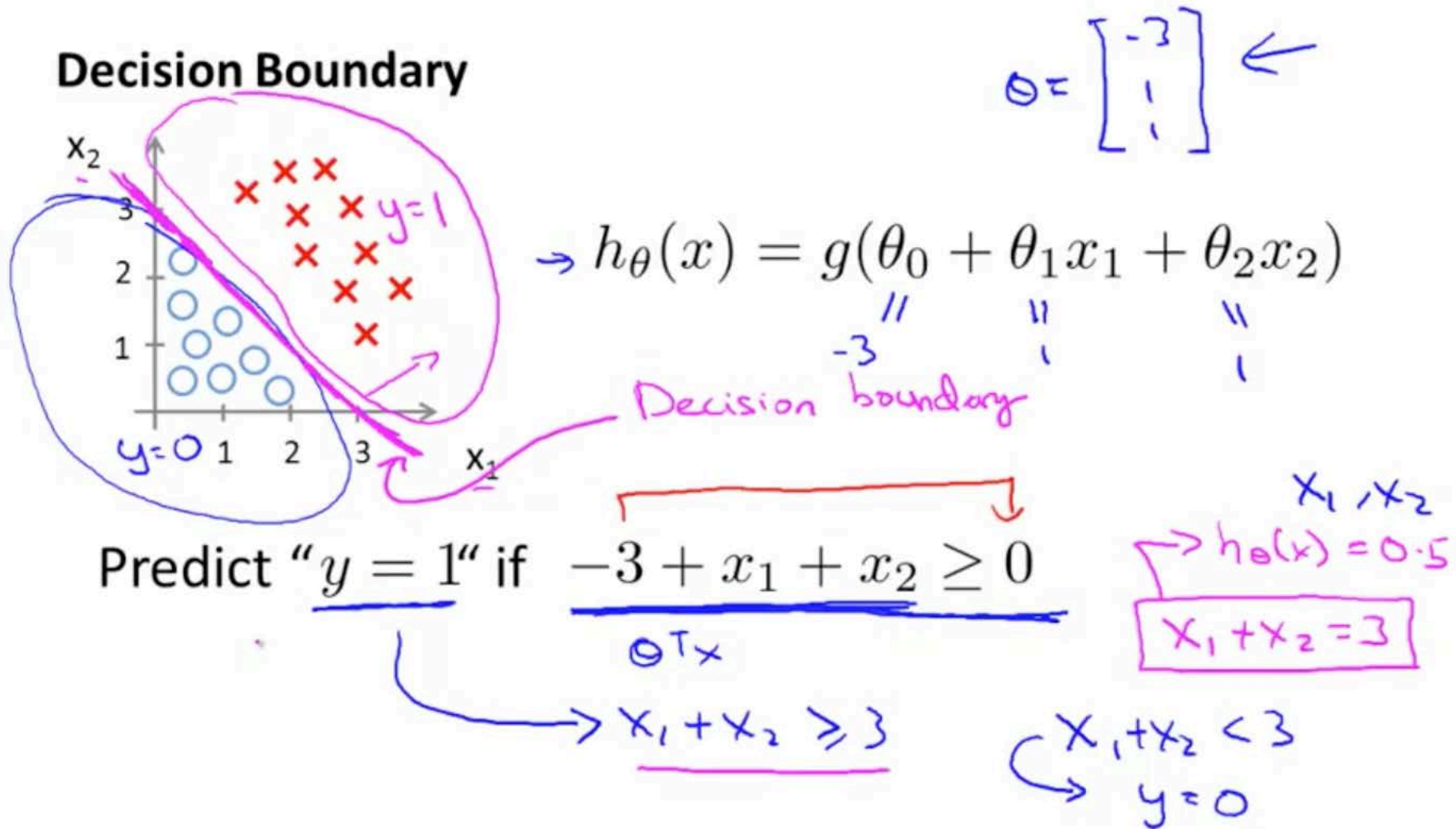
when  $z \geq 0$

$$h_{\theta}(x) = g(\theta^T x) \geq 0.5$$

whenever  $\theta^T x \geq 0$

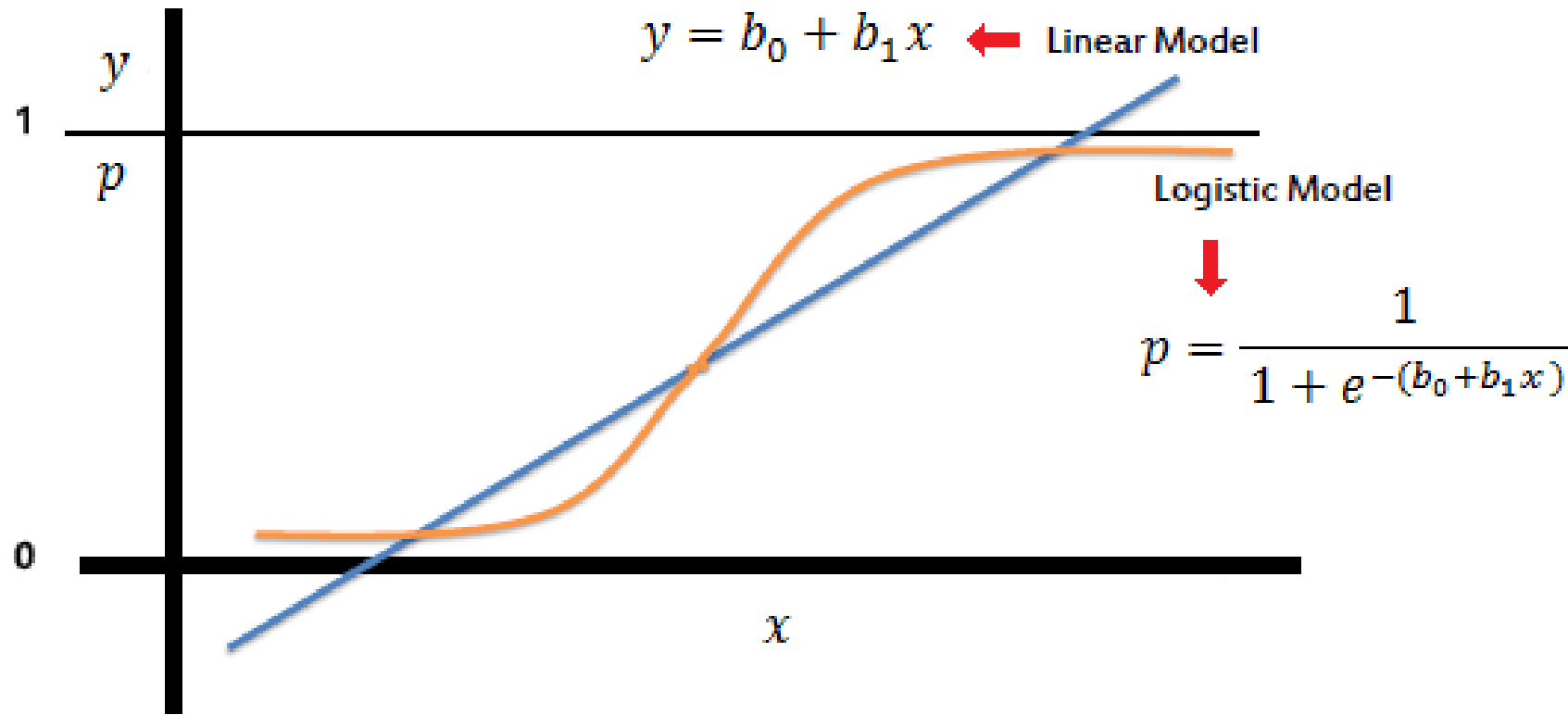
$\uparrow$   
 $z$

# 4 Logistic regression: Decision Boundary



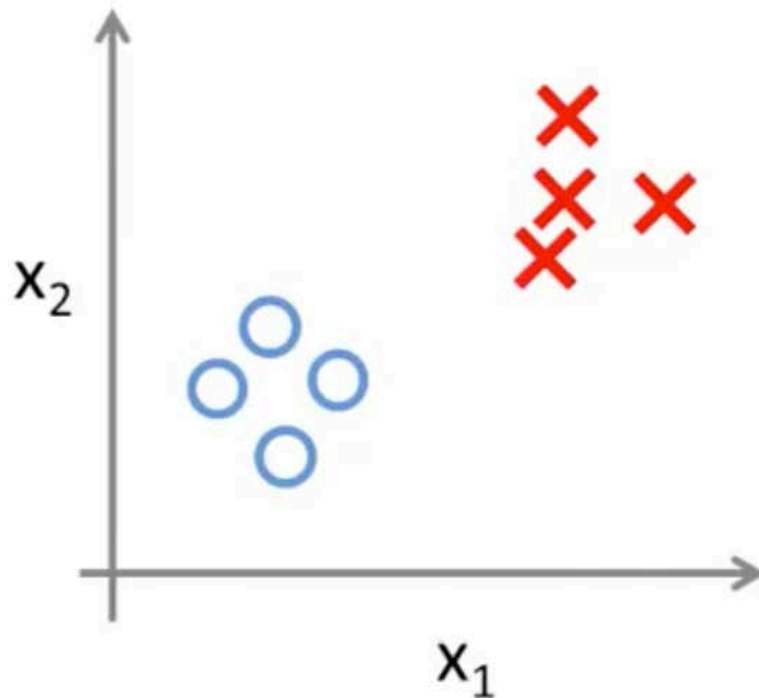


