

Machine learning: grow out of work in AI, have new capacity for computer

#### Range

- image recognition
- recommend system
- Database mining
- Autonomous helicopter
- Handwriting recognition
- Natural language processing
- Self-customizing program: AI

#### Definition

- E: experience: collect what you want to know
- T: task: based on E to do sth
- P: probability: the number of correctly classified result

#### ML algorithms

- Supervised learning
- Unsupervised learning
- Reinforcement learning
- Recommender system

#### Fitting function

- Straight line
- Quadratic function, namely second-order polynomial

Supervised learning: right answers given

Unsupervised learning: predict continuous value output

Classification: A or B: discrete value output (0/1) which several features/attributes impact them

#### Example

- regression: face recognition
- Classification: malignant or benign tumor

Unsupervised learning: approach problems with no idea of the results

Application: organize computer cluster, social network analysis (cohesive people), market segmentation, astronomical data

Octave function: SVD奇异值分解

e.g. SVD (singular value decomposition): linear algebra routine

Clustering application: different genes

Unclustering application: cocktail party algorithm: find structure in a chaotic environment

#### Model representation

1. Training set - learning algorithm - hypothesis

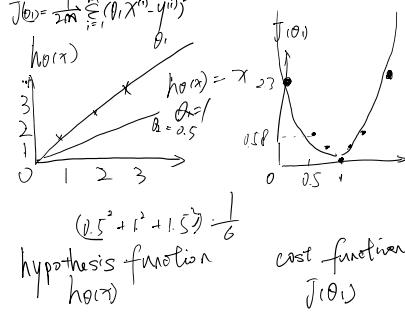
#### Cost function

- Squared error function is commonly used for linear regression

#### Hypothesis function

#### Cost function

#### $h(\theta)$

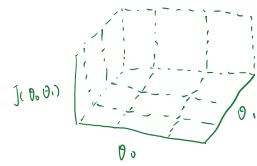


Find a  $\theta$  which make  $J(\theta)$  minimize, in this problem: when  $\theta=1$ , have  $J(\theta)=0$

算出的  $\theta$  而来是让代价函数最小化

高斯的对数似然与最小二乘  
用来得到  $J(\theta)$

Gradient Descent Algorithm 梯度下降  
Apply for linear regression



(=) A denote assignment operator

Repeat until convergence:

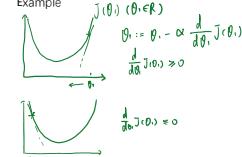
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

learning rate  $\alpha$  gradient derivative term

Simultaneous update - 同步更新 相当于 gradient descent - 梯度下降

$$\begin{aligned} \text{temp0} &:= \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \text{temp1} &:= \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \\ \theta_0 &:= \text{temp0} \\ \theta_1 &:= \text{temp1} \end{aligned}$$

Example



The learning rate:  $\alpha$

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

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

Gradient Descent for Linear Regression

gradient descent algorithm linear regression model

repeat until converge:

$$\begin{aligned} \theta_j &:= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \\ \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) &= \frac{2}{m} \sum_{i=1}^m (\theta_0 * x^{(i)} + \theta_1 * x^{(i)} - y^{(i)}) \\ \Rightarrow \theta_0 &:= \theta_0 - \alpha \frac{2}{m} \sum_{i=1}^m (\theta_0 * x^{(i)} + \theta_1 * x^{(i)} - y^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m (6\theta_0 + 2\theta_1 * x^{(i)} - y^{(i)}) \\ \theta_1 &:= \theta_1 - \alpha \frac{2}{m} \sum_{i=1}^m (\theta_0 * x^{(i)} + \theta_1 * x^{(i)} - y^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m (2\theta_1 * x^{(i)} - y^{(i)}) \\ &= \frac{2}{m} \sum_{i=1}^m (\theta_1 * x^{(i)} - y^{(i)}) \end{aligned}$$

repeat until converge

$$\begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\theta_0 * x^{(i)} - y^{(i)}) \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\theta_1 * x^{(i)} - y^{(i)}) \end{aligned}$$

The cost function of linear regression is called convex function.  $\approx$  bowl-shaped function  
Character: not have local optima, but only one global optimum!

# Linear Algebra Review

2017年8月12日 0:25

Scalar means that an object is a single value, not a vector or matrix. 标量

It is a real number.

Scalar multiplication 标量乘法

Vector: n by 1 or called n-dimensional vector

$$h_1(x) = -40 + 0.25x$$
$$h_2(x) = 200 + 0.1x$$
$$h_3(x) = -130 + 0.4x$$

row data  $\begin{bmatrix} 123 \\ 444 \\ 135 \\ 789 \end{bmatrix}$   $\times \begin{bmatrix} -40 & 200 & -130 \\ 0.25 & 0.1 & 0.4 \end{bmatrix}$

Associative property: YES -  $A * (B * C) = (A * B) * C$

Commutative property: NO -  $A * B \neq B * A$

Identity matrix: denoted I

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

$2 \times 2$

For any matrix A, we have  $A^{m \times n} \cdot I^{n \times n} = I^{m \times m} \cdot A^{m \times n} = A^{m \times n}$

Matrix inverse:

If A is an  $m \times m$  matrix (square matrix), and if it has an inverse

$$A \cdot A^{-1} = I$$

Matrices that don't have an inverse called singular 奇异矩阵 or degenerate 退化矩阵

e.g.  $\begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$

Matrix Transpose

Let A be an  $m \times n$  matrix, and let  $B = A^T$ .

Then B is an  $n \times m$  matrix, and  $B_{ij} = A_{ji}$

## 第二周

2017年10月5日 10:54

### Multivariate linear regression

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

For convenience of notation, define  $x_0 = 1$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$

$$= \theta^T \cdot X$$

Based on inner product of vectors, we have  
theta transfer as theta transpose

$$\theta = [\theta_0 \dots \theta_n] \cdot X \Rightarrow \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

m = the number of training examples  
n = the number of features

### Cost function

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

↓

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\sum_{j=0}^n \theta_j x_j^{(i)}) - y^{(i)})^2$$

### Gradient Descent

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

↓

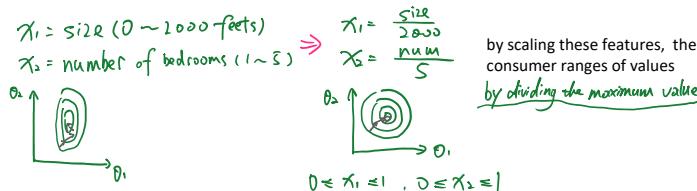
$$\theta_j := \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

通用公式

①

### Feature scaling 特征缩放

idea: make sure features are on a similar scale, then gradient descent will converge more quickly.



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

$$0 \leq x_1 \leq 3 \quad -100 = x_1 \in 100 \times \\ -2 \leq x_2 \leq 0.5 \quad -0.0001 = x_2 \in 0.0001 \times$$

Don't worry if the features are not exactly on the same scale in the same range of values.

