

# ML note of Andrew Ng's course

---

## Table of Contents

ML note of Andrew Ng's course	.....
W1	.....
1. Fundamentals	.....
2. Model Representation	.....
3. Cost Function	.....
a. definition	.....
b. intuition	.....
4. Parameter Learning (gradient descent)	.....
a. introduction	.....
b. GD for linear regression	.....
5. Linear Algebra Fundamentals(Review)	.....
a. matrices and vectors(trivial)	.....
b. addition and scalar multiplication(trivial)	.....
c. matrix matrix multiplication	.....
d. inverse and transpose(trivial)	.....
W2	.....
1. Vectorization(intro)	.....
2. Multivariate Linear Regression	.....
a. hypothesis function of MLR	.....
b. GD for multiple variables	.....
c. practical tricks for GD of MLR	.....
d. Features and Polynomial Regression	.....
e. normal equation(正规方程)	.....
W3	.....
1. Classification Introduction	.....
2. Logistic Regression	.....
a. sigmoid function(logistic function)	.....
b. decision boundary(决策边界)	.....
c. logistic regression cost function & gradient descent	.....
d. other optimization algorithms	.....
3. Multi-class Classification	.....
a. one-vs-all classification	.....

- 4. Overfitting
  - a. introduction
  - b. regularization

W4

- 1. Neural Network: Intro & Representation
  - a. motivation
  - b. model representation
  - c. intuition & application

W5

- 1. Neural Network: Learning
  - a. cost function
  - b. back propagation: compute the derivatives of  $J(\Theta)$
  - c. BP in practice
  - d. summary: how to train a neural network?

W6

- 1. Evaluating a Learning Algorithm
  - a. test set error
  - b. model selection & Train/Validation/Test sets
- 2. Bias(偏差) & Variance(方差)
  - a. detection
  - b. regularization and bias/variance
  - c. learning curves
  - d. review: what to do next?
- 3. ML System Design
  - a. example: build a spam classifier
  - b. handling skewed data(偏斜数据)
  - c. collecting data

W7

- 1. SVM(Support Vector Machine)
  - a. large margin classification(大间距分类)
  - b. kernels(adapting SVM to non-linear classification problems)
  - c. SVM in practice

W8

- 1. Unsupervised learning
  - a. K-Means algorithm
- 2. Dimensionality Reduction

a. motivation	-----
b. PCA(Principle Component Analysis)	-----
c. applying PCA	-----
W9	-----
1. Anomaly Detection	-----
a. density estimation	-----
b. building an anomaly detection system	-----
c. multivariant Gaussian distribution	-----
2. Recommender Systems	-----
a. predicting movie ratings(content-based)	-----
b. collaborative filtering(协同过滤)	-----
c. low rank matrix factorization	-----
W10	-----
1. Stochastic Gradient Descent(large scale GD)	-----
a. batch GD	-----
b. SGD	-----
c. mini-batch GD	-----
2. Advanced Topics	-----
a. online learning	-----
b. map reduce & data parallelism	-----
W11	-----
Case Study: Photo OCR	-----
a. pipeline	-----
b. sliding window technique	-----
c. getting data for the system	-----
d. ceiling analysis(上限分析)	-----
End of Course: Thank you from Andrew Ng	-----

## W1

---

### 1. Fundamentals

- Supervised Learning

Regression (prediction)

Classification

...

- **Unsupervised Learning**

Clustering (grouping)

Non-clustering

...

## 2. Model Representation

### *Glossary #1*

$m$ : the size of training set

$x^{(i)}$  : input of the training set

$y^{(i)}$ : output of the training set

**Linear Regression Problem:**



where  $h_{\theta}(x) = \theta_0 + \theta_1 x$

## 3. Cost Function

### a. definition

Training dataset input  $x$  and output  $y$  known, we need to find  $\theta_0$  and  $\theta_1$  s.t.

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m \left( h_{\theta} \left( x^{(i)} \right) - y^{(i)} \right)^2$$

get the minimum value

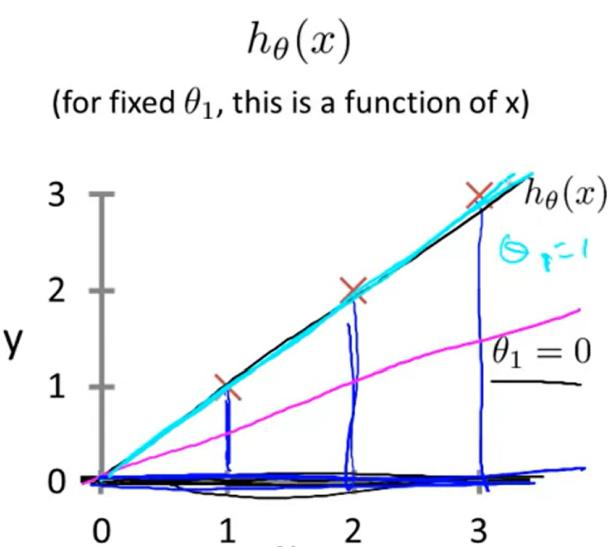
formally, we are going to find

$$\underset{\theta_0, \theta_1}{\operatorname{argmin}} J(\theta_0, \theta_1)$$

where  $J(\theta_0, \theta_1)$  is called **cost function** or **square error cost function**

### b. intuition

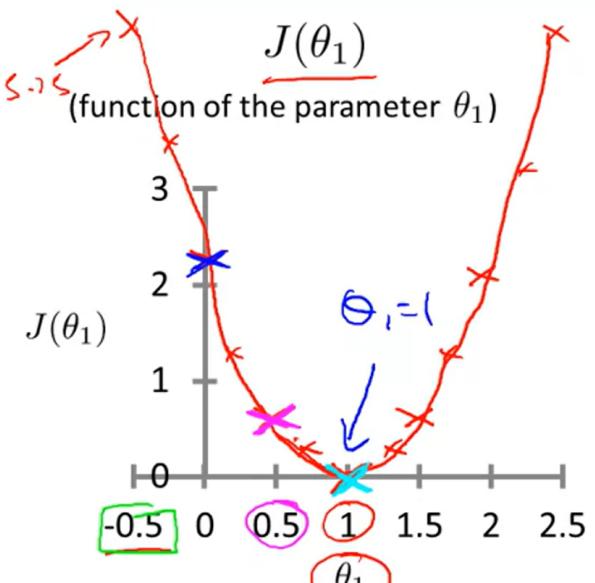
1. simplified problem:



$$J(\theta) = \frac{1}{2m} \left( (1^2 + 2^2 + 3^2) \right)$$

$$= \frac{1}{6} [1 + 4 + 9] = 5$$

visualizations involving both theta zero and theta one. That is without setting theta



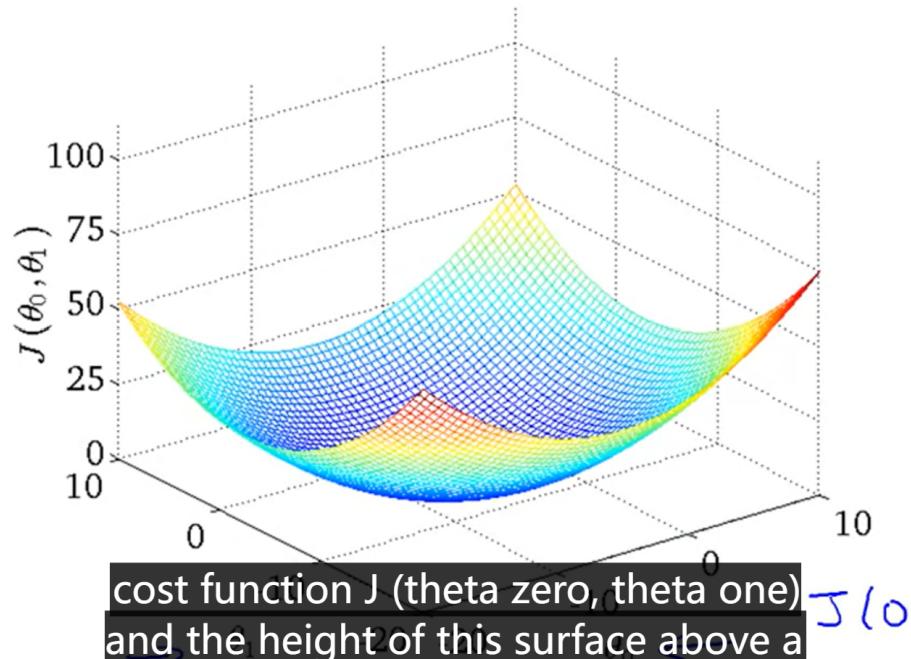
$$\theta_0 = 0$$

dataset: (1,1), (2,2), (3,3)

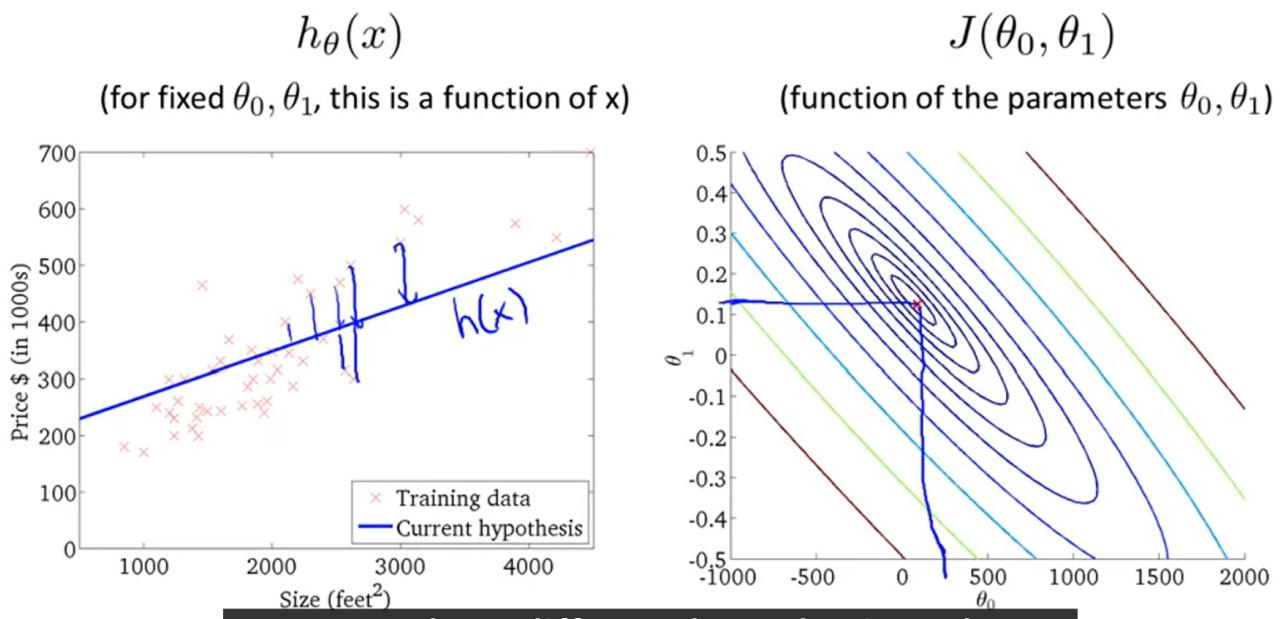
the plot of function  $J$  is on the right. To solve this question, we have to find  $\theta_1$  s.t.  $\frac{d J}{d \theta_1} = 0$

and evidently  $\theta_1 = 1$

2. original problem:



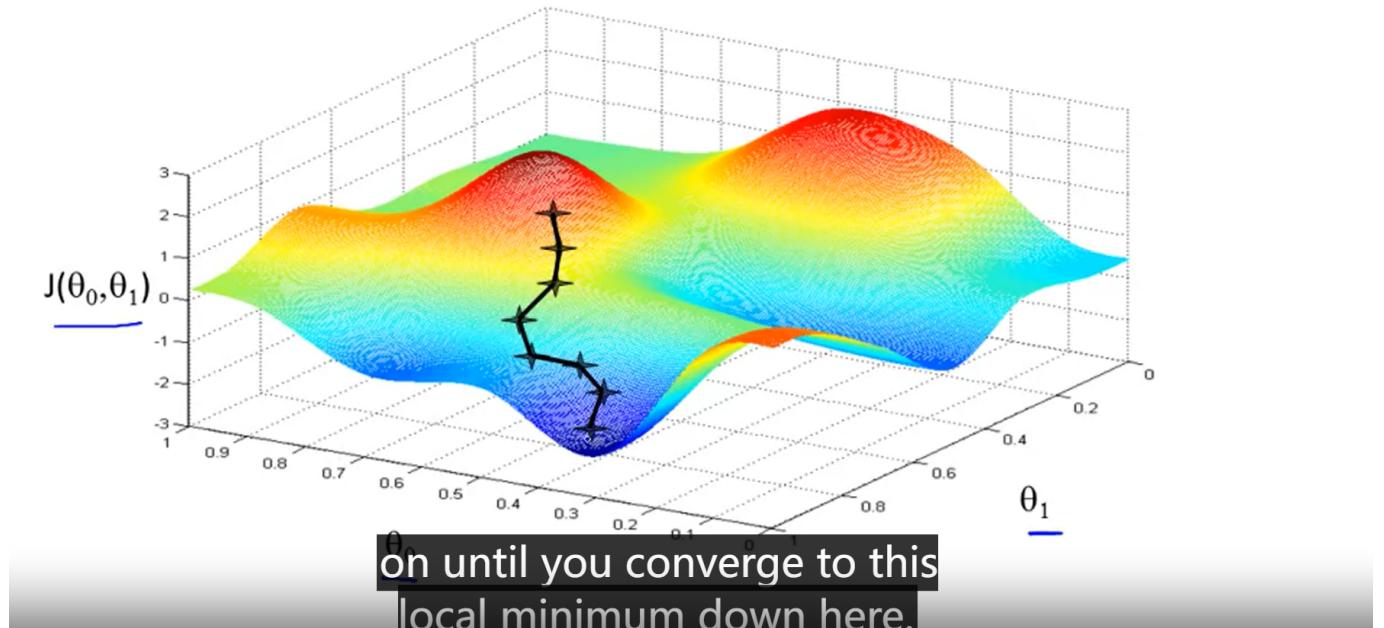
Andrew Ng



3-D figure and it is projected to  $\theta_0 O \theta_1$  as a **contour plot**. At the current point shown in the figure above, we get the minimum of  $J(\theta_0, \theta_1)$ , thus find the best  $h_\theta(x)$  as the hypothesis function.

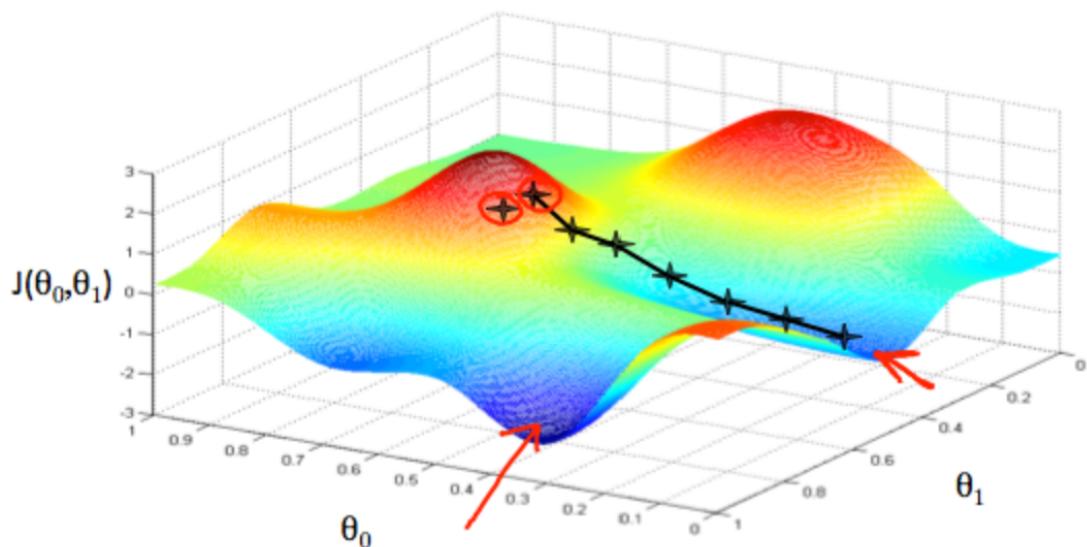
#### 4. Parameter Learning (gradient descent)

## a. introduction



go under the hill on the **steepest way**(partial derivative)

property:



Local optima? Global optima?

## Gradient descent algorithm

repeat until convergence { }  $\theta_0, \theta_1$

$\rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$  (for  $j = 0$  and  $j = 1$ )

learning rate

Simultaneously update  $\theta_0$  and  $\theta_1$

→  $a := b$  ✓  
→  $a := a + 1$  ✗

### Correct: Simultaneous update

→  $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
 →  $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   
 →  $\theta_0 := \text{temp0}$   
 →  $\theta_1 := \text{temp1}$

this is some other algorithm with different properties.

### Incorrect:

→  $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
 →  $\theta_0 := \text{temp0}$   
 $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$

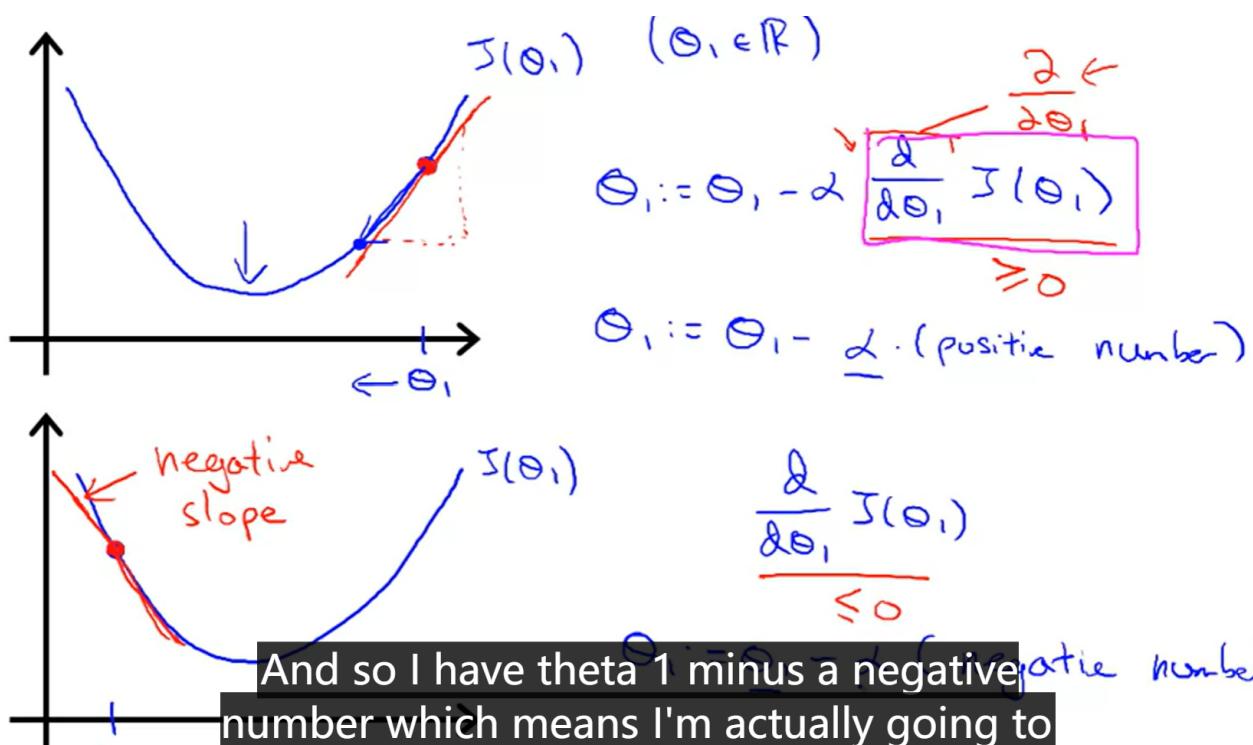
the mathematical implementation of GD

where:

$\alpha$  learning rate: how big a step will be

WARNING: FOCUS ON THE CORRECT UPDATE VERSION

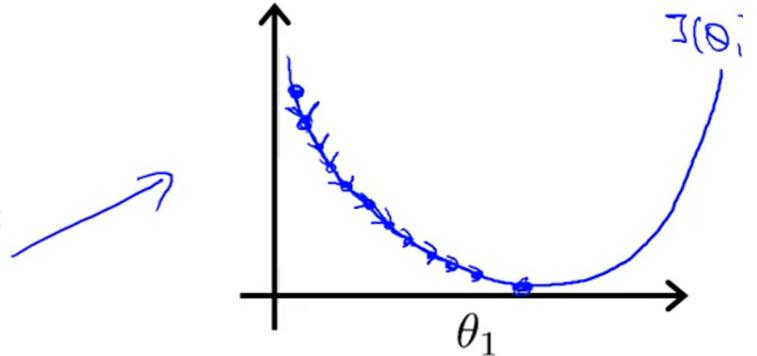
what the derivative does:



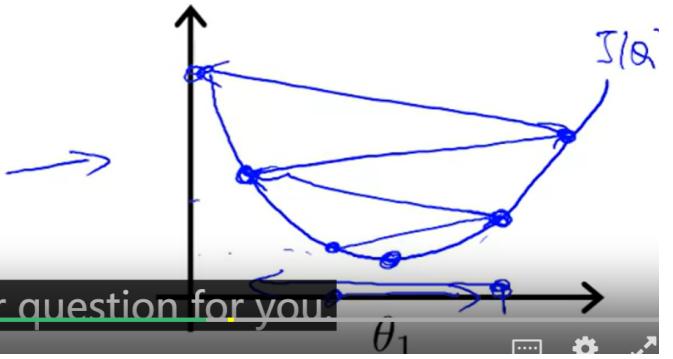
what  $\alpha$  does:

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

If  $\alpha$  is too small, gradient descent can be slow.



If  $\alpha$  is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.



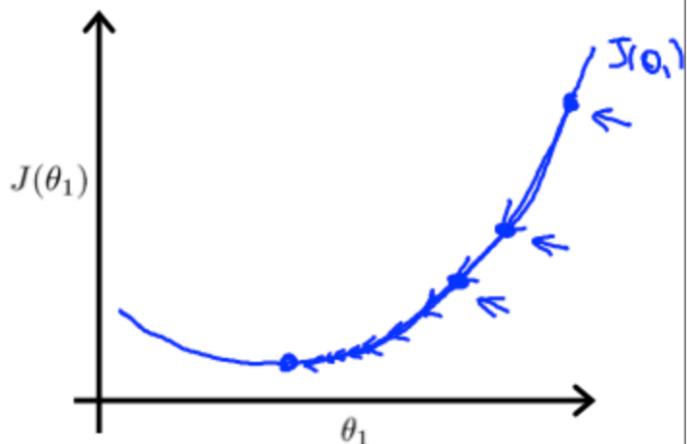
Now, I have another question for you.

if  $\theta_1$  is already the optimal minimum, the algorithm just leave it unchanged, so we get:

Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  fixed.

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease  $\alpha$  over time.



## b. GD for linear regression

the derivative:

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{2}{2m} \cdot \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$= \frac{2}{2m} \cdot \frac{1}{2m} \sum_{i=1}^m (\underline{\theta_0 + \theta_1 x^{(i)}} - y^{(i)})^2$$

$$\Theta_0, j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\Theta_1, j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

So armed with these definitions or  
armed with what we worked out to be

plug back to the GD algorithm:

repeat until convergence {

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

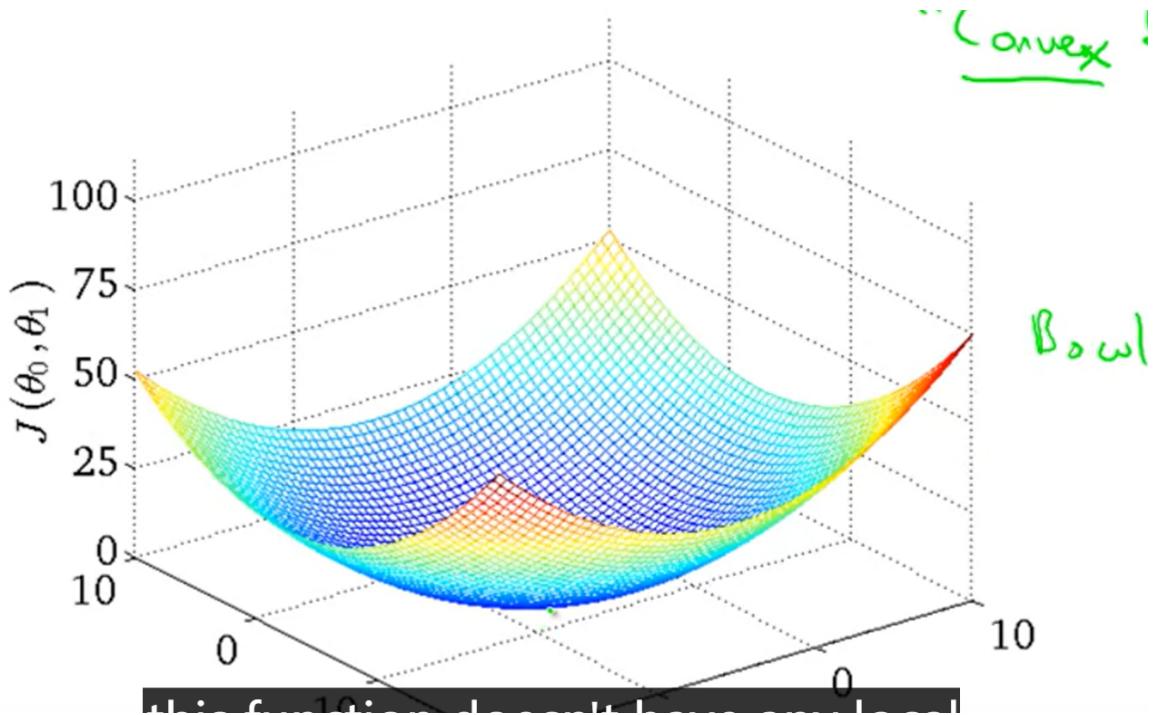
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

An example to calculate this( $\theta_1$ )

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

the cost function of linear regression is always a **convex function** like this



## 5. Linear Algebra Fundamentals(Review)

- a. matrices and vectors(trivial)
- b. addition and scalar multiplication(trivial)
- c. matrix matrix multiplication

neat trick to do prediction by one hypothesis function

House sizes:

- 2104
- 1416
- 1534
- 852

Matrix  $\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}$   $4 \times 2$

$h_{\theta}(x) = -40 + 0.25x$

$h_{\theta}(x)$   $2 \times 1$  Vector  $\begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$

$h_{\theta}(2104)$   $4 \times 1$  matrix  $\begin{bmatrix} -40 \times 1 + 0.25 \times 2104 \\ -40 \times 1 + 0.25 \times 1416 \\ \vdots \\ -40 \times 1 + 0.25 \times 852 \end{bmatrix}$

$h_{\theta}(1416)$

$\boxed{\text{prediction} = \text{DataMatrix} \times \text{Parameters.}}$  for  $i = 1:1000, \text{ prediction}(i) := \dots$

get not only a simpler piece

3:20 13:19 / 13:39 Andrew Ng

another neat trick in the situation of many hypothesis functions

House sizes:

$$\begin{cases} 2104 \\ 1416 \\ 1534 \\ 852 \end{cases}$$

Matrix  $\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}$

Have 3 competing hypotheses:

1.  $h_{\theta}(x) = -40 + 0.25x$
2.  $h_{\theta}(x) = 200 + 0.1x$
3.  $h_{\theta}(x) = -150 + 0.4x$

Matrix  $\begin{bmatrix} -40 \\ 200 \\ -150 \\ 0.25 \\ 0.1 \\ 0.4 \end{bmatrix}$

$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times \begin{bmatrix} -40 & 200 & -150 \\ 0.25 & 0.1 & 0.4 \end{bmatrix} = \begin{bmatrix} 486 & 410 & 692 \\ 314 & 342 & 416 \\ 344 & 353 & 464 \\ 173 & 285 & 191 \end{bmatrix}$

Prediction Predictions

So with just one matrix multiplication!

properties:

1. Not commutative;
2. Associative;
3. Identity matrix is commutative;

d. inverse and transpose(trivial)

## W2

### 1. Vectorization(intro)

trick 1

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

$$= \theta^T x$$

trick 2

$$\begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ \theta_2 &:= \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} \\ (n=2) \end{aligned}$$

(h<sub>θ</sub>)

Vectorized implementation:

$$\Theta := \Theta - \alpha \frac{\delta}{m} X^T \delta$$

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

$$\delta = \begin{bmatrix} \delta_0 \\ \delta_1 \\ \vdots \\ \delta_n \end{bmatrix} \quad \delta_0 = \frac{1}{m} \sum (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

$$X^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

note: treat  $\theta$  as a vector,  $\Sigma$  blah blah  $x$  also a vector.

### 2. Multivariate Linear Regression

*Glossary #2*

$n$ : number of features

$x^{(i)}$ : the  $i^{th}$  entry of the training set

$x_j^{(i)}$ : the value of feature  $j$  of  $x^{(i)}$

### a. hypothesis function of MLR

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

let  $x_0 = 1$

$$\text{let } x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\text{let } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

thus we have

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

### b. GD for multiple variables

just the same thing...

Hypothesis:  $\underline{h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n}$

Parameters:  $\underline{\theta_0, \theta_1, \dots, \theta_n}$   $\underbrace{\quad}_{\text{n+1-dimensional vector}}$

Cost function:

$$\underline{J(\theta_0, \theta_1, \dots, \theta_n)} = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Gradient descent:

Repeat {  
 $\rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$  } J( $\theta$ )  
 simultaneously update for every  $j = 0, \dots, n$ )

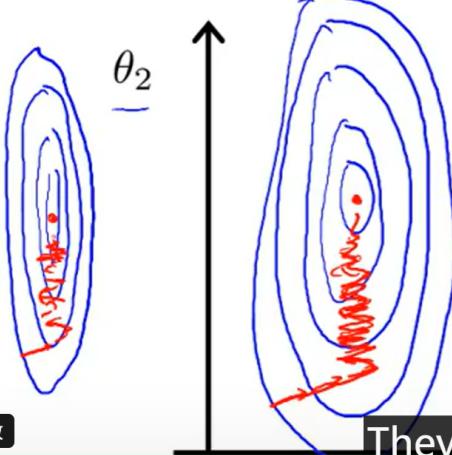
### c. practical tricks for GD of MLR

#### 1. feature scaling (特征缩放)

idea: make sure features are on the similar scale

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

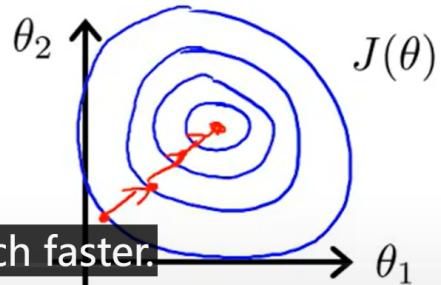
$x_2 = \text{number of bedrooms (1-5)}$   $\leftarrow$



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

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

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



They can convert much faster.

