

Machine Learning

by Sunny

Oct 29th



Week 1

Linear Regression with One Variable

1. Model and Cost Function

1. cost Function

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

$\hat{\theta}_0, \hat{\theta}_1, h_{\theta}(x)$ is close to y .

$$\text{minimize}_{\theta_0, \theta_1} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

m : training examples.

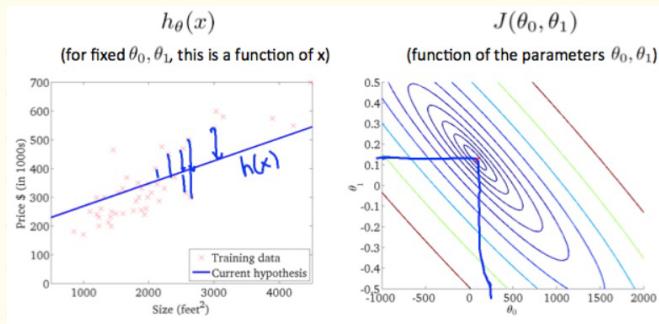
Cost Function

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

goal: minimize $J(\theta_0, \theta_1)$.

Square Cost Function

2. Contour plots.



Week 1

Linear Regression with One Variable

2. Parameter Learning

1. Gradient Descent

outline :

① start with some θ_0, θ_1 ,

② keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$

convergence 45min.

⇒ tips.

$a := b$ assignment eg. $a := a + 1$ ✓

$a = b$. truth assertion. eg. $a = a + 1$ ✗

Gradient descent algorithm

repeat until convergence {

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

}

α : learning rate

下降速率.

Correct: Simultaneous update

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

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

→ $\theta_0 := \text{temp0}$

→ $\theta_1 := \text{temp1}$

Incorrect:

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

→ $\theta_0 := \text{temp0}$

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

→ $\theta_1 := \text{temp1}$

$$\Rightarrow \frac{\partial}{\partial \theta_j} \text{梯度 (微积分)}$$

a) too small : gradient descent slow

b) too big : may overshoot the minimum

越近越小值. 简小 a.

达到 local minimum

达到 minimum.

达到 local optimum

θ_0, θ_1 .

Week 1

Linear Regression with One Variable

2. Gradient Descent For Linear Regression

G-D algorithm

repeat until convergence {

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

(for $j = 0$ and $j = 1$)

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

}

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

统计学

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

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

G-D algorithm

repeat until convergence {

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

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

}

∴ 基于原训练集的梯度

"This method looks at every example in the entire training set, and is called

batch gradient descent"

凸函数.

Convex function

Week 2

Linear Regression with Multiple Variables

1. Multivariate Linear Regression

1. Multiple Features

Notation:

n = number of features

$x^{(i)}$ = input (features) of i^{th} training example.

$x_j^{(i)}$ = value of feature j in i^{th} training example.

Superscript
Subscript

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

方程组 $x_0 = 1$

$$\text{U} h_{\theta}(x) = [\theta_0 \ \theta_1 \ \dots \ \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

transpose 转置

2. Gradient Descent for Multiple Features

Cost function $\rightarrow \hat{\theta}$ (n+1维向量)

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

$$\begin{aligned} \text{变形 } J(\theta) &= \frac{1}{2m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2 \text{ 方} \\ &= \frac{1}{2m} \sum_{i=1}^m \left(\sum_{j=0}^n \theta_j x_j^{(i)} - y^{(i)} \right)^2 \end{aligned}$$

gradient descent

repeat until convergence :

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

Week 2

Linear Regression with Multiple Variables

3. Gradient Descent in Practise

<1> Feature Scaling & Mean Normalization

$$\text{方程} -1 \leq x_{ii} \leq 1 \quad \text{or} \quad 0.5 \leq x_{ii} \leq 0.5$$

$$方程 k_i := \frac{x_i - \mu_i}{s_i}$$

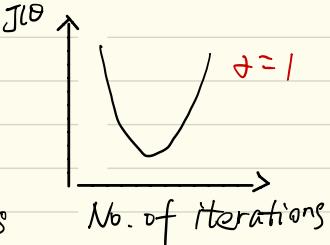
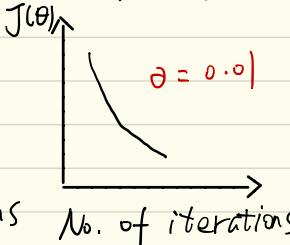
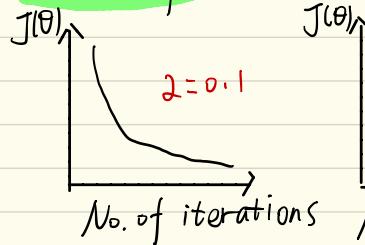
$$\mu_i: \bar{x} \quad s_i: x_i \text{ 标准差}$$

特征缩放的公式 - 方程

give or take a few
方程 -1 到 1.