Mean normalization 均值归一化:

replace  $x_i$  with  $\frac{x_i - \mu}{s_i}$  to make features have approximately zero mean

$$x_1 = \frac{size - 1000}{2000} \Rightarrow \text{average size} = 1000$$

$$x_2 = \frac{bedrooms - 2}{5} \Rightarrow \text{average room number} = 2$$

$$-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$$

$$x_1 \leftarrow \frac{x_1 - \mu_1}{s_1} \leftarrow \text{avg value of } x_1 \text{ in training set}$$

$s_1 \leftarrow \text{range of } x_{1\max} - x_{1\min}$

or

the standard deviation

$$x_2 \leftarrow \frac{x_2 - \mu_2}{s_2}$$

②

### Learning rate $\alpha$

Debugging: make a plot with number of iterations on the x-axis.

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



### Computing parameters analytically

Gradient descent: to minimize the cost function  $J(\theta)$ , with the iterative algorithm that take many steps, multiple iterations of gradient descent to converge the global minimum

**Normal equation:** a method to solve for  $\theta$  analytically (一步求解)

m = the number of training samples  
n = the number of features

$$X^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad y = \begin{bmatrix} y^{(0)} \\ y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \text{ m-dimensional vector}$$

e.g. if  $X^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$   $X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_1^{(m)} \\ 1 & x_2^{(1)} & \dots & x_2^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n^{(1)} & \dots & x_n^{(m)} \end{bmatrix}$   $y = \begin{bmatrix} y^{(0)} \\ y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$

`pinv(X'*X)*X'*y` this code to minimize the cost function of theta for new linear regression

Feature scaling is not necessary of normal equation, but it is important for gradient descent.

Gradient descent	Normal equation
Need to choose $\alpha$	No need to choose $\alpha$
Need many iterations	No need to iterate
Work well even when n is large	Need to compute $(X^T X)^{-1}$
$O(kn^2)$	Slow if n is very large
$N = 1,000,000$	$N = 100 / 1000$

For most sophisticated algorithm, we always choose gradient descent.

the cost of inverting matrix  
 $= O(n^3)$

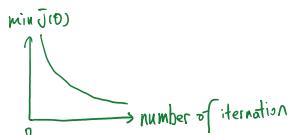
What if  $X^T X$  is non-invertible?

1. Redundant features  
 $x_1$  = size in feet?  
 $x_2$  = sine in m<sup>2</sup>  
In fact  $1m = 3.28$  feet
2. Too many features n (m <= n)

So we should use `pinv()` rather than `inv()`, the `pinv()` function will give you a value of theta if  $X^T X$  is not invertible.

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

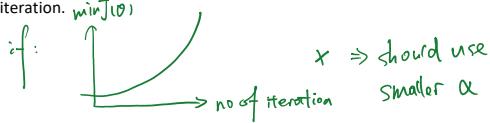
$J(\theta)$  should decrease after every iteration.



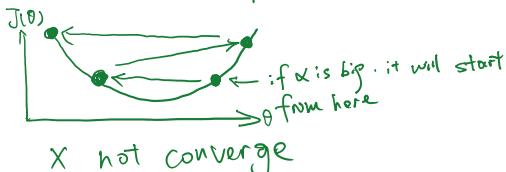
Except plot, automatic convergence test is also ok.

Declare convergence if  $J(\theta)$  decreases by less than  $10^{-3}$  (or smaller) in one iteration.

$\min J(\theta)$

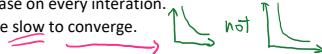


$x \Rightarrow$  should use smaller  $\alpha$



For sufficiently smaller  $\alpha$ ,  $J(\theta)$  should decrease on every iteration.

But if  $\alpha$  is too small, gradient descent can be slow to converge.



In summary:

1. If  $\alpha$  is too small, slow convergence
2. If  $\alpha$  is too large:  $J(\theta)$  may not decrease on every iteration; may not converge

how to test :  $0.001 - 0.003 - 0.01 - 0.03 - 0.1 - 0.3 - 1 \dots$

### Features and polynomial regression

e.g. housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$$



$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot \text{Area} \quad (\text{Area} = \text{frontage} \times \text{depth})$$

We can choose different polynomial regression, quadratic function or cubic function, or else

$\theta_0 + \theta_1 x + \theta_2 x^2$  will come back down

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$  will not go down

now to fit  $(y, x)$  Machinery of multivariate linear regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$= \theta_0 + \theta_1 \text{size} + \theta_2 \text{size}^2 + \theta_3 \text{size}^3$$

$$x_1 = \text{size}, x_2 = \text{size}^2, x_3 = \text{size}^3$$

size :  $1 \sim 10^3$  (e.g.)

$x^2 : 1 \sim 10^6$

$x^3 : 1 \sim 10^9$

Feature scaling:

Other choices of features:

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot \text{size} + \theta_2 \sqrt{\text{size}}$$



### 第三周

2017年8月22日 22:09

#### Classification

##### Example

1. Email: spam / not spam
  2. Online transaction: fraudulent (yes / no)
  3. Tumor: malignant / benign
- $Y \in \{0, 1\}$
- Threshold classifier output  $h_{\text{theta}}(x)$  at 0.5:  
 If  $h_{\text{theta}}(x) \geq 0.5$ , predict  $y = 1$   
 If  $h_{\text{theta}}(x) < 0.5$ , predict  $y = 0$

It is bad to use linear regression for classification problems.

Logistic regression:  $0 \leq h_{\text{theta}}(x) \leq 1$

This is a classification algorithm, not a regression algorithm, it is applied to settings where the label  $y$  is discrete value.

#### Hypothesis representation

A function to represent hypothesis when we have a classification problem.

logistic regression model: want  $0 \leq h_{\text{theta}}(x) \leq 1$

linear regression hypothesis:  $h_{\text{theta}}(x) = \theta^T x$

logistic regression hypothesis:  $h_{\text{theta}}(x) = g(\theta^T x)$

$g(z) = \frac{1}{1 + e^{-z}}$  sigmoid / logistic function

Interpretation of Hypothesis Output  
 $h_{\text{theta}}(x)$  = estimated probability that  $y = 1$  on input  $x$

example if  $x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{undersize} \\ \text{undersize} \end{bmatrix}$

then  $h_{\text{theta}}(x) = \hat{y}$

A patient with features  $x$ , the probability that  $y$  equals 1 is 0.7.

Tell patient that 70% chance of tumor being malignant.

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

Probability that  $y = 1$ , given  $x$ , parameterized by  $\theta$ .

$P(y=1 | x; \theta) \rightarrow P(y=1 | x, \theta) \cdot 1$

#### Decision boundary

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

if  $z \gg 0$ , then  $g(z) \approx 1 \Rightarrow y=1$   
 $z \ll 0$ , then  $g(z) \approx 0 \Rightarrow y=0$   
 so  $\theta^T x \gg 0 \Rightarrow y=1$   
 $\theta^T x \ll 0 \Rightarrow y=0$

So the decision boundary is the line that separates the area where  $y = 0$  and where  $y = 1$ . It is created by our hypothesis function.