Get every feature into approximately a  $-1 \leq x_i \leq 1$  range.

$$x_0 = 1$$

$$\frac{-1}{\nwarrow} \leq \frac{x_i}{\nearrow} \leq \frac{1}{\nwarrow}$$

$$6 \leq x_1 \leq 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.5 \quad \checkmark$$

$$-100 \leq x_3 \leq 100 \quad \times$$

$$-0.0001 \leq x_4 \leq 0.0001 \quad \times$$

## Mean normalization

Replace  $x_i$  with  $\frac{x_i - \mu_i}{\sigma_i}$  to make features have approximately zero mean  
(Do not apply to  $x_0 = 1$ ).

E.g.  $x_1 = \frac{\text{size} - 1000}{2000}$  Average size = 100

$$x_2 = \frac{\#\text{bedrooms} - 2}{5} \quad 1-5 \text{ bedrooms}$$

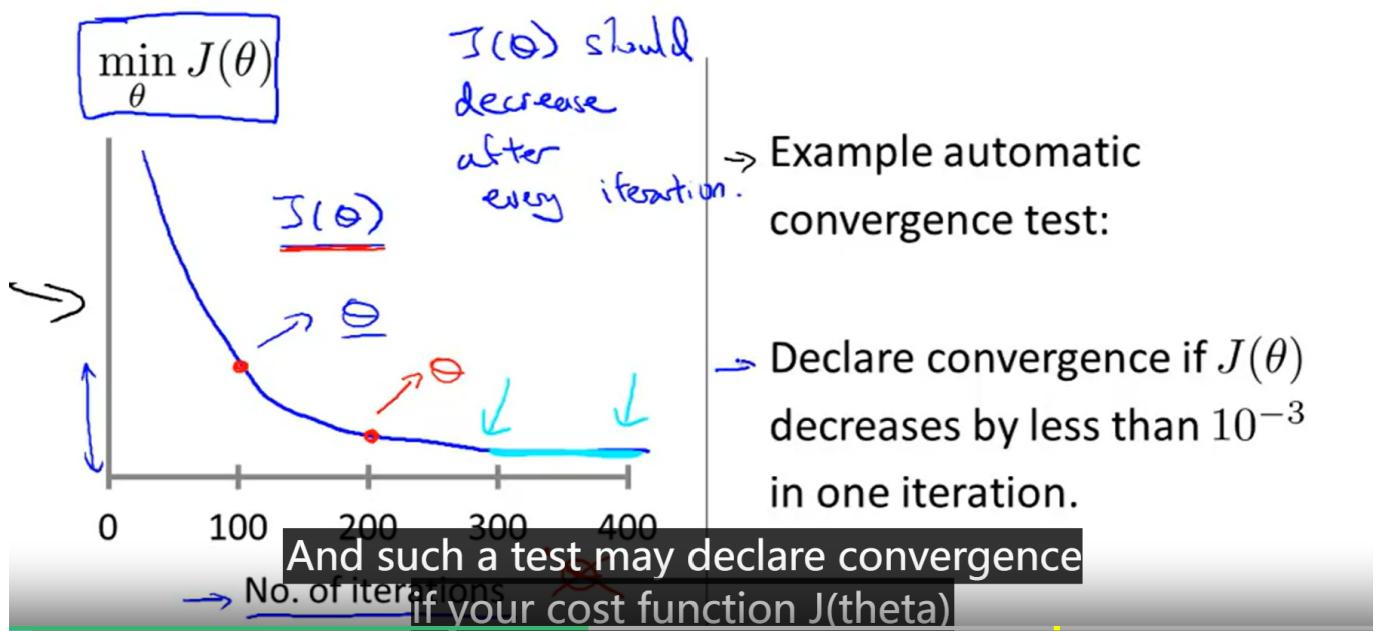
$$[-0.5 \leq x_1 \leq 0.5] \quad [-0.5 \leq x_2 \leq 0.5]$$

$$x_1 \leftarrow \frac{x_1 - \mu_1}{\sigma_1} \quad \begin{array}{l} \text{avg value} \\ \text{of } x_1 \\ \text{in training set} \end{array} \quad \mid \quad x_2 \leftarrow \frac{x_2 - \mu_2}{\sigma_2}$$

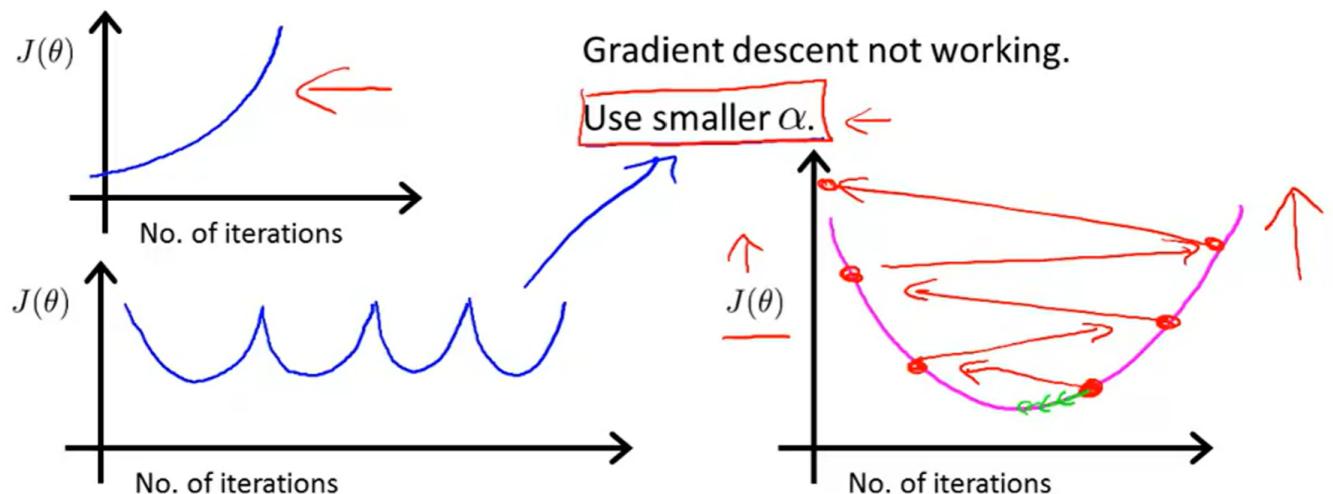
doesn't have to be too exact,  
(or standard deviation)

## 2. learning rate

## Making sure gradient descent is working correctly.



## Making sure gradient descent is working correctly.

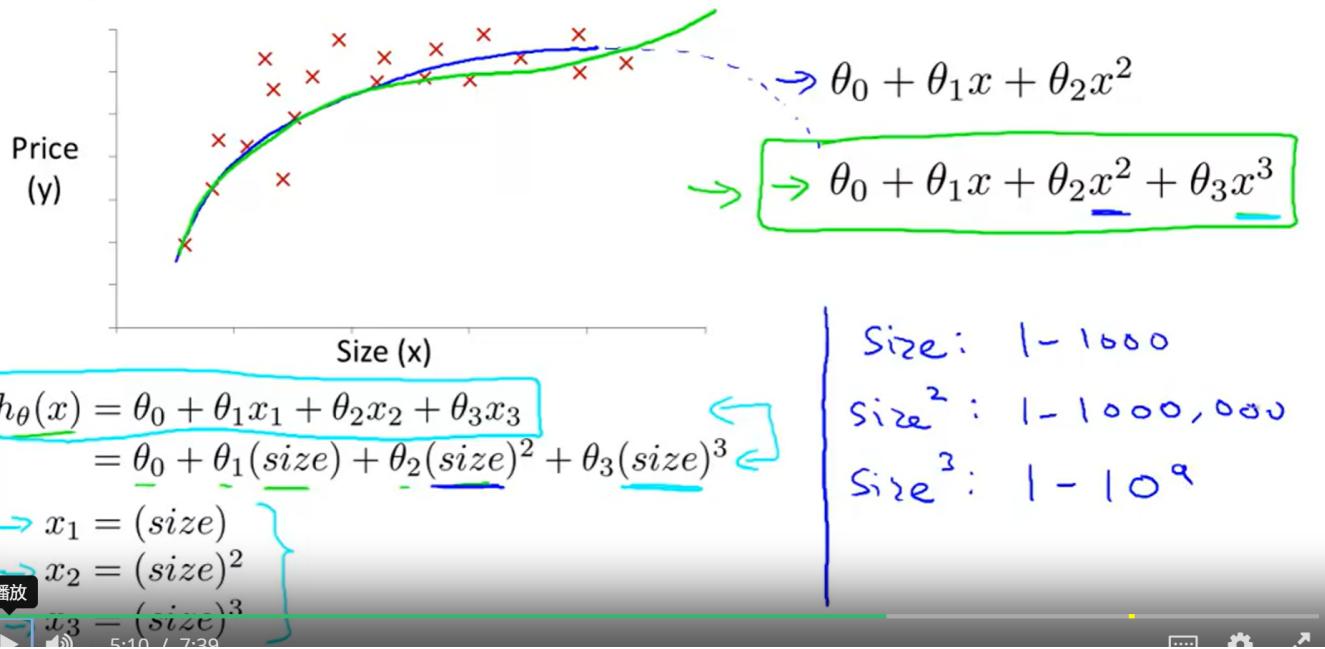


- For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration.
- But if  $\alpha$  is too small, gradient descent can be slow to converge.

$\alpha$  can neither be too small nor too big

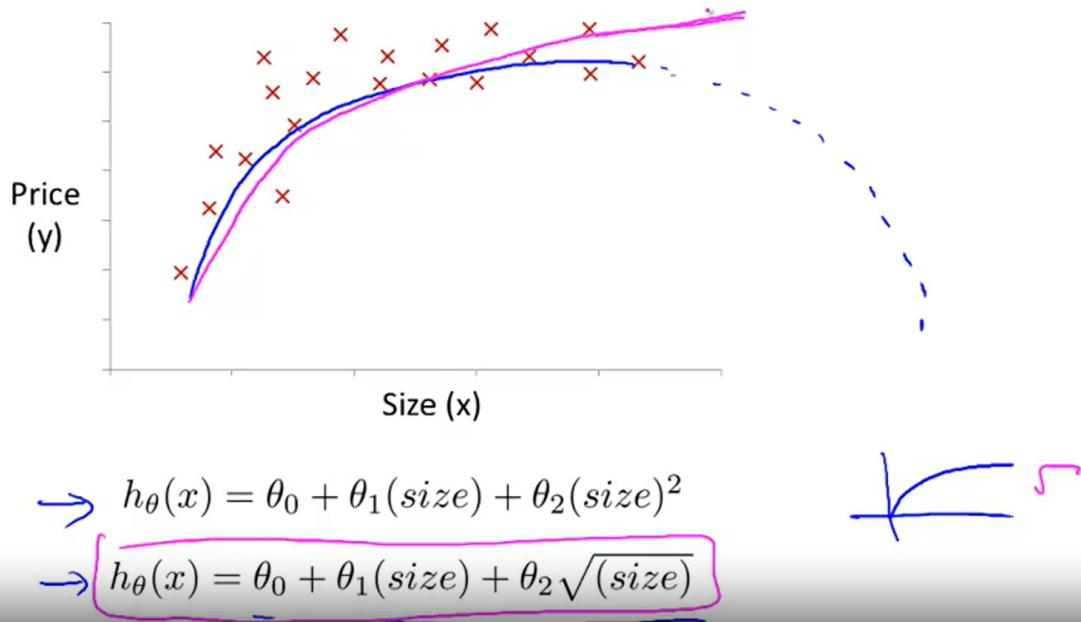
d. Features and Polynomial Regression

## Polynomial regression



note: treat the  $\text{size}$ ,  $\text{size}^2$ ,  $\text{size}^3$  as  $x_1, x_2, x_3$  and the problem is equivalent to MLR

## Choice of features



note: an algorithm will help us to choose a feature

idea:

1. we can combine various features to one feature
2. we can create new feature
3. we have multiple choices of features

### e. normal equation(正规方程)

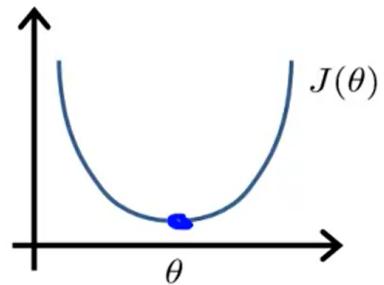
solve  $\theta$  analytically

Intuition: If 1D ( $\theta \in \mathbb{R}$ )

$$\rightarrow J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{\partial}{\partial \theta} J(\theta) = \dots \stackrel{?}{=} 0$$

Solve for  $\theta$



$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad (\text{for every } j)$$

idea: derivation! derivation!

an example:

Examples:  $m = 4$ .

$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$  m x (n+1)

$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$  m-dimensional vector

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

choice of GD & NE

pros and cons of the two methods

Gradient Descent	Normal Equation
Need to choose alpha	No need to choose alpha
Needs many iterations	No need to iterate
$O(kn^2)$	$O(n^3)$ , need to calculate inverse of $X^T X$
Works well when n is large	Slow if n is very large

if  $n$  is large(probably  $n > 10000$ ), then use GD else use NE

use `pinv()` instead of `inv()` in MATLAB in case of  $X^T X$  is non-invertible

## W3

---

### 1. Classification Introduction

binary:- just{0, 1}

multi:- {0, 1, 2...}

using Linear Regression is not a good idea in classification problems

### 2. Logistic Regression

#### a. sigmoid function(logistic function)

in classification problems we have:

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

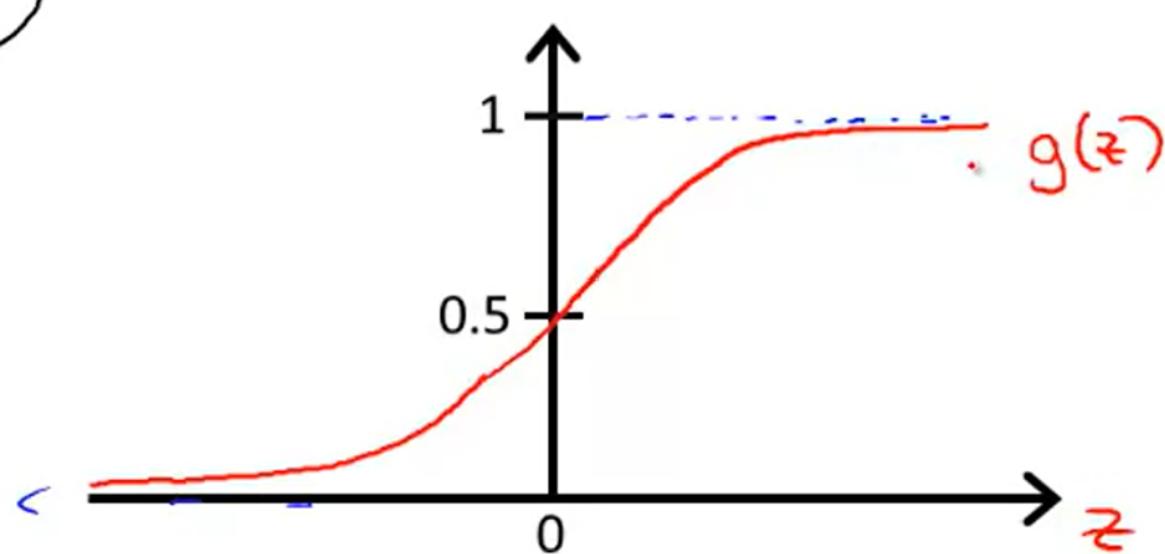
where:

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

thus, we have

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

the diagram of the sigmoid function:



the output of  $h_\theta(x)$  means probability, i.e.  $P(y = 1|x; \theta)$

$h_\theta(x)$  = estimated probability that  $y = 1$  on input  $x$

Example: If  $\underline{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

$$h_\theta(\underline{x}) = 0.7$$

Tell patient that 70% chance of tumor being malignant

and we have:

$$h_\theta(x) = P(y = 1|x; \theta) = 1 - P(y = 0|x; \theta)$$

$$P(y = 0|x; \theta) + P(y = 1|x; \theta) = 1$$

#### b. decision boundary(决策边界)

example of logistic regression prediction:

## Logistic regression

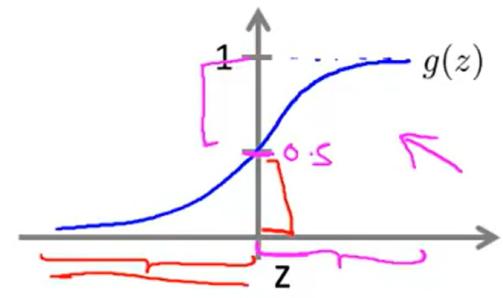
$$\rightarrow h_{\theta}(x) = g(\theta^T x) = P(y=1|x; \theta)$$

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

Suppose predict " $y = 1$ " if  $\underline{h_{\theta}(x) \geq 0.5}$   
 $\underline{\theta^T x \geq 0}$

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

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



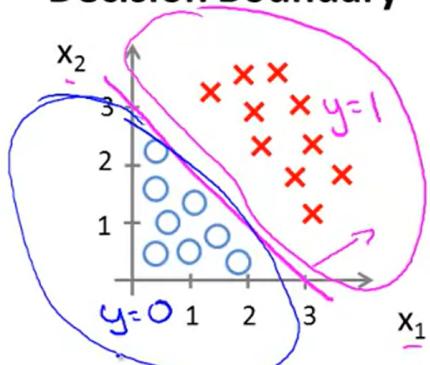
$g(z) \geq 0.5$   
when  $z \geq 0$

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

wherever  $\underline{\theta^T x \geq 0}$

decision boundary example: draw a line(decision boundary) to take different classes apart

## Decision Boundary



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \leftarrow$$

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

Predict " $y = 1$ " if  $\underline{-3 + x_1 + x_2 \geq 0}$   
 $\underline{\theta^T x}$

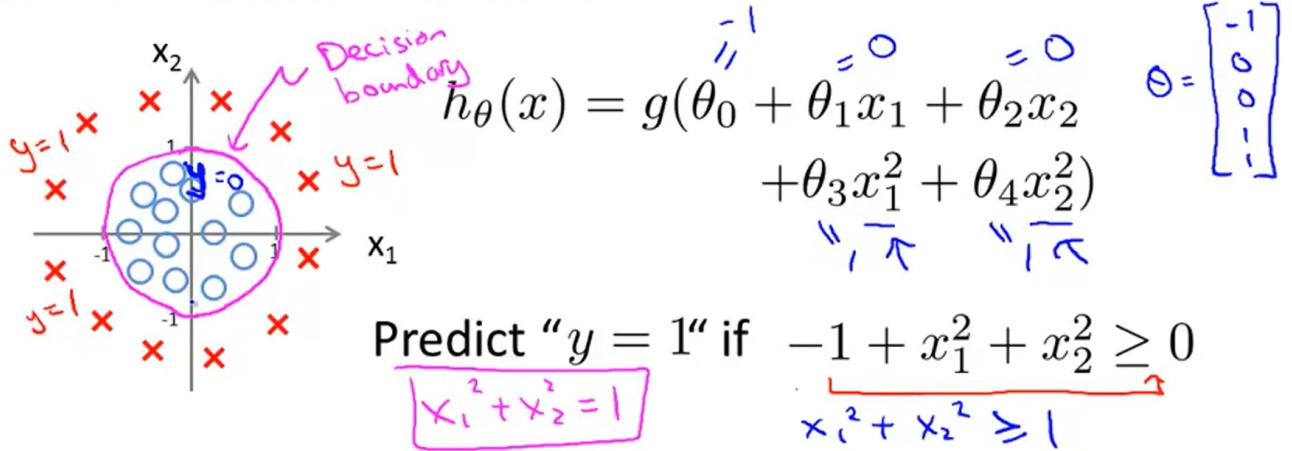
$x_1, x_2$

$$\underline{x_1 + x_2 = 3}$$

that region on the left is the region  $x_2 < 3$   
where our hypothesis will predict  $y = 0$ .

non-linear decision boundary:

## Non-linear decision boundaries



that I can get in this example.

### c. logistic regression cost function & gradient descent

for linear regression we have:

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

rewrite it as:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x), y)$$

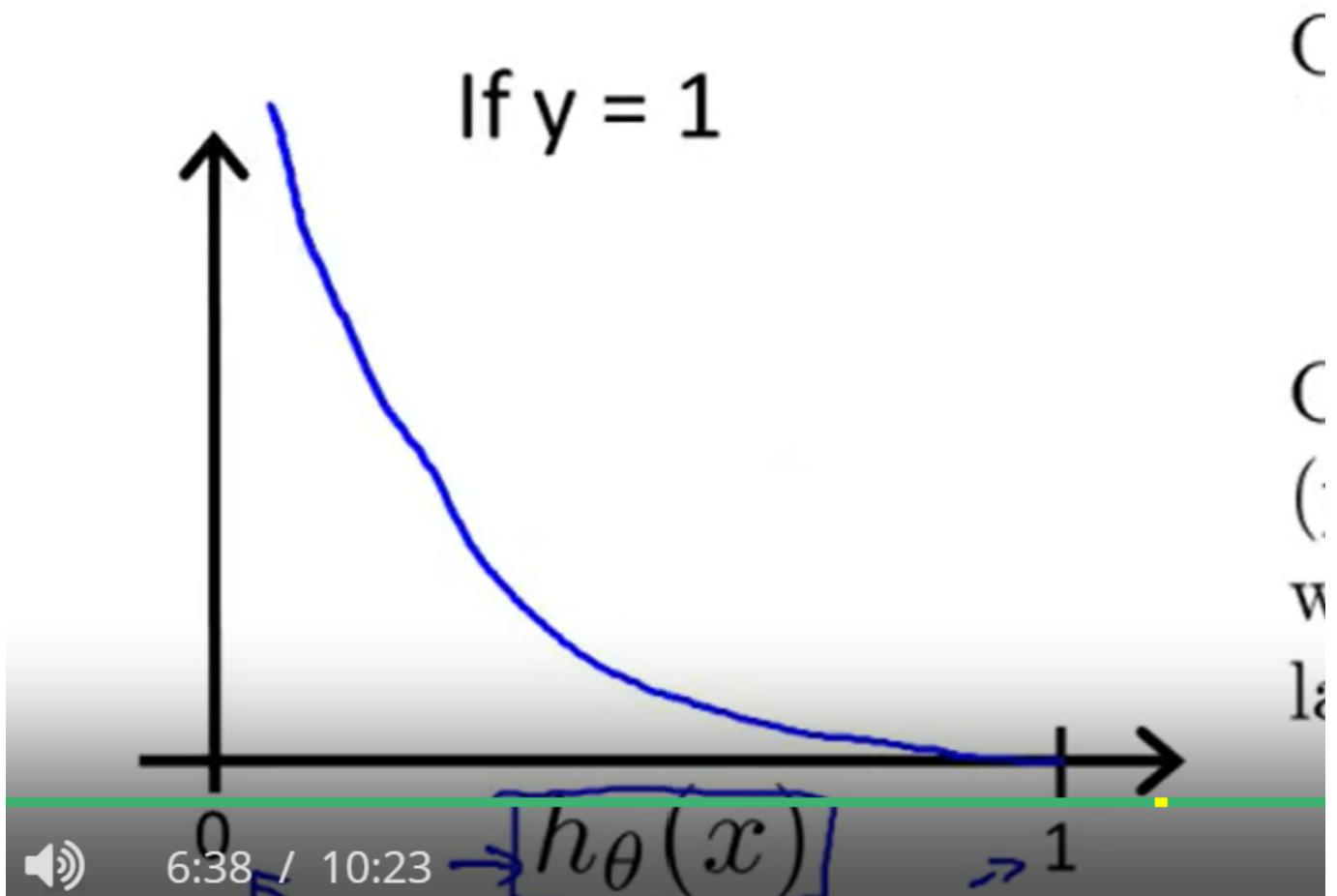
where:

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

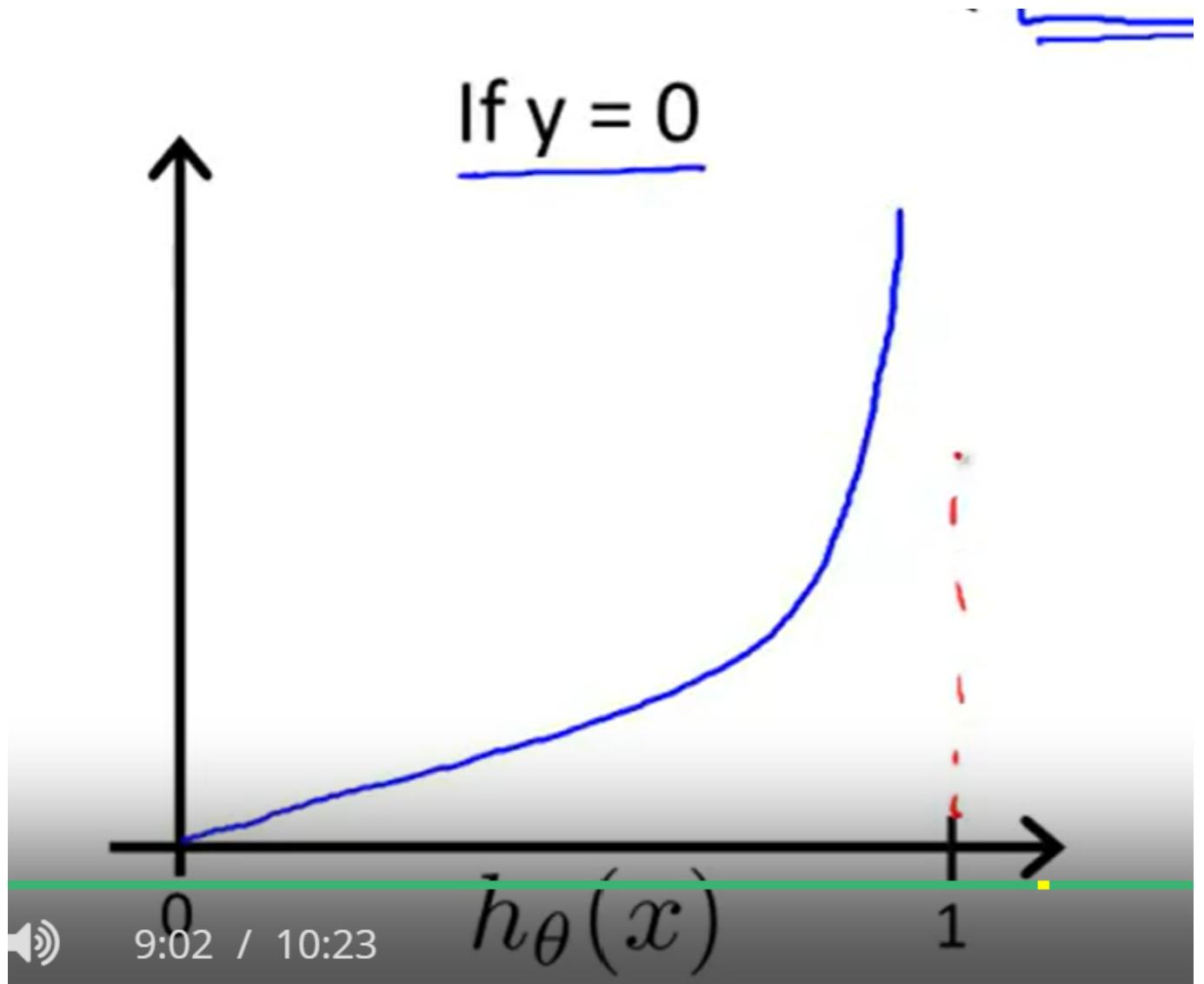
but in logistic regression, we have:

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

$\log$  is actually  $\ln$ , the diagram of the cost function( $y==1$ ) is:



the other part of the diagram is:



simplified cost function:

$$\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1 - y) \log(1 - h_\theta(x))$$

as  $y$  just has two possible values, it is equivalent to the previous equation.

so we have the final version of the cost function of logistic regression:

$$J(\theta) = -\frac{1}{m} \left( \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right)$$

a vectorized version is:

$$h = g(X\theta)$$

$$J(\theta) = \frac{1}{m} \cdot \left( -y^T \log(h) - (1-y)^T \log(1-h) \right)$$

this function is derived from MLE in statistics

and the GD for logistic regression:

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

just like that in linear regression

#### d. other optimization algorithms

Optimization algorithms:

- - Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages:

- No need to manually pick  $\alpha$
- Often faster than gradient descent.

