

# Week 3

- **Classification and Representation:**

- a. Classification:

E.g. email: Spam/not spam?

Online: Transactions: fraudulent or not?

Tumor: malignant/benign?

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

negative class (e.g. benign tumor)  
positive class (e.g. malignant tumor)

How do we develop a classification algorithm?

1. solution: linear regression:

Classification:  $y=0$  or  $1$  but the hypothesis can be larger than  $1$  or smaller than  $0$ .

2. Solution: logistic regression:  $0 <= h_{\theta}(x) <= 1$

- b. Hypothesis Representation:

Our aim:  $0 <= h_{\theta}(x) <= 1$

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

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

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

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

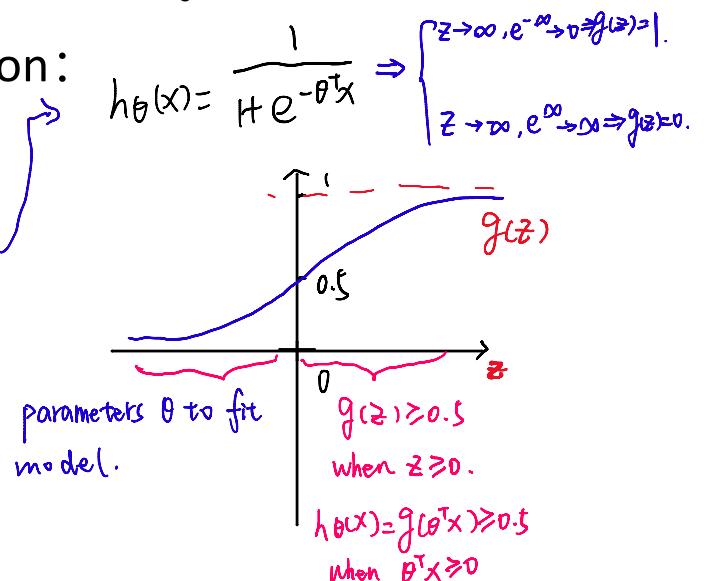
$$h_{\theta}(x) = p(y=1|x; \theta).$$

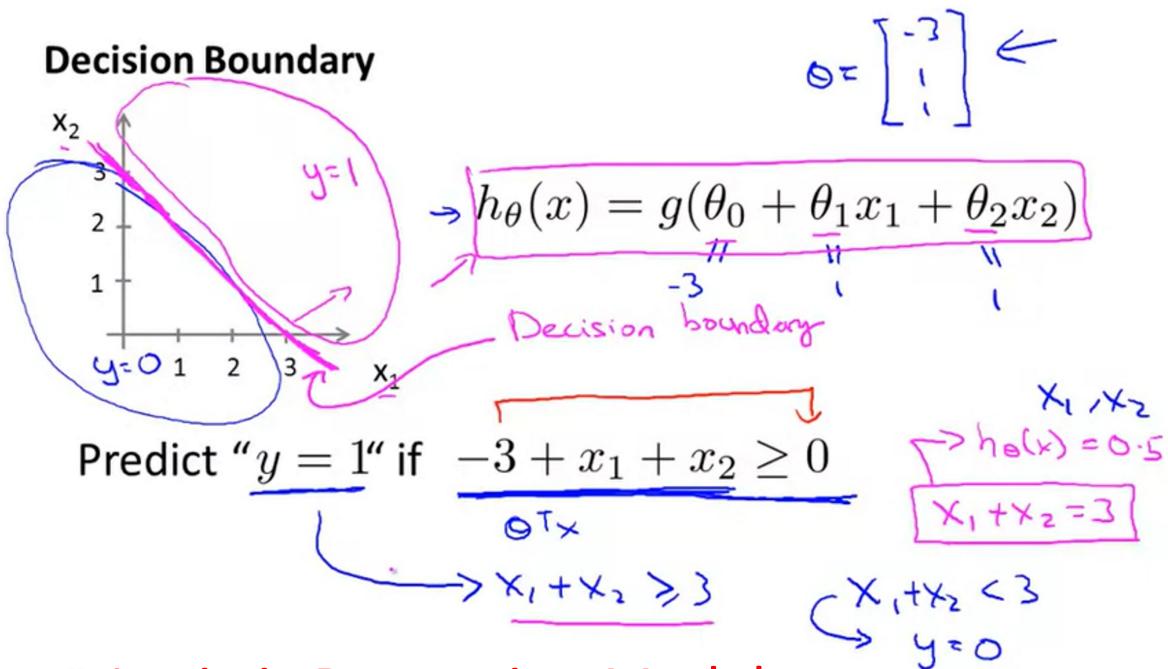
$$P(y=1|x; \theta) + P(y=0|x; \theta) = 1.$$