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

Suppose  $-3 \leq x_1 \leq 1 \leq x_2 \leq 3$

so  $y=1$  if  $0 \leq \theta_0 + \theta_1 x_1 + \theta_2 x_2 \leq 0$

$\therefore -3 \leq x_1 \leq x_2 \geq 0$

$\therefore x_1 \geq -3 \Rightarrow \theta_1 x_1 \geq -3$

$\therefore x_2 \geq -3 \Rightarrow \theta_2 x_2 \geq -3$

$\therefore \theta_1 x_1 + \theta_2 x_2 \geq -3$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq -3$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$

$\therefore \theta_0 + \$

## 第四周

2017年8月28日 23:38

Feature number: n  
Example number: m

### Non-linear hypothesis

#### Neural network

Example: computer vision

$50 \times 50$  pixel images  $\rightarrow 2500$  pixels

$n=2500$  ( $750 \times 300$ )

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n-1} \end{bmatrix}$

with grayscale image

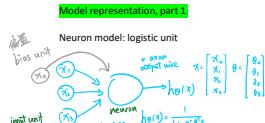
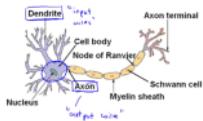
Binary features  $(x_i, x_j) \approx \frac{(1+1+1)}{2} = 2500 \approx 3$  million features

So linear regression is not a good way if the number of features n is very large.

### Neurons and the Brain

Origin: algorithms that try to mimic the brain

#### Neuron in the brain

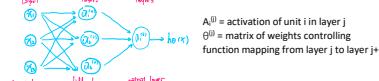


Some terminology with the same meaning for the last videos

Sigmoid (logistic) activation function:  $h(x)$

Theta: weight

Neural network



$A^{(j)}$  = activation of unit i in layer j

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

$$\begin{aligned} a_1^{(2)} &= g(\theta_{10}^{(2)}x_0 + \theta_{11}^{(2)}x_1 + \theta_{12}^{(2)}x_2 + \theta_{13}^{(2)}x_3) \\ a_2^{(2)} &= g(\theta_{20}^{(2)}x_0 + \theta_{21}^{(2)}x_1 + \theta_{22}^{(2)}x_2 + \theta_{23}^{(2)}x_3) \\ a_3^{(2)} &= g(\theta_{30}^{(2)}x_0 + \theta_{31}^{(2)}x_1 + \theta_{32}^{(2)}x_2 + \theta_{33}^{(2)}x_3) \\ h_\theta(x) &= a_1^{(2)} = g(\theta_{10}^{(2)}a_0^{(1)} + \theta_{11}^{(2)}a_1^{(1)} + \theta_{12}^{(2)}a_2^{(1)} + \theta_{13}^{(2)}a_3^{(1)}) \end{aligned}$$

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

#### Model representation, part 2

Vectorization

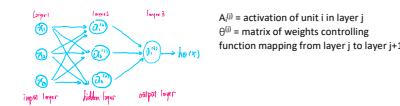
$$a_1^{(2)} = g(\theta_{10}^{(2)}x_0 + \theta_{11}^{(2)}x_1 + \theta_{12}^{(2)}x_2 + \theta_{13}^{(2)}x_3) \geq 0$$

$$a_2^{(2)} = g(\theta_{20}^{(2)}x_0 + \theta_{21}^{(2)}x_1 + \theta_{22}^{(2)}x_2 + \theta_{23}^{(2)}x_3) \geq 0$$

$$a_3^{(2)} = g(\theta_{30}^{(2)}x_0 + \theta_{31}^{(2)}x_1 + \theta_{32}^{(2)}x_2 + \theta_{33}^{(2)}x_3) \geq 0$$

$$h_\theta(x) = a_1^{(2)} = g(\theta_{10}^{(2)}a_0^{(1)} + \theta_{11}^{(2)}a_1^{(1)} + \theta_{12}^{(2)}a_2^{(1)} + \theta_{13}^{(2)}a_3^{(1)})$$

$$\begin{aligned} x &= \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \geq 0 \Rightarrow \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{bmatrix} = g(\begin{bmatrix} \theta_{10}^{(2)} & \theta_{11}^{(2)} & \theta_{12}^{(2)} & \theta_{13}^{(2)} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}) \geq 0 \\ a_0^{(2)} &= X_0^{(2)} \geq 1 \quad \text{similarly } a_1^{(2)}, a_2^{(2)}, a_3^{(2)} \geq 1 \\ \geq a_1^{(2)} &= \theta_{10}^{(2)} a_0^{(1)} + \theta_{11}^{(2)} a_1^{(1)} + \theta_{12}^{(2)} a_2^{(1)} + \theta_{13}^{(2)} a_3^{(1)} \end{aligned}$$



Neural network learning its own features

Other network architecture

Architecture means how the different neurons are connected to each other.

Summary:  $z^{(j)} = \theta^{(j-1)} \cdot X^{(j-1)}$

$$X^{(j)} = g(\theta^{(j)})$$

### Non-linear classification example: XOR / XNOR

$y = x_1 \text{ XOR } x_2$  (当且仅当  $x_1$  或  $x_2$  且只有一个等于 1 时为真)

$y = x_1 \text{ XNOR } x_2 = \text{not}(x_1 \text{ XOR } x_2)$



#### Simple example: AND

$x_1, x_2 \in \{0, 1\}$

$y = x_1 \text{ AND } x_2$

$\begin{cases} 1 & \text{if } x_1 = 1 \text{ and } x_2 = 1 \\ 0 & \text{otherwise} \end{cases}$

$h_\theta(x) = g(-3x_1 + 2x_2 + 2.0)$

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} 2.0 \\ -3.0 \\ 2.0 \end{bmatrix}$

$a_1^{(2)} = g(\theta_{10}^{(2)}x_0 + \theta_{11}^{(2)}x_1 + \theta_{12}^{(2)}x_2) = g(-3x_1 + 2x_2 + 2.0)$

$h_\theta(x) = a_1^{(2)} = g(-3x_1 + 2x_2 + 2.0)$

similar: OR

$\begin{cases} 1 & \text{if } x_1 = 1 \text{ or } x_2 = 1 \\ 0 & \text{otherwise} \end{cases}$

$h_\theta(x) = g(-2x_1 + 2x_2 + 1.0)$

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} 1.0 \\ -2.0 \\ 2.0 \end{bmatrix}$

$a_1^{(2)} = g(\theta_{10}^{(2)}x_0 + \theta_{11}^{(2)}x_1 + \theta_{12}^{(2)}x_2) = g(-2x_1 + 2x_2 + 1.0)$

$h_\theta(x) = a_1^{(2)} = g(-2x_1 + 2x_2 + 1.0)$

#### Simple example: negation (NOT x1)

$x_1 \in \{0, 1\}$

$y = \text{not}(x_1)$

$\begin{cases} 0 & \text{if } x_1 = 1 \\ 1 & \text{if } x_1 = 0 \end{cases}$

$h_\theta(x) = g(-x_1 + 1.0)$

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \begin{bmatrix} 1.0 \\ -1.0 \end{bmatrix}$

$a_1^{(2)} = g(\theta_{10}^{(2)}x_0 + \theta_{11}^{(2)}x_1) = g(-x_1 + 1.0)$

$h_\theta(x) = a_1^{(2)} = g(-x_1 + 1.0)$

#### Simple example: (NOT x1) AND (NOT x2)

$x_1, x_2 \in \{0, 1\}$

$y = x_1 \text{ AND } x_2$

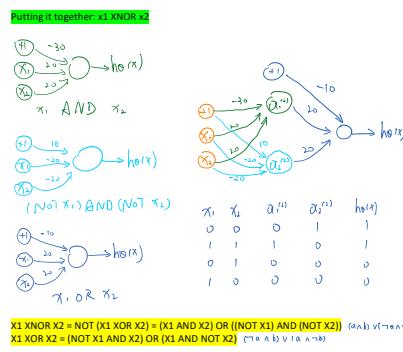
$\begin{cases} 1 & \text{if } x_1 = 1 \text{ and } x_2 = 1 \\ 0 & \text{otherwise} \end{cases}$