<2> Learning Rate

An example, $\alpha = 0.01, 0.1, 1$.



Method

try many α

triple

$$\dots 0.001 \uparrow 0.01 \uparrow 0.1 \uparrow 1 \dots$$

0.003

Automatic convergence test

if θ decrease less than E ✓ (difficult)

To summarize:

If α is too small: slow convergence.

If α is too large: may not decrease on every iteration and thus may not converge.

Week 2

Linear Regression with Multiple Variables

4. Features and Polynomial Regression

(Aware of) combine features into one.

特征 注意形状

$$h_{\theta}(k) = \theta_0 + \theta_1 k_1 + \theta_2 k_2^2 \quad (\text{quadratic function})$$

$$h_{\theta}(k) = \theta_0 + \theta_1 k_1 + \theta_2 k_2^3 + \theta_3 k_3^3 \quad (\text{cubic function})$$

$$h_{\theta}(k) = \theta_0 + \theta_1 k_1 + \theta_2 \sqrt{k_1} \quad (\text{square root function})$$

Feature Scaling is important!

2. Computing Parameters Analytically

1. Normal Equation

gradient descent 和 normal equation

本质都是一种求最优的方法。

~~feature scaling~~

Examples: $m = 4$.

x_0	Size (feet ²) x_1	Number of bedrooms x_2	Number of floors x_3	Age of home (years) x_4	Price (\$1000) 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

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$

$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$

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

$$\text{pinv}(X^T X) * X^T * y$$

m-dimensional vector

comparison of the two method

Week 2

Linear Regression with Multiple Variables

$C-D$	$N-E$
need to choose α	\checkmark
need many iterations	iterate
$O(kn^2)$	$O(n^3)$
work well when $n \uparrow$	slow when $n \uparrow$
generally $n > 10,000$	

*2. Normal Equation Noninvertibility (正常方程不可逆)

normal equation

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

~~若 $X^T X$ 不可逆,~~

例 Octave 里计算 $\text{pinv}(X^T * X) * X^T * y$

Reasons that cause $N-E$ singular

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

- Redundant features (linearly dependent).

E.g. $x_1 = \text{size in feet}^2$ $1m = 3.28 \text{ feet}$
 ~~$x_2 = \text{size in m}^2$~~
 $x_1 = (3.28)^2 x_2$

$$\begin{array}{l} \xrightarrow{m=10} \\ \xrightarrow{n=100} \\ \theta \in \mathbb{R}^{101} \end{array}$$

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

- Delete some features, or use regularization.

\downarrow later

Week 3 Logistic Regression

- Classification and Representation

1. Classification (Logistic)

$$y \in \{0, 1\} \quad 0: \text{"Negative Class"} \\ 1: \text{"Positive Class"}$$

discrete 离散的

⇒ tips

0: absence of something that we want to see

e.g. malignant tumor

1: presence ... , e.g. benign tumor

2. Hypothesis Representation

Logistic Regression Model

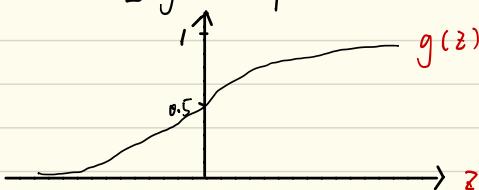
want $0 \leq h_{\theta}(x) \leq 1$

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

$$z: \theta^T x$$

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

called
| Sigmoid function
| Logistic function



$h_{\theta}(x)$: the probability when output is 1.

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

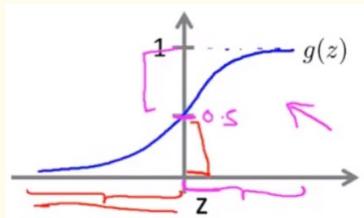
$$P(y=1|x; \theta) + P(y=0|x; \theta) = 1$$

Week 3 Logistic Regression

3. Decision Boundary

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

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



Decision Boundary, to separate the area
where $y=0$ and where $y=1$.