Disadvantages:

- More complex

文

how to use it in MATLAB:

step1: write a function that calculates  $J(\theta)$  and  $\frac{\partial J}{\partial \theta_j}$

```
function [jVal, gradient] = costFunction(theta)  
    jVal = [... code to compute J(theta)...];  
    gradient = [... code to compute derivative of J(theta)...];  
end
```

step2:

use function `fminunc()` to optimize  $\theta$

```

options = optimset('GradObj', 'on', 'MaxIter', 100);
initialTheta = zeros(2,1);
[optTheta, functionVal, exitFlag] = fminunc(@costFunction, initialTheta, options);

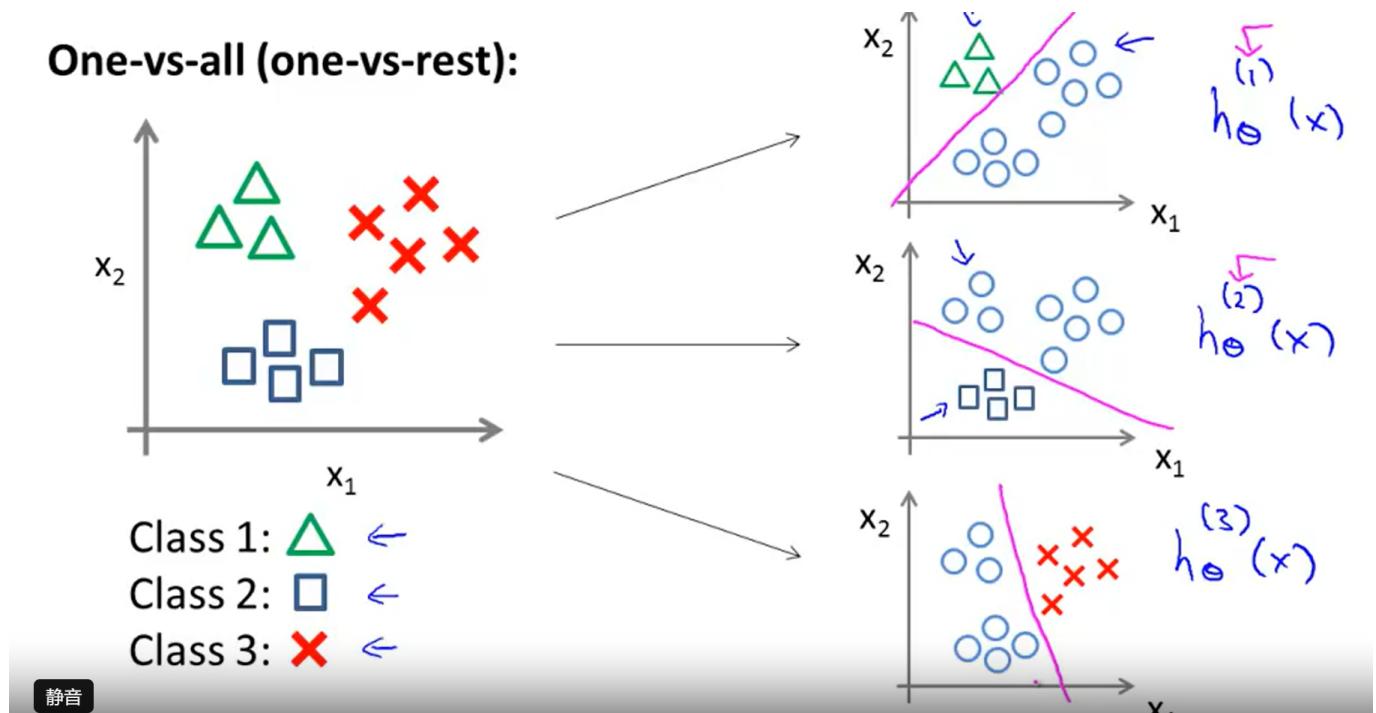
```

### 3. Multi-class Classification

#### a. one-vs-all classification

idea: separate this problem into different binary-classification problems

example:



i.e. we will train logistic regression classifiers as below:

$$h_{\theta}^{(i)}(x) = P(y = i|x; \theta)$$

where  $i = 1, 2, \dots$

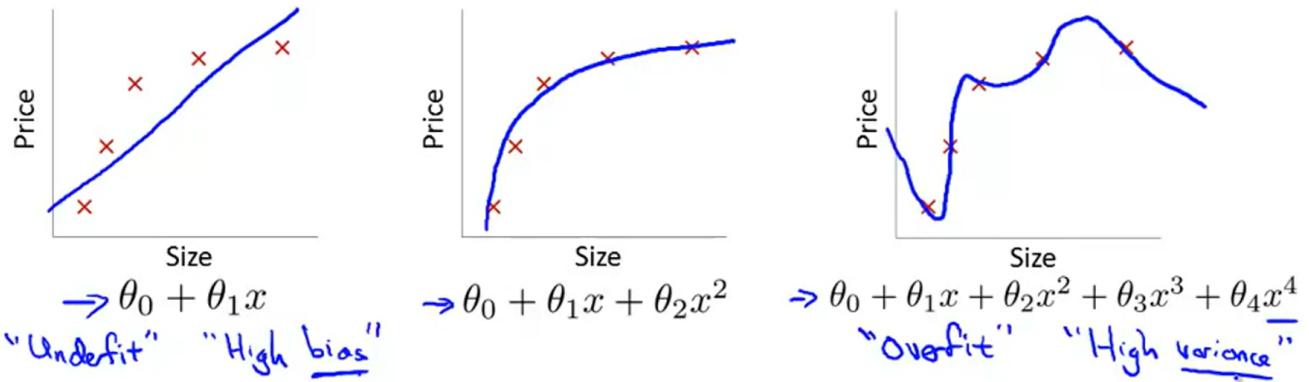
when making predictions, we just have to calculate  $\max_i h_{\theta}^{(i)}(x)$

### 4. Overfitting

#### a. introduction

example:

## Example: Linear regression (housing prices)



2 ways to solve the problem:

1) Reduce the number of features:

- Manually select which features to keep.
- Use a model selection algorithm (studied later in the course).

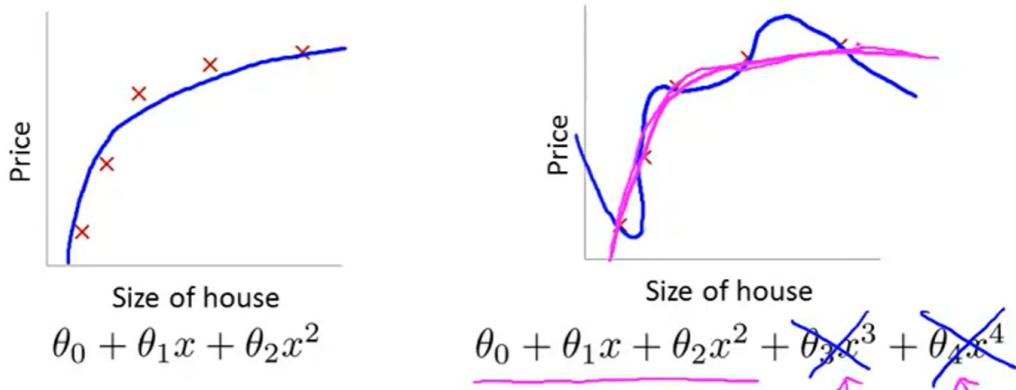
2) Regularization

- Keep all the features, but reduce the magnitude of parameters  $\theta_j$ .
- Regularization works well when we have a lot of slightly useful features.

### b. regularization

1. intuition&idea:

## Intuition



Suppose we penalize and make  $\theta_3, \theta_4$  really small.

$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \underline{\theta_3^2} + 1000 \underline{\theta_4^2}$$

## Regularization.

Small values for parameters  $\theta_0, \theta_1, \dots, \theta_n$

- “Simpler” hypothesis
- Less prone to overfitting

$$\rightarrow \underline{\theta_3, \theta_4} \approx 0$$

### 2. cost function

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

where

$$\lambda \sum_{j=1}^n \theta_j^2$$

is the regularization term.

$\lambda$  cannot be too large, otherwise it may cause underfitting

### 3. regularized linear regression

gradient descent:

## Gradient descent

$$\frac{\theta_0}{\uparrow} \quad \theta_1, \theta_2, \dots, \theta_n$$

Repeat {

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

$$\rightarrow \theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j = \cancel{X}, 1, 2, 3, \dots, n)$$

}

normal equation:

Suppose  $m \leq n$ ,

(#examples) (#features)

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

If  $\lambda > 0$ ,

$$\theta = \left( X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{bmatrix} \right)^{-1} X^T y$$

So finally I want to just quickly  
describe the issue of non-invertibility.

note that  $X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{bmatrix}$  is always invertible

4. regularized logistic regression

gradient descent:

just like that in linear regression, the only difference is the definition of  $h_\theta(x)$

Repeat {

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

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

$(j = \cancel{0}, 1, 2, 3, \dots, n)$

$\theta_0, \dots, \theta_n$

}

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

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

advanced algorithms in MATLAB:

### Advanced optimization

```

function [jVal, gradient] = costFunction(theta)
    jVal = [ code to compute J(theta) ];
    
$$J(\theta) = \left[ -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

    gradient(1) = [ code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$  ];
    
$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

    gradient(2) = [ code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$  ];
    
$$\left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \right) + \frac{\lambda}{m} \theta_1$$

    gradient(3) = [ code to compute  $\frac{\partial}{\partial \theta_2} J(\theta)$  ];
    
$$\vdots \quad \left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)} \right) + \frac{\lambda}{m} \theta_2$$

    gradient(n+1) = [ code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$  ];

```

just alike.

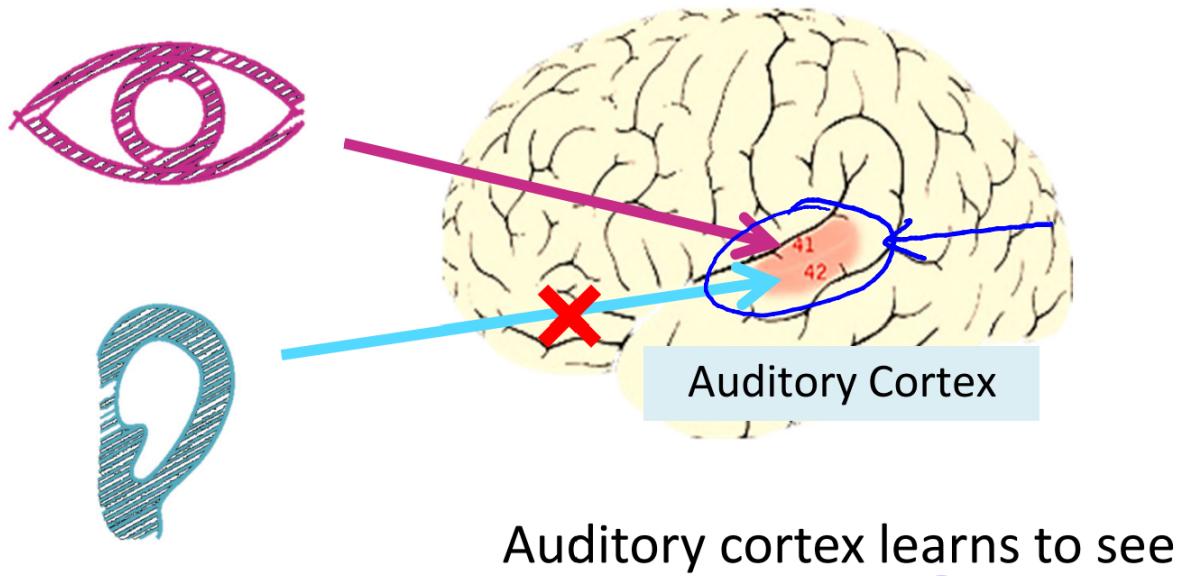
## W4

### 1. Neural Network: Intro & Representation

#### a. motivation

too many features if we want to train a real classifier like a car detector

brains can actually learn to handle different kinds of signals on the same part once they are connected with different sensors.



b. model representation

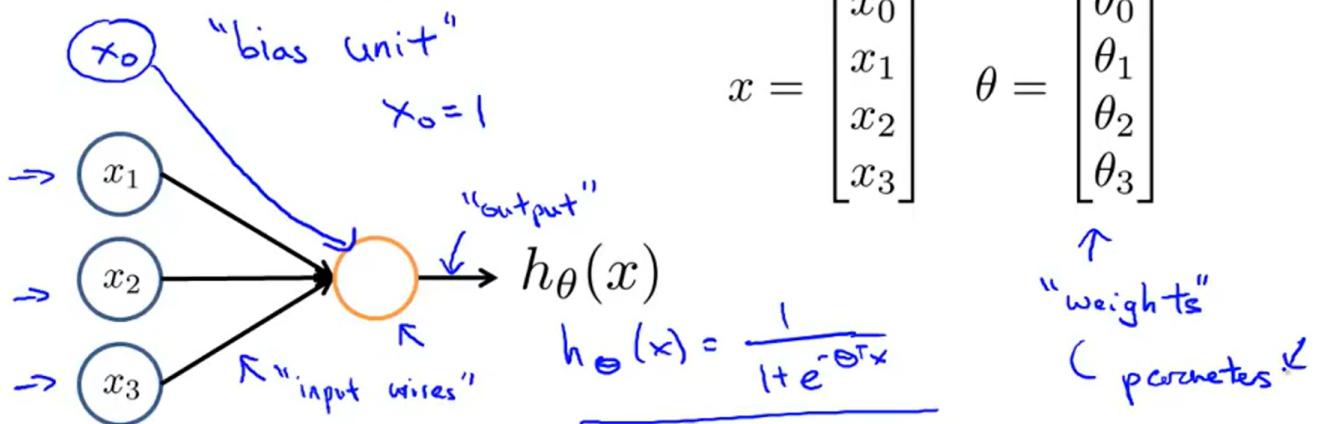
*Glossary #3*

bias unit: input  $x_0$

sigmoid activation function:  $g(x) = \frac{1}{1+e^{-x}}$

weights: vector  $\theta$

## Neuron model: Logistic unit

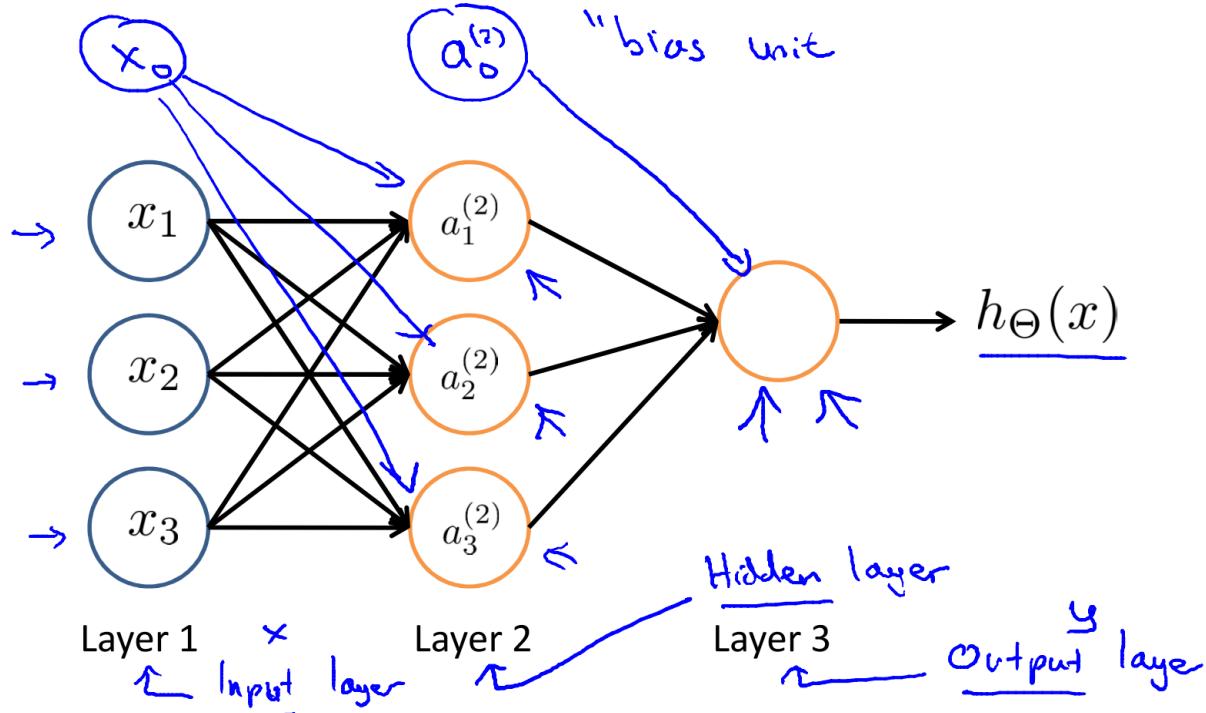


Sigmoid (logistic) activation function.

~~But I'll mostly continue to use the~~

input layer, hidden layer; output layer: as shown in the figure below:

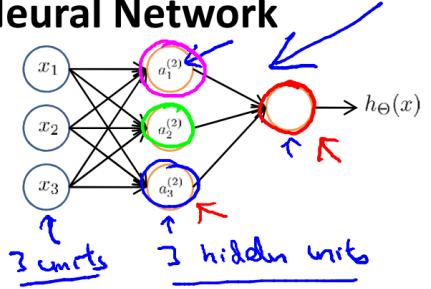
## Neural Network



the computational steps:

Andri

## Neural Network



$\rightarrow a_i^{(j)}$  = "activation" of unit  $i$  in layer  $j$

$\rightarrow \Theta^{(j)}$  = matrix of weights controlling function mapping from layer  $j$  to layer  $j + 1$

$$h_{\Theta}(x)$$

$$\rightarrow a_1^{(2)} = g(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3)$$

$$\rightarrow a_2^{(2)} = g(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3)$$

$$\rightarrow a_3^{(2)} = g(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3)$$

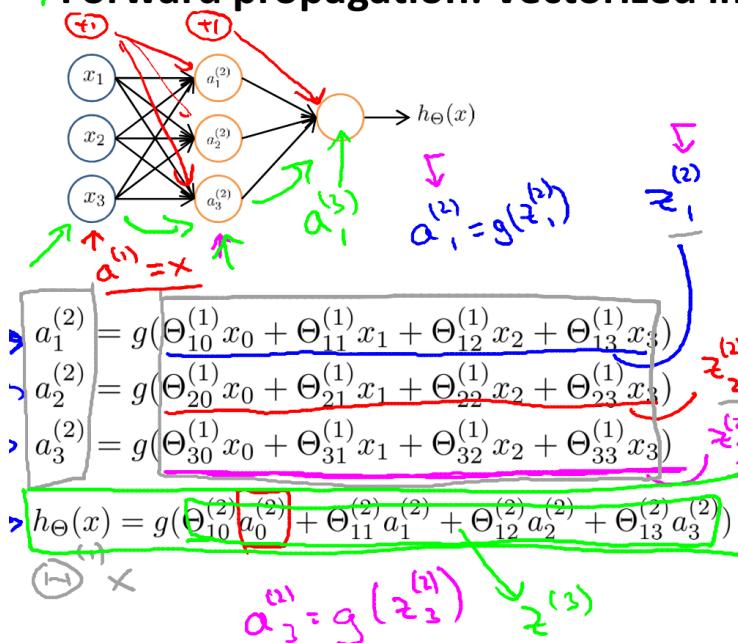
$$\rightarrow h_{\Theta}(x) = a_1^{(3)} = g(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)})$$

$\rightarrow$  If network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j + 1$ , then  $\Theta^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$ .

$$s_{j+1} \times (s_j + 1)$$

vectorized steps: forward propagation

## Forward propagation: Vectorized implementation



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$\boxed{\begin{aligned} z^{(2)} &= \Theta^{(1)} \times a^{(1)} \\ a^{(2)} &= g(z^{(2)}) \\ z^{(3)} &= \Theta^{(2)} a^{(2)} \\ h_{\Theta}(x) &= g(z^{(3)}) \end{aligned}}$$

$a^{(1)} \in \mathbb{R}^3$   
 $a^{(2)} \in \mathbb{R}^3$   
 $a^{(3)} \in \mathbb{R}^1$

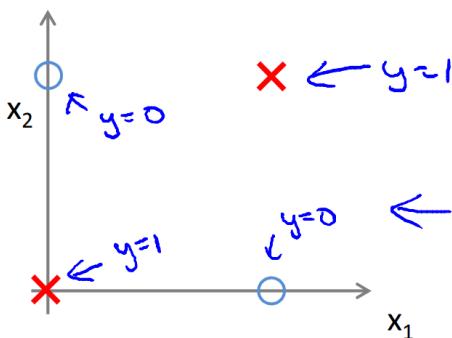
NN just do things like logistic regression except that the features fed to  $h_{\theta}(x)$  are from the hidden layers. NN learn its own features

### c. intuition & application

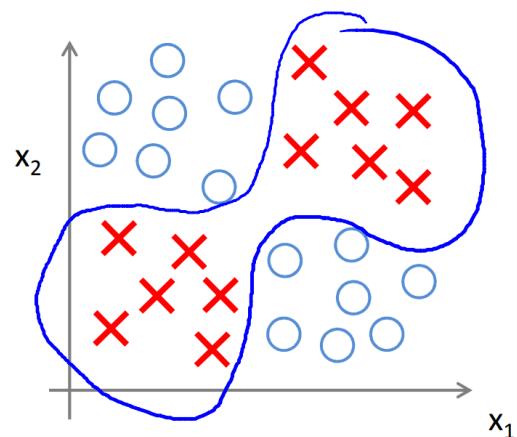
intuition: XNOR network

## Non-linear classification example: XOR/XNOR

→  $x_1, x_2$  are binary (0 or 1).



$$y = \underbrace{x_1 \text{ XOR } x_2}_{\substack{\rightarrow x_1 \text{ XNOR } x_2 \\ \rightarrow \text{NOT } (x_1 \text{ XOR } x_2)}}$$

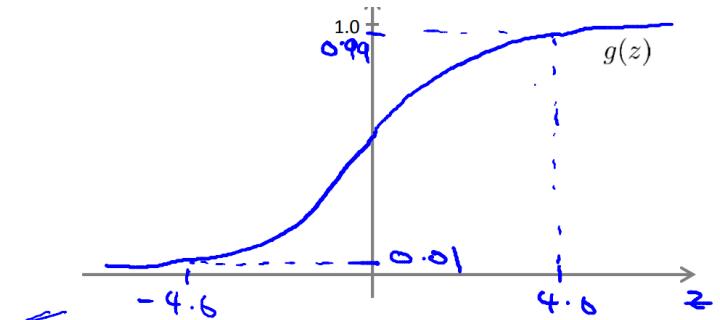
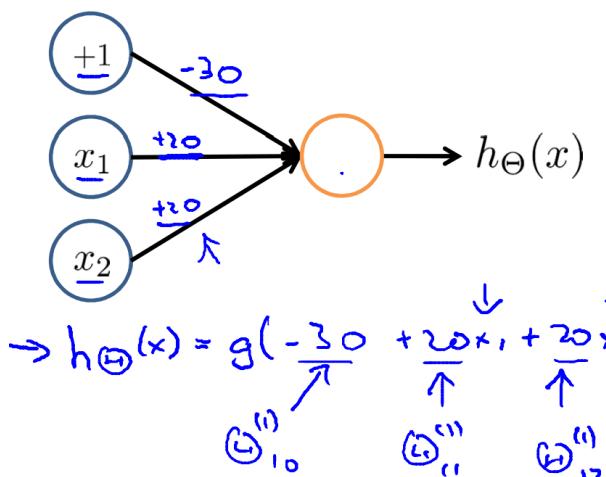


starting from AND

## Simple example: AND

→  $x_1, x_2 \in \{0, 1\}$

→  $y = x_1 \text{ AND } x_2$

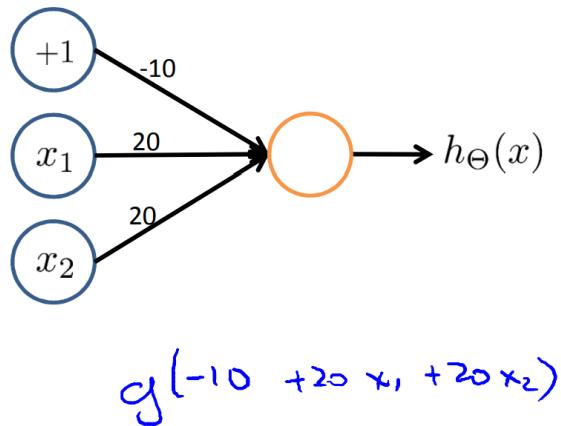


$x_1$	$x_2$	$h_{\Theta}(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

$h_{\Theta}(x) \approx x_1 \text{ AND } x_2$

and OR

## Example: OR function



$x_1$	$x_2$	$h_\Theta(x)$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	$\approx 1$
1	1	$\approx 1$

and NOT

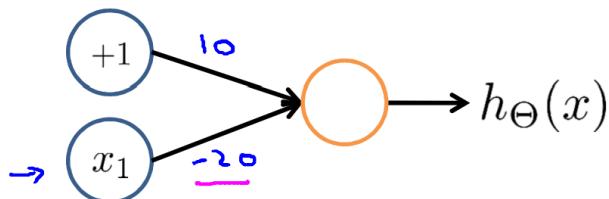
$\rightarrow x_1 \text{ AND } x_2$

$\rightarrow x_1 \text{ OR } x_2$

$\{0, 1\}$

**Negation:**

NOT  $x_1$



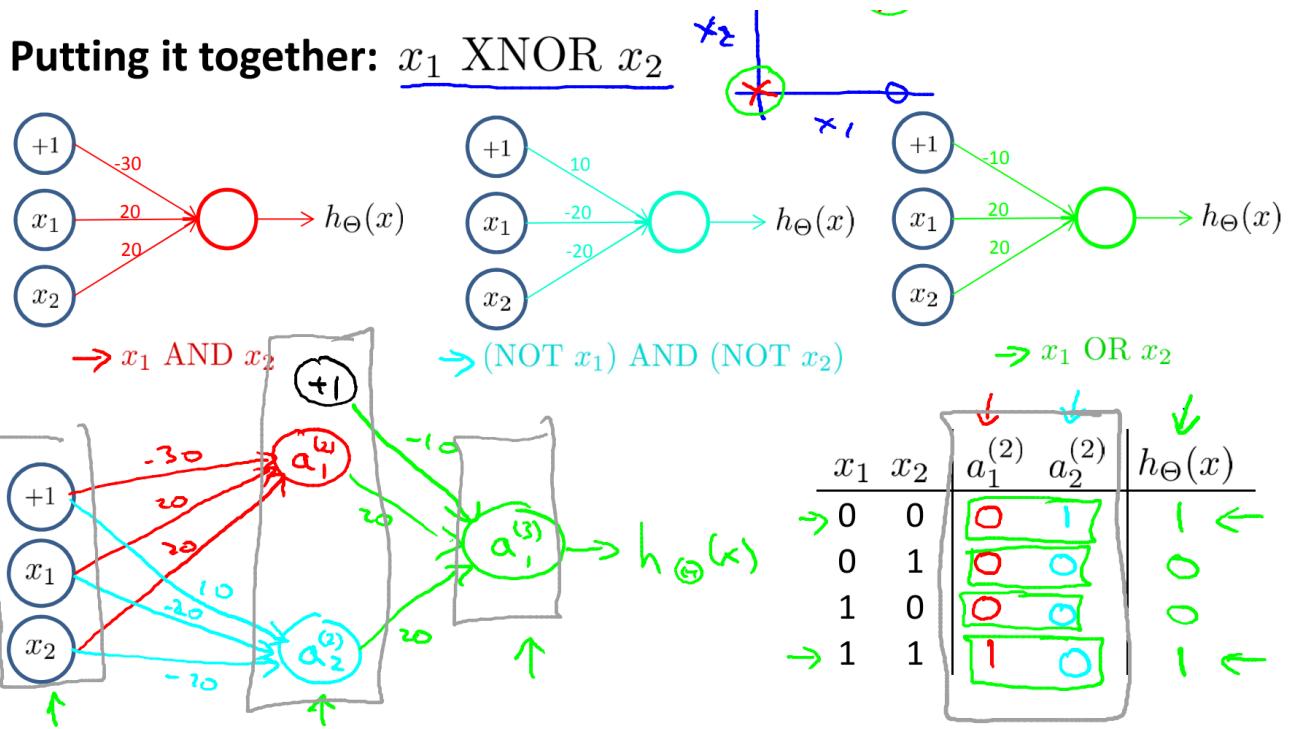
$x_1$	$h_\Theta(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

$$h_\Theta(x) = g(10 - 20x_1)$$

$\rightarrow (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$

$= 1$  if and only if  
 $\rightarrow x_1 = x_2 = 0$

finally we can get:



application: handwritten digits classification

## W5

### 1. Neural Network: Learning

a. cost function

**Glossary #4**

$L$ : the number of layers of a neural network

$s_l$ : the number of units of each layer (without bias unit)

$K$ : the dimension of the output layer

cost function of NN in classification problems:

## Neural network:

$$h_{\Theta}(x) \in \mathbb{R}^K \quad (h_{\Theta}(x))_i = i^{th} \text{ output}$$

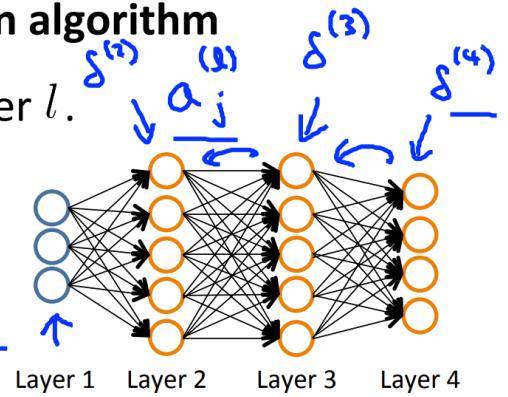
$$\begin{aligned} J(\Theta) = & -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\Theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_k) \right] \\ & + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2 \end{aligned}$$

b. back propagation: compute the derivatives of  $J(\Theta)$

1. compute  $\delta^{(j)}$

### Gradient computation: Backpropagation algorithm

Intuition:  $\underline{\delta_j^{(l)}}$  = "error" of node  $j$  in layer  $l$ .