$h_\theta(x) = g(-x_1 + 1.0) \cdot g(-x_2 + 1.0)$

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} 1.0 \\ -1.0 \\ -1.0 \end{bmatrix}$

$a_1^{(2)} = g(\theta_{10}^{(2)}x_0 + \theta_{11}^{(2)}x_1 + \theta_{12}^{(2)}x_2) = g(-x_1 + 1.0) \cdot g(-x_2 + 1.0)$

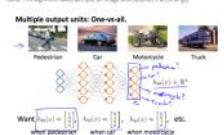
$h_\theta(x) = a_1^{(2)} = g(-x_1 + 1.0) \cdot g(-x_2 + 1.0)$



### Multi-class classification

#### Multiclass Classification

To classify 3-class multiclass classes we let our hypothesis function return a vector of values. Say we want to classify an image as either a person, a car, a motorcycle, or a truck. We can do this by defining four output nodes:



We can define a set of training classes as:

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Each  $y_i \in \{0, 1\}$  represents a different label corresponding to either a car, person, truck, or motorcycle. The input layer,  $x$ , is combined with some new information which leads to our final hypothesis function. The output layer:

$$h_\theta(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_n(x) \end{bmatrix}$$

Our resulting hypothesis for one set of inputs may look like:

$$h_\theta(x) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

## 第五周

2017年9月4日 1:01

### Cost function and backpropagation

$L$  = total no. of layers in network  
 $S_i$  = no. of units (not counting bias unit) in layer  $i$   
 $K$  = number of output units / classes

① Binary classification  
 $y = 0 \text{ or } 1 \rightarrow 1 \text{ output unit} \rightarrow \theta \rightarrow h(\theta)$

② Multi-class classification ( $K$  classes)  
 $y \in \mathbb{R}^K \quad \exists [1] [0] [1] \quad K=3$   
 car bike plane  
 $K$  output units.  $h(\theta) \in \mathbb{R}^K$  (generally  $K > 3$ )

Cost function  
 Logistic regression:

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

*do not regularize the bias term  $\theta_0$*

Neural network:

$$h_\theta(x) \in \mathbb{R}^K, h_\theta(x)_j = \text{ith output} \quad \theta_0 \text{ bias}$$

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)})_k) + (1-y_k^{(i)}) \log(1-h_\theta(x^{(i)})_k) \right] + \frac{\lambda}{2m} \sum_{j=1}^{L-1} \sum_{k=1}^{S_j} \sum_{l=1}^{S_{j+1}} (\theta_{jkl})^2$$

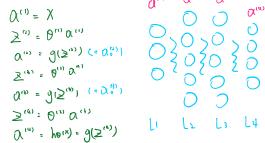
Backpropagation algorithm: a neural-network terminology for minimizing our cost function

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)})_k) + (1-y_k^{(i)}) \log(1-h_\theta(x^{(i)})_k) \right] + \frac{\lambda}{2m} \sum_{j=1}^{L-1} \sum_{k=1}^{S_j} \sum_{l=1}^{S_{j+1}} (\theta_{jkl})^2$$

Gradient computation

Given one training example  $(x, y)$

Forward propagation:



Intuition:  $\delta_j^{(l)}$  = 'error' of code  $j$  in layer  $l$

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

$$\begin{aligned} \delta_j^{(4)} &= y_j - h_\theta(x) \\ \delta_j^{(4)} &= (y_j - h_\theta(x)) \cdot (1 - h_\theta(x)) \\ \delta_j^{(4)} &= (y_j - h_\theta(x)) \cdot (1 - h_\theta(x)) \\ \delta_j^{(4)} &= (y_j - h_\theta(x)) \cdot (1 - h_\theta(x)) \\ \delta_j^{(4)} &= (y_j - h_\theta(x)) \cdot g(z^{(4)}) \end{aligned}$$

*Intuition:  $\delta_j^{(l)}$  =  $\delta_j^{(l+1)} \cdot g'(z^{(l)})$*

no  $\delta^{(1)}$  because 1st layer correspond to input layer. observed as training set, so have no error

$$\frac{\partial}{\partial \theta_j} J(\theta) = \alpha_j^{(t)} \delta_j^{(t+1)} \quad (\text{when ignore } \lambda \text{ regularization})$$

### Backpropagation algorithm

Training set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$   
 Set  $\Delta_j^{(l)} = 0$  for all  $i, i_l$  (use  $\Delta$  to compute  $\frac{\partial}{\partial \theta_j} J(\theta)$ )

Upper case:  $\Delta$  Lower case:  $\delta$

For  $i = 1 \text{ to } m$   
 Set  $a^{(i)} = x^{(i)}$   
 Perform forward propagation to compute  $a^{(1)}$  for  $l = 2, 3, \dots, L$   
 Using  $y^{(i)}$ , compute  $\delta^{(L)} = a^{(L)} - y^{(i)}$   
 Compute  $\delta^{(1)}, \delta^{(2)}, \dots, \delta^{(L)}$

$$\Delta_j^{(l)} = \Delta_j^{(l)} + a_j^{(l)} \delta_j^{(l+1)}$$

(vectorwise)  $\Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)}(a^{(l)})^\top$

$$\begin{aligned} \Delta_j^{(0)} &= \frac{1}{m} \Delta_j^{(0)} + \lambda \theta_j^{(0)} \text{ if } j \neq 0 \\ \Delta_j^{(0)} &= \frac{1}{m} \Delta_j^{(0)} + \lambda \theta_j^{(0)} \text{ if } j = 0 \end{aligned}$$

The capital-delta matrix  $D$  is used as an accumulator to add up our values as we go along and eventually compute our partial derivative.

Suppose you have two training examples  $(x^{(1)}, y^{(1)})$  and  $(x^{(2)}, y^{(2)})$ . Which of the following is a correct sequence of operations for computing the gradient? (Below, FP = forward propagation, BP = back propagation).

- ① FP using  $x^{(1)}$  followed by BP using  $x^{(2)}$ . Then BP using  $y^{(1)}$  followed by BP using  $y^{(2)}$ .
- ② FP using  $x^{(1)}$  followed by BP using  $y^{(2)}$ . Then FP using  $x^{(2)}$  followed by BP using  $y^{(1)}$ .
- ③ BP using  $y^{(1)}$  followed by FP using  $x^{(1)}$ . Then BP using  $y^{(2)}$  followed by FP using  $x^{(2)}$ .
- ④ FP using  $x^{(1)}$  followed by BP using  $y^{(1)}$ . Then FP using  $x^{(2)}$  followed by BP using  $y^{(2)}$ .