# 4 Logistic regression vs Linear regression

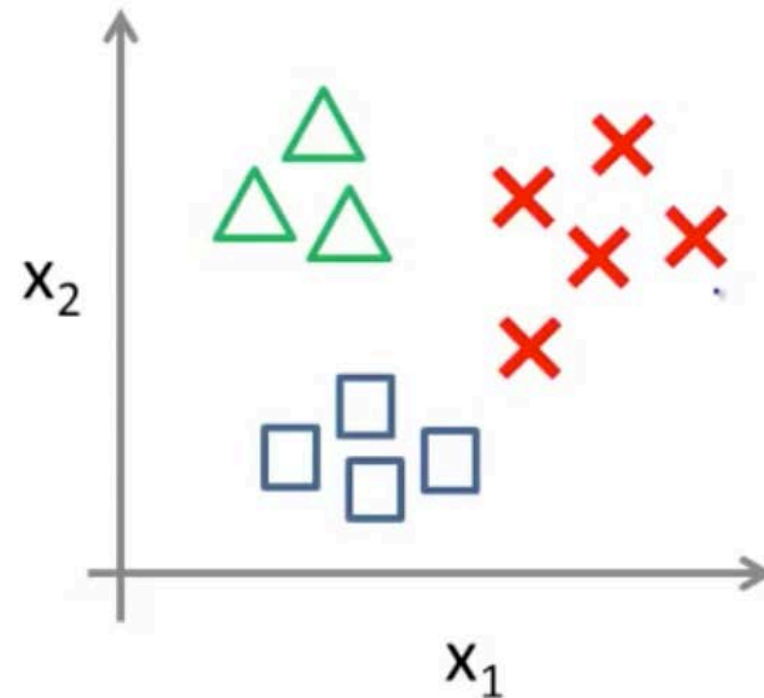


# 4 Logistic regression: Multiclass Classification

Binary classification:

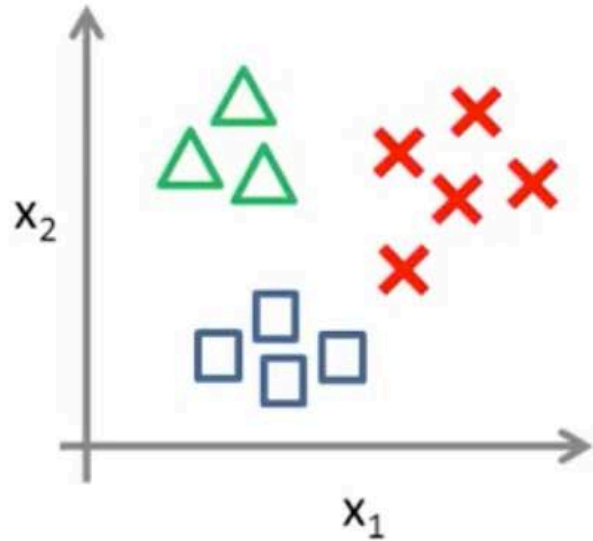


Multi-class classification:



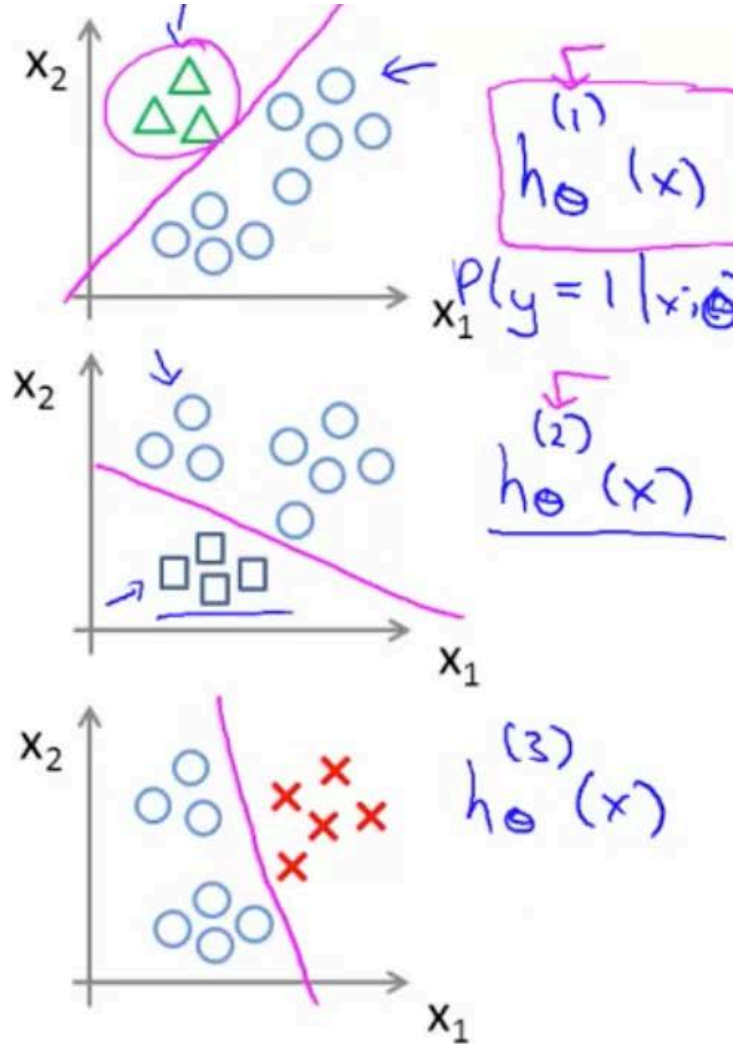
# 4 Logistic regression: Multiclass Classification

One-vs-all (one-vs-rest):



Class 1:   $\leftarrow$   
Class 2:   $\leftarrow$   
Class 3:   $\leftarrow$

$$h_{\theta}^{(i)}(x) = P(\underline{y = i} | x; \theta) \quad (\underline{i = 1, 2, 3})$$



# 4 Logistic regression: Multiclass Classification

## One-vs-all

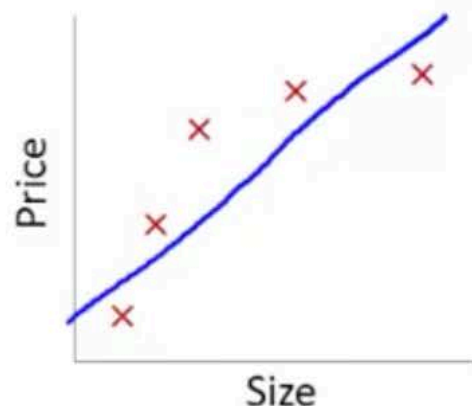
Train a logistic regression classifier  $h_{\theta}^{(i)}(x)$  for each class  $i$  to predict the probability that  $y = i$ .

On a new input  $x$ , to make a prediction, pick the class  $i$  that maximizes

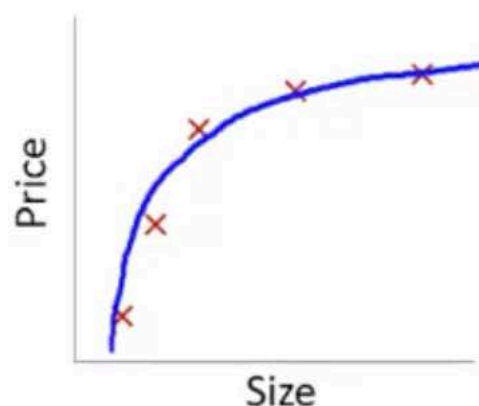
$$\max_i h_{\theta}^{(i)}(x)$$

# 5 Regularization: Overfitting

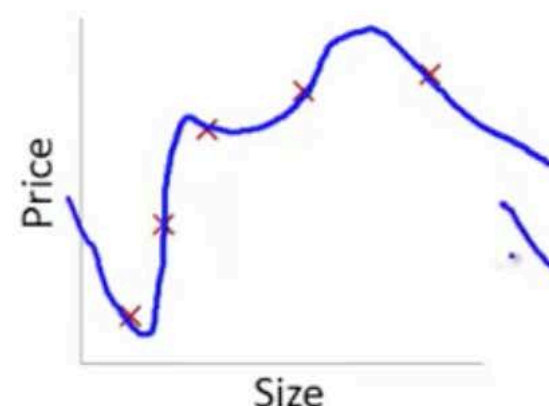
Example: Linear regression (housing prices)



$\rightarrow \theta_0 + \theta_1 x$   
"Underfit" "High bias"



$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$   
"Just right"

