1. 分类问题。首先考虑简单情况，即y取值0,1。y={0,1}。Negative Class/Positive Class

## Classification

→ Email: Spam / Not Spam?

→ Online Transactions: Fraudulent (Yes / No)?
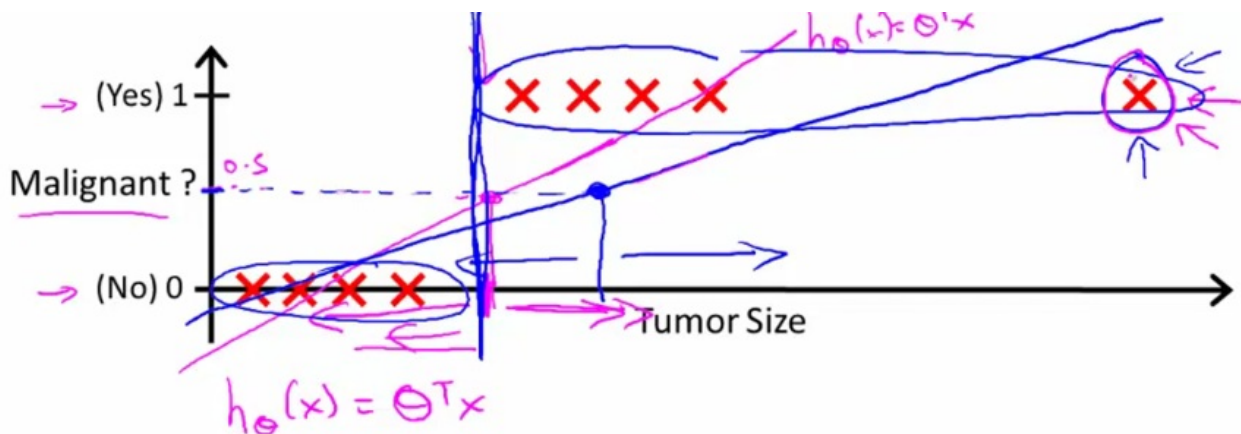
→ Tumor: Malignant / Benign ?

→ $y \in \{0, 1\}$　　0: "Negative Class" (e.g., benign tumor)

　　　　　　　　　　1: "Positive Class" (e.g., malignant tumor)

→ $y \in \{0, 1, 2, 3\}$

我们以后再关心多类的问题

考虑用线性回归来解决分类问题，阈值分类器，根据阈值来划分：



$h_\theta(x) = \Theta^T x$

→ Threshold classifier output $h_\theta(x)$ at 0.5:

→ If $h_\theta(x) \geq 0.5$, predict "y = 1"

If $h_\theta(x) < 0.5$, predict "y = 0"

线性回归算法来解决分类问题

Andrew N

但这显然不是一个好主意。另外，h(x)的范围有可能也不止在[0, 1]。

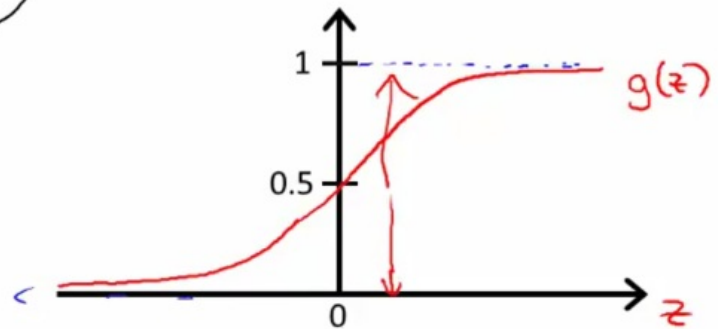2. 逻辑回归（Logistic Regression）。S型函数/逻辑函数。

## Logistic Regression Model

Want $\boxed{0 \le h_\theta(x) \le 1}$

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

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

$\theta^T x$

$\to$ Sigmoid function
$\hookrightarrow$ Logistic function

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



$g(z)$

Parameters $\theta$.