### Backpropagation intuition

Forward propagation



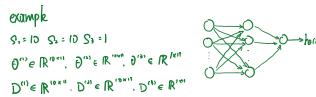
$\Sigma$ : the weighted sum of inputs of the input units  
 then apply the sigmoid function  $\Rightarrow \alpha$

# Backpropagation in practice

### Unrolling parameters from matrices into vectors

```
Function [jVal, gradientvec] = costFunction(thetaVec)
optTheta = fminunc(@costFunction, initialTheta, options)
```

Neural Network (L=4)  
 $\theta(1), \theta(2), \theta(3)$  = matrices(Theta1, Theta2, Theta3)  
 $D(1), D(2), D(3)$  = matrices(D1, D2, D3)  
 "unroll" into vectors

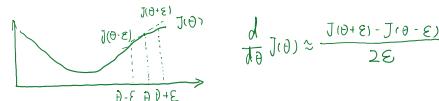


```
thetaVec = [Theta1(:); Theta2(:); Theta3(:)];
Dvec = [D1(:); D2(:); D3(:)];
Theta1 = reshape(thetaVec(1:110), 10, 11);
Theta2 = reshape(thetaVec(111:220), 10, 11);
Theta3 = reshape(thetaVec(221:231), 1, 11);
```

```
Function [jVal, gradientvec] = costFunction(thetaVec)
From thevec, get  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$ 
Use forward prop/backprop to compute D1, D2, D3 and J(theta)
Unroll D1, D2, D3 to get gradientvec
```

### Gradient checking

Numerical estimation of gradients



Implement:

gradapprox = (J(theta + epsilon) - J(theta - epsilon)) / (2 \* epsilon)

$\theta \in \mathbb{R}^n$   $\theta$  is unrolled version of  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$

$\theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{J(\theta + \epsilon, \theta_1, \dots, \theta_j, \dots, \theta_n) - J(\theta - \epsilon, \theta_1, \dots, \theta_j, \dots, \theta_n)}{2\epsilon}$$

Epsilon = 1e-4

For  $i = 1 \dots n$  is the dimension of our parameter of vector theta  
 $\text{Thetaplus} = \text{theta};$   
 $\text{Thetaplus}(i) = \text{thetaplus}(i) + \text{epsilon};$   
 $\text{Thetaminus} = \text{theta};$   
 $\text{Thetaminus}(i) = \text{thetaminus}(i) - \text{epsilon};$   
 $\text{Gradapprox}(i) = (J(\text{thetaplus}) - J(\text{thetaminus})) / (2 * \text{epsilon});$

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{i-1} \\ \theta_i \\ \theta_{i+1} \\ \vdots \\ \theta_n \end{bmatrix}$$

check that gradApprox  $\approx$  Dvec

Implementation note:

1. Implement backprop to compute Dvec (unrolled D1, D2, D3)
2. Implement numerical gradient check to compute gradapprox
3. Make sure they give similar values ( $\approx$ )
4. Turn off gradient checking, using backprop code for learning

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.

### Random initialization - initial value of $\theta$

For gradient descent and advanced optimization method, need initial value of  $\theta$

OptTheta = fminunc(@costfunction, initialTheta, options)

Consider gradient descent

Zero initialization

$$\begin{aligned} \theta_{ij}^{(l)} &= 0 \quad \text{for all } i, j, l \quad \theta_{ij}^{(l)} \\ \alpha_i^{(l)} &= \alpha_i^{(l)}, \text{ also } \delta_i^{(l)} = \delta_i^{(l)} \\ \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) &\approx \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) \end{aligned}$$

After each update, parameters corresponding to inputs going into each of two hidden units are identical, so one way to avoid it is to break it

Symmetry breaking - 对称性破缺 (one way for random initialization)  
 Initialize each  $\theta_{ij}^{(l)}$  to a random value in  $[-\epsilon, \epsilon]$

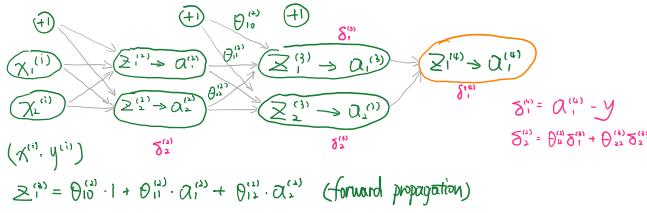
```
Theta1 = rand(10,11) * (2 * init_epsilon) - init_epsilon
Theta2 = rand(10,11) * (2 * init_epsilon) - init_epsilon
Theta3 = rand(1,11) * (2 * init_epsilon) - init_epsilon
```

*(pick a random number  $r = \text{rand}(1,1) + (2 * \text{init\_epsilon}) - \text{init\_epsilon}$   
 set  $\theta_{ij}^{(l)} = r$  for all  $i, j, l$ )*  
*because this fails to break symmetry X*

### Put it together

Pick a network architecture (connectivity pattern between neurons)

### Forward propagation



When the algorithm is applying forward propagation, what is backpropagation doing?

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{j=1}^{n_l} \sum_{k=1}^{n_{l+1}} (\theta_{jk}^{(l)})^2$$

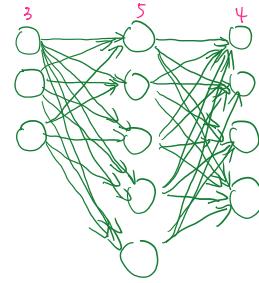
Focus on a single example  $x^{(0)}, y^{(0)}$ , the case of 1 output unit and ignoring regularization ( $\lambda = 0$ )

$$\delta_j^{(l)} = \text{error of cost for } a_j^{(l)} \text{ (unit } j \text{ in layer } l\text{). formally } \delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(x^{(0)}) \text{ (for } j \geq 0\text{)}$$

where cost( $\hat{y} = y^{(0)} \log h_\theta(x^{(0)}) + (1-y^{(0)}) \log(1-h_\theta(x^{(0)}))$ )

### Put it together

Pick a network architecture (connectivity pattern between neurons)



No. of input units: dimension of features  $x^{(0)}$

No. of 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)

### Training a neural network

1. Randomly initialize weight
2. Implement forward propagation to get  $h_\theta(x^{(0)})$  for any  $x^{(0)}$
3. Implement code to compute cost function  $J(\theta)$
4. Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$

To calculate backprop, we use loop

For  $i = 1:m$   $\underbrace{(x^{(i)}, y^{(i)})}_{\text{Perform forward propagation and backpropagation using example } (x^{(i)}, y^{(i)})} \cdots (x^{(m)}, y^{(m)})$

get activations  $a^{(l)}$ , and delta terms  $\delta^{(l)}$  for  $l = 2, \dots, L$

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

$$\text{compute } \frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$$

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  $\theta$



{+}

it's who has

ation to predict y

d no other features

ently predict y?

ers: logistic/linear

i units.