由 θ 決定 而不是 x 的結果.

$$h_{\theta}(x) \geq 0.5 \rightarrow y=1$$

$$h_{\theta}(x) < 0.5 \rightarrow y=0$$

Remember

$$z = 0, e^0 = 1 \Rightarrow g(z) = 1/2$$

$$z \rightarrow \infty, e^{-\infty} \rightarrow 0 \Rightarrow g(z) = 1$$

$$z \rightarrow -\infty, e^{\infty} \rightarrow \infty \Rightarrow g(z) = 0$$

二. Logistic Regression Model

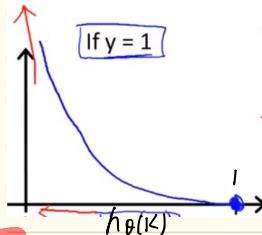
1. Cost Function

Week 3 Logistic Regression

Logistic regression cost function

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

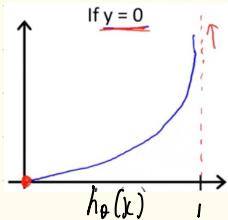
$y=1$



→ Cost = 0 if $y=1, h_\theta(x)=1$
But as $h_\theta(x) \rightarrow 0$
 $\text{Cost} \rightarrow \infty$

→ Captures intuition that if $h_\theta(x)=0$,
(predict $P(y=1|x;\theta)=0$), but $y=1$,
we'll penalize learning algorithm by a very
large cost.

$y=0$



cost = 0 if $y=0, h_\theta(x)=0$
But as $h_\theta(x) \rightarrow 1$

$\text{Cost} \rightarrow \infty$
if $h_\theta(x)=1$, (predict $P(y=0|x;\theta)=1$)
but $y=0 \Rightarrow$ a large cost

To Summarize

$\text{Cost}(h_\theta(x), y) = 0$ if $h_\theta(x) = y$

$\text{Cost}(h_\theta(x), y) \rightarrow \infty$ if $y=0$ and $h_\theta(x) \rightarrow 1$

$\text{Cost}(h_\theta(x), y) \rightarrow \infty$ if $y=1$ and $h_\theta(x) \rightarrow 0$

This cost function guarantees that $J(\theta)$ is convex for Logistic regression.

2. Simplified Cost Function and Gradient Descent

L-R cost function

Week 3

Logistic Regression

Logistic regression cost function

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

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

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

$y = 0 \text{ or } 1 \text{ always}$

Gradient Descent

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

repeat }

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

simultaneously update

同线性回归

→ 向量实现方式

$$\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \vec{y})$$

3. Advanced Optimization

Optimization algorithms:

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

Advantages:

- No need to manually pick α
- Often faster than gradient descent.

Disadvantages:

- More complex

→ 调用内置函数即可，算法过程见后图。

eg.

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

Week 3 Logistic Regression

$$\underline{\text{theta}} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \begin{array}{l} \text{theta}(1) \\ \text{theta}(2) \\ \vdots \\ \text{theta}(n+1) \end{array}$$

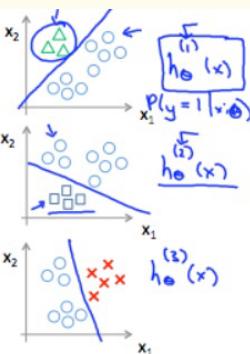
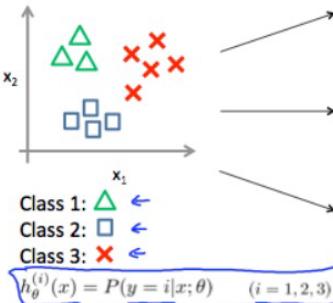
[θ₀ θ₁ ... θ_n] 值

```
function [jVal, gradient] = costFunction(theta)
    jVal = [code to compute  $J(\theta)$ ];
    gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ];
    gradient(2) = [code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ];
    :
    gradient(n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ];
```

三. Multiclass Classification

1. Multiclass Classification : One vs all

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



$$y \in \{0, 1, \dots, n\}, h_{\theta}^{(0)}(x) = P(y=0|x; \theta)$$

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

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

即，给定 x 预测 y 值。取 $h_{\theta}^{(i)}(x)$ 最大者

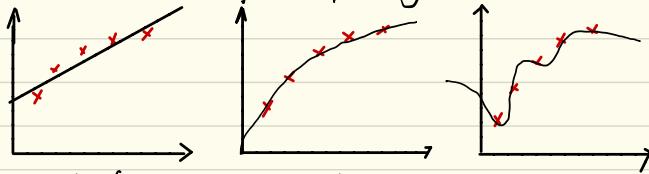
one vs rest

Week 3

Logistic Regression

④ Solving the Problem of Overfitting

1. The Problem of Overfitting



underfit ✓ overfit

function is too simple, too few features

overfitting (caused by a complicated function)
work well for training sets, but predict bad.