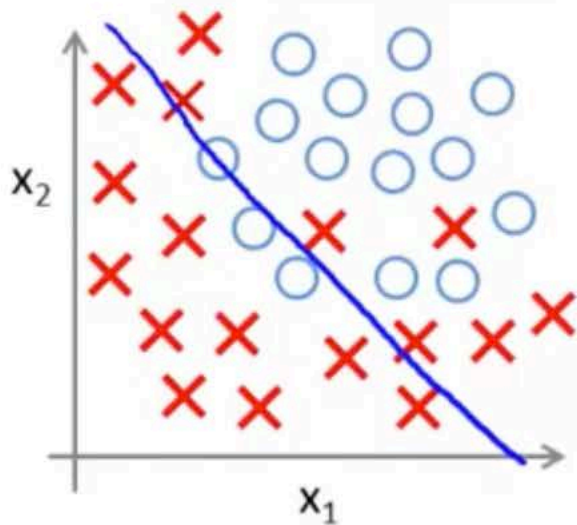


$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$   
"Overfit" "High variance"

**Overfitting:** If we have too many features, the learned hypothesis may fit the training set very well ( $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0$ ), but fail to generalize to new examples (predict prices on new examples).

# 5 Regularization: Overfitting

Example: Logistic regression

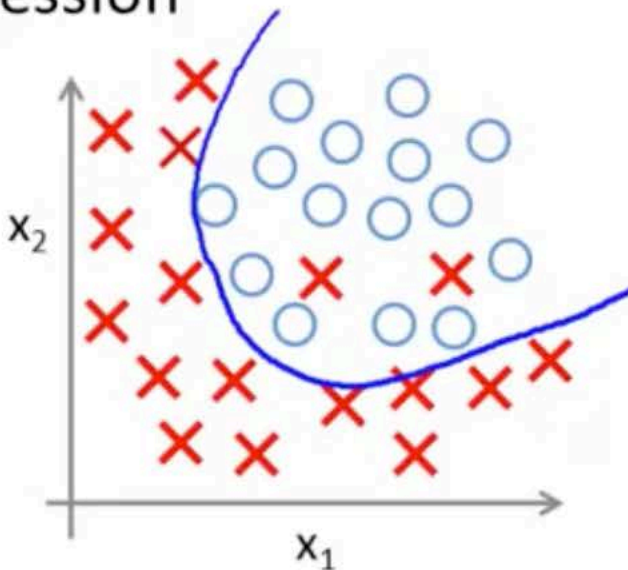


$$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

( $g$  = sigmoid function)

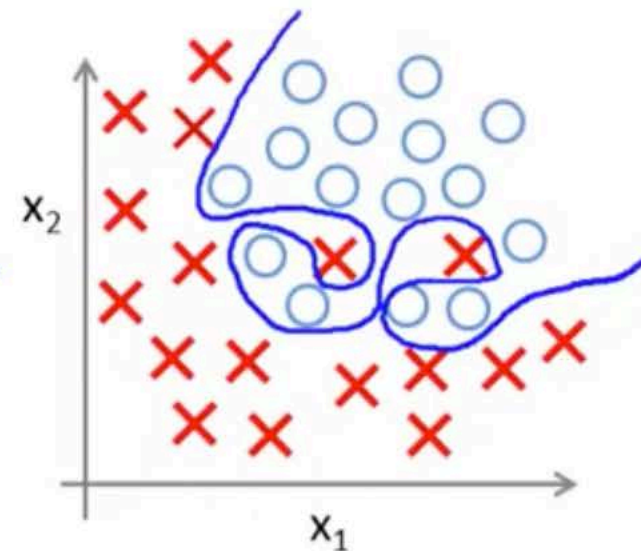
↩

"Underfit"



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

↗



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$

↖

# 5 Regularization: Overfitting

## Addressing overfitting:

### Options:

1. Reduce number of features.
  - — Manually select which features to keep.
  - — Model selection algorithm (later in course).
2. Regularization.
  - — Keep all the features, but reduce magnitude/values of parameters  $\theta_j$ .
  - Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .



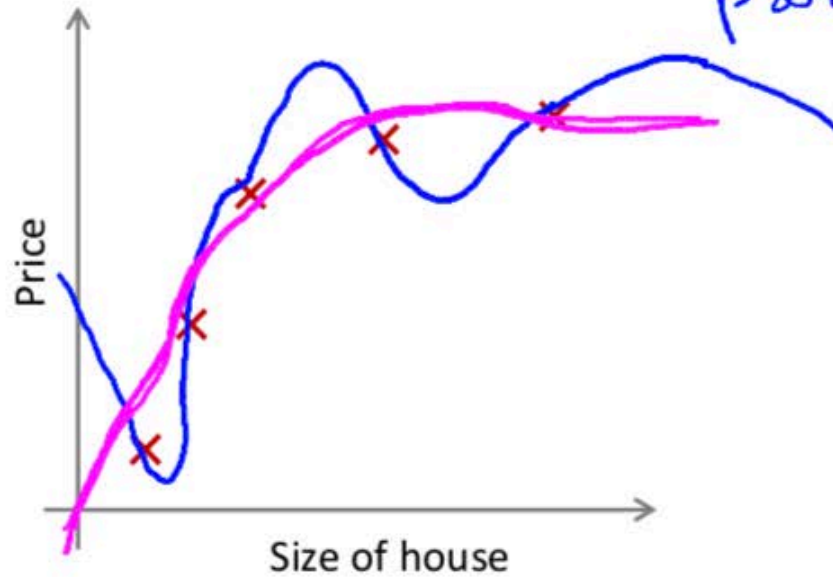
# 5 Regularization

Regularization.

$$\rightarrow J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$\min_{\theta} J(\theta)$

regularization parameter





# 5 Regularization: Linear Regression

In regularized linear regression, we choose  $\theta$  to minimize

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

What if  $\lambda$  is set to an extremely large value (perhaps far too large for our problem, say  $\lambda = 10^{10}$ )?