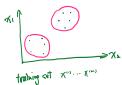


# 第八周

2017年10月2日 0:06

## Clustering

### Unsupervised learning



- Application of clustering
1. Marketing segmentation
  2. Social network analysis
  3. Organize computing clusters
  4. Astronomical data analysis

### K-means algorithm



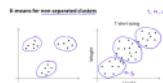
Input:

- K (number of clusters)
- Training set  $(x_1^m, x_2^m)$
- $x_i^m \in \mathbb{R}^n$  (only last coordinate)

Randomly initialize K clusters centroids  $\mu_1, \dots, \mu_K \in \mathbb{R}^n$

Repeat:  
 cluster assignment For  $i = 1$  to  $m$   $c(i) :=$  index (From 1 to K) of cluster centroid closest to  $x(i)$   
 move centroid For  $k = 1$  to  $K$   $\mu_k =$  average(mean) of the points assigned to cluster k  
 $\lambda^m = \frac{1}{m} [x_1^m, x_2^m, \dots, x_m^m] \in \mathbb{R}^n$ , n维向量

### K-means for non-separated clusters



Optimization objective:  
 $c^0 =$  index of cluster  $(1, \dots, K)$  to which example  $x^0$  is currently assigned  
 $\mu_k =$  cluster centroid  $\mu_k \in \mathbb{R}^n$

$\mu_k$  = cluster centroid of cluster to which example  $x^0$  has been assigned

Optimization objective:  
 $J(C^0, \dots, C^m, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x_i^0 - \mu_{c(i)}\|^2$   
 min  $J(C^0, \dots, C^m, \mu_1, \dots, \mu_K)$  —> related distortion 生直感

### Random initialization

should have  $K < m$

randomly pick K training examples

set  $\mu_1, \dots, \mu_K$  equal to those K examples

### Local optima

try more times for initialization

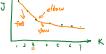
```
For i = 1 to 100 {
    randomly initialize K-means
    Run K-means. Get  $C^0, \dots, C^m, \mu_1, \dots, \mu_K$  (initial)
    Compute cost function (distortion)
     $J(C^0, \dots, \mu_1, \dots, \mu_K)$ 
}
```

pick the best one based on  $J \rightarrow \downarrow$

### Choosing the numbers of clusters

What is the right value of K?

### Elbow method



## Dimensionality reduction - data compression

### Data compression

reduce data from 2D to 1D  
 $\mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n \times \mathbb{R}$   
 $\mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n \times \mathbb{R}$   
 $\mathbb{R}^n \rightarrow \mathbb{R}^m$

### Reduce data from 3D to 2D

Data Compression  
 Reduce data from 3D to 2D  
 Reduce data from 2D to 1D

A lower dimensional object ( $\mathbb{R}^n \times \mathbb{R}^m$ ) of m examples where  $\mathbb{R}^n \times \mathbb{R}^m$  for some value of n and m < n  
 original  $(x_1^m, x_2^m, \dots, x_n^m) \in \mathbb{R}^n$

## Principal component analysis (PCA)

### PCA problem formulation

length of data blue line  
 segment call projection vector  
 投影矢量

PCA is to find a surface onto which to project the data.

### Principal Component Analysis (PCA) problem formulation

$3D \rightarrow 2D$   
 $R = 2$

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

### PCA algorithm

#### Data preprocessing

Training set  $X^0, X^1, \dots, X^m$   
 pre-processing: feature scaling — mean normalization  
 $\mu_j = \frac{1}{m} \sum_{i=1}^m x_{ij}$   
 replace each  $x_{ij}$  with  $x_{ij} - \mu_j$

#### Principal Component Analysis (PCA) algorithm

Reduce data from 2D to 1D  
 Reduce data from 3D to 2D  
 Reduce data from n-dimensions to k-dimensions

Compute "covariance matrix" 协方差  
 $\Sigma = \frac{1}{m} \sum_{i=1}^m (X^0)(X^0)^T$  Sigma

Compute "eigenvectors" of matrix  $\Sigma$ :  
 $[U, S, V] = svd(\Sigma)$   
 singular value decomposition 特征值分解  
 eigenvalues  
 covariance under diagonal matrix  
 properly symmetric positive definite (Hessian)  
 $U = [u^1 | u^2 | \dots | u^k] \in \mathbb{R}^{n \times k}$   
 $U = LD \Rightarrow U^T = U^{-1}$   
 $X \in \mathbb{R}^n \Rightarrow Z \in \mathbb{R}^k$   
 $Z = U \Sigma^{\frac{1}{2}} V^T X \in \mathbb{R}^{k \times 1}$

### Summary

After mean normalization (ensure every feature has zero mean and optional feature scaling):

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (X^0)(X^0)^T \quad X = \begin{bmatrix} (X^0)^T \\ \vdots \\ (X^m)^T \end{bmatrix}$$

$$\Sigma = I / m + X^T X$$

$$[U, S, V] = svd(\Sigma)$$

$$Underline = U(1, :k)$$

$$Z = Underline \times X$$

## Applying PCA

### Reconstruction from compressed representation

Reconstruction from compressed representation  
 Applying PCA  
 $X \approx Underline \times Z$  (Underline may not be invertible)

### Choosing the number of principal component analysis

Average squared projection error 平均平方映射误差  
 Total variation in the data  $\frac{1}{m} \sum_{i=1}^m \|x_i^m\|^2$

Typically, choose k to be smallest value so that:

$$\frac{\sum_{i=1}^m \|x_i^m - Underline_i\|^2}{\sum_{i=1}^m \|x_i^m\|^2} \leq 0.01 \text{ or } 0.05$$

99% of variance is retained 99%的方差被保留

### Choosing k

#### Algorithm:

Try PCA with  $k=1, \dots, 2, \dots, 3, \dots$

compute Underline  $\Sigma^{\frac{1}{2}}, \dots, \Sigma^{\frac{1}{2}}, \Sigma^{\frac{1}{2}}, \dots, \Sigma^{\frac{1}{2}}$

check if  $\frac{\sum_{i=1}^m \|x_i^m - Underline_i\|^2}{\sum_{i=1}^m \|x_i^m\|^2} \leq 0.01$

$[U, S, V] = svd(\Sigma)$   
 $S = \begin{bmatrix} s_1 & & & \\ & s_2 & & \\ & & \ddots & \\ & & & s_m \end{bmatrix}$   
 for given k:  
 $1 - \frac{s_1^2}{s_1^2 + s_2^2 + \dots + s_k^2} \geq 0.95$

### Summary

$$[U, S, V] = svd(\Sigma)$$

pick the smallest value of k for which

$$\frac{s_1^2}{s_1^2 + s_2^2 + \dots + s_k^2} \geq 0.95 \quad (95\% \text{ of variance retained})$$

### Advice for applying PCA

#### Supervised learning speedup

$$(X^0, y^0), (X^1, y^1), \dots, (X^n, y^n) \quad X^i \in \mathbb{R}^{n \times p}, y^i \in \mathbb{R}^{n \times 1}$$

#### Extract inputs:

Unlabeled dataset:  $X^0, \dots, X^m \in \mathbb{R}^{n \times p}$   
 $y^0, \dots, y^m \in \mathbb{R}^{n \times 1}$

New training set:  
 $(X^0, y^0), \dots, (X^n, y^n) \quad \text{for } i = 1 \dots n$

Note: mapping  $x^0 \rightarrow y^0$  should be defined by running PCA only on the training set. This mapping can be applied as well to the examples  $x_0^0$  and  $x_{test}^0$  in the cross validation and test sets.

#### Compression

- Reduce memory/disk needed to store data
- Speed up learning algorithm

#### Visualization

Bad use of PCA: to prevent overfitting  
 This might work OK, but is not a good way to address overfitting. Use regularization instead.

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m \|y^i - \theta^T x^i\|^2 + \frac{\lambda}{2} \|\theta\|^2$$

Bad use of PCA: use it where it shouldn't be

#### PCA is sometimes used where it shouldn't be

Design of ML system

→ Get training set  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

→ Run PCA to reduce  $x^{(i)}$  in dimension to get  $x^{(i)}$

→ Train logistic regression on  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

→ Run k-means to  $x^{(i)}$ . Run kNN on  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

→ How about doing the whole thing without using PCA?