Address overfitting

Options:

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

2. Cost Function

若解决上图中的 overfitting, 则需 $\theta_3 \rightarrow 0, \theta_4 \rightarrow 0$

Ridge cost function

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 + 1000\theta_3^2 + 1000\theta_4^2$$

$$\Rightarrow \theta_3 \approx 0, \theta_4 \approx 0$$

当不知道哪项从何处减小 θ_i 时。

修正 Ridge cost function 如下：

使得向 2th 层发展

Week 3

Logistic Regression

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

λ too large \rightarrow underfit
 $(\theta_1, \dots, \theta_n) \perp, \approx 0$

$$h_{\theta}(k) = \theta_0$$

λ : regularization parameter

3. Regularized Linear Regression

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

① Gradient Descent

repeat

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

\downarrow 每一步缩小
 $\theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$

② Normal Equation

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

where

$$L = \begin{bmatrix} 0 & & & \\ 1 & 1 & & \\ & 1 & \ddots & \\ & & \ddots & 1 \end{bmatrix} \quad (n+1) \times (n+1) \text{ dimension}$$

note: if $m < n$, $X^T X$ 不可逆.

但 $+ \lambda L$ 一定可逆

4. Regularized Logistic Regression

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

Week 3 Logistic Regression

Gradient descent

Repeat {

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

separately

$$\Rightarrow \theta_j := \theta_j - \alpha \left[\underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\substack{(j = \text{X } 1, 2, 3, \dots, n) \\ \theta_0, \dots, \theta_n}} + \frac{\lambda}{m} \theta_j \right] \Leftarrow$$

}

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

Week 4

Neural Networks : Representation

- Motivations

1. Non-linear Hypotheses

too many features

2. Neurons and the Brain

Neural Network origins :

Algorithms that try to mimic the brain.

二. Neural Networks

1. Model Representation I

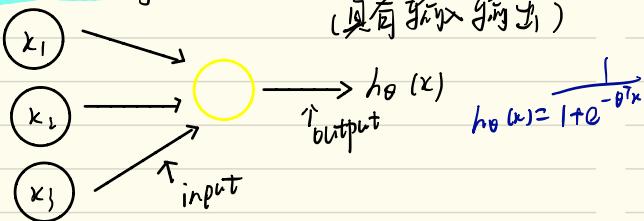
① Neurons in the brain

输入 / 通过神经元传播到输出。

抽象. 传出. (传递信息)

② Neural Network

model : Logistic unit



Sigmoid (logistic) activation function.

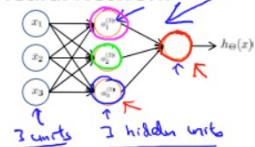
$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} = \text{"weight"}$$

Neural Network

Week 4

Neural Networks : Representation

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 $j+1$

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

Subscript 表示第几层，
i 层

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

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

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

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

(2)

↓

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

2. Model Representation II

Vectorized implementation

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \end{bmatrix}$$

$$z^{(1)} = \Theta^{(0)}x \Rightarrow a^{(1)} = g(z^{(1)})$$

Setting $x = a^{(0)}$

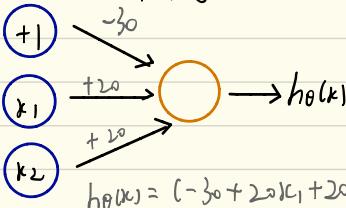
$$\Rightarrow z^{(j)} = \Theta^{(j-1)}a^{(j-1)} \quad a^{(j)} = g(z^{(j)})$$

3. Applications

1. Examples and Intuitions 1

Simple example : AND

$x_1, x_2 \in \{0, 1\}$, $y = x_1 \text{ AND } x_2$

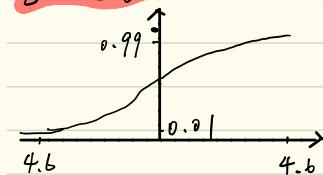


$$h_{\Theta}(x) = (-30 + 20x_1 + 20x_2)$$

Week 4

Neural Networks: Representation

g(z) 回答



x_1	x_2	$h_{\Theta}(x)$
0	0	$g(-3.0) \approx 0$
0	1	$g(-1.0) \approx 0$
1	0	$g(-1.0) \approx 0$
1	1	$g(1.0) \approx 1$