c. Decision Boundary: → property of hypothesis under parameters.

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

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





## • Logistic Regression Model:

### a. Cost Function:

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

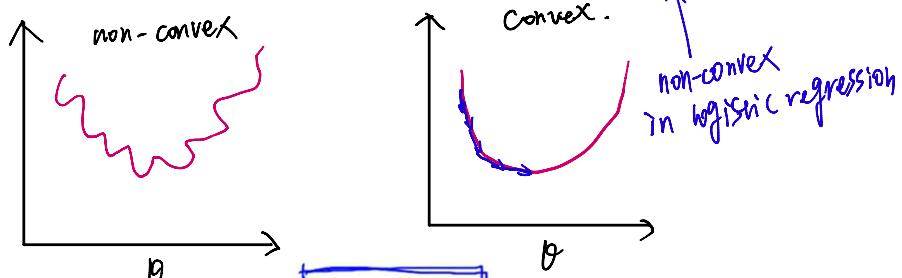
m examples

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

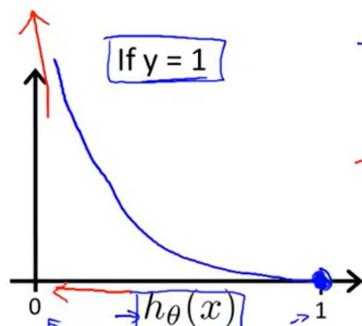
$$h_\theta = \frac{1}{1 + e^{-\theta^T x}}$$

Linear regression:  $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} \text{cost}(h_\theta(x^{(i)}), y^{(i)})^2$

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

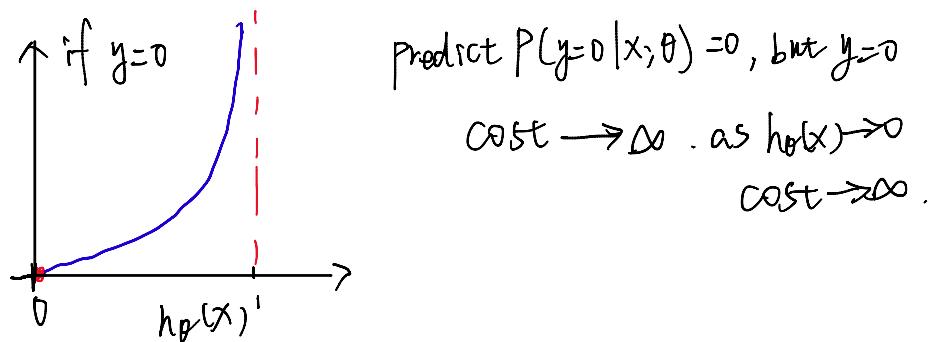


$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$



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

$\rightarrow$  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.



## b. Simplified Cost Function and Gradient Descent:

*Simplified form:*

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

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)})) \right] \end{aligned}$$

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

Repeat {

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

(simultaneously update all  $\theta_j$ )

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

linear regression

$$\begin{aligned} h_\theta(x) &= \theta^T x \\ \rightarrow h_\theta(x) &= \frac{1}{1+e^{-\theta^T x}} \end{aligned}$$

logistic regression.

Algorithm looks identical to linear regression!