→ Before implementing PCA, first try running whatever you want to do with the original/raw data. Only if that doesn't do what you want, then implement PCA and consider using it.



Anomaly detection examples:

**Fraud detection:**  
 $x^T \in \text{features of users' activities}$   
Model  $p(x)$  from data  
Identify unusual users by checking which have  $p(x) < \epsilon$

**Manufacturing**

Monitoring computers in a data center:  
 $x^T \in \text{features of machine } i$   
 $x_1 = \text{memory use}, x_2 = \text{number of disk accesses}$   
 $x_3 = \text{CPU load} / \text{network traffic}$



Parameter estimation 参数估计

$$\text{Dataset: } (x^{(1)}, y^{(1)})^T, (x^{(2)}, y^{(2)})^T, \dots, (x^{(m)}, y^{(m)})^T$$

$$\text{Density: } f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\text{Probability: } p(x) = \int_{\text{anomalous}} f(x) dx$$

$$\hat{\mu} = \bar{x}, \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

**Algorithm:**

- Choose features  $x_i$  that you think might be indicative of anomalous examples
- Fit parameters  $\mu_1, \dots, \mu_n, \dots, \mu_d$
- Given new examples  $x$ , compute  $p(x)$ :

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x_1 - \mu_1)^2}{2\sigma_1^2}} \times \dots \times \frac{1}{\sqrt{2\pi}\sigma_d} e^{-\frac{(x_d - \mu_d)^2}{2\sigma_d^2}}$$

Assume  $\sigma_1 = \dots = \sigma_d$ **Building an anomaly detection system****Developing and evaluating an anomaly detection system**

The importance of real-number evaluation

When developing a learning algorithm (choosing features), making decisions is much easier if we have a way of evaluating our learning algorithm.

Assume we have some labeled data, of anomalous and non-anomalous examples.  $y=0$  normal,  $y=1$  anomalous.

Training set:  $X^T \times Y^T$  ( $X^T$  rows are  $x$ )Cross validation:  $(X^T_{cv} \times Y^T_{cv}) \times (X^T_{test} \times Y^T_{test})$ Test set:  $(X^T_{test} \times Y^T_{test})$ 

Aircraft engines motivating example

From model  $p(x)$ 

CV: 2000 good engines, 10 anomalous engines

Test: 2000 good engines, 10 anomalous engines

Algorithm evaluation

Fit model  $p(x)$  on training set:  $X^T \times Y^T$ On a cross-validation/test example  $x$ , predict $\Pr[\text{anomalous}] = p(x)$ 

Possible evaluation metrics:

- True positive, false positive, false negative, true negative
- Precision:  $\Pr[\text{positive} | \text{actual positive}]$
- Recall:  $\Pr[\text{actual positive} | \text{predicted positive}]$

Can also use cross-validation set to choose parameters

**Anomaly detection vs supervised learning****Anomaly detection**

Very small number of positive examples ( $y=1$ ) ( $0-20$  is common)

Large number of negative ( $y=0$ ) examples

Many different "types" of anomalies. Hard for any algorithm to learn from positive examples what anomalies look like

Future anomalies may look nothing like any of the anomalous examples we've seen so far

**Application**

- Fraud detection
- Manufacturing
- Monitor machines in data center

**Choosing what features to use**

Gaussian curve - bell shaped curve - normal distribution

Non-Gaussian features

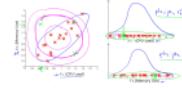
Error analysis for anomaly detection

Want  $p(x)$  large for normal examples $p(x)$  small for anomalous examplesMost common problem:  $p(x)$  is comparable (both large) for normal and anomalous examples

**Question:** suppose anomaly detection algorithm is comparing the probability of  $p(x)$  for many normal examples and for anomalous examples in your cross validation dataset. The way to help change the algorithm is to try coming up with more features to distinguish between the normal and the anomalous examples.

**Multivariate Gaussian Distribution  $\mathcal{N}(\mu, \Sigma)$** 

Motivating example: Multivariate Gaussians in a data center



700 R^2, model poor at w/sgd

Parameter:  $p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ (2.1) want to decrease  $\Pr[\text{actual } y=1 | \text{predicted } y=1]$  (precision)(2.2) want to increase  $\Pr[\text{actual } y=0 | \text{predicted } y=0]$  (recall)

Multivariate Gaussian (Normal) examples

700 R^2, model poor at w/sgd

Parameter:  $p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ (2.1) want to decrease  $\Pr[\text{actual } y=1 | \text{predicted } y=1]$  (precision)(2.2) want to increase  $\Pr[\text{actual } y=0 | \text{predicted } y=0]$  (recall)

Multivariate Gaussian Distribution

parameters:  $\mu, \Sigma$ 

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

Parameters:  $\mu$ ,  $\Sigma$ 

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

Anomaly detection with the multivariate Gaussian

1. Fit model  $p(x)$  by  $\text{sgd}$ 

$$\hat{x} = \arg \min_{x \in \mathbb{R}^d} \Pr[\text{actual } y=1 | \text{predicted } y=1]$$

2. Given a new example  $x$ , compute

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

3. Given  $x$  and  $y$ ,  $\Pr[\text{actual } y=1 | \text{predicted } y=1]$ 

Relationship to original model

The contours of the Gaussian are always axis aligned

Algorithm:

1. Fit model  $p(x)$  by  $\text{sgd}$ 2. Given a new example  $x$ , compute

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

3. Given  $x$  and  $y$ ,  $\Pr[\text{actual } y=1 | \text{predicted } y=1]$ 

Contour plot of Gaussian

Algorithm:

1. Fit model  $p(x)$  by  $\text{sgd}$ 2. Given a new example  $x$ , compute

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

3. Given  $x$  and  $y$ ,  $\Pr[\text{actual } y=1 | \text{predicted } y=1]$ 

Contour plot of Gaussian

**Collaborative Filtering****Problem formulation**

(the user rates item zero to five stars)

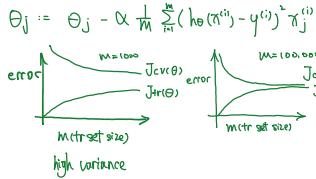
User: Movie: Rating: User: Movie: Rating:

# 第十周

2017年10月18日 18:01

## Gradient descent with large datasets

### Learning with large datasets



## Stochastic gradient descent 随机梯度下降

### Linear regression with gradient descent

$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$   
 $J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Repeat {  
 $\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$   
 for every  $j=0, \dots, n$   
}

→ called batch gradient descent (批量梯度下降)

### Stochastic gradient descent

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

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

1. randomly shuffle dataset
  2. repeat { for  $i=1 \dots m$  {
  $\theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$  just use one at a time  
 (for every  $j=0 \dots n$ )
 }
- $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$

**Stochastic gradient descent**

1. Randomly shuffle (reorder) training examples
2. Repeat { outer loop: 1 ~ 10 times  
 for  $i = 1, \dots, m$  {
  $\Rightarrow \theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$   
 (for every  $j = 0, \dots, n$ )
 } depend on number of dataset
 }

continuously wander around the optimum

## Mini-Batch gradient descent - 小批量梯度下降

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

$b = \text{mini-batch size}$  e.g.  $b=10$ , range 2 ~ 100

Get 10 examples  $(x^{(1)}, y^{(1)}), \dots, (x^{(10)}, y^{(10)})$

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

### Mini-batch gradient descent

Say  $b=10, m=1000$ . problem: must have good vectorized implementation.