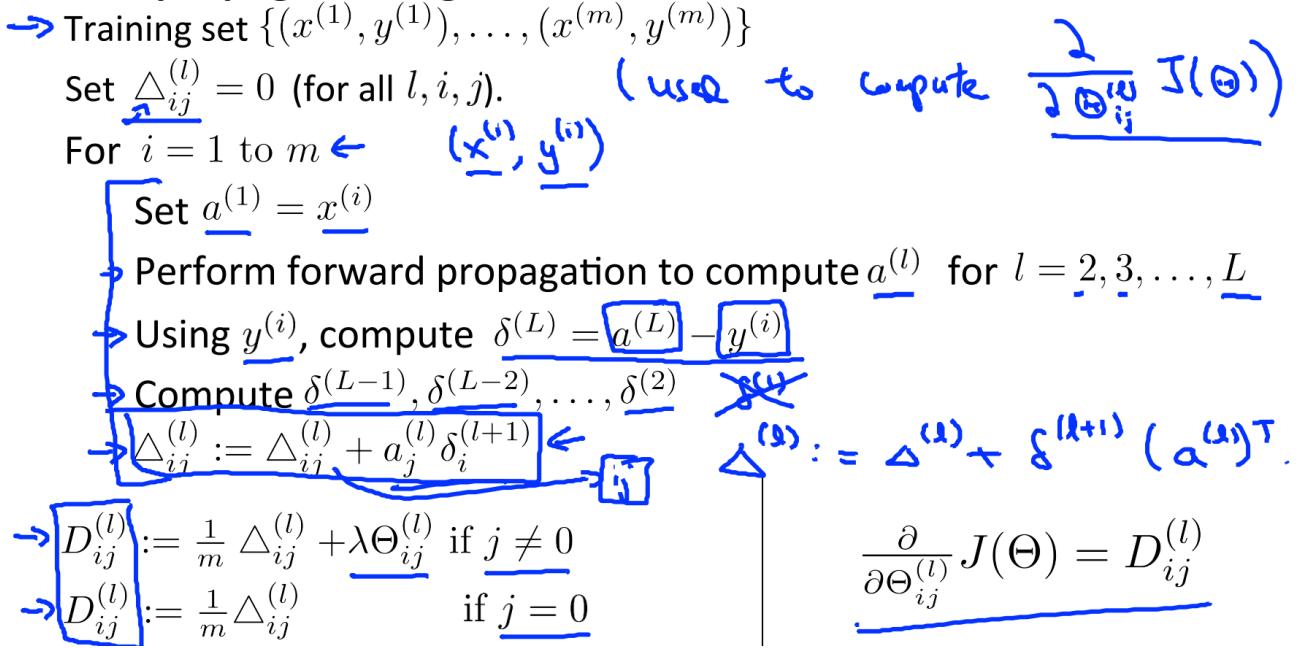
For each output unit (layer  $L = 4$ )

$$\delta_j^{(4)} = \underline{a_j^{(4)} - y_j} \quad (h_{\Theta}(x))_j \quad \underline{\delta^{(4)}} = \underline{a^{(4)} - y}$$

$$\begin{aligned} \delta^{(3)} &= (\Theta^{(3)})^T \delta^{(4)} \cdot * g'(z^{(3)}) \\ \delta^{(2)} &= (\Theta^{(2)})^T \delta^{(3)} \cdot * g'(z^{(2)}) \\ \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) &= a_i^{(l)} \delta_i^{(l+1)} \quad (\text{ignoring } \lambda; \text{ if } \lambda = 0) \end{aligned}$$

2. compute derivatives

## Backpropagation algorithm



steps of BP

Given training set  $\{(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})\}$

1. Set  $\Delta_{ij}^{(l)} := 0$  for all  $(l, i, j)$ , (hence you end up having a matrix full of zeros)

for training example  $t = 1$  to  $m$ :

1. Set  $a^{(1)} := x^{(t)}$

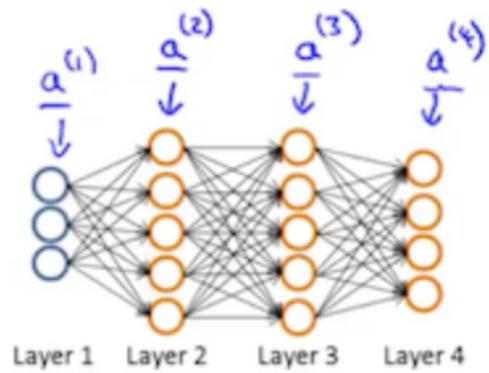
2. Perform forward propagation to compute  $a^{(l)}$  for  $l=2,3,\dots,L$ :

## Gradient computation

Given one training example ( $x, y$ ):

Forward propagation:

$$\begin{aligned}
 \underline{a^{(1)}} &= \underline{x} \\
 \rightarrow z^{(2)} &= \Theta^{(1)} a^{(1)} \\
 \rightarrow a^{(2)} &= g(z^{(2)}) \quad (\text{add } a_0^{(2)}) \\
 \rightarrow z^{(3)} &= \Theta^{(2)} a^{(2)} \\
 \rightarrow a^{(3)} &= g(z^{(3)}) \quad (\text{add } a_0^{(3)}) \\
 \rightarrow z^{(4)} &= \Theta^{(3)} a^{(3)} \\
 \rightarrow a^{(4)} &= \underline{h_{\Theta}(x)} = g(z^{(4)})
 \end{aligned}$$



3. Using  $y^{(t)}$ , compute  $\delta^{(L)} = a^{(L)} - y^{(t)}$

4. Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$  using  $\delta^{(l)} = ((\Theta^{(l)})^T \delta^{(l+1)}) . * a^{(l)} . * (1 - a^{(l)})$

where  $a^{(l)} . * (1 - a^{(l)}) = g'(z^{(l)})$  namely the derivative of the sigmoid function

5.  $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

end of `for` loop

1.  $D_{i,j}^{(l)} := \frac{1}{m}(\Delta_{i,j}^{(l)} + \lambda \Theta_{i,j}^{(l)})$  if  $j \neq 0$ .

2.  $D_{i,j}^{(l)} := \frac{1}{m} \Delta_{i,j}^{(l)}$  if  $j = 0$

then we get the partial derivative  $\frac{\partial}{\partial \Theta_{i,j}^{(l)}} J(\Theta) = D_{i,j}^{(l)}$

### c. BP in practice

#### 1. parameter unrolling from matrices to vectors(系数展开)

reason: advanced optimizing algorithms treat the inputs as vectors but in NN the inputs are matrices

**idea:** unroll matrices to vectors

example:

## Example

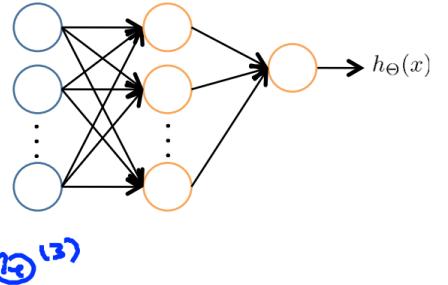
$$s_1 = 10, s_2 = 10, s_3 = 1$$

$\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$

$D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$

```
> thetaVec = [ Theta1(:); Theta2(:); Theta3(:) ];
> DVec = [D1(:); D2(:); D3(:)];

Theta1 = reshape(thetaVec(1:110), 10, 11);
> Theta2 = reshape(thetaVec(111:220), 10, 11);
> Theta3 = reshape(thetaVec(221:231), 1, 11);
```



usage in advanced learning algorithm

## Advanced optimization

```
function [jVal, gradient] = costFunction(theta)
    ...
optTheta = fminunc(@costFunction, initialTheta, options)
```

Neural Network ( $L=4$ ):

$\rightarrow \Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$  - matrices (Theta1, Theta2, Theta3)

$\rightarrow D^{(1)}, D^{(2)}, D^{(3)}$  - matrices (D1, D2, D3)

“Unroll” into vectors

## 2. gradient checking

reason: to debug

idea: check that whether the gradient given by  $D^{(l)}$  is similar to the gradient given by its definition(using the formula  $\frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$ )

implementation:

```
for i = 1:n, ←  
    thetaPlus = theta;  
    thetaPlus(i) = thetaPlus(i) + EPSILON;  
    thetaMinus = theta;  
    thetaMinus(i) = thetaMinus(i) - EPSILON;  
    gradApprox(i) = (J(thetaPlus) - J(thetaMinus))  
                    / (2*EPSILON) ;  
end;
```

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i + \epsilon \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \theta_i + \epsilon$$

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

Check that gradApprox  $\approx$  DVec ←  
From backprop.

notice that we have to re-initialize thetaPlus and thetaMinus because we just calculate the partial derivative of each  $\theta^{(i)}$

advice on how to use gradient check:

### Implementation Note:

- - Implement backprop to compute DVec (unrolled  $D^{(1)}, D^{(2)}, D^{(3)}$ ).
- - Implement numerical gradient check to compute gradApprox.
- - Make sure they give similar values.
- - Turn off gradient checking. Using backprop code for learning.

$$\begin{array}{c} \downarrow \quad \downarrow \quad \downarrow \\ g^{(1)}, g^{(2)}, g^{(3)} \end{array} \xrightarrow{\text{DVec}}$$

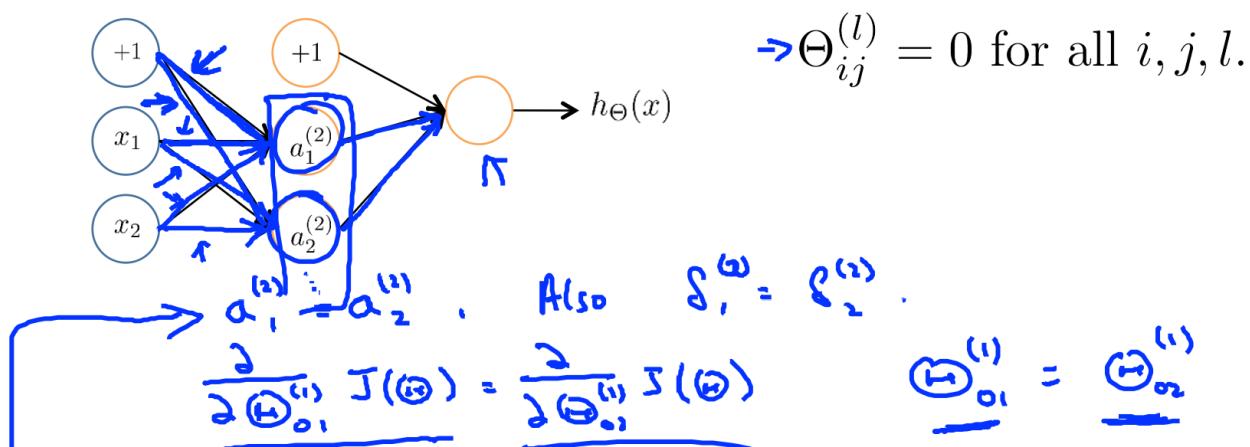
### Important:

- - Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of costFunction(...)) your code will be very slow.

### 3. random initialization

reason:

## Zero initialization



After each update, parameters corresponding to inputs going into each of two hidden units are identical.

$$\underline{a_1^{(1)} = a_2^{(1)}}$$

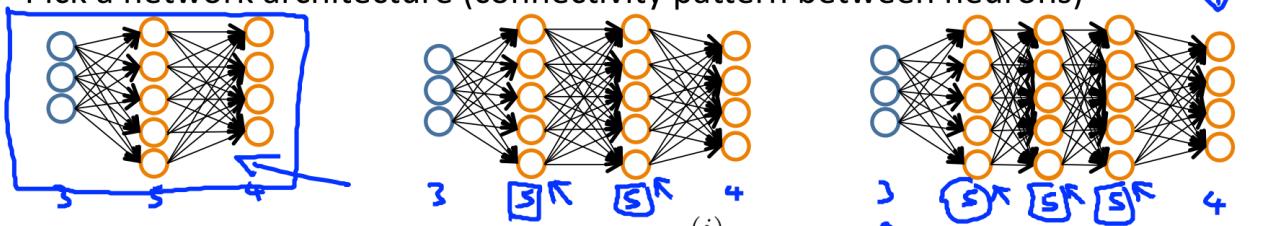
Ar

d. summary: how to train a neural network?

1. pick a network architecture

### Training a neural network

Pick a network architecture (connectivity pattern between neurons)



→ No. of input units: Dimension of features  $x^{(i)}$

→ No. output units: Number of classes

[Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)]

$$y \in \{1, 2, 3, \dots, 10\}$$

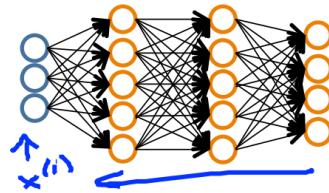
~~yes~~

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow$$

2. training through 6 steps

## Training a neural network

- 1. Randomly initialize weights
- 2. Implement forward propagation to get  $h_{\Theta}(x^{(i)})$  for any  $x^{(i)}$
- 3. Implement code to compute cost function  $J(\Theta)$
- 4. Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$
- for  $i = 1:m$  {  $(x^{(i)}, y^{(i)})$      $(x^{(i)}, y^{(i)})$ , ....  $(x^{(i)}, y^{(i)})$  }
  - Perform forward propagation and backpropagation using example  $(x^{(i)}, y^{(i)})$
  - (Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for  $l = 2, \dots, L$ ).
  - $\Delta^{(2)} := \Delta^{(1)} + \delta^{(1)} (a^{(1)})^T$
  - ...  
}
- compute  $\frac{\partial}{\partial \Theta_{jk}^{(m)}} J(\Theta)$ .



Andre

## Training a neural network

- 5. Use gradient checking to compare  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$  computed using backpropagation vs. using numerical estimate of gradient of  $J(\Theta)$ .
- Then disable gradient checking code.
- 6. Use gradient descent or advanced optimization method with backpropagation to try to minimize  $J(\Theta)$  as a function of parameters  $\Theta$

$$\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta) \quad \nwarrow$$

$J(\Theta)$  — non-convex.

## W6

### 1. Evaluating a Learning Algorithm

a. test set error

## Evaluating your hypothesis

Dataset:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
<hr/>	
1427	199
1380	212
1494	243

Handwritten annotations:

- A bracket on the left side of the table is labeled "70%".
- A bracket on the right side of the first seven rows is labeled "Training set".
- A bracket on the right side of the last three rows is labeled "Test set".
- To the right of the "Training set" bracket, there is an arrow pointing to a box containing the formula  $(x^{(1)}, y^{(1)})$ ,  $(x^{(2)}, y^{(2)})$ , ...,  $(x^{(m)}, y^{(m)})$ .
- To the right of the "Test set" bracket, there is an arrow pointing to a box containing the formula  $(x_{test}^{(1)}, y_{test}^{(1)})$ ,  $(x_{test}^{(2)}, y_{test}^{(2)})$ , ...,  $(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$ .
- Next to the "Test set" box, handwritten text says "m<sub>test</sub> = no. of test example" and "(x<sub>test</sub><sup>(i)</sup>, y<sub>test</sub><sup>(i)</sup>)".

Andrew

training : test  $\approx 7 : 3$

learn through the training set; compute  $J_{test}$  of the test set by the formulas below (varies between linear regression and classification):

- For linear regression:  $J_{test}(\Theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_\Theta(x_{test}^{(i)}) - y_{test}^{(i)})^2$

2. For classification ~ Misclassification error (aka 0/1 misclassification error):

$$err(h_\Theta(x), y) = \begin{cases} 1 & \text{if } h_\Theta(x) \geq 0.5 \text{ and } y = 0 \text{ or } h_\Theta(x) < 0.5 \text{ and } y = 1 \\ 0 & \text{otherwise} \end{cases}$$

This gives us a binary 0 or 1 error result based on a misclassification. The average test error for the test set is:

$$\text{Test Error} = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} err(h_\Theta(x_{test}^{(i)}), y_{test}^{(i)})$$

### b. model selection & Train/Validation/Test sets

#### Glossary #6

$d$ : the degree of the current model

$\theta^{(d)}$ : the parameters learnt by the model of degree  $d$

The problem is:

using test set:

Problem:  $J_{test}(\theta^{(5)})$  is likely to be an optimistic estimate of generalization error. I.e. our extra parameter (d = degree of polynomial) is fit to test set.

and to solve it, we split the dataset into 3 sets: Training/(Cross) Validation/Test sets instead of 2 training : validation : test  $\approx 6 : 2 : 2$

and compute the error of each set

#### steps of selecting the model using TVT sets

1. learn parameters of each model using the training set
2. compute  $J_{cv}$  using the validation set and choose the degree that minimizes  $J_{cv}$  and select the model
3. estimate generalized performance using the test set

## 2. Bias(偏差) & Variance(方差)

### a. detection

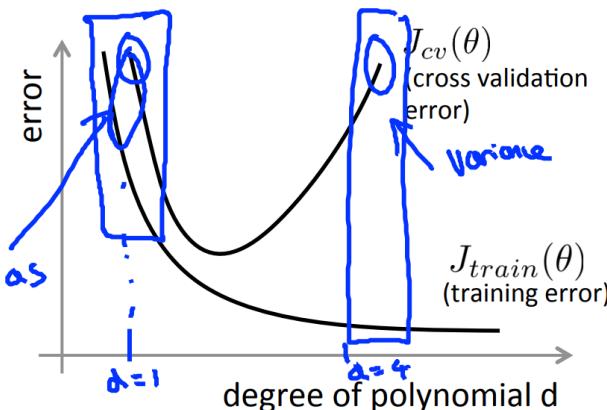
high bias: underfit:  $J_{train}(\theta)$

high variance: overfit:  $J_{cv}(\theta)$

how to detect:

## Diagnosing bias vs. variance

Suppose your learning algorithm is performing less well than you were hoping. ( $J_{cv}(\theta)$  or  $J_{test}(\theta)$  is high.) Is it a bias problem or a variance problem?



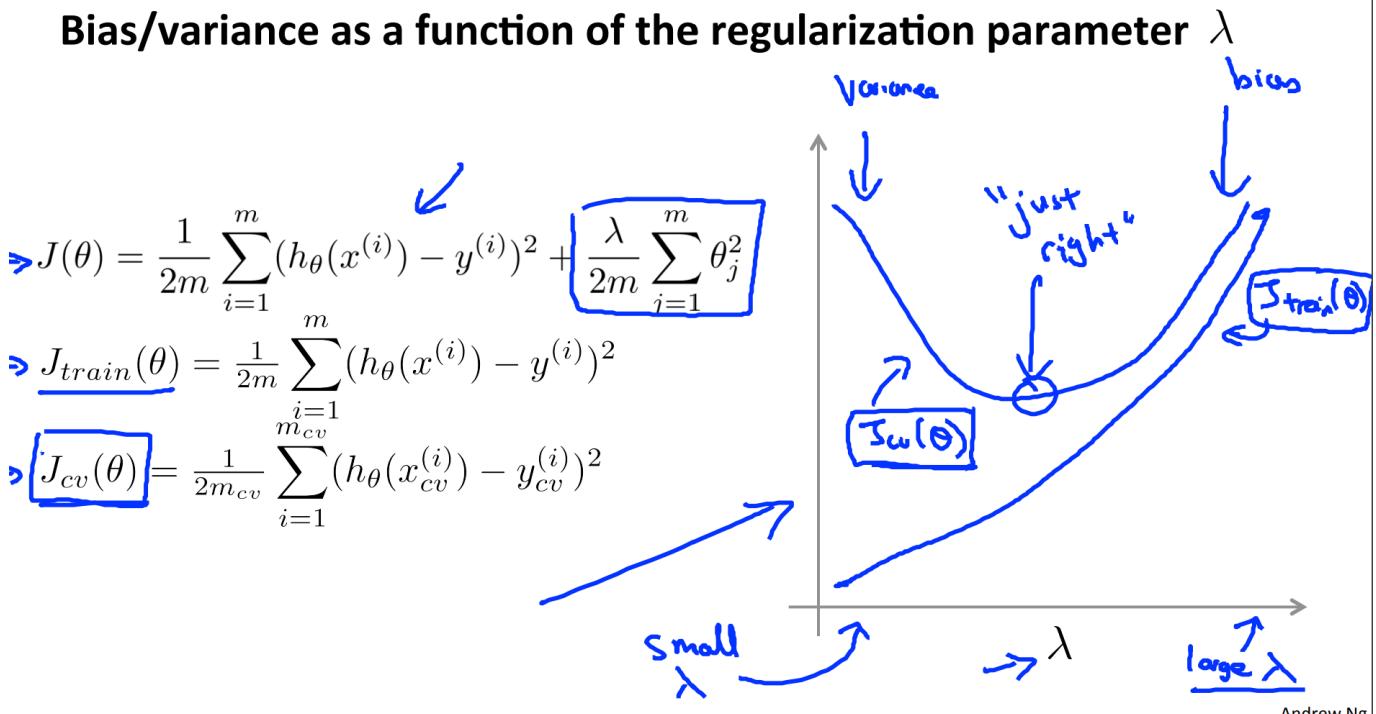
Bias (underfit):  
 $\rightarrow J_{train}(\theta)$  will be high  
 $J_{cv}(\theta) \approx J_{train}(\theta)$

Variance (overfit):  
 $\rightarrow J_{train}(\theta)$  will be low  
 $J_{cv}(\theta) \gg J_{train}(\theta)$   
 $\gg$

Andrew Ng

### b. regularization and bias/variance

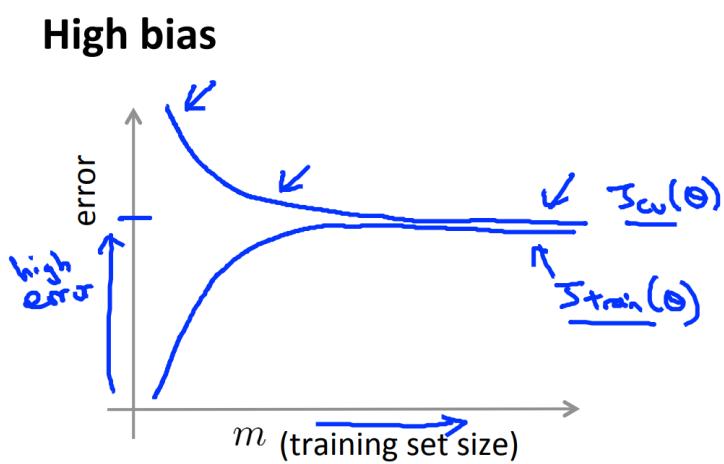
how to choose the regularization parameter  $\lambda$ :



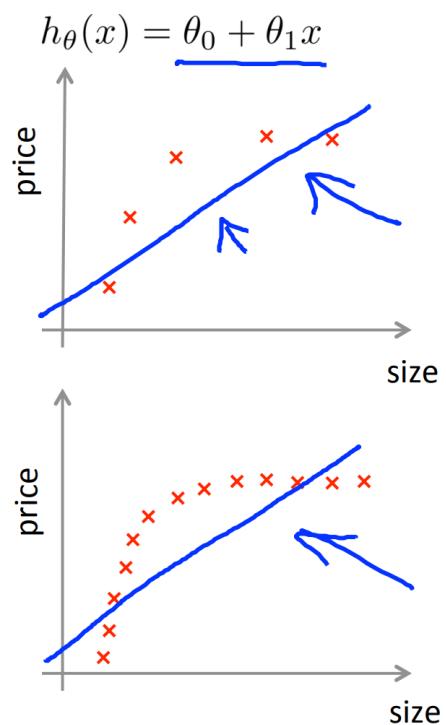
note that all the definition of  $J_{train}$ ,  $J_{cv}$  and  $J_{test}$  are without the regularization term

### c. learning curves

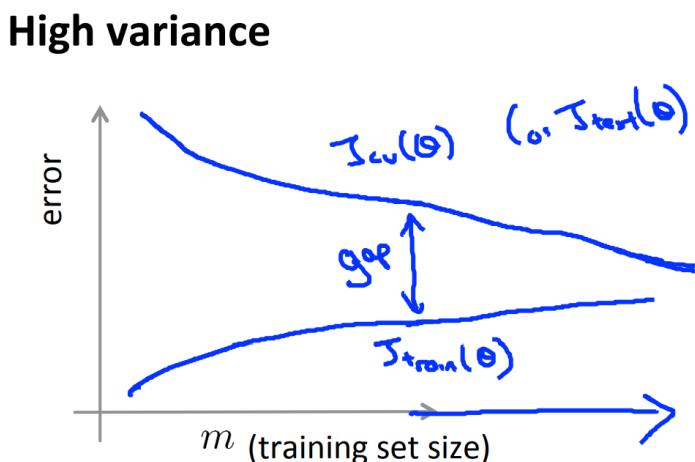
high bias



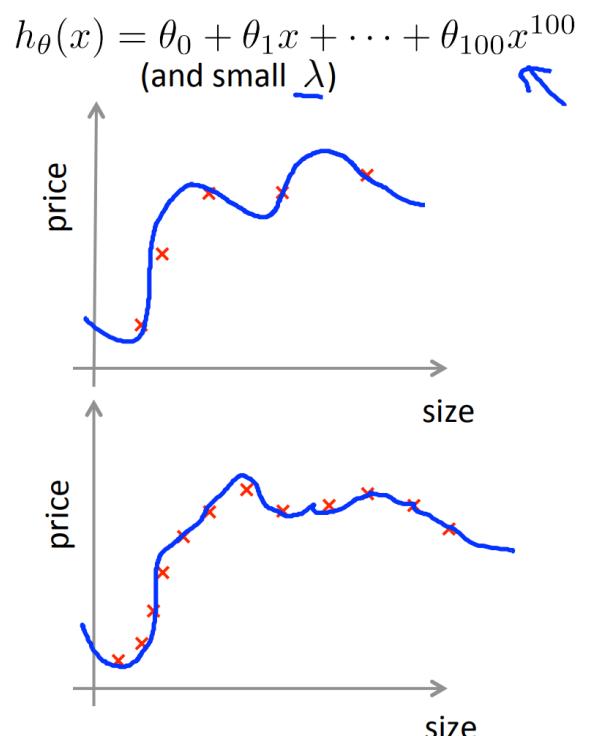
If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.



high variance



If a learning algorithm is suffering from high variance, getting more training data is likely to help. ↙



d. review: what to do next?

## Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

- Get more training examples → fixes high variance
- Try smaller sets of features → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features ( $x_1^2, x_2^2, x_1x_2$ , etc) → fixes high bias.
- Try decreasing  $\lambda$  → fixes high bias
- Try increasing  $\lambda$  → fixes high variance

## 3. ML System Design

a. example: build a spam classifier

prioritizing what to do next

### Building a spam classifier

Supervised learning.  $x$  = features of email.  $y$  = spam (1) or not spam (0).

Features  $x$ : Choose 100 words indicative of spam/not spam.

E.g. deal, buy, discount, andrew, now, ...

$$x = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \begin{array}{l} \text{andrew} \\ \text{buy} \\ \text{deal} \\ \text{discount} \\ \vdots \\ \text{now} \end{array} x \in \mathbb{R}^{100}$$

$$x_j = \begin{cases} 1 & \text{if word } j \text{ appear} \\ 0 & \text{otherwise.} \end{cases}$$

From: cheapsales@buystufffromme.com  
To: ang@cs.stanford.edu  
Subject: Buy now!

Deal of the week! Buy now!

Note: In practice, take most frequently occurring  $n$  words (10,000 to 50,000) in training set, rather than manually pick 100 words.

## Building a spam classifier

How to spend your time to make it have low error?

- Collect lots of data
  - E.g. "honeypot" project.
- Develop sophisticated features based on email routing information (from email header).
- Develop sophisticated features for message body, e.g. should "discount" and "discounts" be treated as the same word? How about "deal" and "Dealer"? Features about punctuation?
- Develop sophisticated algorithm to detect misspellings (e.g. m0rtgage, med1cine, w4tches.)

### error analysis

recommended approach when building a machine learning system

1. implement a simple algorithm
2. cross validation
3. list possible approaches to improve it

recommended approach when doing error analysis

1. find the specific type of the mistaken test cases

e.g.(spam classifier)

Pharma: 12  
Replica/fake: 4  
Steal passwords: 53  
Other: 31

→ Deliberate misspellings: 5  
(m0rgage, med1cine, etc.)  
→ Unusual email routing: 16  
→ Unusual (spamming) punctuation: 32

2. list possible cues to improve the performance

doing numerical evaluation is important

It is very important to get error results as a single, numerical value. Otherwise it is difficult to assess your algorithm's performance. I

### b. handling skewed data(偏斜数据)

definition

the number of data in one class is much greater than that in another class (e.g. only 0.05% of humanities have cancer) and causes overfit

precision(查准率) and recall(召回率)

## Precision/Recall

$y = 1$  in presence of rare class that we want to detect

		Actual class	
		1	0
Predicted class	1	True positive False positive	
	0	False negative True negative	

$$y=0 \\ \text{recall} = 0$$

### Precision

(Of all patients where we predicted  $y = 1$ , what fraction actually has cancer?)

$$\frac{\text{True positives}}{\# \text{predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$$

### Recall

(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

$$\frac{\text{True positives}}{\# \text{actual positives}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}}$$

tradeoff between precision and recall

## Trading off precision and recall

→ Logistic regression:  $0 \leq h_\theta(x) \leq 1$

Predict 1 if  $h_\theta(x) \geq 0.3$  ~~0.7~~ ~~0.9~~ ~~0.3~~ ←

Predict 0 if  $h_\theta(x) < 0.3$  ~~0.7~~ ~~0.9~~ ~~0.3~~

→ Suppose we want to predict  $y = 1$  (cancer) only if very confident.

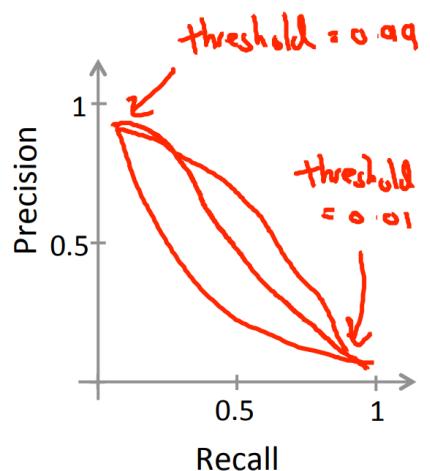
→ Higher precision, lower recall

→ Suppose we want to avoid missing too many cases of cancer (avoid false negatives).

→ Higher recall, lower precision.

$$\rightarrow \text{precision} = \frac{\text{true positives}}{\text{no. of predicted positive}}$$

$$\rightarrow \text{recall} = \frac{\text{true positives}}{\text{no. of actual positive}}$$



More generally: Predict 1 if  $h_\theta(x) \geq \text{threshold}$  ←

## F score

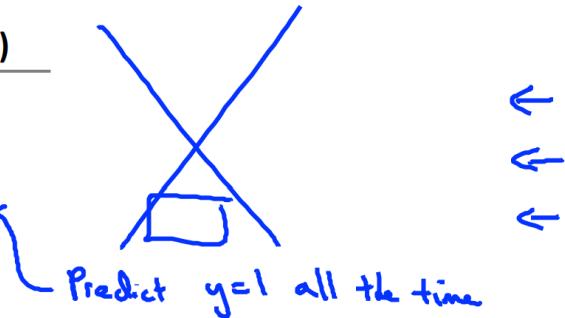
### $F_1$ Score (F score)

How to compare precision/recall numbers?