## c. Advanced Optimization:

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

### Advantages:

- no need to manually pick  $\alpha$
- Often faster than GD

### Disadvantages:

- more complex

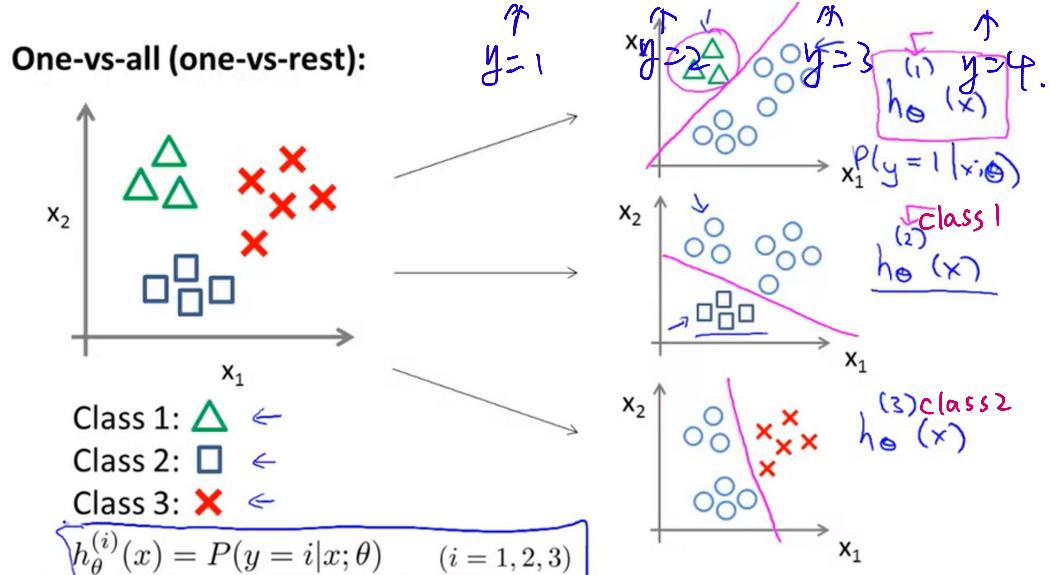
## • Multiclass Classification:

a. One-vs-all:

Email foldering/tagging: work, friends, family, hobby

Medical diagrams: not ill, cold, flu.

Weather: sunny, cloudy, rain, snow



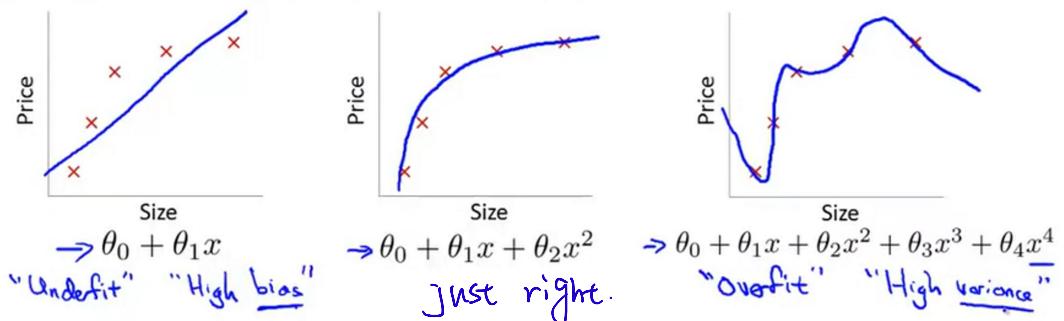
On a new input  $x$ , to make a prediction, pick the class  $i$  that maximizes