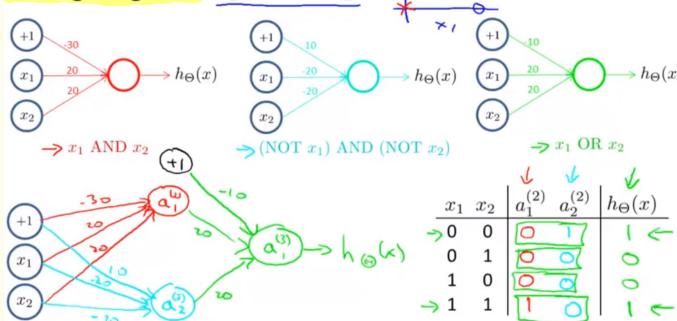
化简逻辑表达式。
取常数 0, 1 观察

Q: $\hat{y}_{\Theta}^0 \Theta = [-10 \ 20 \ 20]$, 答什么? (x_1 or x_2)

2. Examples and Intuition II

Example: XNOR (异或加. 反之为 1)

Putting it together: x_1 XNOR x_2



3. Multiclass Classification

Multiple output units: One-vs-all.



Want $h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, etc.
when pedestrian when car when motorcycle

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

$y^{(i)}$ one of $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$ 表示第 i 个类

类别

输出层神经元 n ↑
表示第 i 个类

Week 4

Neural Networks : Representation

课后题

Which of the following statements are true? Check all that apply.

- A two layer (one input layer, one output layer; no hidden layer) neural network can represent the XOR function.

需进行叠不可延年

- Any logical function over binary-valued (0 or 1) inputs x_1 and x_2 can be (approximately) represented using some neural network.

- Suppose you have a multi-class classification problem with three classes, trained with a 3 layer network. Let $a_1^{(3)} = (h_\theta(x))_1$ be the activation of the first output unit, and similarly $a_2^{(3)} = (h_\theta(x))_2$ and $a_3^{(3)} = (h_\theta(x))_3$. Then for any input x , it must be the case that $\underline{a_1^{(3)} + a_2^{(3)} + a_3^{(3)} = 1}$.

由前层计算所得向量各元素之和

- The activation values of the hidden units in a neural network, with the sigmoid activation function applied at every layer, are always in the range $(0, 1)$.

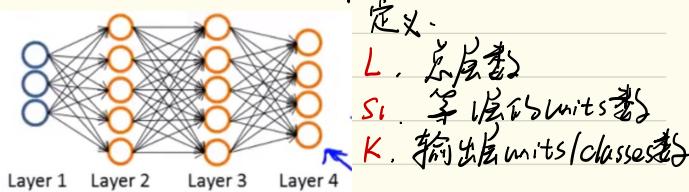
Week 5

Neural Networks : Learning

1. Cost Function and Backpropagation

1. Cost Function

Neural Network (Classification)



k 分类模型是 k 个 0-1 矢量，则有 k 分类代价函数 $J_{\theta}(x)$ [输出层有 k 个 ...]

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K (y_k^{(i)} \log(h_{\theta}(x^{(i)})_k) + (1-y_k^{(i)}) \log(1-h_{\theta}(x^{(i)})_k)) \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{j,i}^{(l)})^2$$

注：日知每列数为当前层 node 个数。

行数为下一层 node 个数。

Note

- the double sum simply adds up the logistic regression costs calculated for each cell in the output layer
- the triple sum simply adds up the squares of all the individual θ s in the entire network.
- the i in the triple sum does not refer to training example i

2. Backpropagation

Gradient computation

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

cost function

$$\min_{\theta} (J(\theta)) \quad \text{need code to compute} \begin{cases} J(\theta) \\ \frac{\partial}{\partial \theta_{j,i}^{(l)}} J(\theta) \end{cases}$$

Week 5

Neural Networks : Learning

Backpropagation algorithm

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

for total L layer. $\delta_j^{(L)} = a_j^{(L)} - y_j$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)}) \rightarrow a^{(3)} \cdot (1-a^{(3)})$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \cdot g'(z^{(2)}) \rightarrow a^{(2)} \cdot (1-a^{(2)})$$

$$\Rightarrow \delta^{(1)} = (\Theta^{(1)})^T \delta^{(2)} \cdot a^{(1)} \cdot (1-a^{(1)})$$

Algorithm

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

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j). (use to compute $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$)

For $i = 1$ to m $\leftarrow (x^{(i)}, y^{(i)})$.