	Precision(P)	Recall (R)
Algorithm 1	0.5	0.4
Algorithm 2	0.7	0.1
Algorithm 3	0.02	1.0

Average:  ~~$\frac{P+R}{2}$~~

$F_1$  Score:  $2 \frac{PR}{P+R}$



$$P=0 \text{ or } R=0 \Rightarrow F\text{-score} = 0$$

$$P=1 \text{ and } R=1 \Rightarrow F\text{-score} = 1$$

note: choose a proper threshold using cross validation set to compute the F score

c. collecting data

**"It's not who has the best algorithm that wins.**

**It's who has the most data."**

rationale

1. sufficient to predict/classify

## Large data rationale

→ Assume feature  $x \in \mathbb{R}^{n+1}$  has sufficient information to predict  $y$  accurately.

Example: For breakfast I ate two eggs.

Counterexample: Predict housing price from only size (feet<sup>2</sup>) and no other features.

Useful test: Given the input  $x$ , can a human expert confidently predict  $y$ ?

## 2. parameters of the learning algorithm

### Large data rationale

→ Use a learning algorithm with many parameters (e.g. logistic regression/linear regression with many features; neural network with many hidden units). low bias algorithms.

→  $J_{\text{train}}(\theta)$  will be small.

Use a very large training set (unlikely to overfit) low variance

→  $J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$

→  $J_{\text{test}}(\theta)$  will be small

a large scale of data can help in both situations

## W7

### 1. SVM(Support Vector Machine)

#### a. large margin classification(大间距分类)

## an alternative viewpoint of logistic regression

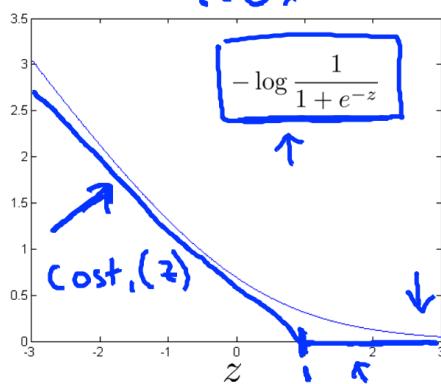
a new definition for the cost function of logistic regression  $\text{cost}_1(z)$  ( $y=1$ ) and  $\text{cost}_0(z)$  ( $y=0$ ) as the blue line shows :

### Alternative view of logistic regression

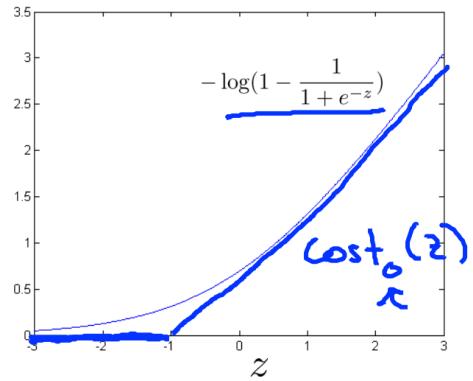
Cost of example:  $-(y \log h_\theta(x) + (1 - y) \log(1 - h_\theta(x)))$

$$= \left[ -y \log \frac{1}{1 + e^{-\theta^T x}} \right] - \left[ (1 - y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right) \right]$$

If  $y = 1$  (want  $\theta^T x \gg 0$ ):



If  $y = 0$  (want  $\theta^T x \ll 0$ ):



definition of SVM:

$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

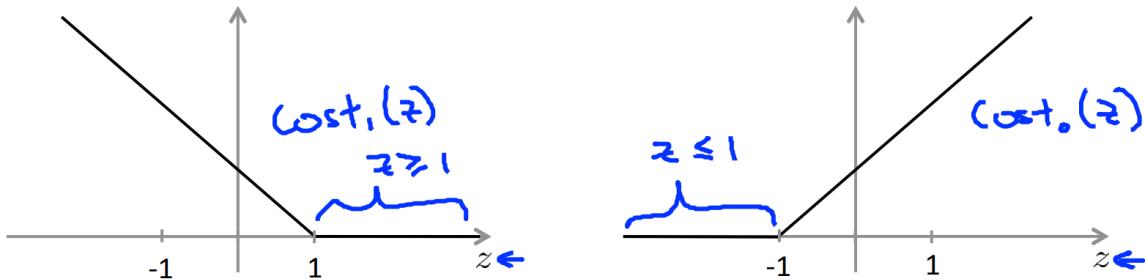
$m$  is removed from the cost function and  $C = \frac{1}{\lambda}$ , and  $h_\theta(x)$  is replaced by  $\text{cost}_{0,1}(\theta^T x^{(i)})$

### intuition

We want  $\theta^T x \geq 1$  instead of 0 if  $y = 1$  otherwise we want  $\theta^T x \leq -1$  instead of 0 because we want to minimize the cost function by setting the first sum 0:

## Support Vector Machine

$$\rightarrow \min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \underline{\text{cost}_1(\theta^T x^{(i)})} + (1 - y^{(i)}) \underline{\text{cost}_0(\theta^T x^{(i)})} \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$



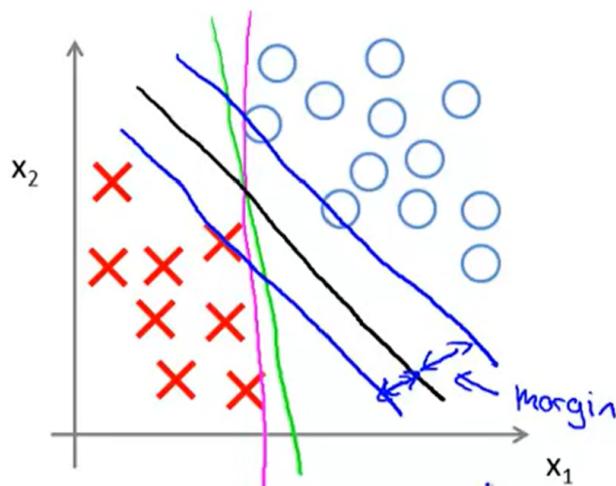
→ If  $y = 1$ , we want  $\underline{\theta^T x \geq 1}$  (not just  $\geq 0$ )

→ If  $y = 0$ , we want  $\underline{\theta^T x \leq -1}$  (not just  $< 0$ )

$$C = 100,000$$

SVM choose the best decision boundary with the largest margin (the distance between the boundary and the nearest (to the boundary) points of each class):

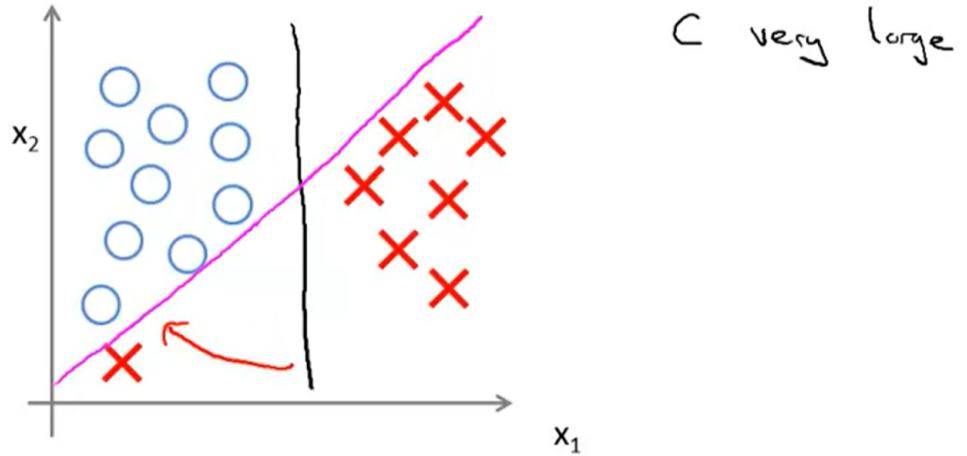
### SVM Decision Boundary: Linearly separable case



blue lines are margins and SVM will choose the black one

but in datasets with outliers like this:

## Large margin classifier in presence of outliers



if  $C$  is very large the SVM will choose the magenta line as the boundary; otherwise it will remain choosing the black one.

### math behind LMC

the SVM tries to solve:

$$\begin{aligned}
 \rightarrow \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 &= \frac{1}{2} \|\theta\|^2 \leftarrow \\
 \text{s.t. } p^{(i)} \cdot \|\theta\| &\geq 1 \quad \text{if } y^{(i)} = 1 \\
 p^{(i)} \cdot \|\theta\| &\leq -1 \quad \text{if } y^{(i)} = -1 \\
 \text{where } p^{(i)} \text{ is the projection of } r^{(i)} \text{ on }
 \end{aligned}$$

i.e. (simplified version)

## SVM Decision Boundary

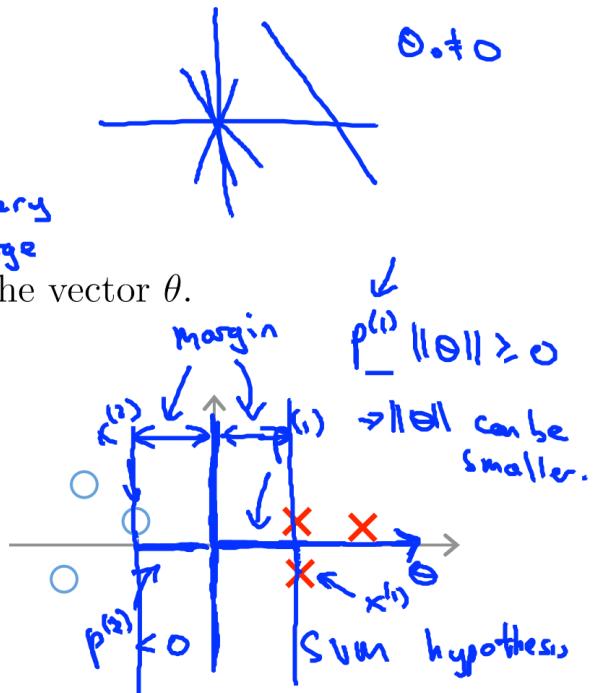
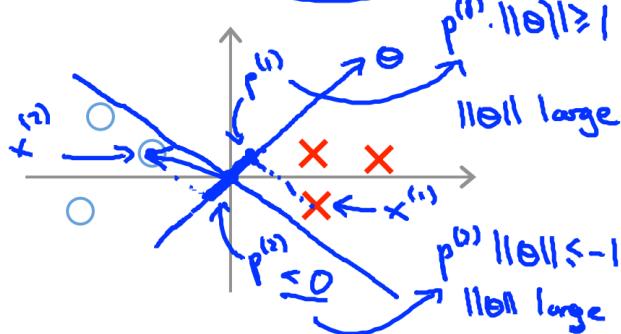
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2 \leftarrow$$

$$\text{s.t. } \begin{cases} p^{(i)} \cdot \|\theta\| \geq 1 & \text{if } y^{(i)} = 1 \\ p^{(i)} \cdot \|\theta\| \leq -1 & \text{if } y^{(i)} = -1 \end{cases}$$

*C very large*

where  $p^{(i)}$  is the projection of  $x^{(i)}$  onto the vector  $\theta$ .

Simplification:  $\theta_0 = 0$



Andrew Ng

the reason why SVM doesn't choose the boundary on the left (seems less 'natural' than that on the left) is:  $|p^{(i)}|$  is small so that  $\|\theta\|$  shall be large but that is **contradictory to the first minimizing condition**

here is why  $\theta$  is vertical to the boundary:

<http://cs229.stanford.edu/materials/smo.pdf>

**Q3) How do we know that the theta is perpendicular to the decision boundary?**

(thanks to Mentor Chirag for this derivation)

We know that  $\theta' * x = 0$  for any  $x$  on the boundary since the boundary is where sigmoid = 0.5.

$$\frac{1}{1+e^{-z}} = \frac{1}{2} \Rightarrow e^{-z} = 1 \Rightarrow -z = 0 \Rightarrow \theta' * x = 0$$

So pick two random points on the boundary a and b, then

$$\theta' * a = 0 \text{ and } \theta' * b = 0$$

$$\Rightarrow \theta' * (a - b) = 0 \Rightarrow \theta \cdot (a - b) = 0$$

and we know that when the dot product of two vectors is 0, the angle between them is 90 degrees. And since the vector (a-b) is on the decision boundary,  $\theta$  is perpendicular to the decision boundary.

## b. kernels(adapting SVM to non-linear classification problems)

Gaussian kernels

### Kernels and Similarity

$$f_1 = \text{similarity}(x, \underline{l}^{(1)}) = \exp\left(-\frac{\|x - \underline{l}^{(1)}\|^2}{2\sigma^2}\right)$$

$\downarrow \downarrow$

If  $x \approx \underline{l}^{(1)}$  :

$$f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$$

$\underline{l}^{(1)} \rightarrow f_1$   
 $\underline{l}^{(2)} \rightarrow f_2$   
 $\underline{l}^{(3)} \rightarrow f_3$   
 $\uparrow \quad \uparrow \quad \uparrow$   
 $x$

If  $x$  if far from  $\underline{l}^{(1)}$  :

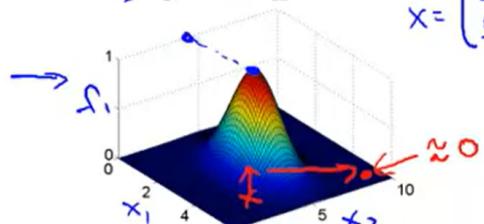
$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0.$$

example

### Example:

$$\rightarrow \underline{l}^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - \underline{l}^{(1)}\|^2}{2\sigma^2}\right)$$

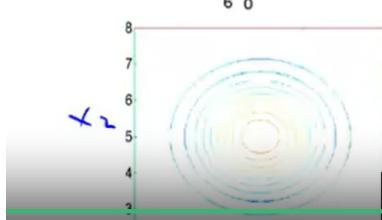
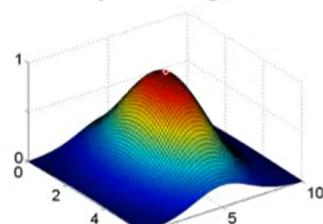
$$\rightarrow \sigma^2 = 1$$



$$x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

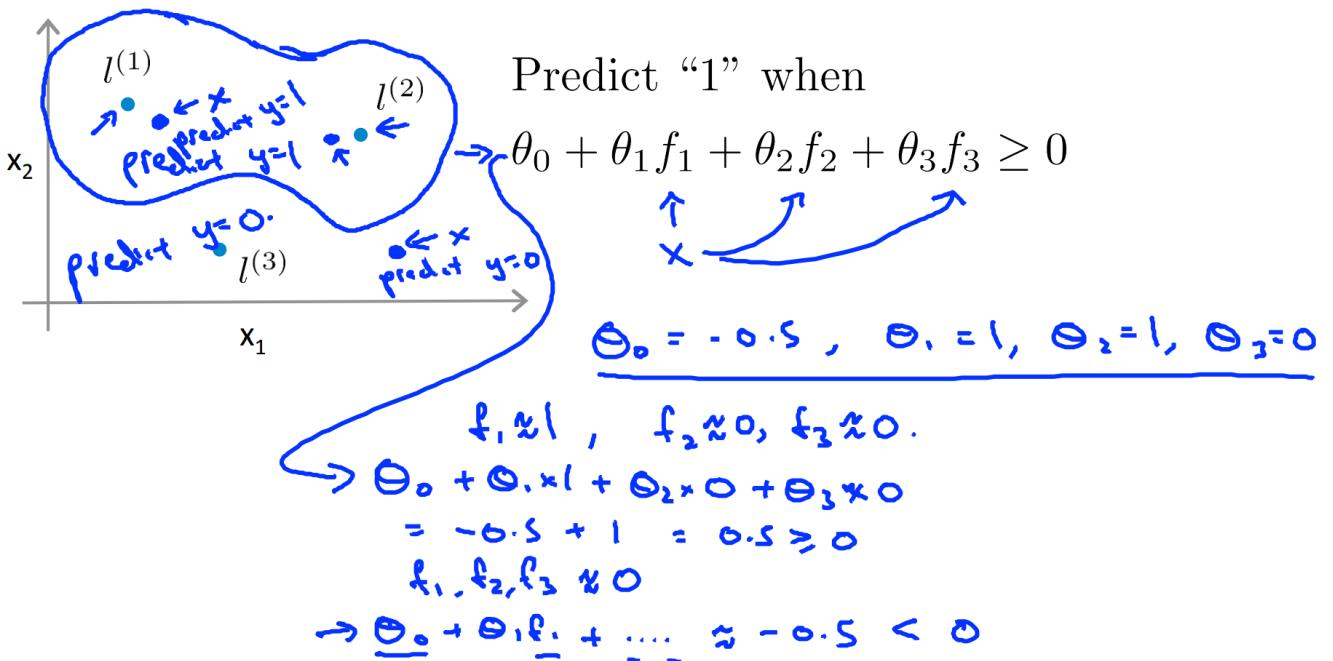
$$\sigma^2 = 0.5$$

$$\sigma^2 = 3$$



move away from, you know! So,

how to predict



choose the landmarks

just choose  $l^{(i)} = x^{(i)}$

### SVM with Kernels

- Given  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$ ,
- choose  $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$ .

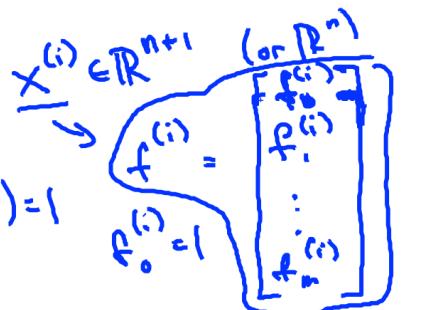
Given example  $x$ :

$$\begin{aligned} \rightarrow f_1 &= \text{similarity}(x, l^{(1)}) \\ \rightarrow f_2 &= \text{similarity}(x, l^{(2)}) \\ \vdots \\ \rightarrow f_m &= \text{similarity}(x, l^{(m)}) \end{aligned}$$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1$$

For training example  $(x^{(i)}, y^{(i)})$ :

$$\begin{aligned} x^{(i)} \rightarrow f_1^{(i)} &= \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} &= \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} &= \text{sim}(x^{(i)}, l^{(m)}) = \exp\left(-\frac{\|x^{(i)} - l^{(m)}\|^2}{2\sigma^2}\right) = 1 \end{aligned}$$



combining SVM with kernels

just modify  $\theta^T x^{(i)}$  to  $\theta^T f^{(i)}$  like this:

## SVM with Kernels

Hypothesis: Given  $x$ , compute features  $f \in \mathbb{R}^{m+1}$

$$\rightarrow \text{Predict "y=1" if } \theta^T f \geq 0$$

$$\theta \in \mathbb{R}^{n+1}$$

$$\theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$$

Training:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$\rightarrow \theta_0$

$\rightarrow \sum_j \theta_j^2 = \theta^T \theta$  (ignoring  $\theta_0$ )

$\rightarrow \theta^T M \theta$  ( $M = I_{n+1}$ )

Andrew Ng

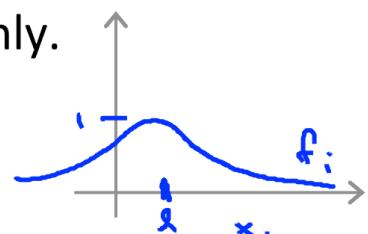
parameters

## SVM parameters:

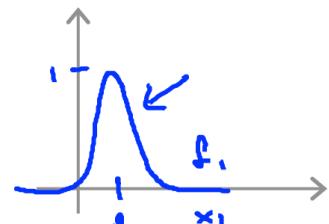
$C$  ( $= \frac{1}{\lambda}$ ).  $\rightarrow$  Large  $C$ : Lower bias, high variance. (small  $\lambda$ )  
 $\rightarrow$  Small  $C$ : Higher bias, low variance. (large  $\lambda$ )

$\sigma^2$  Large  $\sigma^2$ : Features  $f_i$  vary more smoothly.  
 $\rightarrow$  Higher bias, lower variance.

$$\exp\left(-\frac{\|x - \mu^{(i)}\|^2}{2\sigma^2}\right)$$



Small  $\sigma^2$ : Features  $f_i$  vary less smoothly.  
Lower bias, higher variance.



## c. SVM in practice

in programming we don't need to write the SVM functions ourselves but we have to specify two things below:

- the choice of parameter  $C$

- the choice of the kernel(similarity function)
- (when using Gaussian kernel) implement the kernel function

### choice of the kernel

1. no kernel(i.e. linear kernel): gives a standard linear classifier i.e. predict  $y = 1$  if  $\theta^T x \geq 0$ .  
one situation to use it is there are many features but few training data
2. Gaussian kernel: few features and many training data, **do perform feature scaling!**
3. other choices:

### Other choices of kernel

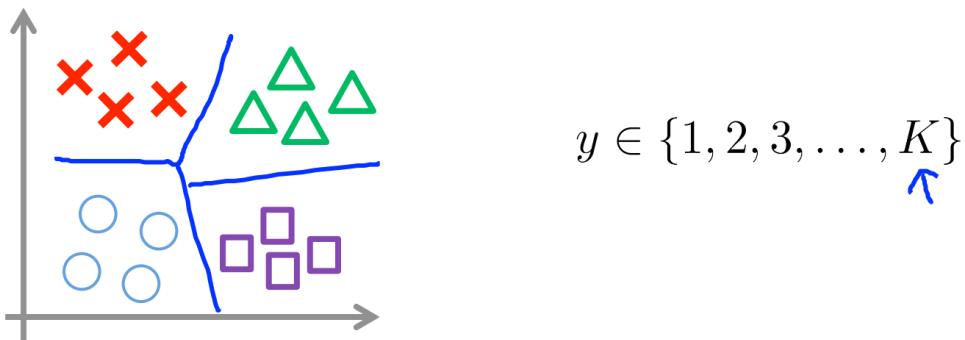
Note: Not all similarity functions  $\text{similarity}(x, l)$  make valid kernels.

→ (Need to satisfy technical condition called “Mercer’s Theorem” to make sure SVM packages’ optimizations run correctly, and do not diverge).

- Many off-the-shelf kernels available:
- Polynomial kernel:  $k(x, l) = \underbrace{(x^T l)}_{\text{constant}} + \underbrace{(x^T l)^2}_{\text{degree 2}}, \underbrace{(x^T l)^3}_{\text{degree 3}}, \underbrace{(x^T l + 1)^1}_{\text{degree 1}}, \underbrace{(x^T l + 5)^0}_{\text{degree 0}}$
  - More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...  $\text{sim}(x, l)$

### multi-classification using SVM

## Multi-class classification



Many SVM packages already have built-in multi-class classification functionality.

- Otherwise, use one-vs.-all method. (Train  $K$  SVMs, one to distinguish  $y = i$  from the rest, for  $i = 1, 2, \dots, K$ ), get  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$   
Pick class  $i$  with largest  $(\theta^{(i)})^T x$

$$y=1 \quad y=2 \quad \dots \quad \theta^{(k)}$$

### logistic regression vs. SVM

$n$  is the number of features;  $m$  is the number of training examples

- $n$  is large relative to  $m$ , use logistic regression or SVM without kernels
- $n$  small while  $m$  intermediate: SVM with Gaussian kernel
- $n$  small while  $m$  large: add more features, use logistic regression or SVM without kernels

### Logistic regression vs. SVMs

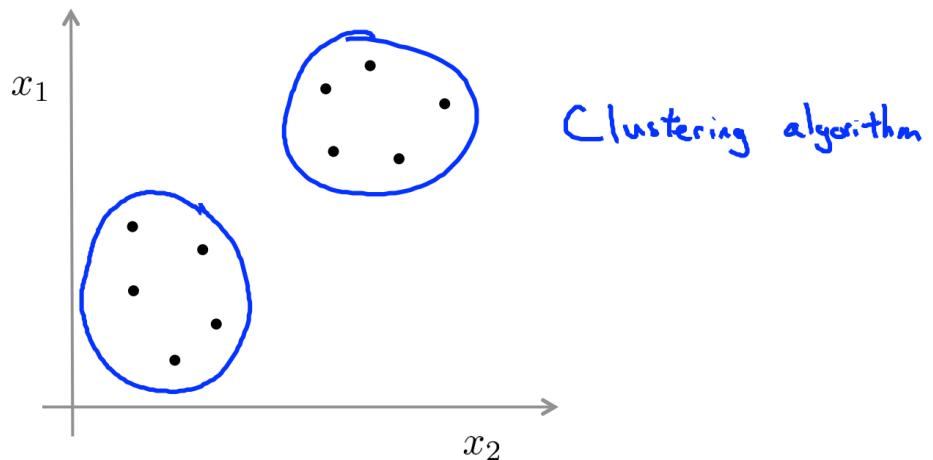
- $n$  = number of features ( $x \in \mathbb{R}^{n+1}$ ),  $m$  = number of training examples
- If  $n$  is large (relative to  $m$ ): (e.g.  $n \geq m$ ,  $n = 10,000$ ,  $m = 10 \dots 1000$ )
- Use logistic regression, or SVM without a kernel ("linear kernel")
- If  $n$  is small,  $m$  is intermediate: ( $n = 1-1000$ ,  $m = 10 - 10,000$ ) ←
  - Use SVM with Gaussian kernel
- If  $n$  is small,  $m$  is large: ( $n = 1-1000$ ,  $m = 50,000+$ )
  - Create/add more features, then use logistic regression or SVM without a kernel
- Neural network likely to work well for most of these settings, but may be slower to train.



## 1. Unsupervised learning

what it looks like:

### Unsupervised learning



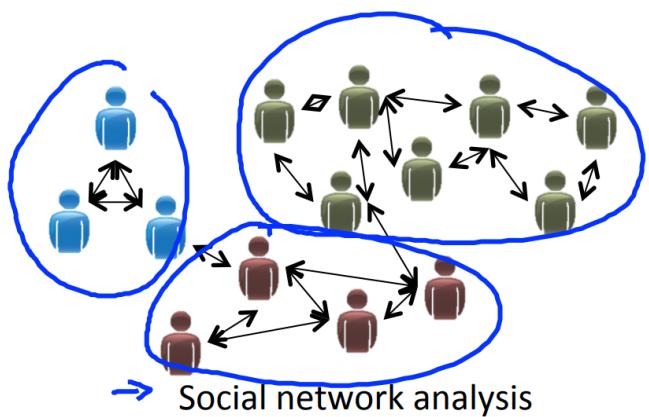
Training set:  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$  ←

where it is applied:

### Applications of clustering



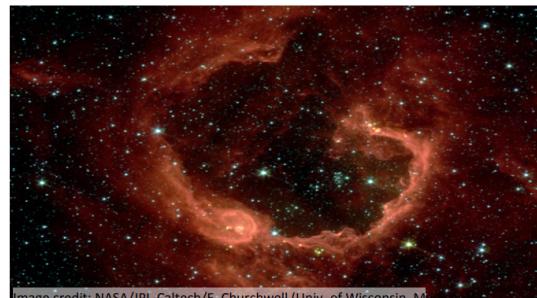
→ Market segmentation



→ Social network analysis



Organize computing clusters

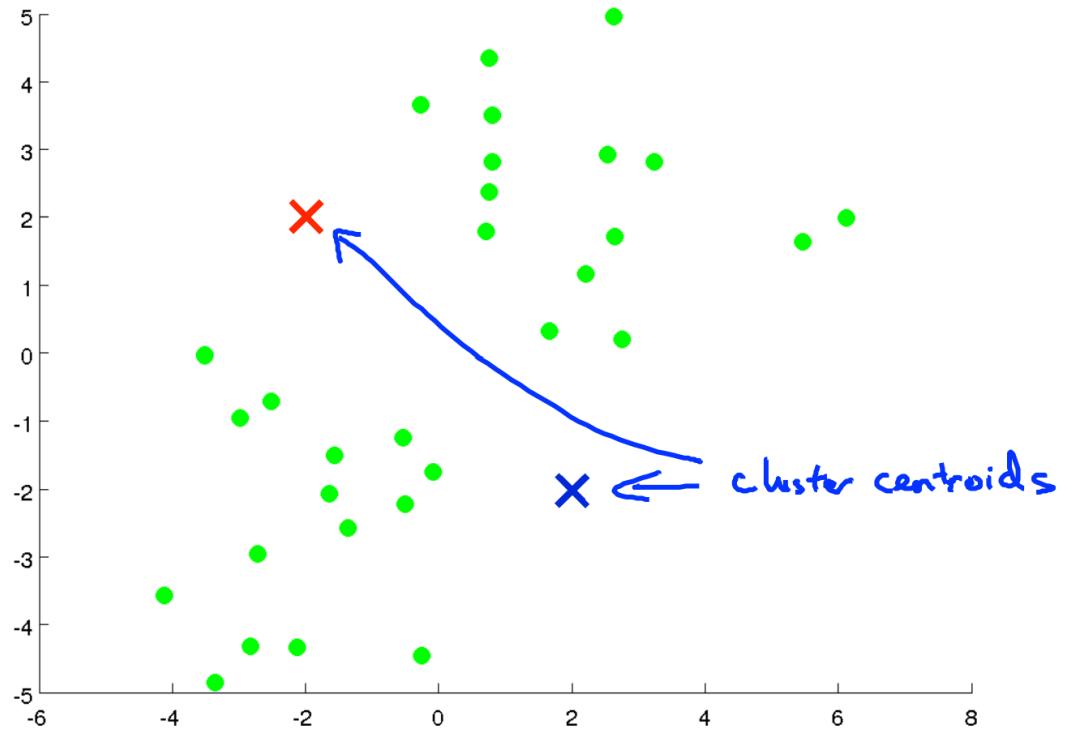


→ Astronomical data analysis

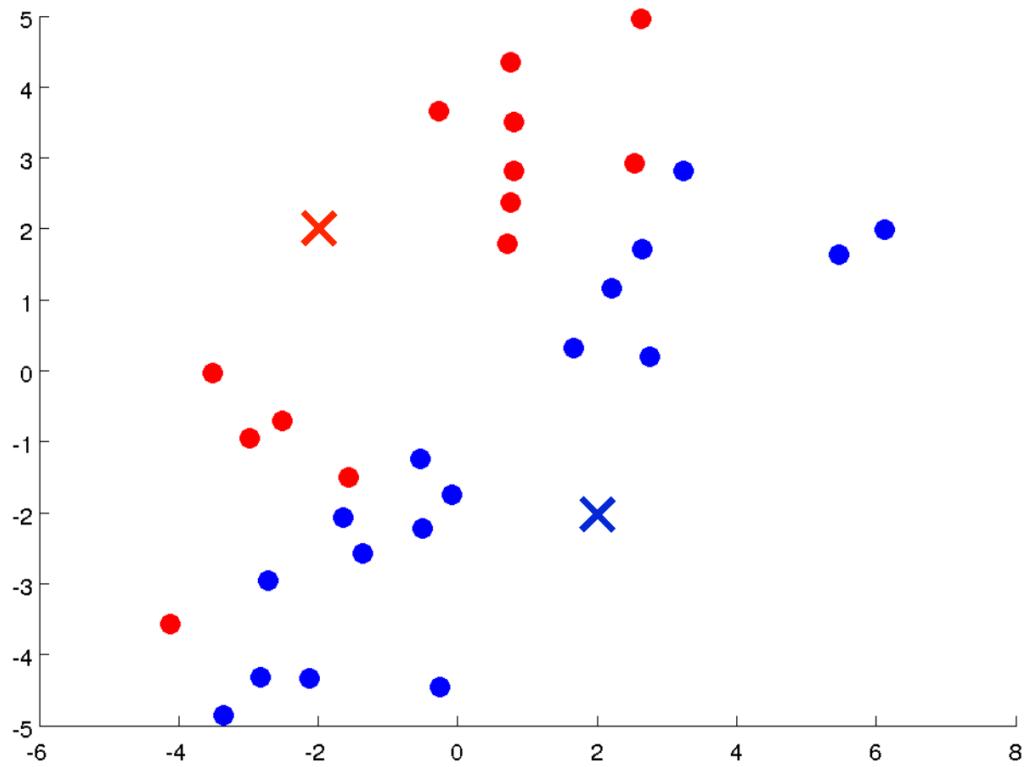
## a. K-Means algorithm

steps

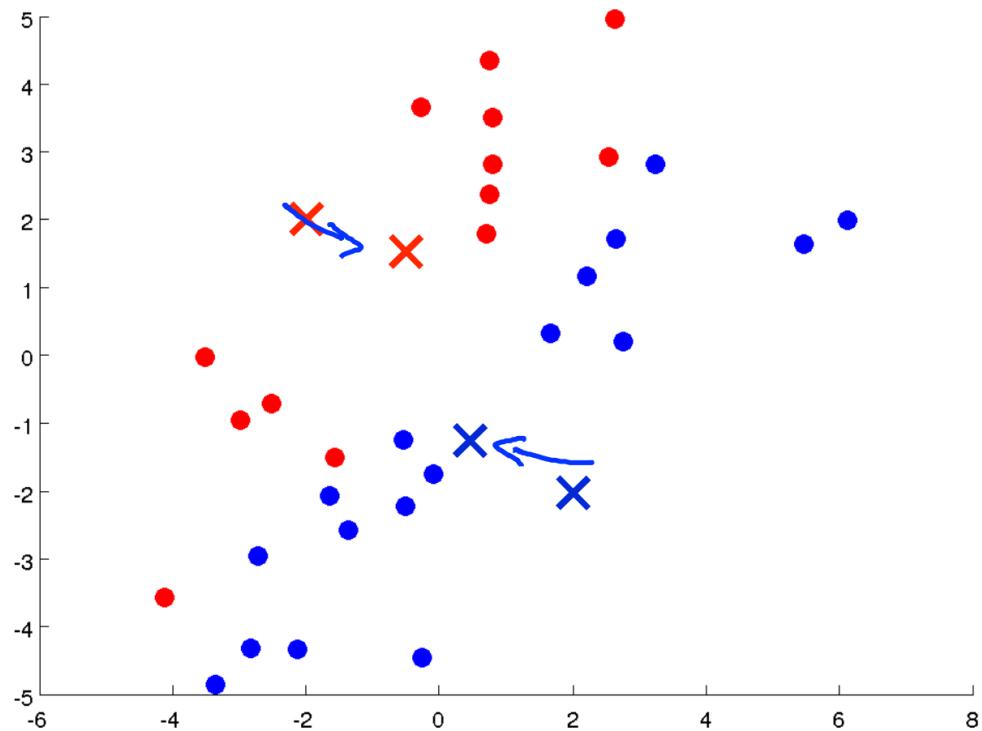
1. randomly choose cluster centroids(聚类中心)



2. assign the training data into clusters by the distance between the data and the cluster centroids



3. move the cluster centroid to the 'average' of each cluster



4. repeat until the centroid remains unchange

more generally, the algorithm goes like this:

## Input:

- $K$  (number of clusters) 
- Training set  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$  

$x^{(i)} \in \mathbb{R}^n$  (drop  $x_0 = 1$  convention)

## K-means algorithm

$$\mu_1 \quad \mu_2$$


Randomly initialize  $K$  cluster centroids  $\underline{\mu}_1, \underline{\mu}_2, \dots, \underline{\mu}_K \in \mathbb{R}^n$

Repeat {

*Cluster assignment step*

for  $i = 1$  to  $m$   
 $c^{(i)}$  := index (from 1 to  $K$ ) of cluster centroid  
closest to  $x^{(i)}$

min  $\|x^{(i)} - \underline{\mu}_k\|^2$

for  $k = 1$  to  $K$   
 $\rightarrow \underline{\mu}_k$  := average (mean) of points assigned to cluster  $k$

$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$        $\rightarrow c^{(1)}=2, c^{(2)}=2, c^{(3)}=2, c^{(4)}=2$

$\underline{\mu}_2 = \frac{1}{4} \left[ \underline{x}^{(1)} + \underline{x}^{(2)} + \underline{x}^{(3)} + \underline{x}^{(4)} \right] \in \mathbb{R}^n$

Note that if no points were assigned to a cluster, a common way to handle this is eliminating the cluster and modify  $K$  to  $K - 1$ . If you really need  $K$  clusters, just reinitialize the cluster centroid.

optimization objective

Glossary #7

$c^{(i)}$ : index of current cluster to which example  $x^{(i)}$  is assigned

$K$ : number of clusters

$\mu_k$ : cluster centroid  $k$

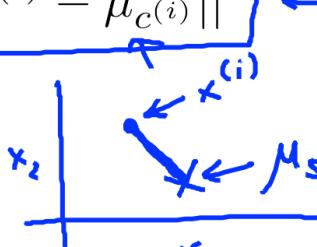
$\mu_{c^{(i)}}$ : centroid of cluster to which example  $x^{(i)}$  is assigned

the optimization objective(**distortion cost function(失真代价函数)**):

$$\rightarrow J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$\min_{\substack{c^{(1)}, \dots, c^{(m)}, \\ \mu_1, \dots, \mu_K}} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

*Distortion*



the cluster assignment step is actually optimizing  $c^{(i)}$

the cluster centroid step is actually optimizing  $\mu_k$

function  $J$  won't sometimes increase!

**random initialization**

principles:

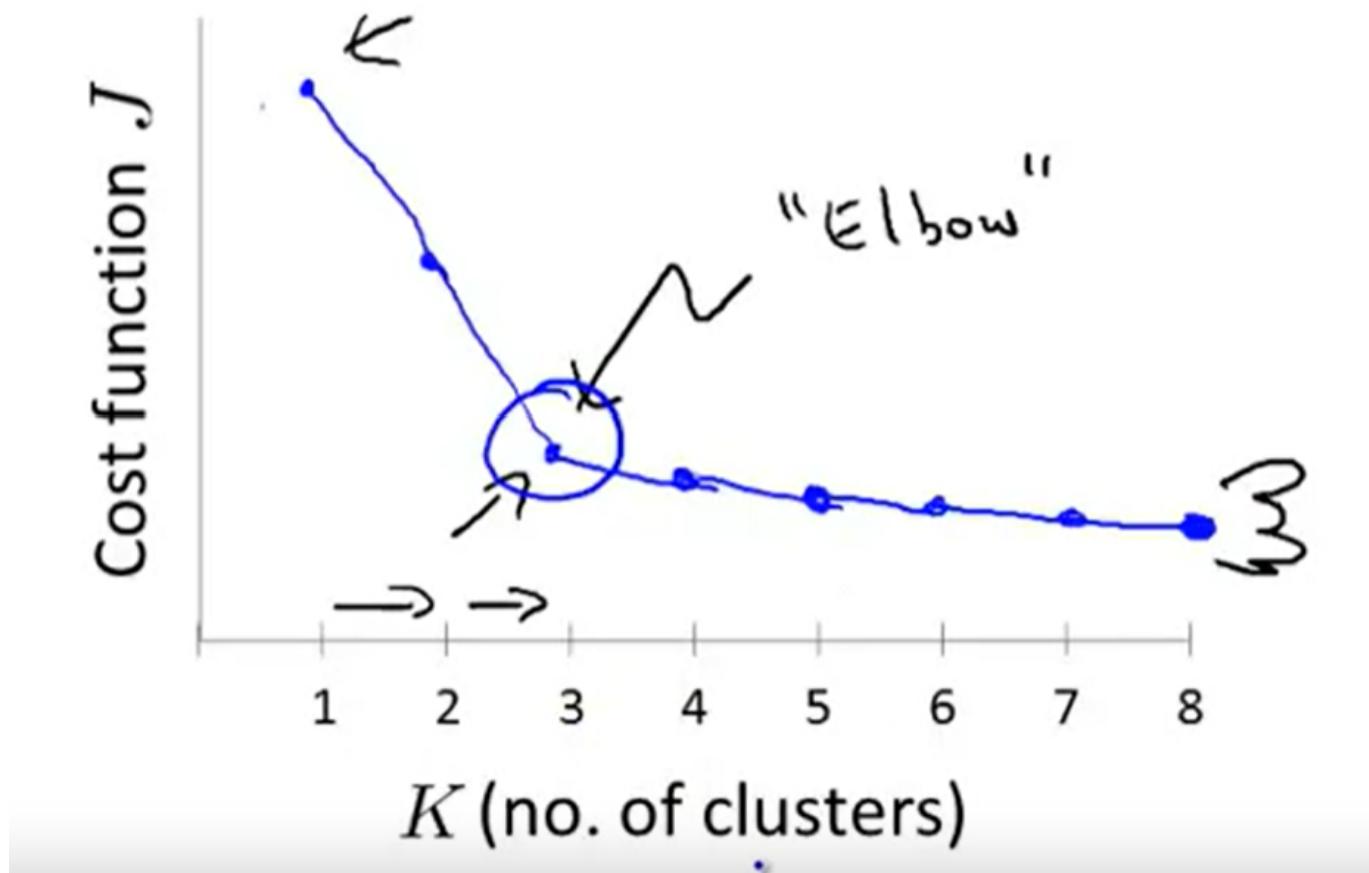


to avoid local optima, we can run K-means many times and select the one with the lowest error as the result.

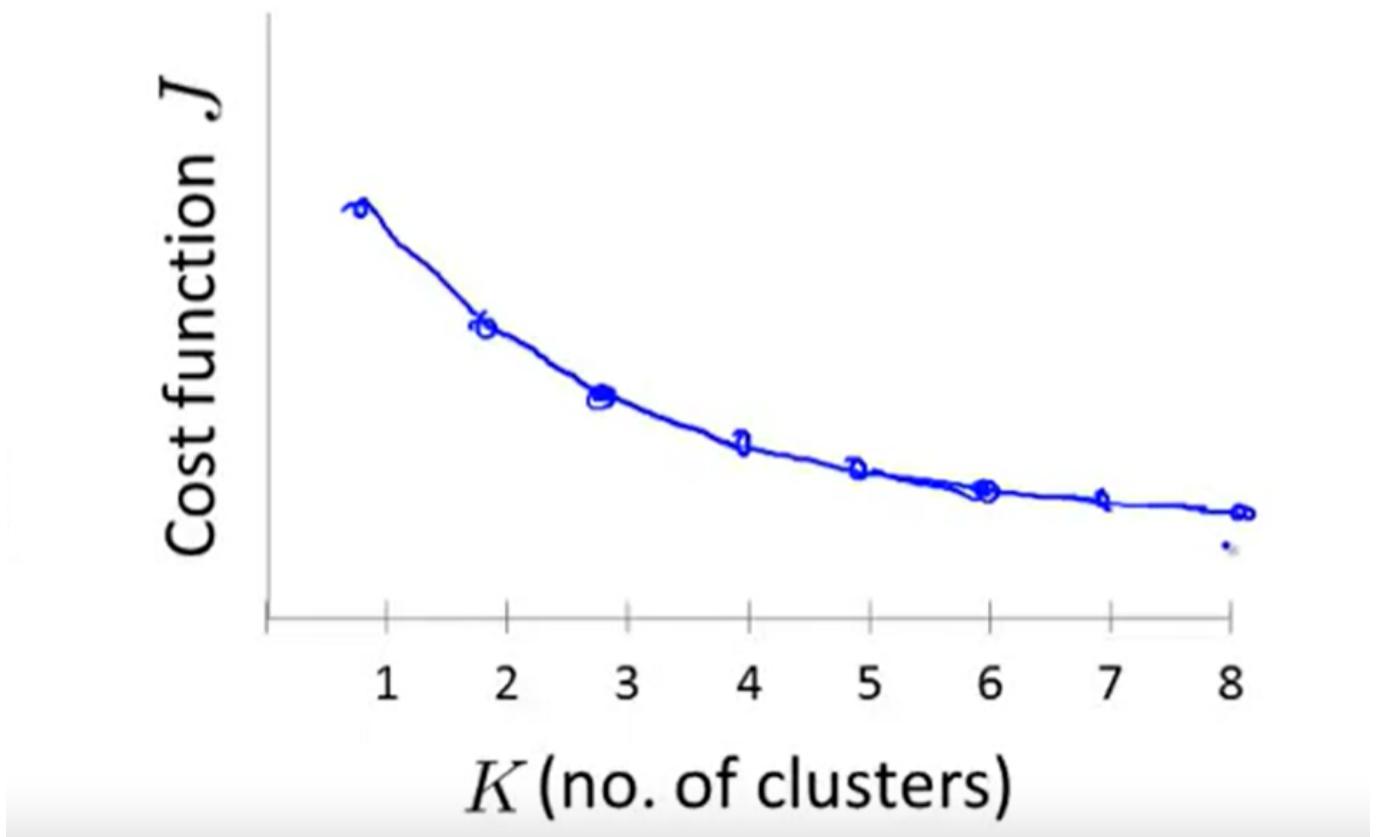
**choosing  $K$**

*Elbow method:*

for  $k=1$  to  $m$ (exclusive) run K-Means and compute  $J$ , correct  $K$  shall be the value at the "elbow":



but if the plot is like this:



It is hard to find the "elbow"

## 2. Dimensionality Reduction

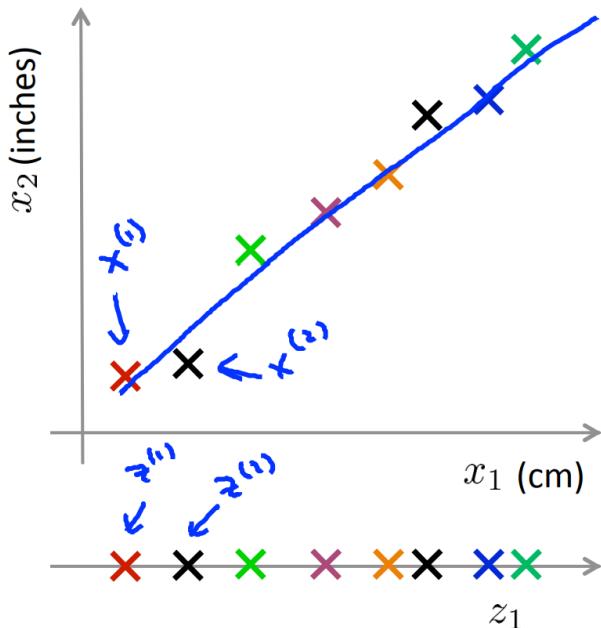
a. motivation

data compression

make the algorithm run faster

(2->1)

# Data Compression



Reduce data from  
2D to 1D

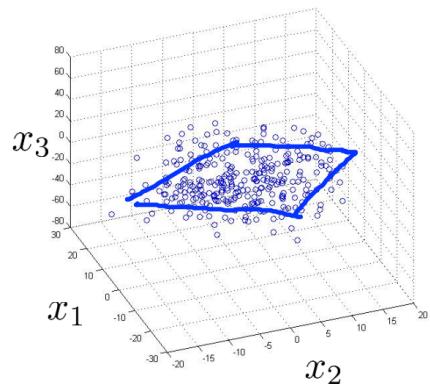
$$\begin{aligned} x^{(1)} \in \mathbb{R}^2 &\rightarrow z^{(1)} \in \mathbb{R} \\ x^{(2)} \in \mathbb{R}^2 &\rightarrow z^{(2)} \in \mathbb{R} \\ &\vdots \\ x^{(m)} \in \mathbb{R}^2 &\rightarrow z^{(m)} \in \mathbb{R} \end{aligned}$$

(3->2)

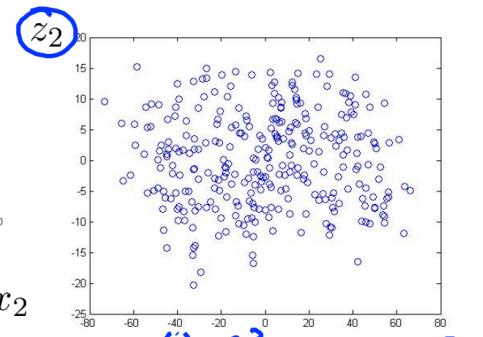
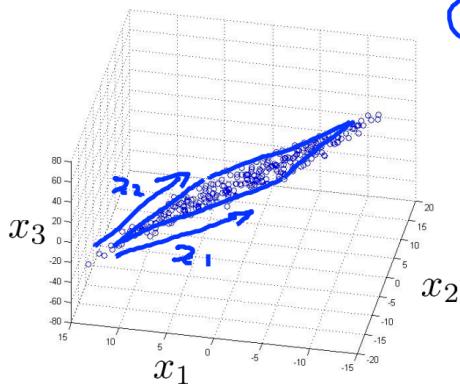
# Data Compression

$1000 \rightarrow 100$

Reduce data from 3D to 2D



$$x^{(i)} \in \mathbb{R}^3$$



$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad z^{(i)} \in \mathbb{R}^2$$

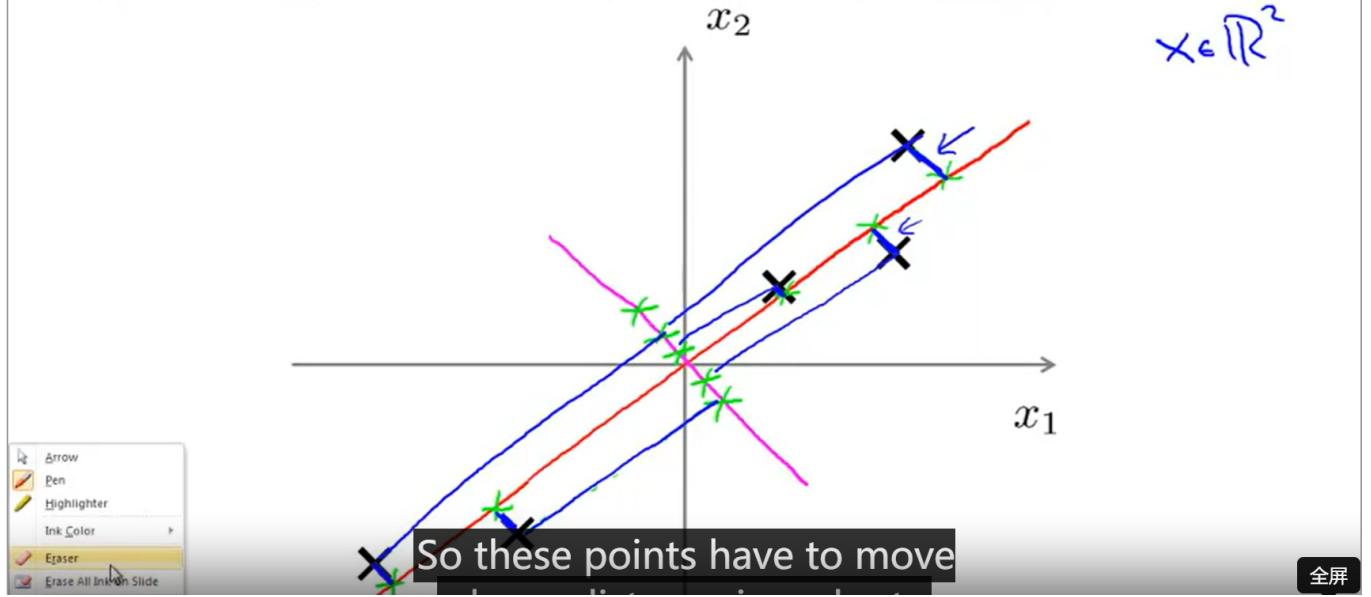
visualization

helps visualize the data easily

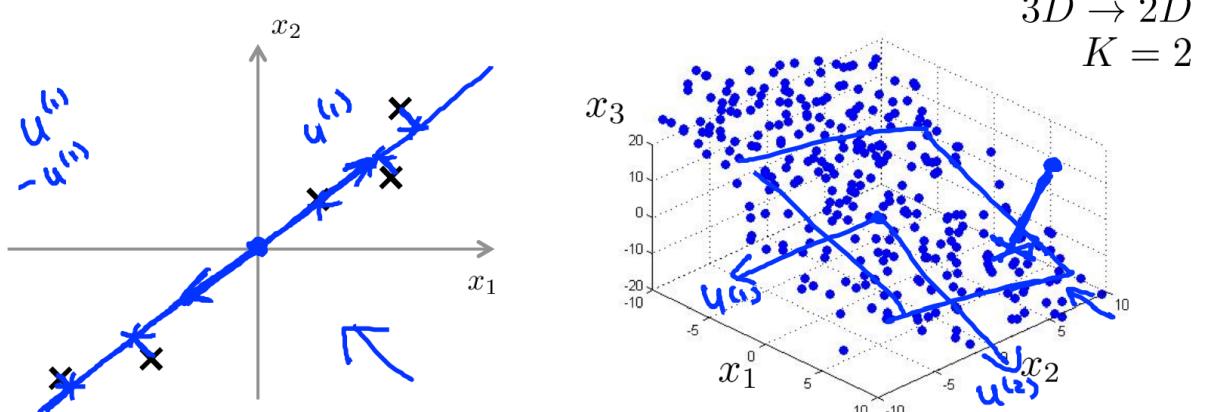
b. PCA(Principle Component Analysis)

PCA tries to find a line(sup-plain, i.e.,  $k$  vectors) onto which data is projected where the sum of the distance is minimized

### Principal Component Analysis (PCA) problem formulation



### Principal Component Analysis (PCA) problem formulation



Reduce from 2-dimension to 1-dimension: Find a direction (a vector  $u^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to minimize the projection error.

Reduce from  $n$ -dimension to  $k$ -dimension: Find  $k$  vectors  $u^{(1)}, u^{(2)}, \dots, u^{(k)}$  onto which to project the data, so as to minimize the projection error.

PCA is NOT linear regression

- PCA aims to minimize the distance but LR aims to minimize the y-direction square error
- PCA is not for prediction but data dimension reduce

steps of the algorithm

1. data pre-processing: apply mean normalization(must); feature scaling(optional) to all features

## Data preprocessing

Training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$  ←

Preprocessing (feature scaling/mean normalization):

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

Replace each  $x_j^{(i)}$  with  $x_j^{(i)} - \mu_j$ .

If different features on different scales (e.g.,  $x_1$  = size of house,  $x_2$  = number of bedrooms), scale features to have comparable range of values.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

2. compute covariance matrix and the eigenvectors of the covariance matrix using SVD(Singular-Value Decomposition)

## Principal Component Analysis (PCA) algorithm

Reduce data from  $n$ -dimensions to  $k$ -dimensions

Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)}) (x^{(i)})^T \quad \text{Sigma}$$

Compute "eigenvectors" of matrix  $\Sigma$ :

$$\rightarrow [U, S, V] = \underline{\text{svd}}(\text{Sigma}) ;$$

*Singular value decomposition  
eig (Sigma)*

The first  $k$  columns of the  $U$  matrix are exactly  $u^{(1)}, u^{(2)}, \dots, u^{(k)}$  and  $z = [u^{(1)} \ u^{(2)} \ \dots \ u^{(k)}]^T x$

c. applying PCA

data reconstruction

$$\text{as } z = [u^{(1)} \ u^{(2)} \ \dots \ u^{(k)}]^T x$$

$$x = ([u^{(1)} \ u^{(2)} \ \dots \ u^{(k)}]^T)^{-1} z = [u^{(1)} \ u^{(2)} \ \dots \ u^{(k)}] z \text{ (pseudo inverse here)}$$

$$\text{Here is why } [u^{(1)} \ u^{(2)} \ \dots \ u^{(k)}]^T = [u^{(1)} \ u^{(2)} \ \dots \ u^{(k)}]^{-1}$$



Neil Ostrove 社区助教

Both are because the eigenvectors returned by Singular Value Decomposition (SVD) form an orthonormal basis.

喜欢 2个回复 ⋮ 更多

choosing  $k$

principle:

$$\begin{aligned} &\rightarrow \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2 \leq 0.01 \quad (1\%) \\ &\rightarrow \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2 \end{aligned}$$

algorithm:

## Choosing $k$ (number of principal components)

Algorithm:

Try PCA with  $k = 1$   ~~$k=2$~~   ~~$k=4$~~

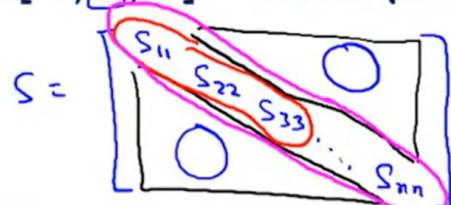
Compute  $U_{reduce}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{approx}^{(1)}, \dots, x_{approx}^{(m)}$

Check if

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01?$$

$$\underline{k=17}$$

$$\rightarrow [U, S, V] = \text{svd}(Sigma)$$



For given  $k$   $\underline{k=3}$

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01$$

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

application

## 1. supervised learning speedup

**Supervised learning speedup**  
→  $(\underline{x}^{(1)}, \underline{y}^{(1)}), (\underline{x}^{(2)}, \underline{y}^{(2)}), \dots, (\underline{x}^{(m)}, \underline{y}^{(m)})$

Extract inputs:

Unlabeled dataset:  $\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(m)} \in \mathbb{R}^{10000}$  ←  
 $\downarrow PCA$

$\underline{z}^{(1)}, \underline{z}^{(2)}, \dots, \underline{z}^{(m)} \in \mathbb{R}^{1000}$  ←

New training set:

$(\underline{z}^{(1)}, \underline{y}^{(1)}), (\underline{z}^{(2)}, \underline{y}^{(2)}), \dots, (\underline{z}^{(m)}, \underline{y}^{(m)})$

Note: Mapping  $x^{(i)} \rightarrow z^{(i)}$  should be defined by running PCA  
only on the training set. This mapping can be applied as well to  
the examples  $x_{cv}^{(i)}$  and  $x_{test}^{(i)}$  in the cross validation and test sets.

Andrew Ng

2. DO NOT use PCA to prevent overfitting, apply regularization instead

3. DO NOT use PCA if unnecessary

## W9

### 1. Anomaly Detection

#### a. density estimation

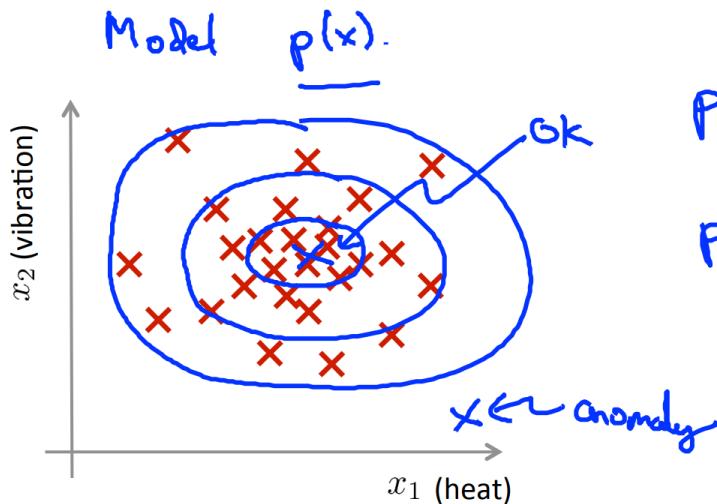
motivation

find out the point that is anomalous:

## Density estimation

→ Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

→ Is  $x_{test}$  anomalous?



$p(x_{test}) < \varepsilon \rightarrow$  flag anomaly

$p(x_{test}) \geq \varepsilon \rightarrow$  OK

Gaussian distribution(trivial)

algorithm

1. choose probable anomalous features  $x_i$

2. fit  $\mu_j, \sigma_j^2$  using:

3. compute  $p(x)$  for a new example  $x$  using:

3. Given new example  $x$ , compute  $p(x)$ :

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

4. anomaly if  $p(x) < \varepsilon$

b. building an anomaly detection system

evaluation metrics

## Algorithm evaluation

- Fit model  $p(x)$  on training set  $\{x^{(1)}, \dots, x^{(m)}\}$
- On a cross validation/test example  $x$ , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (normal)} \end{cases}$$

Possible evaluation metrics:

- True positive, false positive, false negative, true negative
- Precision/Recall
- $F_1$ -score

Can also use cross validation set to choose parameter  $\varepsilon$

precision/recall/F score (see W6)

Note that classification accuracy is not used as a metric in case of skewed data

anomaly detection vs. supervised learning

different features:

### Anomaly detection

- Very small number of positive examples ( $y = 1$ ). (0-20 is common).
- Large number of negative ( $y = 0$ ) examples.  $p(x) \leq$
- Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like;
- future anomalies may look nothing like any of the anomalous examples we’ve seen so far.

vs.