- Algorithm works fine; setting  $\lambda$  to be very large can't hurt it
- Algorithm fails to eliminate overfitting.
- Algorithm results in underfitting. (Fails to fit even training data well).
- Gradient descent will fail to converge.

# 5 Regularization: Linear Regression

Gradient descent

Repeat {

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$\theta_0$   
↑

$\theta_1, \theta_2, \dots, \theta_n$

$\frac{\partial}{\partial \theta_0} J(\theta)$

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j \right]$$

$(j = \cancel{x}, 1, 2, 3, \dots, n)$

}

$$\theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$\rightarrow J(\theta)$

$\theta_j^2$

$$1 - \alpha \frac{\lambda}{m} < 1$$

0.99

$\theta_j \times 0.99$

# 5 Regularization: Linear Regression

Normal equation

$$\underline{X} = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

$m \times (n+1)$

$$\underset{\uparrow}{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \mathbb{R}^m$$

$$\rightarrow \min_{\theta} \underline{J(\theta)}$$

$$\frac{\partial J(\theta)}{\partial \theta_j} \stackrel{\text{set}}{=} 0$$

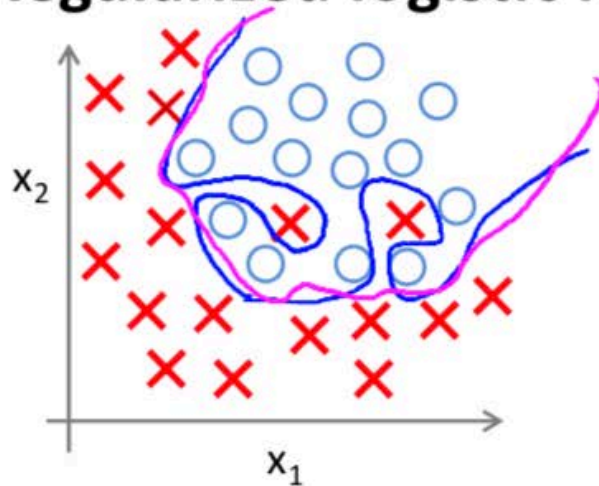
$$\Rightarrow \Theta = \left( X^T X + \lambda \underbrace{\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{(n+1) \times (n+1)} \right)^{-1} X^T y$$

$\in \mathbb{R}^n \quad n=2$

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

# 5 Regularization: Logistic Regression

Regularized logistic regression.



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

Handwritten annotations: A blue arrow points down to  $\theta_2 x_1^2$ , a blue arrow points to the entire expression inside the  $g(\dots)$  function, and a pink arrow points up to  $\theta_5 x_1^2 x_2^3$ .

Cost function:

$$\rightarrow J(\theta) = - \left[ \frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Handwritten annotations: A pink box highlights the parameters  $\theta_1, \theta_2, \dots, \theta_n$  in the regularization term.

# 5 Regularization: Logistic Regression

## Gradient descent

Repeat {

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{n} \theta_j \right] \leftarrow$$

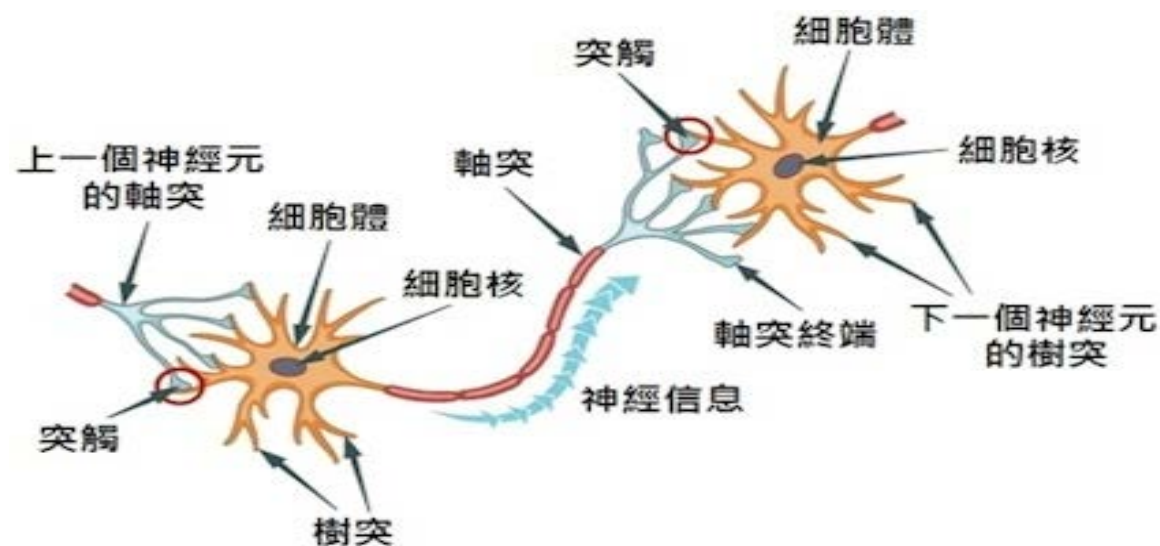
$(j = \text{red X}, 1, 2, 3, \dots, n)$   
 $\theta_1, \dots, \theta_n$