Repeat {

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

$\Rightarrow \theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_{\theta}(x^{(k)}) - y^{(k)}) x_j^{(k)}$   
 (for every  $j = 0, \dots, n$ )

}

## Stochastic gradient descent convergence

Checking for convergence

Batch gradient descent:

## Advanced topics

### Online learning - 在线学习机制

#### 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  $(x, y)$ : use once  $(x^{(i)}, y^{(i)})$ , then push it

$$\theta_j := \theta_j - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} (j=0, \dots, n)$$

} Can adapt to change user preference

Other online learning example:

#### Product search (learning to search)

User searches for "iPhone 6"

Have 100 phones in store. Will return 10 results.

→  $x$  = 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

→ Learn  $p(y=1|x;\theta)$  predict CTR (click-through rate)

Use to show user the 10 phones they are most likely

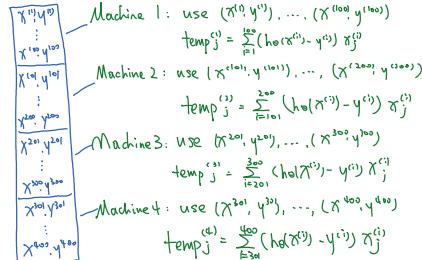
to click on

Other examples:

- choosing special offers to show user;
- customized selection of news articles;
- product recommendation

## Map reduce and data parallelism 映射约减

$$\text{Batch gradient descent } \theta_j := \theta_j - \alpha \frac{1}{400} \sum_{i=1}^{400} (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$



Combine:

$$\theta_j := \theta_j - \alpha \frac{1}{400} (\text{temp}^{(1)} + \text{temp}^{(2)} + \text{temp}^{(3)} + \text{temp}^{(4)})$$

## Map reduce and summation over the training set

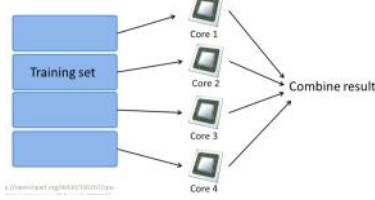
Many learning algorithms can be expressed as computing sums of functions over the training set.

e.g. for advanced optimization, with logistic regression, need:

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

$$\frac{\partial}{\partial \theta_j} J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

#### Multi-core machines



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

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

Stochastic gradient descent:

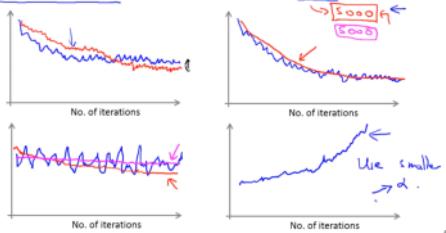
$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

During learning, compute  $\text{cost}(\theta, (x^{(i)}, y^{(i)})$ ) before updating  $\theta$  using  $(x^{(i)}, y^{(i)})$

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

### Checking for convergence

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



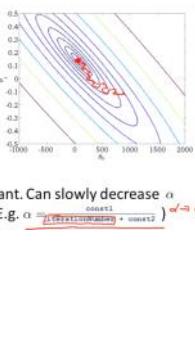
### Stochastic gradient descent

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

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

1. Randomly shuffle dataset.

2. Repeat {  
 for  $i := 1, \dots, m$  {  
 $\theta_j := \theta_j - \alpha(h_{\theta}(x^{(i)}) - y^{(i)})x_j^{(i)}$   
 for  $j = 0, \dots, n$   
 }  
 }



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

# 第十一周

2017年10月19日 1:51

Photo OCR – Photo optical character recognition 照片光学字符识别

Problem description and pipeline 机器学习流水线

Photo OCR pipeline

1. Text detection



2. Character segmentation

ANTIQUES

3. Character classification

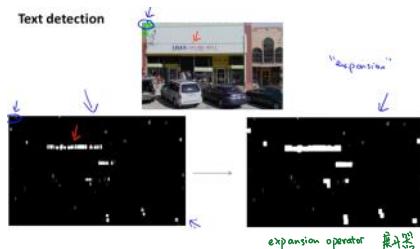


Sliding windows classifier

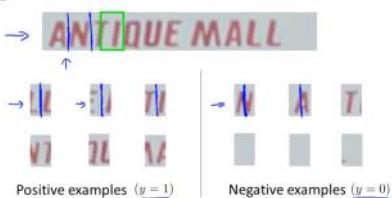
Sliding window detection



Text detection

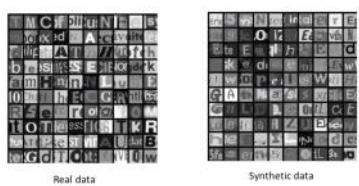


1D Sliding window for character segmentation

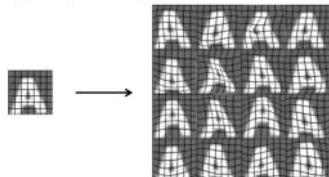


Getting lots of data and artificial data

Artificial data synthesis for photo OCR



Synthesizing data by introducing distortions



Summary: main topic

1. Supervised learning

- Linear regression, logistic regression, neural network, SVM

2. Unsupervised learning

- K-means, PCA, Anomaly detection

3. Special applications / special topics

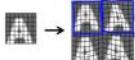
- Recommender systems, large scale machine learning

4. Advice on building a machine learning system

- Bias / variance, regularization; deciding what to work on next: evaluation of learning algorithms, learning curves, error analysis, ceiling analysis

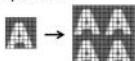
### Synthesizing data by introducing distortions

→ Distortion introduced should be representation of the type of 要具有代表性的噪声  
noise/distortions in the test set.



→ Audio:  
Background noise,  
bad cellphone connection

→ Usually does not help to add purely random/meaningless noise to your data.



$x_i$  = intensity (brightness) of pixel  $i$   
 $x_i \leftarrow x_i + \text{random noise}$

### Discussion on getting more data

1. Firstly, 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
  - "Crows source" (e.g. Amazon mechanical turk)

### Ceiling analysis: what part of the pipeline to work on next - 上限分析

Image → text detection → character segmentation → character recognition

Q: which part is the most important?

Component Accuracy

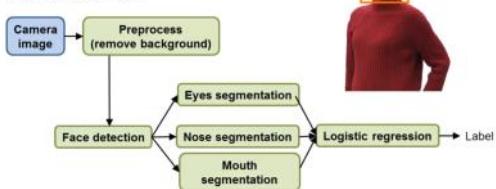
overall system	72%
text detection	89
character seg	90
character rec	100

### Another ceiling analysis example 上限分析

#### Face recognition from images (artificial example)

##### Another ceiling analysis example

Face recognition from images  
(Artificial example)



Component	Accuracy
Overall system	85%
Preprocess (remove background)	85.1
Face detection	91
Eyes segmentation	95
Nose segmentation	96
Mouth segmentation	97
Logistic regression	100