一下这个模型的解释

逻辑函数h(x)值的意义为：y=1的概率。即：

## Interpretation of Hypothesis Output

$h_\theta(x)$

$h_\theta(x)$ = estimated probability that $\boxed{y = 1}$ on input x $\leftarrow$

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

$$h_\theta(x) = 0.7 \qquad y=1$$

Tell patient that 70% chance of tumor being malignant

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

"probability that y = 1, given x, parameterized by $\theta$"

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

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

这就是说

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

3. 决策边界（Decision Boundary）。g(z)>=0 when z>=0; whenever theta'x >=0。

## Logistic regression
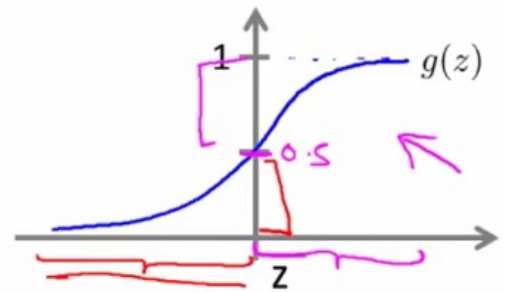
$$h_\theta(x) = g(\theta^T x) = P(y=1 \mid x; \theta)$$

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

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

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

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

$$h_\theta(x) = g(\theta^T x)$$
$$\theta^T x < 0$$

我们就预测y等于0

$$g(z) \geq 0.5$$
$$\text{when } z \geq 0$$
$$h_\theta(x) = g(\theta^T x) \geq 0.5$$
$$\text{whenever } \theta^T x \geq 0$$
$$z$$

$$g(z) < 0.5$$

决策边界是参数的性质而不是训练集的性质。事实上，决策边界就是：theta'x = 0。在决策边界右边被预测为1，左边被预测为0。

## Decision Boundary

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

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

-3    1    1

Decision boundary

Predict "$y = 1$" if $-3 + x_1 + x_2 \geq 0$

$$\theta^T x$$

$$x_1 + x_2 \geq 3$$
$$x_1, x_2$$
$$h_\theta(x) = 0.5$$
$$x_1 + x_2 = 3$$

$$x_1 + x_2 < 3$$
$$y = 0$$

在这幅图中 我画了一个训练集

非线性决策边界：

## Non-linear decision boundaries



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$

Predict "$y = 1$" if $\;-1 + x_1^2 + x_2^2 \geq 0$

$x_1^2 + x_2^2 = 1$

$x_1^2 + x_2^2 \geq 1$

## Non-linear decision boundaries



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ \vdots \end{bmatrix}$$

Predict "$y = 1$" if $\;-1 + x_1^2 + x_2^2 \geq 0$

$x_1^2 + x_2^2 = 1$

$x_1^2 + x_2^2 \geq 1$

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_6 x_1^3 x_2 + \ldots)$$

因此 通过这些可视化图形

但是，如何自动选择参数theta，以便给定一个数据集，可以根据数据自动拟合参数呢？

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

m examples $\quad x \in \begin{bmatrix} x_0 \\ x_1 \\ \cdots \\ x_n \end{bmatrix} \mathbb{R}^{n+1} \quad x_0 = 1, y \in \{0, 1\}$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$
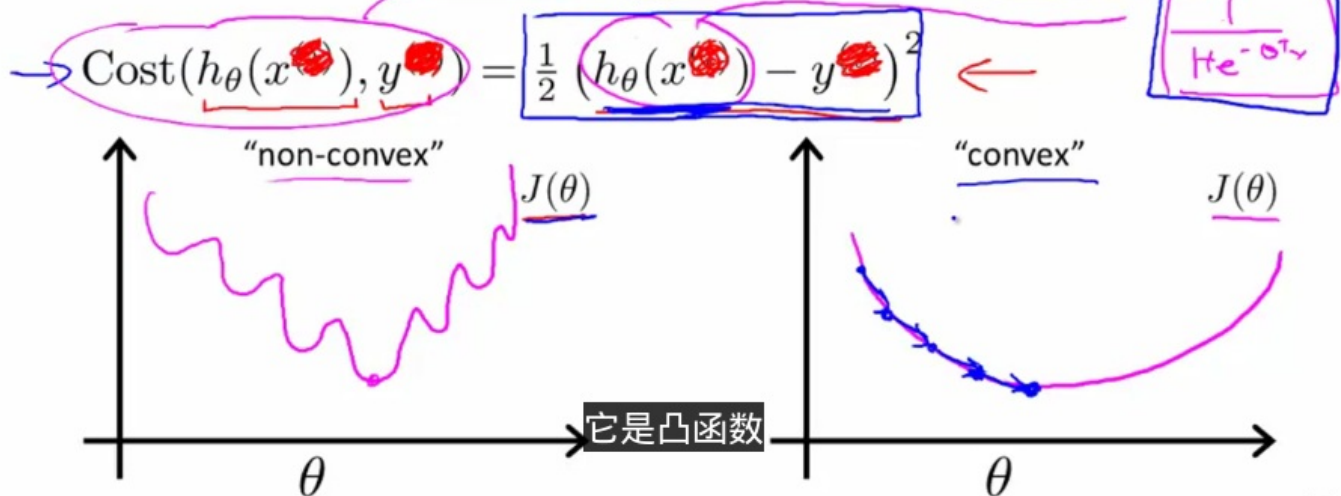
## How to choose parameters $\theta$ ?