}

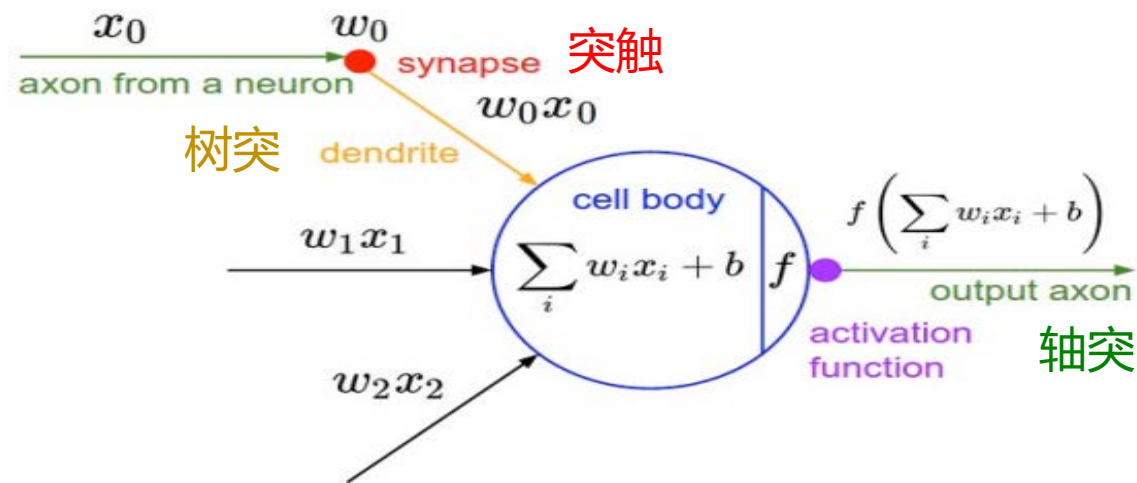
$$\frac{2}{2\theta_j} \underline{J(\theta)}$$

$$\underline{h_{\theta}(x)} = \frac{1}{1 + e^{-\theta^T x}}$$

# 6 Perceptron: Artificial Neuron



生物神经元的基本结构

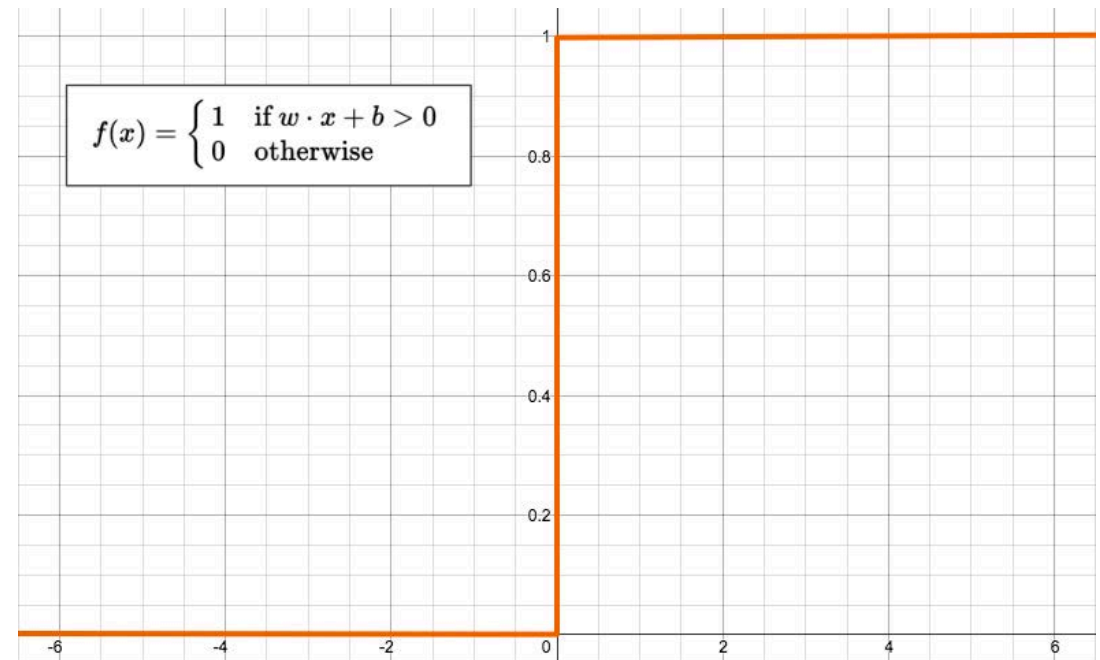
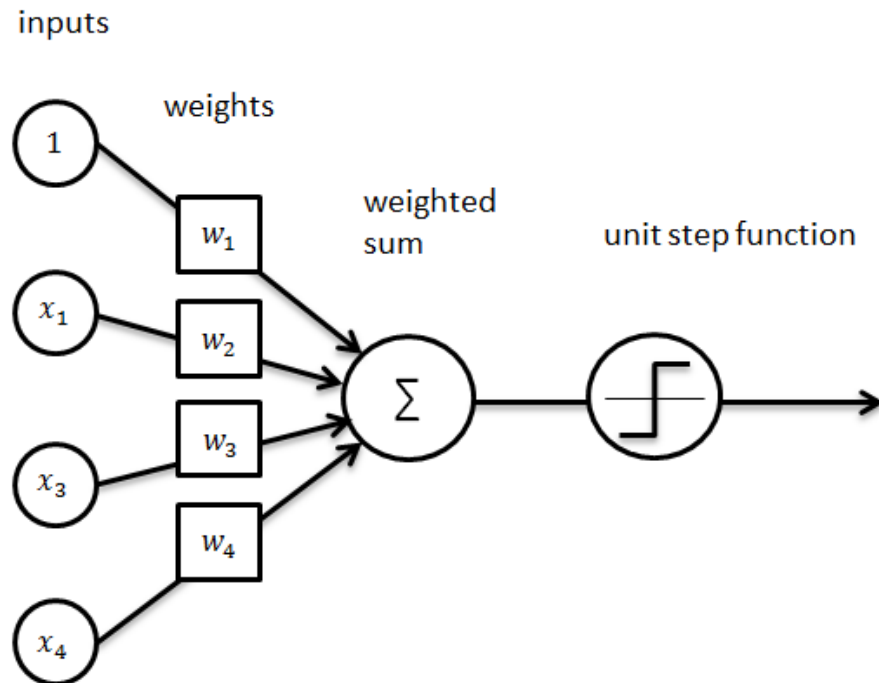


人工神经元的结构模型

# 6 Perceptron: Model

- **Perceptron**

- First compute a weighted sum of the inputs from other neurons.
- Then output a 1 if the weighted sum exceeds the threshold.