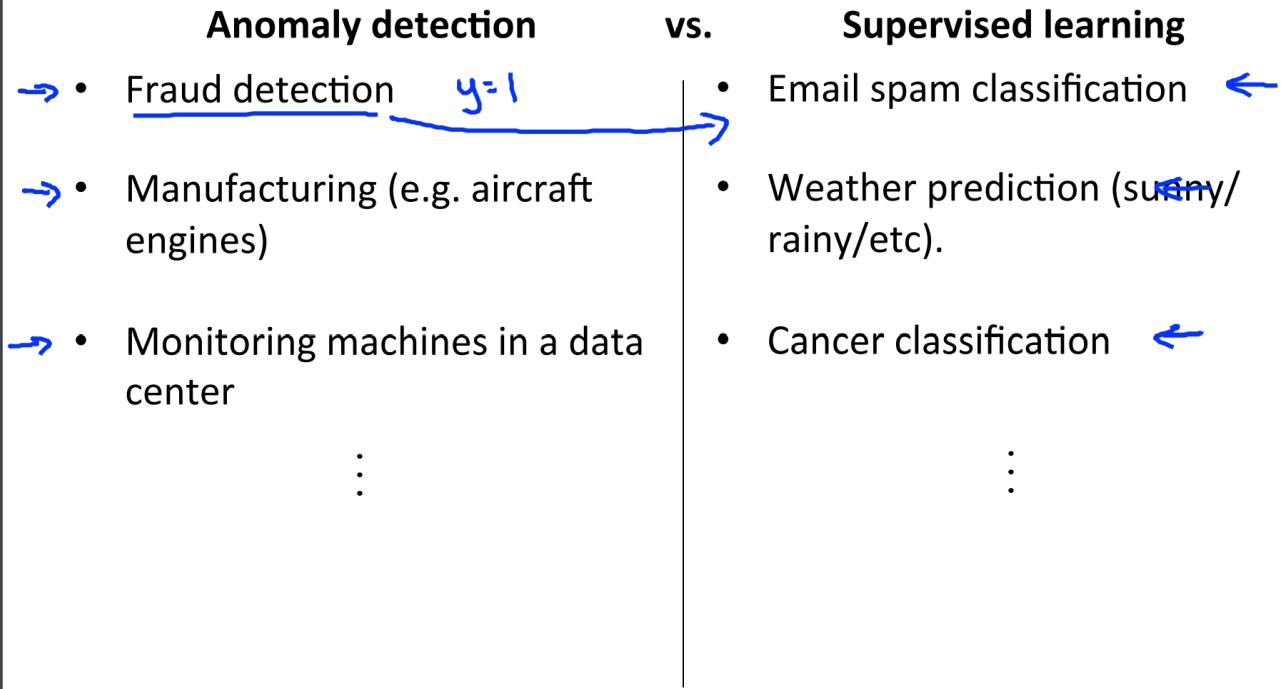
### Supervised learning

Large number of positive and negative examples. ←

Enough positive examples for ← algorithm to get a sense of what positive examples are like, future ← positive examples likely to be similar to ones in training set. ←

Spam ←

different usages:

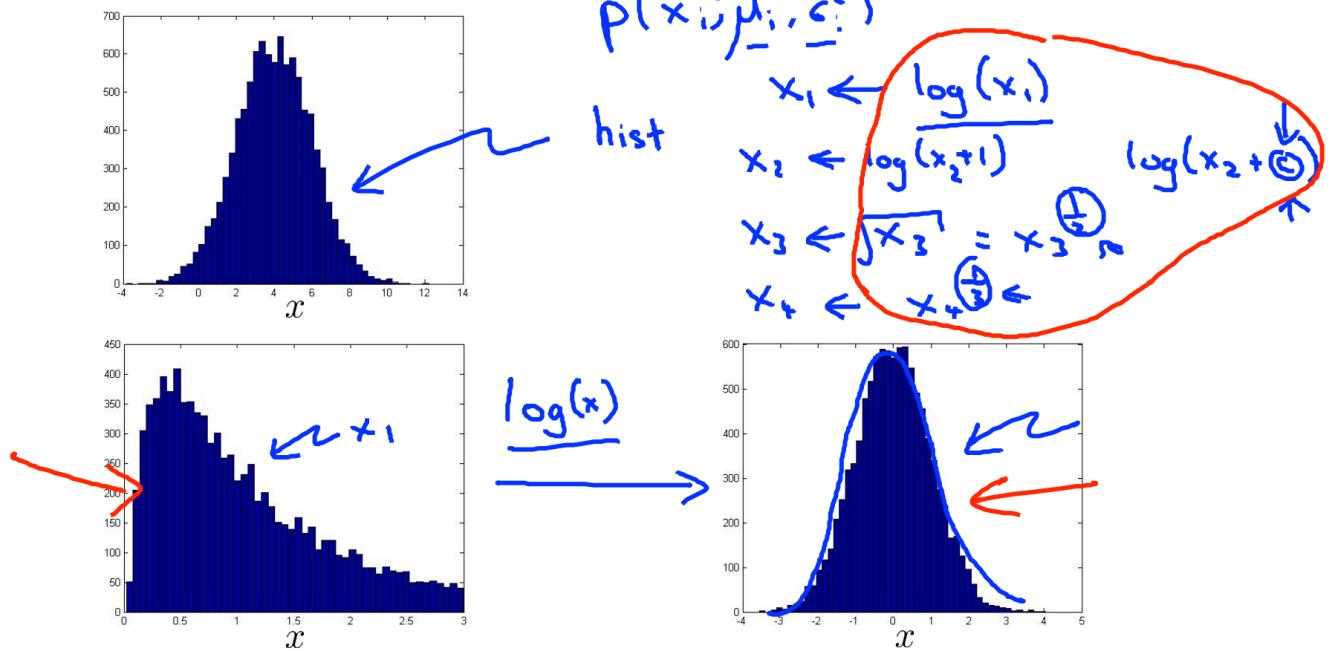


Anomaly detection

choosing what features to use

handling non-Gaussian features:

## Non-gaussian features



do some transformations to make it more Gaussian:

1.  $x := \log(x + c)$
2.  $x := x^c$

the constant  $c$  is what we adjust.

use `hist(x)` to plot the histogram of the data

error analysis:

what we want:

## → Error analysis for anomaly detection

- Want  $p(x)$  large for normal examples  $x$ .
- $p(x)$  small for anomalous examples  $x$ .

the problem is:

Most common problem:

- $p(x)$  is comparable (say, both large) for normal and anomalous examples

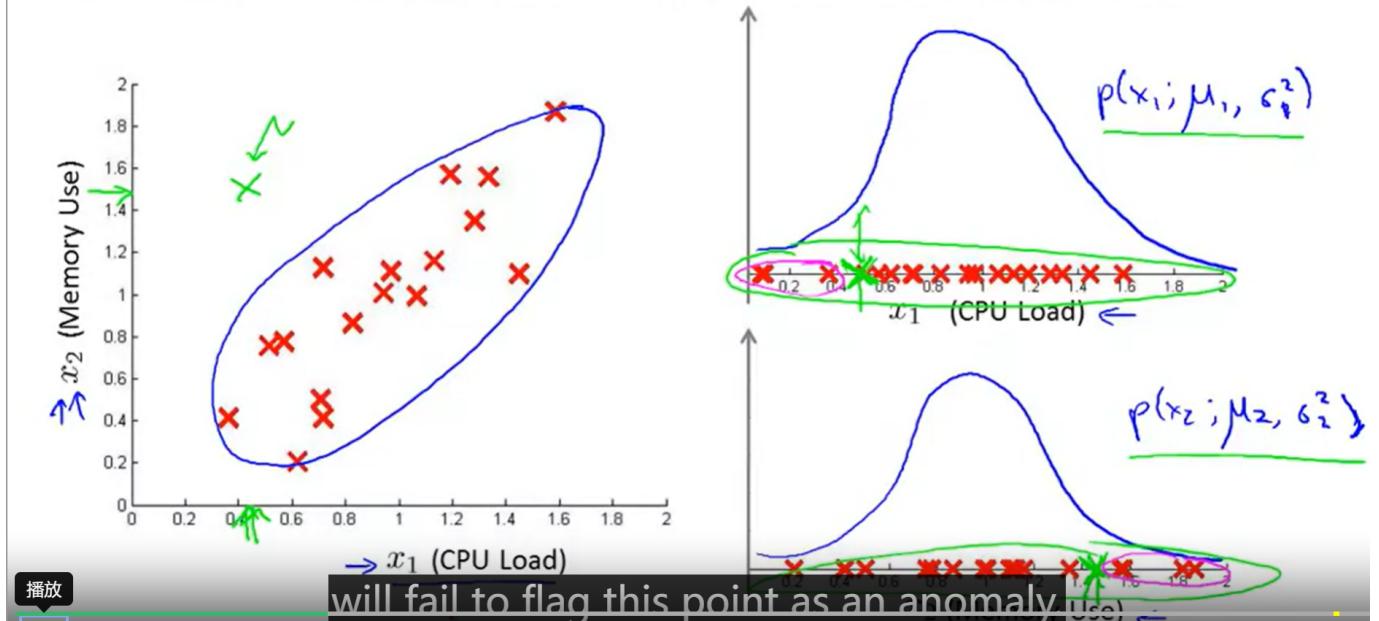
↑

solution: combine features that will distinguish between normal cases and anomalous cases

c. multivariate Gaussian distribution

an example when univariate Gaussian fails to detect the anomaly:

## Motivating example: Monitoring machines in a data center



the green point is anomalous but the Gaussian values of both features are relatively high.

### MGD

Instead of model  $p(x_i)$  separately, we model  $p(x)$  in one go.

parameters:  $\mu \in \mathbb{R}^n$   $\Sigma \in \mathbb{R}^{n \times n}$ , where  $\Sigma$  is the covariance matrix

the PDF for MGD:

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{1}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

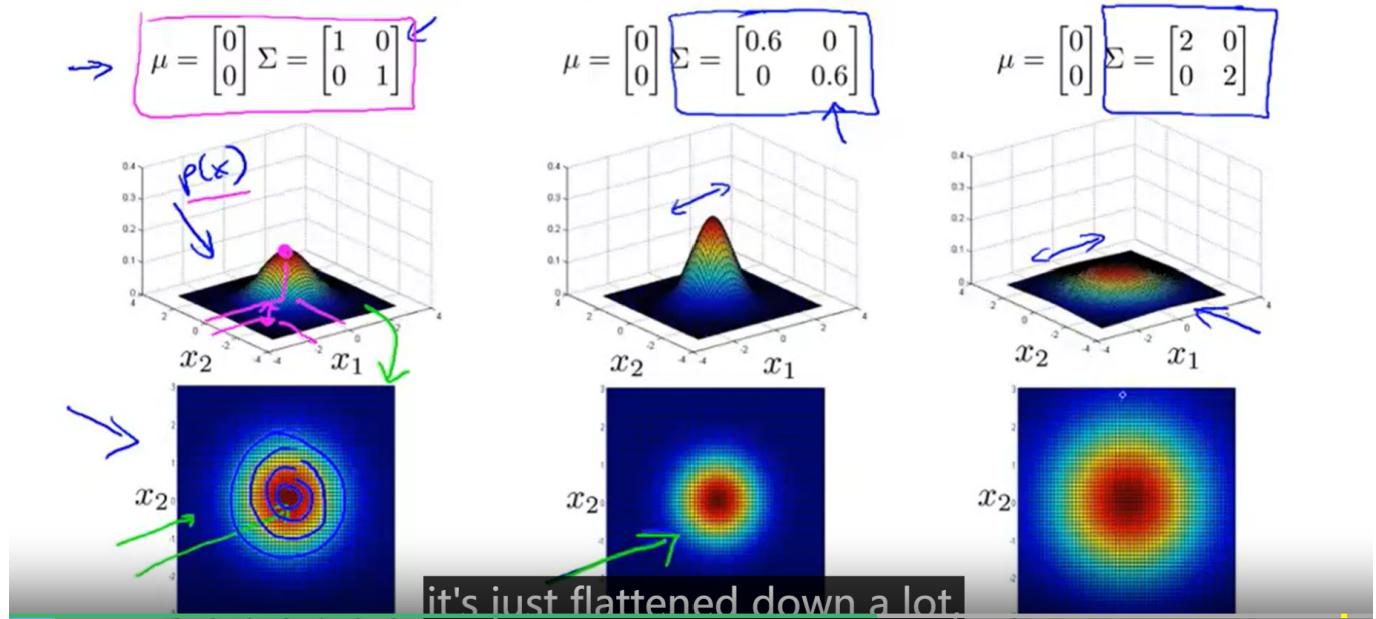
where  $x$  is a vector

### examples

1. varying the diagonal elements of  $\Sigma$ (the range of each feature):

synchronized:

## Multivariate Gaussian (Normal) examples



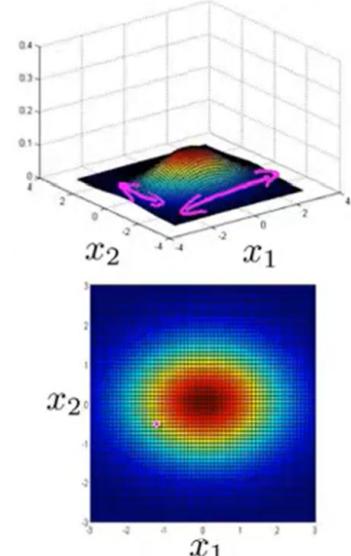
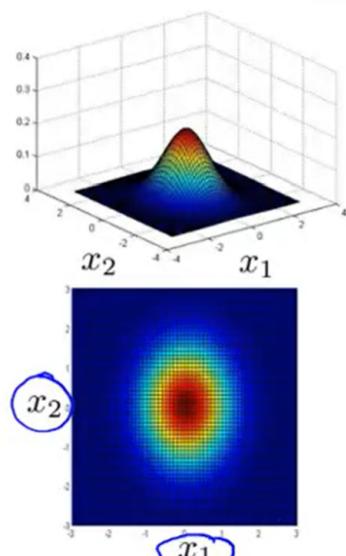
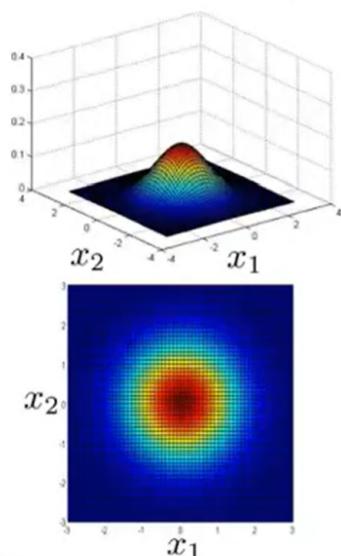
asynchronous:

## Multivariate Gaussian (Normal) examples

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

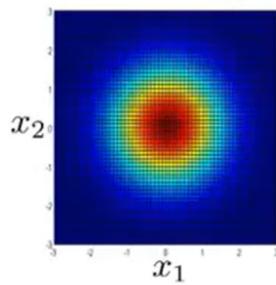
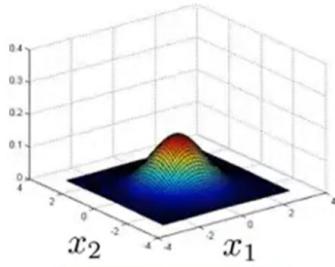
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

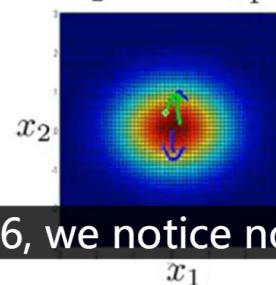
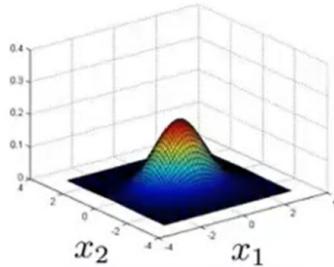


## Multivariate Gaussian (Normal) examples

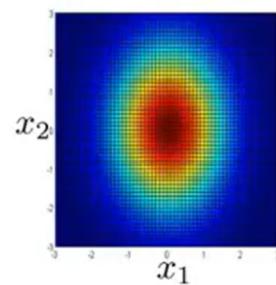
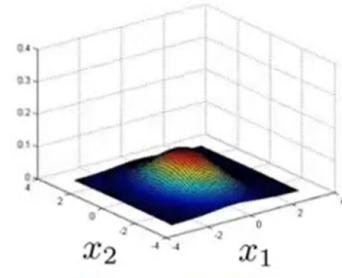
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 0.6 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$



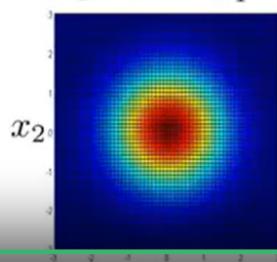
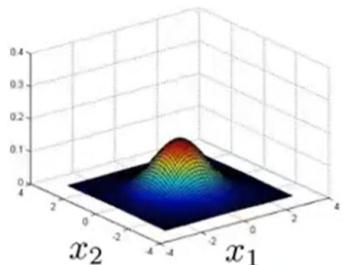
is 0.6, we notice now X2

Andrew

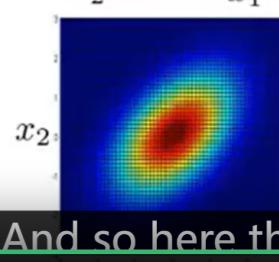
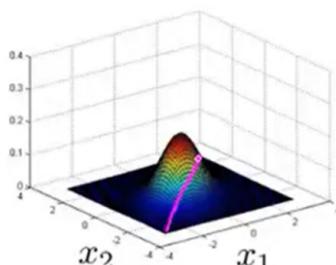
2. varying off-diagonal elements of  $\Sigma$  (the correlation between each feature):

## Multivariate Gaussian (Normal) examples

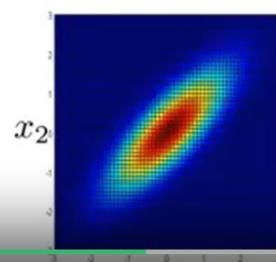
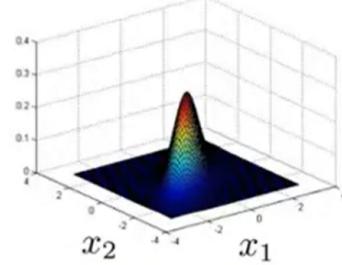
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



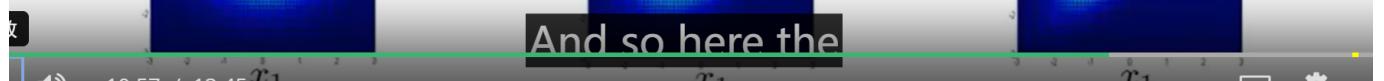
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



And so here the

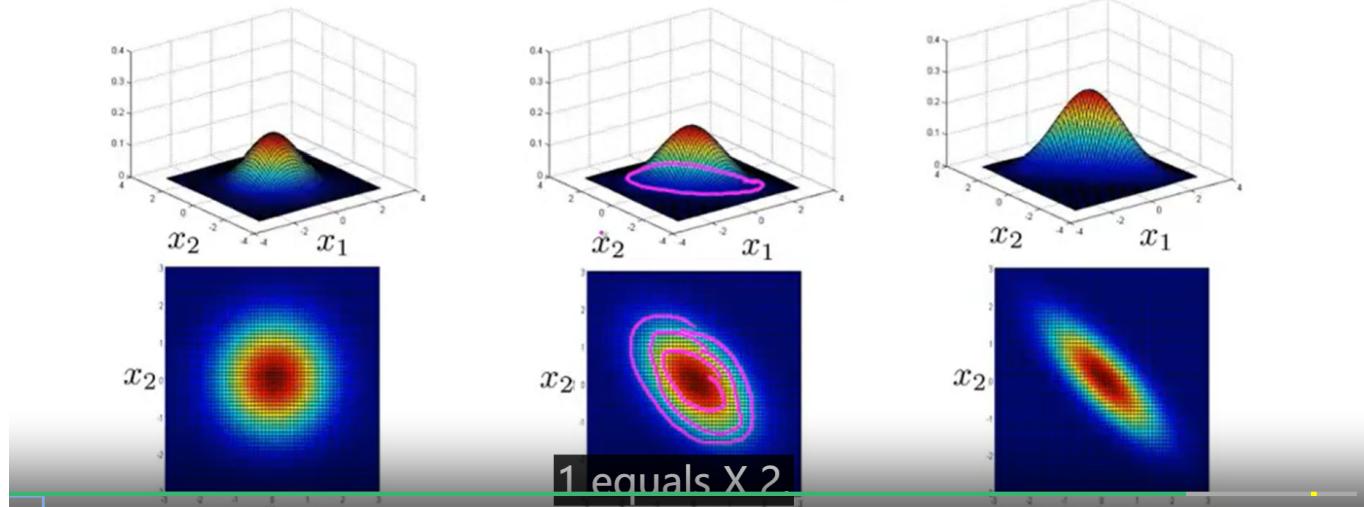


## Multivariate Gaussian (Normal) examples

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$



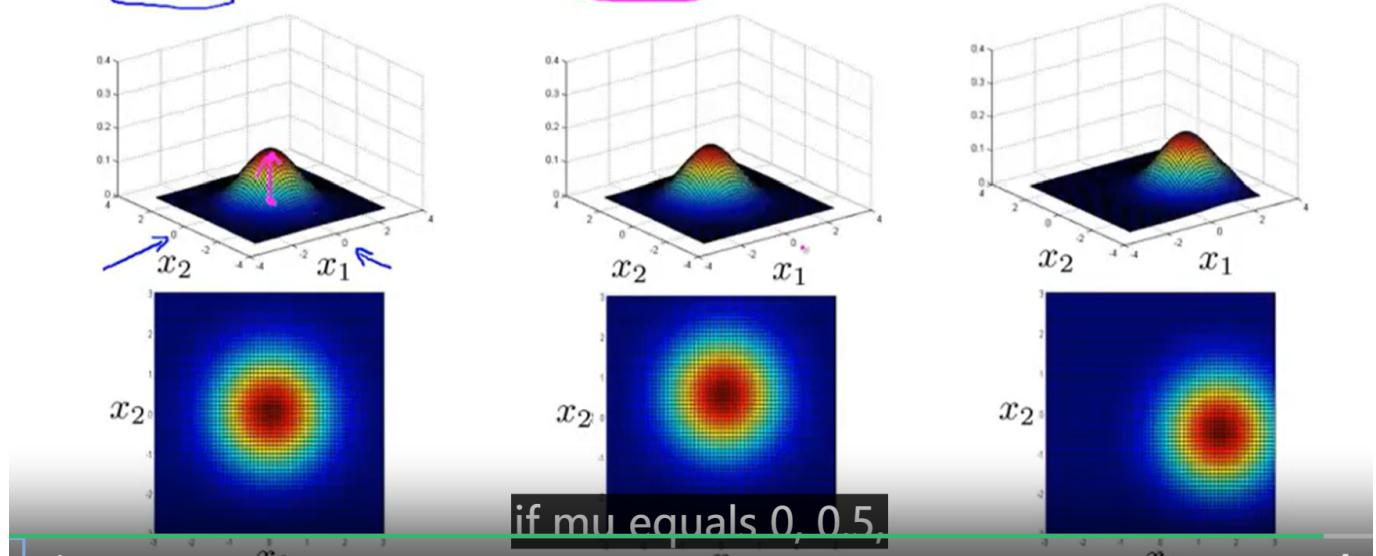
3. varying  $\mu$ (the location of the peak of the function)

## Multivariate Gaussian (Normal) examples

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 1.5 \\ -0.5 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



anomaly detection using MGD

steps:

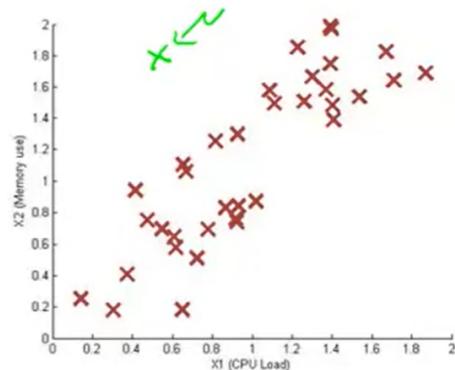
1. fitting parameters using the same approach as PCA.
2. compute the PDF of MGD of the new example:

## Anomaly detection with the multivariate Gaussian

1. Fit model  $p(x)$  by setting

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$



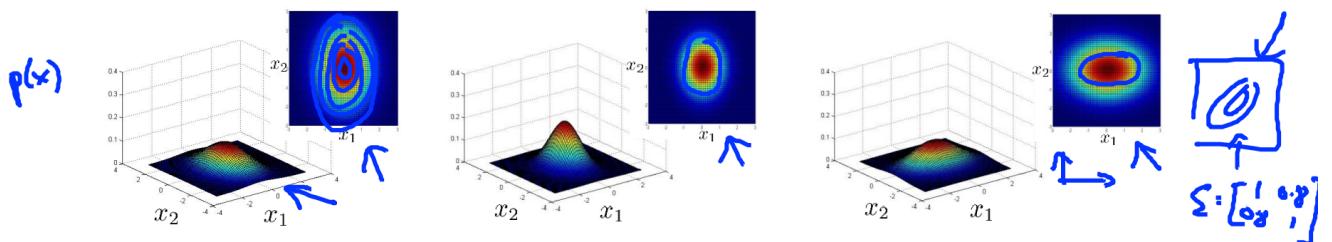
2. Given a new example  $x$ , compute

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

relationship to the original model:

### Relationship to original model

Original model:  $p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$



Corresponds to multivariate Gaussian

$$\rightarrow p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

where

$$\Sigma = \begin{bmatrix} \dots & \dots \\ \dots & \dots \end{bmatrix}$$

Andrew Ng

constraints: the off-diagonal elements shall all be 0s

The MGD automatically take the correlations between features

here is more concrete comparisons:

## → Original model

$$p(x_1; \mu_1, \sigma_1^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$$

Manually create features to capture anomalies where  $x_1, x_2$  take unusual combinations of values.

$$\rightarrow X_3 = \frac{x_1}{x_2} = \frac{\text{CPU load}}{\text{memory}}$$

- Computationally cheaper (alternatively, scales better to large  $n=10,000$ ,  $n=100,000$ )  
OK even if  $m$  (training set size) is small

## vs. → Multivariate Gaussian

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left( -\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

→ Automatically captures correlations between features

$$\Sigma \in \mathbb{R}^{n \times n}$$

$$\Sigma^{-1}$$

Computationally more expensive

$$\rightarrow \Sigma \sim \frac{n^2}{2}$$

$$\begin{cases} x_1 = x_2 \\ x_2 = x_4 + x_5 \end{cases}$$

Must have  $m > n$  or else  $\Sigma$  is non-invertible.

$$m \geq 10n$$

Andrew Ng

if  $\Sigma$  is non-invertible, then there are 2 possible reasons:

- $m \leq n$
- redundant features

## 2. Recommender Systems

### a. predicting movie ratings(content-based)

#### Glossary #8

$n_u$  : number of users

$n_m$  : number of movies

$r(i, j)$ : 1 if user  $j$  has rated movie  $i$

$y^{(i,j)}$  the rate given by user  $j$  to movie  $i$  (defined only when  $r(i, j) = 1$ )

recommender system tries to fill in the blanks marked '?' and find out what a user likes

## Example: Predicting movie ratings

→ User rates movies using one to five stars  
 zero

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	5	0	6
Romance forever	5	?	0	0
Cute puppies of love	?	5	0	?
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	?

$n_u = 4$        $n_m = 5$

new movies to that user to watch



→  $n_u = \text{no. users}$   
 →  $n_m = \text{no. movies}$   
 $r(i, j) = 1$  if user  $j$  has rated movie  $i$   
 $y^{(i,j)}$  = rating given by user  $j$  to movie  $i$  (defined only if  $r(i, j) = 1$ )

content based recommendations

idea:

**Content-based recommender systems**

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$x_0 = 1$	$x_1$ (romance)	$x_2$ (action)
Love at last	1	5	0	0	1	0.9	0
Romance forever	2	5	?	0	1	1.0	0.01
Cute puppies of love	3	?	4	0	1	0.99	0
Nonstop car chases	4	0	5	4	1	0.1	1.0
Swords vs. karate	5	0	5	?	1	0	0.9

→ For each user  $j$ , learn a parameter  $\theta^{(j)} \in \mathbb{R}^3$ . Predict user  $j$  as rating movie  $i$  with  $(\theta^{(j)})^T x^{(i)}$  stars.

problem formulation

## Problem formulation

- $r(i, j) = 1$  if user  $j$  has rated movie  $i$  (0 otherwise)
  - $y^{(i,j)}$  = rating by user  $j$  on movie  $i$  (if defined)
  - $\theta^{(j)}$  = parameter vector for user  $j$
  - $x^{(i)}$  = feature vector for movie  $i$
  - For user  $j$ , movie  $i$ , predicted rating:  $\underline{(\theta^{(j)})^T(x^{(i)})}$        $\theta^{(j)} \in \mathbb{R}^{n+1}$
  - $\underline{m^{(j)}}$  = no. of movies rated by user  $j$
- To learn  $\underline{\theta^{(j)}}$ :

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: r(i,j)=1} \left( (\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2m^{(j)}} \cdot \sum_{k=1}^n (\theta_k^{(j)})^2$$

to simplify,  $m^{(j)}$  is ignored:

the optimization objective:

## Optimization objective:

To learn  $\underline{\theta^{(j)}}$  (parameter for user  $j$ ):

$$\rightarrow \min_{\theta^{(j)}} \frac{1}{2} \sum_{i: r(i,j)=1} \left( (\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{k=1}^n (\theta_k^{(j)})^2$$

To learn  $\underline{\theta^{(1)}}, \underline{\theta^{(2)}}, \dots, \underline{\theta^{(n_u)}}$ :

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: r(i,j)=1} \left( (\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

And I can then use that to make

algorithm: gradient descent again:

## Optimization algorithm:

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} \left( (\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$\mathcal{J}(\theta^{(1)}, \dots, \theta^{(n_u)})$

Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k = 0)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left( \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0)$$

*these 1 over m terms,*

b. collaborative filtering(协同过滤)

idea

assumption #1: we have no idea about the features of the movie(this is more realistic)

## Problem motivation

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$x_1$ (romance)	$x_2$ (action)
Love at last	5	5	0	0	?	?
Romance forever	5	?	?	0	?	?
Cute puppies of love	?	4	0	?	?	?
Nonstop car chases	0	0	5	4	?	?
Swords vs. karate	0	0	5	?	?	?

assumption #2: users have told us what kind of movie they like. i.e., the  $\theta$  parameters are given

## Problem motivation

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$x_1$ (romance)	$x_2$ (action)
Love at last	5	5	0	0	?	?
Romance forever	5	?	?	0	?	?
Cute puppies of love	?	4	0	?	?	?
Nonstop car chases	0	0	5	4	?	?
Swords vs. karate	0	0	5	?	?	?

$$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(2)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}, \theta^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}, \theta^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$$

that she really

we can infer the value of features from the  $\theta$  values and the rates:

## Optimization algorithm

Given  $\theta^{(1)}, \dots, \theta^{(n_u)}$ , to learn  $x^{(i)}$ :

$$\min_{x^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

Given  $\theta^{(1)}, \dots, \theta^{(n_u)}$ , to learn  $x^{(1)}, \dots, x^{(n_m)}$ :

$$\min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

this is essentially the same algorithm as the content-based one.