- Set $a^{(1)} = x^{(i)}$
- Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \dots, L$
- Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$
- Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ ~~Take~~
- $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ ~~Take~~

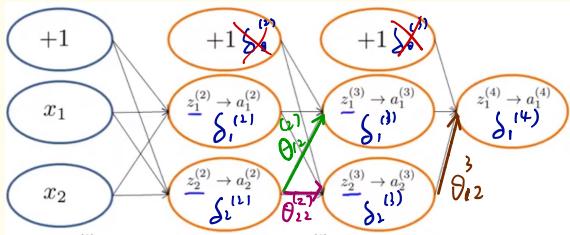
$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \text{ if } j \neq 0$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \text{ if } j = 0$$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$$

$\Delta = \delta \text{ is } \text{大} \text{ 倍}$

3. Backpropagation Intuition.



$$\delta_1^{(4)} = y^{(1)} - a_1^{(4)} \quad (\text{for example } i)$$

$$\delta_2^{(2)} = \theta_{12}^{(2)} \delta_1^{(3)} + \theta_{22}^{(2)} \delta_2^{(3)}$$

$$\delta_2^{(1)} = \theta_{12}^{(1)} \delta_1^{(2)}$$

:

Week 5

Neural Networks : Learning

$K = 1$ 分类问题, for example +
 $\text{cost}(t) = y^{(t)} \log(h_\theta(x^{(t)})) + (1 - y^{(t)}) \log(1 - h_\theta(x^{(t)}))$
 $\Rightarrow \delta_j^{(t)} = \frac{\partial}{\partial z_j^{(t)}} \text{cost}(t)$? 不太懂

2. Backpropagation in Practice

1. Implementation Note: Unrolling Parameters

Unroll into a long vector

e.g. $\text{thetaVector} = [\Theta_1(:); \Theta_2(:); \Theta_3(:)]$

get back: $\Theta_1 = \text{reshape}(\text{thetaVector}(1:10), 10, 1)$

取1到10到组成10x1矩阵

$\hat{A}_k - \hat{y}_j$

To summarize:

Learning Algorithm

- Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.
- Unroll to get initialTheta to pass to
- fminunc(@costFunction, initialTheta, options)

```
function [jval, gradientVec] = costFunction(thetaVec)
    From thetaVec, get  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ .
    Use forward prop/back prop to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\Theta)$ .
    Unroll  $D^{(1)}, D^{(2)}, D^{(3)}$  to get gradientVec.
```

2. Gradient Checking

从直觉方法估测梯度. 与 backprop 计算出的对比.

发现 subtle bugs

$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \text{采用中值法估计梯度} \\ = \text{gradApprox}$$

算法过程

```
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; check gradApprox ≈ Dvec
```

↑ from backprop

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_i + \epsilon \\ \theta_n \end{bmatrix}$$

Week 5

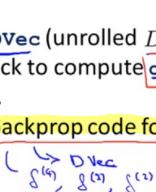
Neural Networks : Learning

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.

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

Random initialization: Symmetry breaking

Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$
(i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)

E.g. random 10x11 matrix between $-\epsilon$ and ϵ

```
Theta1 = rand(10,11)*(2*INIT_EPSILON)  
- INIT_EPSILON;
```

$[-\epsilon, \epsilon]$

```
Theta2 = rand(1,11)*(2*INIT_EPSILON)  
- INIT_EPSILON;
```

Why zeros(n, 1)?

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

4. Putting It Together

Training a neural network (a network archt)

- ① No. of input units : Dimension of features $x^{(i)}$
- ② No. of output units : Number of classes
- ③ Reasonable default : 1 hidden layer

or > 1 hidden & have same no. of hidden units in every layer.

no initial weight

Week 5

Neural Networks : Learning

Training a Neural Network

1. Randomly initialize the weights
2. Implement forward propagation to get $h_{\Theta}(x^{(i)})$ for any $x^{(i)}$
3. Implement the cost function
4. Implement backpropagation to compute partial derivatives
5. Use gradient checking to confirm that your backpropagation works. Then disable gradient checking.
6. Use gradient descent or a built-in optimization function to minimize the cost function with the weights in theta.

III. Application of Neural Networks

1. Autonomous Driving Example