# 6 Perceptron: Minimizing Squared Errors

$$h_{\theta}(x) = g(\theta^T x) = g\left(\sum_{i=0}^n \theta_i x_i\right) \quad g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

- The squared error for a single training example with input  $x$  and true output  $y$  is:

$$E = \frac{1}{2} (y^{(i)} - h_{\theta}(x^{(i)}))^2$$

- Gradient Descent

undefined and omitted

$$\begin{aligned} \frac{\partial E}{\partial \theta_j} &= (y^{(i)} - h_{\theta}(x^{(i)})) \times (-1) \times \boxed{g'(in)} \times x_j^{(i)} \\ &= -(y^{(i)} - h_{\theta}(x^{(i)})) \times x_j^{(i)} \end{aligned}$$

- Update Rule

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$



# Homework (1)

- Task: Regression
- Dataset: Housing

(<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/regression.html#housing>)

- Algorithms: Linear Regression

## housing

- Source: [UCI](#) / Housing (Boston)
- # of data: 506
- # of features: 13
- Files:
  - [housing](#)
  - [housing\\_scale](#) (scaled to [-1,1])

```
24.0 1:0.00632 2:18.00 3:2.310 4:0 5:0.5380 6:6.5750 7:65.20 8:4.0900 9:1 10:296.0 11:15.30 12:396.90 13:4.98
21.6 1:0.02731 2:0.00 3:7.070 4:0 5:0.4690 6:6.4210 7:78.90 8:4.9671 9:2 10:242.0 11:17.80 12:396.90 13:9.14
34.7 1:0.02729 2:0.00 3:7.070 4:0 5:0.4690 6:7.1850 7:61.10 8:4.9671 9:2 10:242.0 11:17.80 12:392.83 13:4.03
33.4 1:0.03237 2:0.00 3:2.180 4:0 5:0.4580 6:6.9980 7:45.80 8:6.0622 9:3 10:222.0 11:18.70 12:394.63 13:2.94
36.2 1:0.06905 2:0.00 3:2.180 4:0 5:0.4580 6:7.1470 7:54.20 8:6.0622 9:3 10:222.0 11:18.70 12:396.90 13:5.33
28.7 1:0.02985 2:0.00 3:2.180 4:0 5:0.4580 6:6.4300 7:58.70 8:6.0622 9:3 10:222.0 11:18.70 12:394.12 13:5.21
```

# Homework (2)

- Task: Classification
- Dataset: Mnist (<http://yann.lecun.com/exdb/mnist/>). The MNIST database of handwritten digits, available from this page, has a training set of 60,000 examples, and a test set of 10,000 examples.
- Algorithms: Logistic Regression

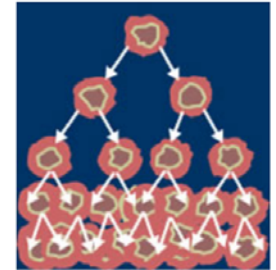
# Homework (3)

- Task: Classification
- Dataset: Breast Cancer Wisconsin (Diagnostic) Data Set  
(<https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29>)
- Algorithms: Perceptron

## Breast Cancer Wisconsin (Diagnostic) Data Set

Download: [Data Folder](#), [Data Set Description](#)

**Abstract:** Diagnostic Wisconsin Breast Cancer Database



Data Set Characteristics:	Multivariate	Number of Instances:	569	Area:	Life
Attribute Characteristics:	Real	Number of Attributes:	32	Date Donated	1995-11-01
Associated Tasks:	Classification	Missing Values?	No	Number of Web Hits:	692187

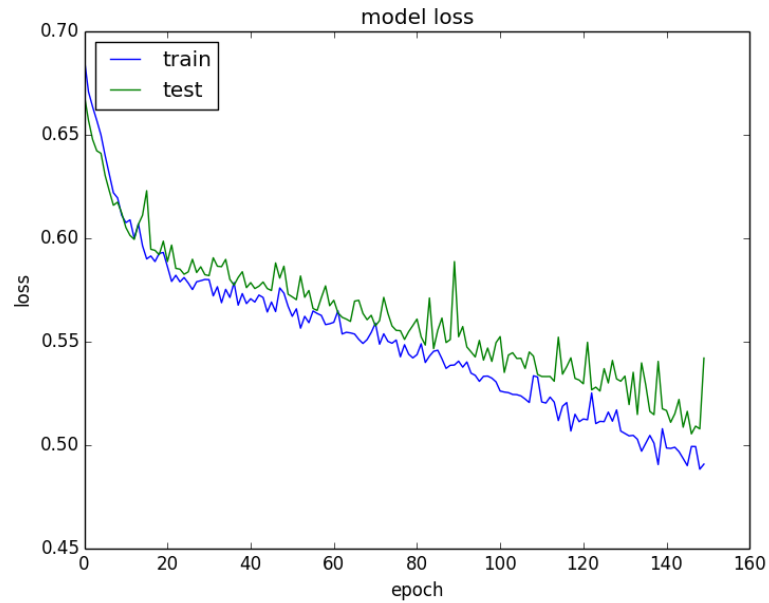
# Tips for Homework

- **Evaluation Metrics:**

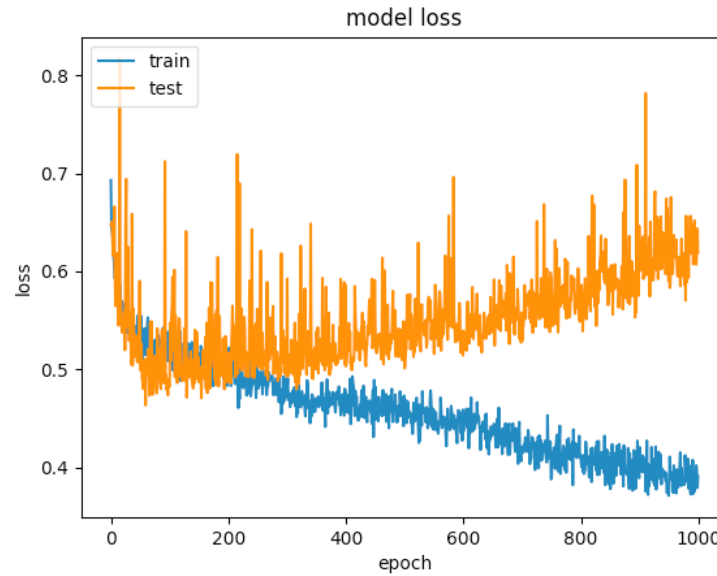
- ✓ Regression: Root Mean Squared Error
- ✓ Classification: Accuracy