4. 代价函数。分类问题代价函数是非凸函数。

## Cost function

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

logistic

$\rightarrow \text{Cost}(h_\theta(x^{(i)}), y)$

$\rightarrow \text{Cost}(h_\theta(x), y) = \frac{1}{2}\left(h_\theta(x) - y\right)^2 \leftarrow$

$\frac{1}{1+e^{-\theta^T x}}$

"non-convex"

"convex"

$J(\theta)$

$J(\theta)$

它是凸函数

$\theta$

$\theta$
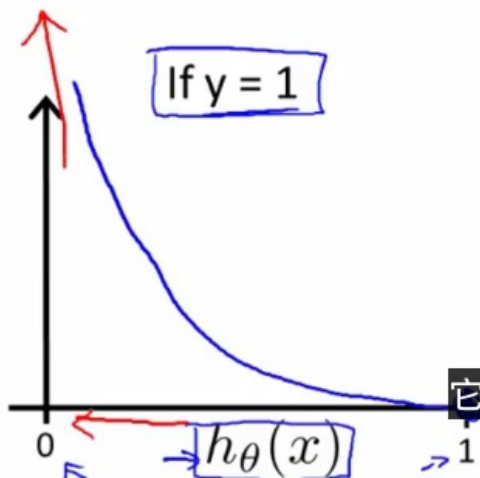
寻找代价函数为凸函数：

在y=1的时

## Logistic regression cost function

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

If y = 1

→ Cost = 0 if $y = 1, h_\theta(x) = 1$
But as $h_\theta(x) \to 0$
$Cost \to \infty$

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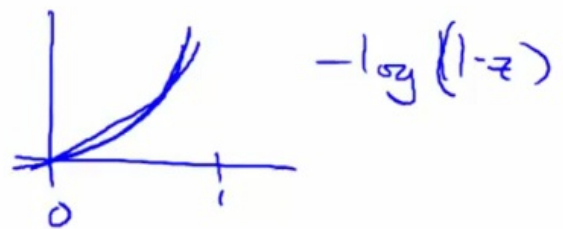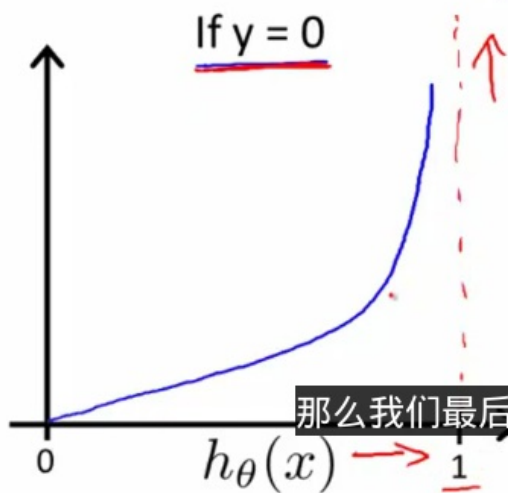we'll penalize learning algorithm by a very