## Collaborative filtering

$$\begin{matrix} r^{(i,j)} \\ y^{(i,j)} \end{matrix}$$

[ Given  $x^{(1)}, \dots, x^{(n_m)}$  (and movie ratings),  
can estimate  $\theta^{(1)}, \dots, \theta^{(n_u)}$  ]

[ Given  $\theta^{(1)}, \dots, \theta^{(n_u)}$ ,  
can estimate  $x^{(1)}, \dots, x^{(n_m)}$  ]

Guess  $\Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \dots$

the chain is interesting isn't it? But it is not as efficient as it can be. we can solve  $\theta$  and  $x$  simultaneously by combining the cost functions together.

optimization objective

Minimizing  $x^{(1)}, \dots, x^{(n_m)}$  and  $\theta^{(1)}, \dots, \theta^{(n_u)}$  simultaneously:

$$J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$\min_{\substack{x^{(1)}, \dots, x^{(n_m)} \\ \theta^{(1)}, \dots, \theta^{(n_u)}}} J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$

Andrew Ng

we ignore the convention that  $x_0 = 1$  because the algorithm can do this itself.

algorithm

steps:

1. initialize  $x$  and  $\theta$  to small random values (to ensure that  $x^{(1)}$  to  $x^{(n_m)}$  can be learnt separately)
2. use GD to minimize  $J$

2. Minimize  $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$  using gradient descent (or an advanced optimization algorithm). E.g. for every  $j = 1, \dots, n_u, i = 1, \dots, n_m$  :

$$x_k^{(i)} := x_k^{(i)} - \alpha \left( \sum_{j:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left( \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right)$$

$\frac{\partial J}{\partial x_k^{(i)}}$

3. predict!

3. For a user with parameters  $\underline{\theta}$  and a movie with (learned) features  $\underline{x}$ , predict a star rating of  $\underline{\theta^T x}$ .

停

( $\Theta^{(i)}$ )<sup>T</sup> user J is going to

c. low rank matrix factorization

vectorization: low rank matrix factorization

**Collaborative filtering**

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

$\times (\Theta)^T$

$(\Theta^{(i)})^T (x^{(i)})$

Predicted ratings:  $(i, j) \rightarrow$

$$\begin{bmatrix} (\theta^{(1)})^T (x^{(1)}) & (\theta^{(2)})^T (x^{(1)}) & \dots & (\theta^{(n_u)})^T (x^{(1)}) \\ (\theta^{(1)})^T (x^{(2)}) & (\theta^{(2)})^T (x^{(2)}) & \dots & (\theta^{(n_u)})^T (x^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ (\theta^{(1)})^T (x^{(n_m)}) & (\theta^{(2)})^T (x^{(n_m)}) & \dots & (\theta^{(n_u)})^T (x^{(n_m)}) \end{bmatrix}$$

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

$\rightarrow \Theta = \begin{bmatrix} -(\Theta^{(1)})^T \\ -(\Theta^{(2)})^T \\ \vdots \\ -(\Theta^{(n_u)})^T \end{bmatrix}$

is also called low rank

finding related movies

## Finding related movies

For each product  $i$ , we learn a feature vector  $\underline{x^{(i)}} \in \mathbb{R}^n$ .

$\rightarrow x_1 = \text{romance}, x_2 = \text{action}, x_3 = \text{comedy}, x_4 = \dots$

How to find movies  $j$  related to movie  $i$ ?

small  $\|\underline{x^{(i)}} - \underline{x^{(j)}}\| \rightarrow$  movie  $j$  and  $i$  are "similar"

5 most similar movies to movie  $i$ :

Find the 5 movies  $j$  with the smallest  $\|\underline{x^{(i)}} - \underline{x^{(j)}}\|$ .

mean normalization

if a user hasn't rated any movies, the  $\theta$  parameters will eventually go to all zeros because of the regularization term and it is not possible to recommend any movies

## Users who have not rated any movies

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)	$\downarrow$
Love at last	5	5	0	0	?	0
Romance forever	5	?	?	0	?	0
Cute puppies of love	?	4	0	?	?	0
Nonstop car chases	0	0	5	4	?	0
Swords vs. karate	0	0	5	?	?	0

$\rightarrow$   $\min_{x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{(i,j): r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$

$n=2 \quad \underline{\Theta^{(s)}} \in \mathbb{R}^2 \quad \underline{\Theta^{(s)}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$\frac{\lambda}{2} \left[ (\Theta_1^{(s)})^2 + (\Theta_2^{(s)})^2 \right] \leftarrow$

what we do:

add the mean rate of each movie to the prediction:

## Mean Normalization:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 \\ 2.5 & ? & ? & -2.5 \\ ? & 2 & -2 & ? \\ 2.25 & -2.25 & 2.75 & 1.75 \\ -1.25 & -1.25 & 3.75 & -1.25 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

For user  $j$ , on movie  $i$  predict:

$$\rightarrow (\Theta^{(s)})^\top (x^{(i)}) + \mu_i$$

$\downarrow$   
learn  $\underline{\Theta}^{(i)}, \underline{x}^{(i)}$

User 5 (Eve):

$$\Theta^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad (\Theta^{(s)})^\top (x^{(i)}) + \boxed{\mu_i}$$

W10

## 1. Stochastic Gradient Descent(large scale GD)

### a. batch GD

just the same as the algorithm in Week 2

## Batch gradient descent

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

Repeat {

$$\rightarrow \theta_j := \theta_j - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\frac{\partial}{\partial \theta_j} J_{train}(\theta)}$$

(for every  $j = 0, \dots, n$ )

}

b. SGD

algorithm

## Stochastic gradient descent

$\rightarrow$  1. Randomly shuffle (reorder) training examples

$\rightarrow$  2. Repeat {

for  $i := 1, \dots, m$  {

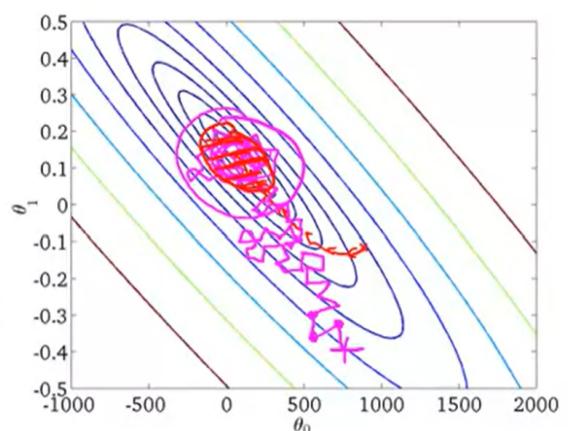
$$\rightarrow \theta_j := \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(for every  $j = 0, \dots, n$ )

}

}

we had this outer loop repeat which says to do this inner loop multiple



Observe one training example every time and the model eventually wanders around the global minimum.

convergence problem

convergence checking:

## Checking for convergence

→ Batch gradient descent:

→ Plot  $J_{train}(\theta)$  as a function of the number of iterations of gradient descent.

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

M = 300, 500, 800

→ Stochastic gradient descent:

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

→  $(x^{(i)}, y^{(i)})$ ,  $(x^{(i+1)}, y^{(i+1)})$ , ...

→ During learning, compute  $\underset{\uparrow \uparrow}{cost(\theta, (x^{(i)}, y^{(i)}))}$  before updating  $\theta$  using  $(x^{(i)}, y^{(i)})$ .

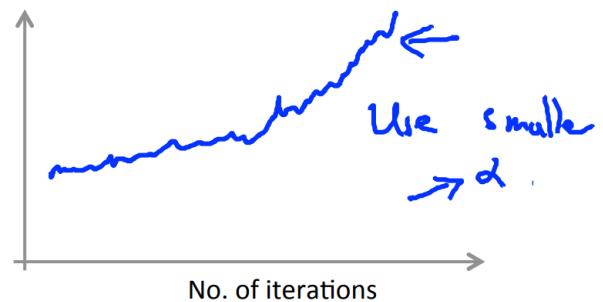
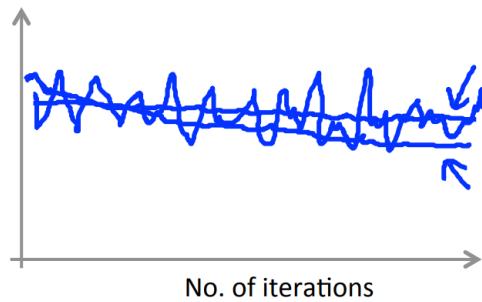
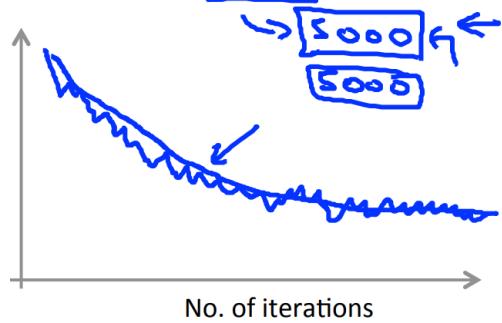
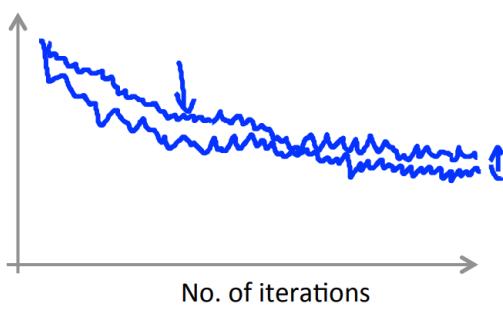
→ Every 1000 iterations (say), plot cost( $\theta, (x^{(i)}, y^{(i)})$ ) averaged over the last 1000 examples processed by algorithm.

Andrew Ng

4 possible situations:

## Checking for convergence

Plot  $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ , averaged over the last 1000 (say) examples



Andri

1 and 2: no bugs

3: may have to change something in the algorithm or set the average number larger(say, 5000)

4: use smaller learning rate  $\alpha$

*an approach to make SGD converge to the global minimum*

slowly decrease  $\alpha$  in each epoch.  $\alpha$  can be set to  $\alpha = \frac{\text{const1}}{\text{iteration number} + \text{const2}}$

the meandering will be smaller and smaller while approaching the global minimum:

## Stochastic gradient descent

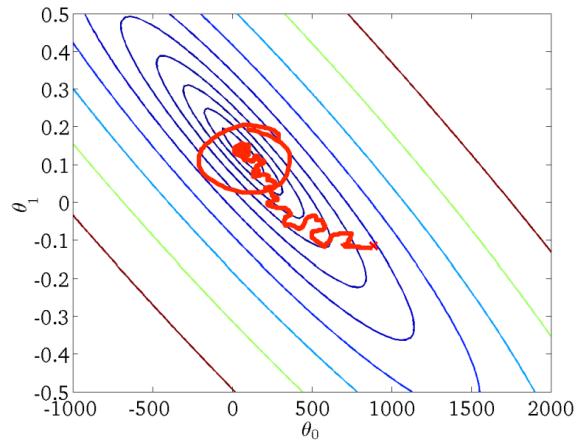
$$cost(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2$$

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m cost(\theta, (x^{(i)}, y^{(i)}))$$

```

1. Randomly shuffle dataset.
2. Repeat {
    for i := 1, ..., m {
        θ_j := θ_j - α(h_θ(x^{(i)}) - y^{(i)})x_j^{(i)}
            (for j = 0, ..., n)
    }
}

```



Learning rate  $\alpha$  is typically held constant. Can slowly decrease  $\alpha$  over time if we want  $\theta$  to converge. (E.g.  $\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}$ )  $\alpha \rightarrow 0$

### c. mini-batch GD

idea

## Mini-batch gradient descent

→ Batch gradient descent: Use all  $m$  examples in each iteration

→ Stochastic gradient descent: Use 1 example in each iteration

Mini-batch gradient descent: Use  $b$  examples in each iteration

$b = \text{mini-batch size}$ .  $b = 10$ . 2-100

algorithm

## Mini-batch gradient descent

Say  $b = 10, m = 1000$ .

Repeat {

for  $i = 1, 11, 21, 31, \dots, 991$  {

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

(for every  $j = 0, \dots, n$ )

}

}

## 2. Advanced Topics

### a. online learning

learn from a data stream. In each epoch we just consider one training example:

#### Online learning

Shipping service website where user comes, specifies origin and destination, you offer to ship their package for some asking price, and users sometimes choose to use your shipping service ( $y = 1$ ), sometimes not ( $y = 0$ ).

Features  $x$  capture properties of user, of origin/destination and asking price. We want to learn  $p(y = 1|x; \theta)$  to optimize price.

Repeat forever {  
Get  $(x, y)$  corresponding to user.  
Update  $\theta$  using  $(x, y)$ :  
 $\rightarrow \theta_j := \theta_j - \alpha (h_\theta(x) - y) \cdot x_j \quad (j=0, \dots, n)$   
}  $\rightarrow$  Can adapt to changing user preference.

another example: information retrieval

## Other online learning example:

### Product search (learning to search)

User searches for "Android phone 1080p camera"  $\leftarrow$

Have 100 phones in store. Will return 10 results.

→  $x = \text{features of phone, how many words in user query match name of phone, how many words in query match description of phone, etc.}$

→  $y = 1$  if user clicks on link.  $y = 0$  otherwise.

$(x, y) \leftarrow$   
↑      ↓

→ Learn  $p(y = 1|x; \theta)$ .  $\leftarrow$  predicted CTR

→ Use to show user the 10 phones they're most likely to click on.

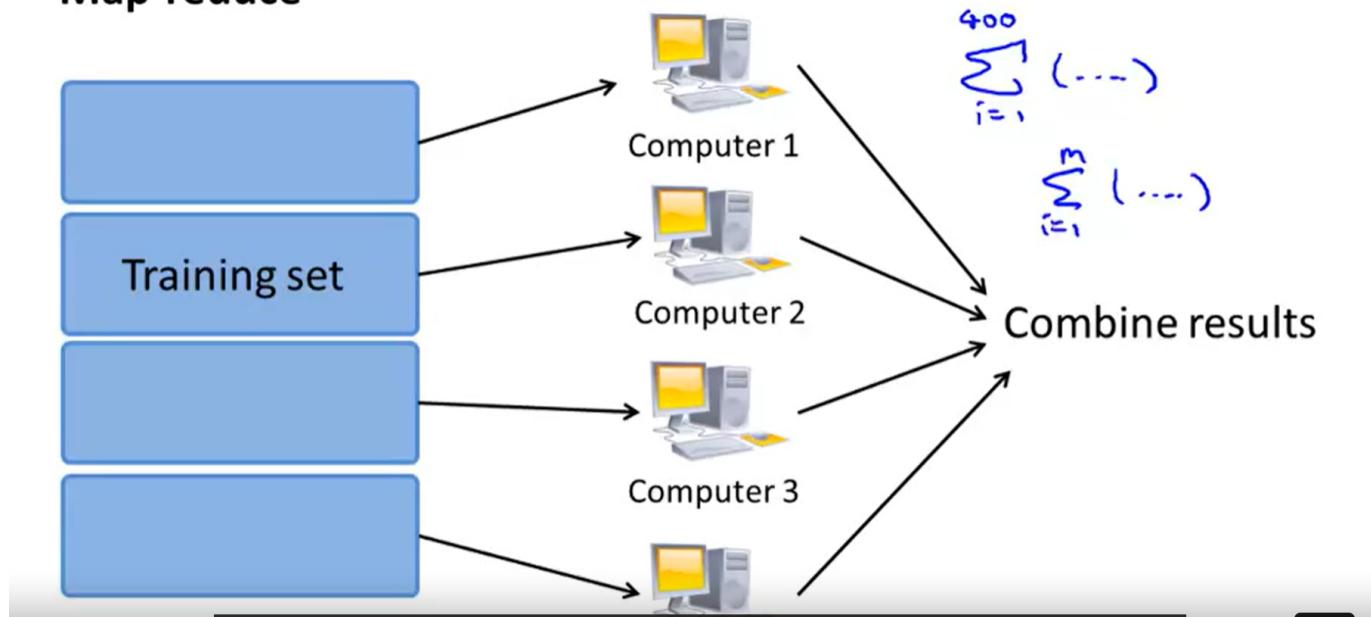
Other examples: Choosing special offers to show user; customized selection of news articles; product recommendation; ...  
Maybe, you can run your

## b. map reduce & data parallelism

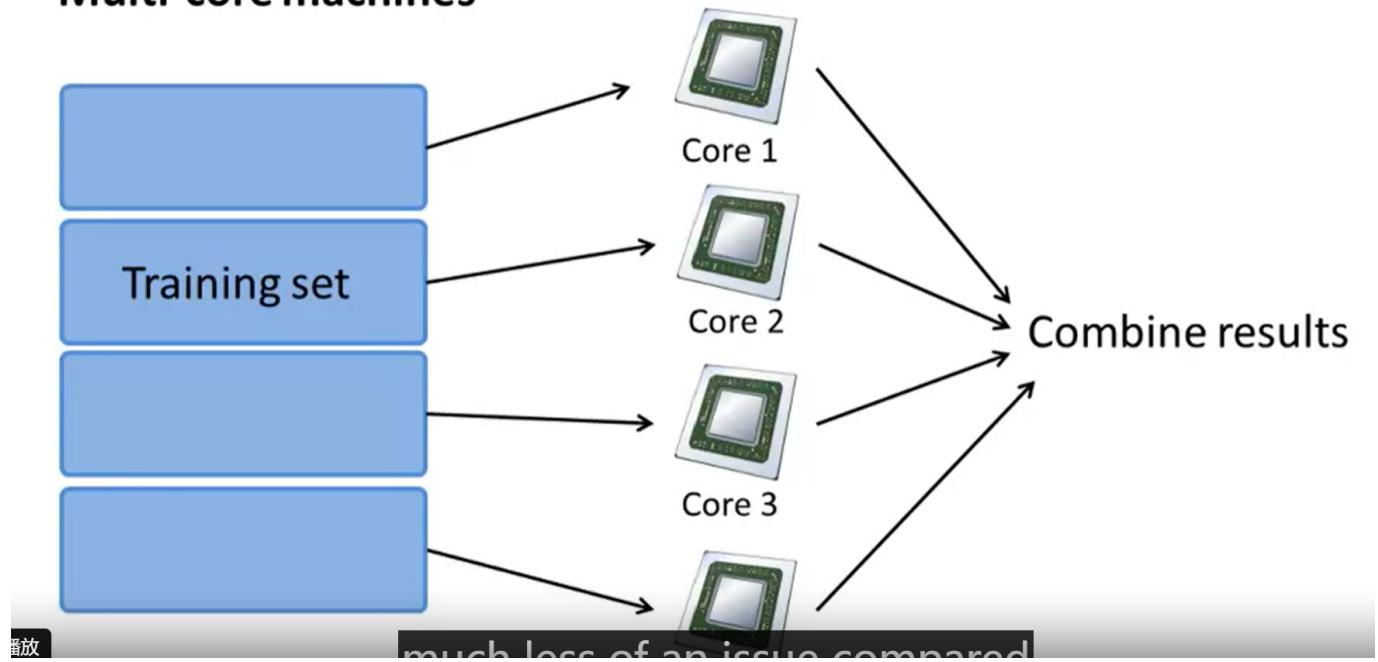
map reduce idea:

compute the gradient on different computers(or different cores) and combine them in a master server

### Map-reduce



## Multi-core machines



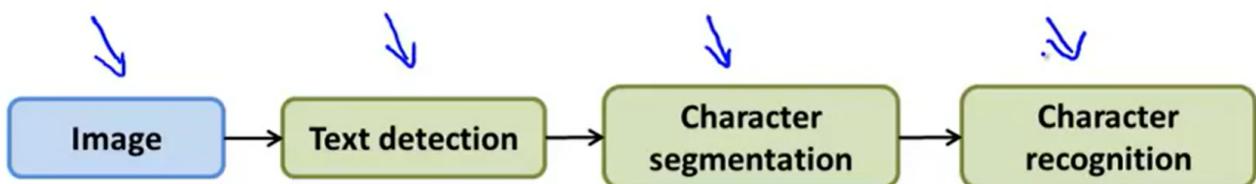
## W11

### Case Study: Photo OCR

#### a. pipeline

1. text detection
2. character segmentation
3. character classification
4. (optional) spelling correction

### Photo OCR pipeline



#### b. sliding window technique

## pedestrian detection problem

1. determine the size of the window
2. slide the window from the left top to the right bottom to detect whether there is a pedestrian
3. choose a larger window size(resize it to the image size of the training set, say 80\*40) and repeat 1,2.

## text detection problem

1. use slide window to detect all the text regions in an image:



2. the result may be like this



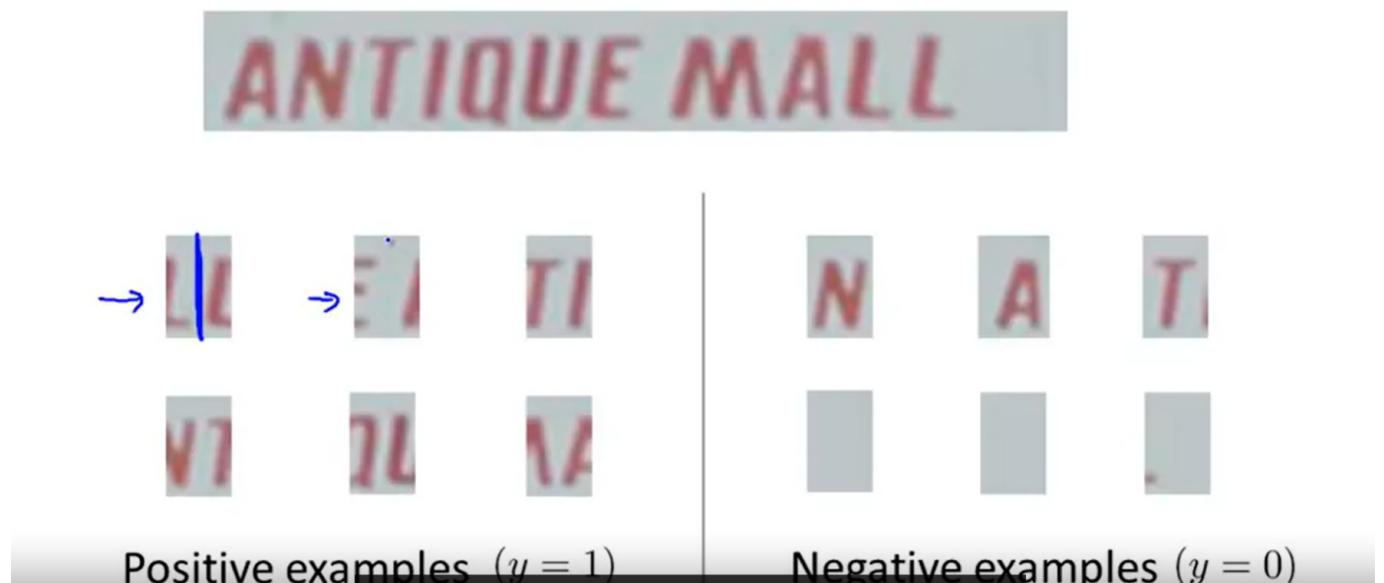
3. apply expansion operator:



character segmentation

train a classifier to determine whether there is a split between characters

## 1D Sliding window for character segmentation



applying SW on this:



c. getting data for the system

1. artificial data synthesis

## Artificial data synthesis for photo OCR



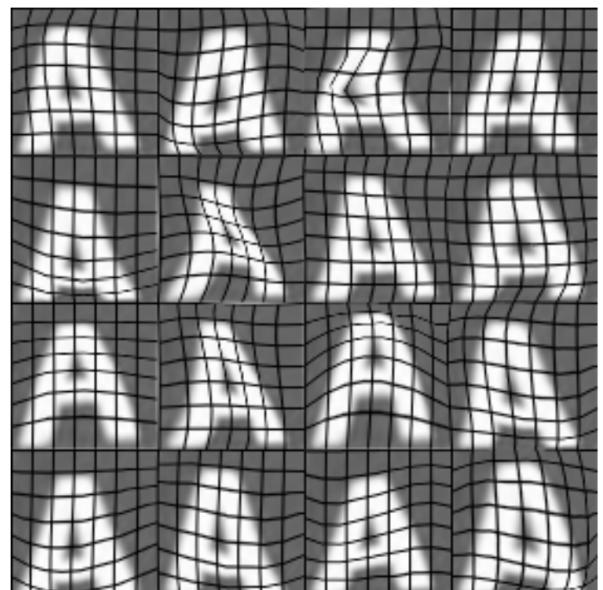
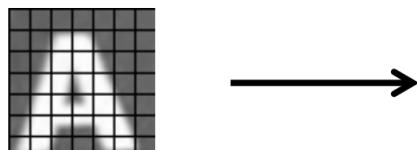
Real data



Synthetic data

2. introducing distortions:

### Synthesizing data by introducing distortions



advice on getting more data

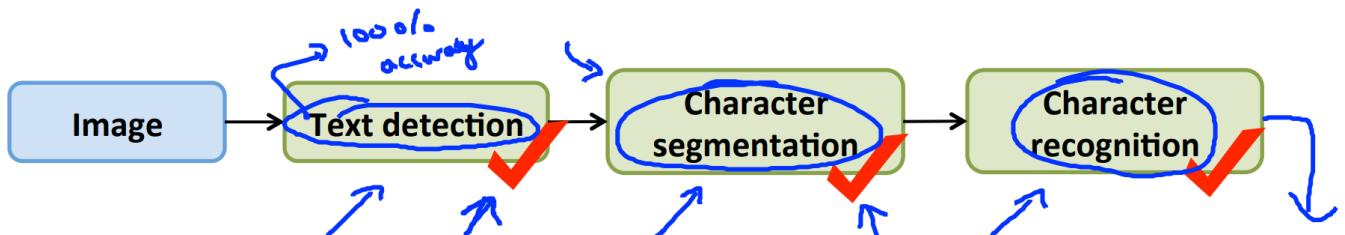
## Discussion on getting more data

1. Make sure you have a low bias classifier before expending the effort. (Plot learning curves). E.g. keep increasing the number of features/number of hidden units in neural network until you have a low bias classifier.
2. "How much work would it be to get 10x as much data as we currently have?"
  - Artificial data synthesis
  - Collect/label it yourself
  - "Crowd source" (E.g. Amazon Mechanical Turk)

d. ceiling analysis(上限分析)

examples:

### Estimating the errors due to each component (ceiling analysis)



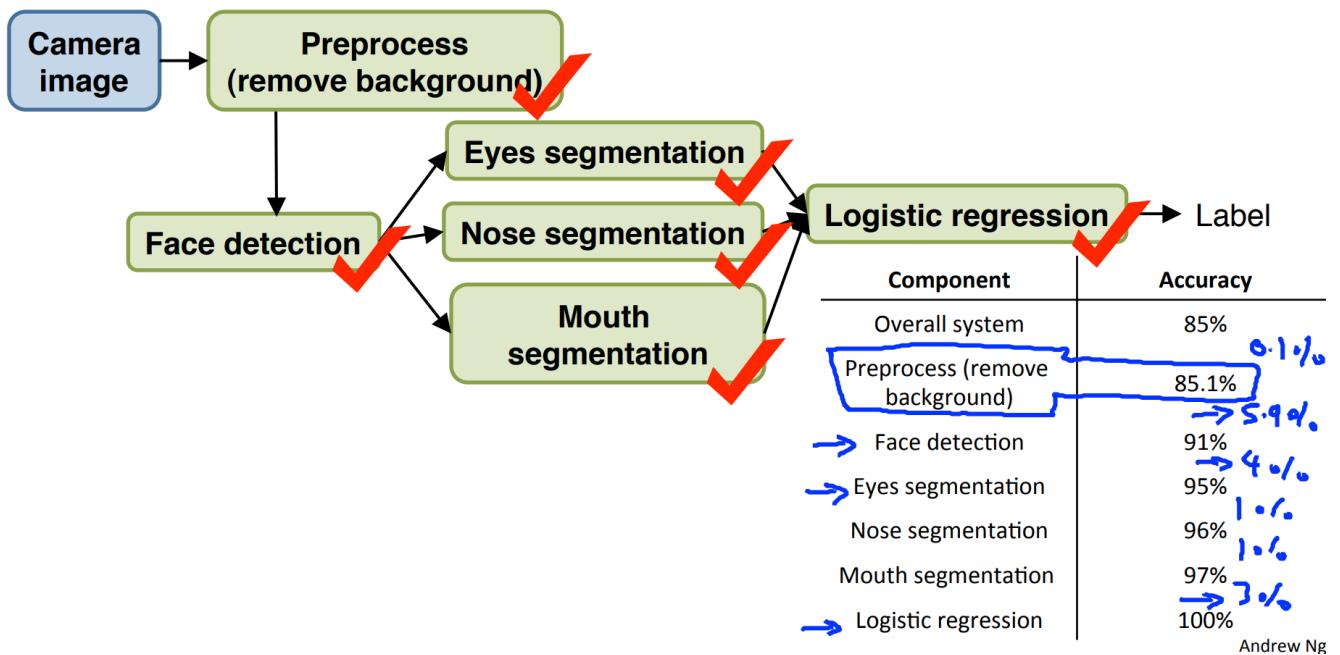
What part of the pipeline should you spend the most time trying to improve?

Component	Accuracy
Overall system	72%
→ Text detection	89%
Character segmentation	90%
Character recognition	100%

Handwritten annotations show the following improvements:  
- Overall system: 72% → 89% (17%)  
- Text detection: 89% → 90% (1%)  
- Character segmentation: 90% → 100% (10%)

Andrew

## Another ceiling analysis example



Manually make each step of the pipeline to the ground truth and check how much the accuracy improved.

End of Course: Thank you from Andrew Ng

Thank you.

- Andrew Ng