# Tips for Homework

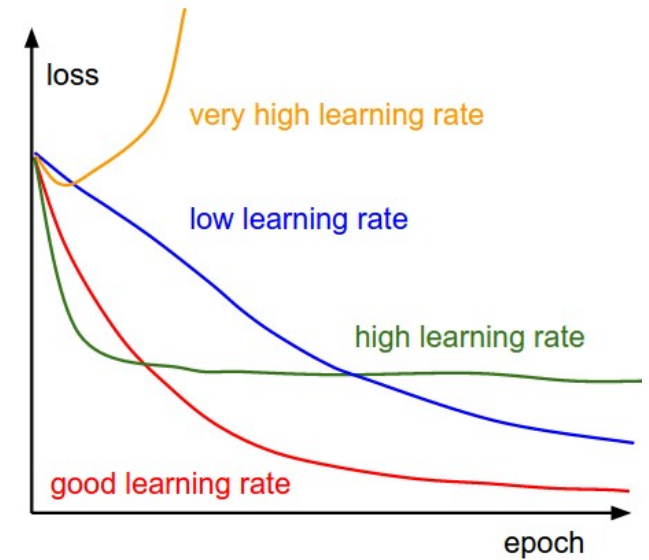
- Plot Loss Function Figure



**Normal**



**Overfitting**

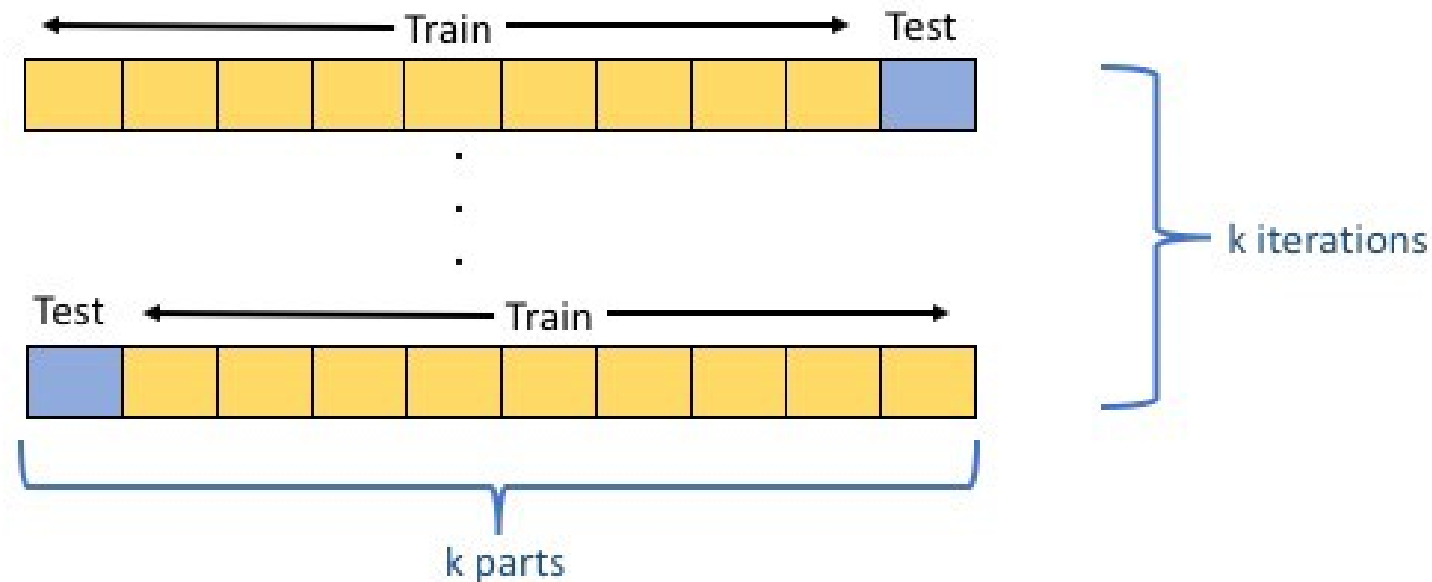


**Learning Rate**

# Tips for Homework

- **K-Folds Cross Validation**

1. Divide the sample data into  $k$  parts.
2. Use  $k-1$  of the parts for training, and 1 for testing.
3. Repeat the procedure  $k$  times, rotating the test set.
4. Determine an expected performance metric (mean square error, misclassification error rate, confidence interval, or other appropriate metric) based on the results across the iterations



# More Materials

- 1. Deep Learning, deeplearning.ai (Andrew Ng et al.) @ Coursera 网易云课堂
- 2. Andrew Ng. Machine Learning. Coursera, 网易公开课
- 3. 机器学习基石, Hsuan-Tien Lin, 林軒田, 台湾大学 @ Coursera
- 4. 机器学习技法, Hsuan-Tien Lin, 林軒田, 台湾大学 @ Coursera
- 5. CS231n: Convolutional Neural Networks for Visual Recognition Feifei Li et al., University of Stanford
- 6. Ian Goodfellow, Yoshua Bengio, and Aaron Courville, Deep Learning, Book in preparation for MIT Press, 2016, <http://www.deeplearningbook.org/>
- 7. 台大, 李宏毅, [http://speech.ee.ntu.edu.tw/~tlkagk/courses\\_ML16.html](http://speech.ee.ntu.edu.tw/~tlkagk/courses_ML16.html)
- MIT公开课, Gilbert Strang, 《线性代数》

# Programming Skills

- **Python**
- **Ubuntu**
- **Keras**
- **Pytorch**
- **Tensorflow**





暑期培训

END

# CS229 Lecture 2

Machine Learning

顾晓玲

[guxl@hdu.edu.cn](mailto:guxl@hdu.edu.cn)