它是被这样体现出来

$h_\theta(x)$

0        1

候：

在y=0的时候：

## Logistic regression cost function

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

If y = 0

$-\log(1-z)$

那么我们最后就要付出非常大的代价值

$h_\theta(x) \longrightarrow 1$

0

Andrew

显而易见，这样的代价函数是凸函数，没有局部最优值。

5. 实现逻辑回归。

## Logistic regression cost function

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

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

To fit parameters $\theta$:

$$\min_\theta J(\theta)$$ Get $\theta$

To make a prediction given new $x$:

Output $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$  $p(y=1 \mid x; \theta)$

你就把这个想成

求偏导，计算代价函数，发现和线性回归梯度下降更新规则相同！

## Gradient Descent

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

Want $\min_\theta J(\theta)$:

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

Repeat {

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

(simultaneously update all $\theta_j$)

}

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

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

Algorithm looks identical to linear regression!

在先前的视频中

6. 优化算法。

## Optimization algorithm

Given $\theta$, we have code that can compute

- $J(\theta)$ ←
- $\frac{\partial}{\partial \theta_j} J(\theta)$ ←  (for $j = 0, 1, \ldots, n$)

Optimization algorithms:
→ - Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

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

Disadvantages:
- More complex

称为线性搜索(line search)算法 它可以自动

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

→ $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$   $\theta_1 = 5, \theta_2 = 5.$

→ $J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$

↪ $\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$

↪ $\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$

```
function [jVal, gradient]
            = costFunction(theta)
jVal = (theta(1)-5)^2 + ...
            (theta(2)-5)^2;
gradient = zeros(2,1);
gradient(1) = 2*(theta(1)-5);
gradient(2) = 2*(theta(2)-5);
```

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

```
1.    function [jVal, gradient] = costFunction(theta)
2.    jVal = (theta(1)-5)^2+(theta(2)-5)^2;
3.    gradient = zeros(2,1);
4.    gradient(1) = 2*(theta(1)-5);
5.    gradient(2) = 2*(theta(2)-5);
6.
7.    octave:1> options = optimset('GradObj', 'on', 'MaxIter', '100');
8.    octave:2> initialTheta = zeros(2,1);
9.    octave:3> [optTheta, functionVal, exitFlag]=fminunc(@costFunction, initialTheta, options)
10.   optTheta =
11.
12.       5.0000
13.       5.0000
14.
15.   functionVal =    7.8886e-31
16.   exitFlag =  1
```

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

theta(1) ←

theta(2)

theta(n+1)

```
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 $\frac{\partial}{\partial \theta_1} J(\theta)$] ;

$\vdots$

gradient(n+1) = 因为 Octave 的标号 te $\frac{\partial}{\partial \theta_n} J(\theta)$ ] ;

Suppose you want to use an advanced optimization algorithm to minimize the cost function for logistic regression with parameters $\theta_0$ and $\theta_1$. You write the following code:

```
function [jVal, gradient] = costFunction(theta)
  jVal = % code to compute J(theta)
  gradient(1) = CODE#1 % derivative for theta_0
  gradient(2) = CODE#2 % derivative for theta_1
```

What should CODE#1 and CODE#2 above compute?

○ CODE#1 and CODE#2 should compute $J(\theta)$.

○ CODE#1 should be theta(1) and CODE#2 should be theta(2).

◉ CODE#1 should compute $\frac{1}{m} \sum_{i=1}^{m} [(h_\theta(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}] (= \frac{\partial}{\partial \theta_0} J(\theta))$, and

CODE#2 should compute $\frac{1}{m} \sum_{i=1}^{m} [(h_\theta(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}] (= \frac{\partial}{\partial \theta_1} J(\theta))$.

**Correct Response**

○ None of the above.

因此可以自己写代价函数，从而使用octave高级优化的方法。

7. 逻辑回归的多分类方法。y取值不再是0,1,而可以是离散的多个。

## Multiclass classification

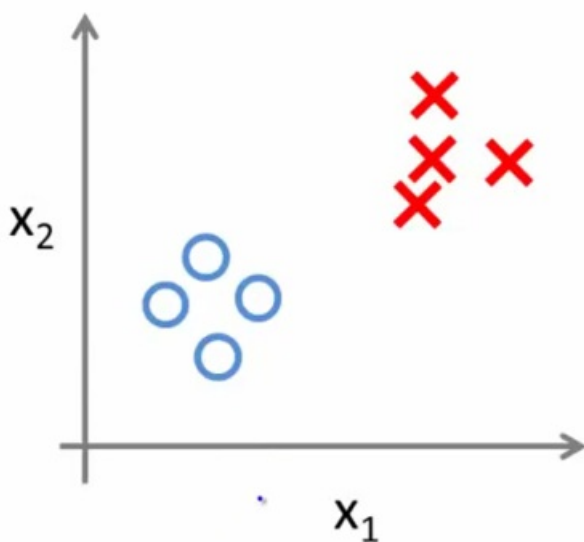Email foldering/tagging: Work, Friends, Family, Hobby

$$y=1 \quad y=2 \quad y=3 \quad y=4$$
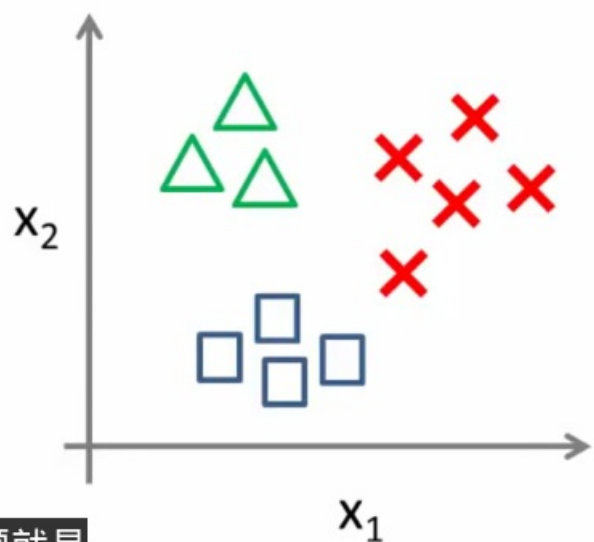
Medical diagrams: Not ill, Cold, Flu

$$y=1 \quad 2 \quad 3$$

Weather: Sunny, Cloudy, Rain, Snow

$$y=1 \quad 2 \quad 3 \quad 4 \quad \Leftarrow$$

Binary classification:

Multi-class classification:

$x_2$     $x_1$

问题就是

$x_2$     $x_1$

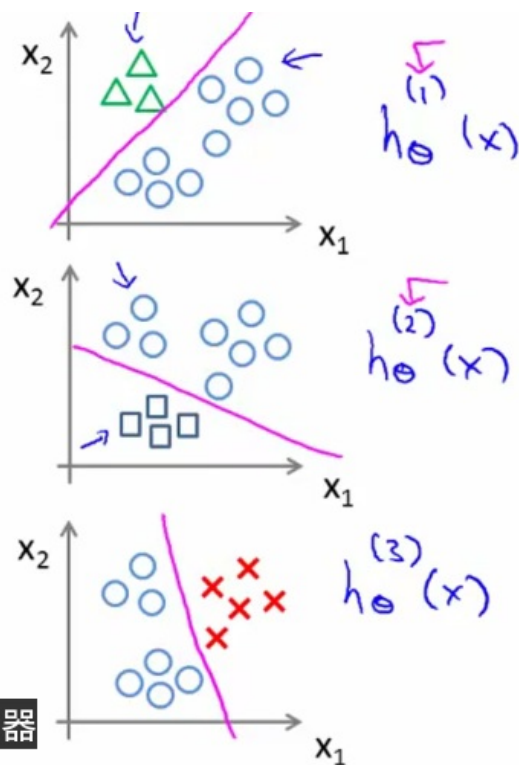这种问题可以转化成y取二值的逻辑回归分析：即从one-vs-all -> one-vs-rest。

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



Class 1: △ ←
Class 2: □ ←
Class 3: ✗ ←

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

我们都找到了一个分类器

从而将一个三元分类问题转换成三次一元分类。

**One-vs-all**

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

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

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