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

## • Solving the Problem of Overfitting:

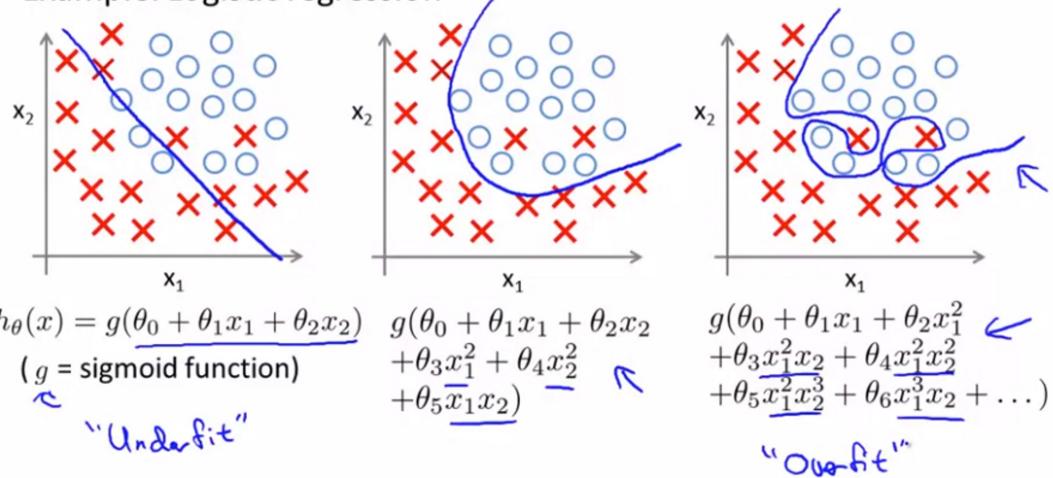
a. The problem of Overfitting:

Example: Linear regression (housing prices)



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

### Example: Logistic regression



A lot of features and little training examples may result in Overfitting.

Options:

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

## b. Cost Function:

### Regularization.

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

- "Simpler" hypothesis
- Less prone to overfitting

$$\theta_3, \theta_4 \rightarrow \approx 0$$

Housing:

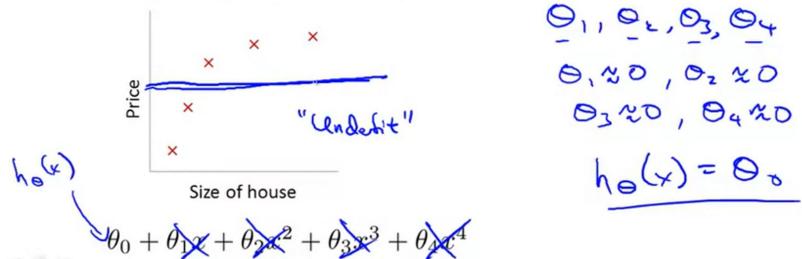
- Features:  $x_1, x_2, \dots, x_{100}$
- Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

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

regularization term.  
 regularization parameter.  
 fit training data  
 keep  $\theta$  small  
 control a trade off between two different goals.

modify cost function to shrink all parameters

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



### c. Regularized Linear Regression:

#### Gradient descent

Repeat {

we don't penalize  $\theta_0$

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

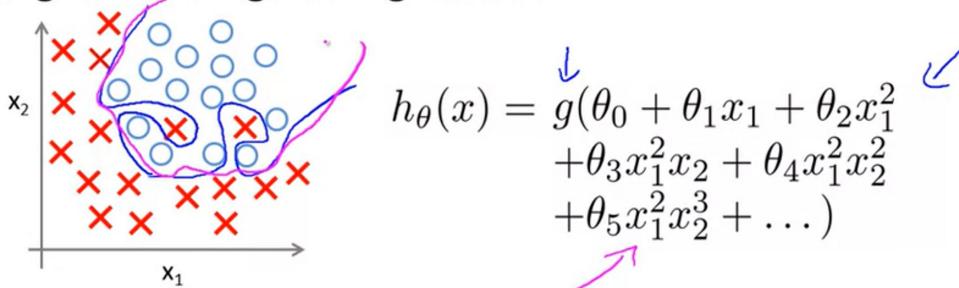
$$\begin{aligned} \theta_j &:= \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \\ &\quad \text{update as usual.} \\ \theta_j &:= \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \\ &\quad \text{shrink } \theta_j \text{ a little bit} \end{aligned}$$

#### Normal equation

$$\begin{aligned} X &= \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \in \mathbb{R}^{m \times (n+1)} & y &= \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m \\ \rightarrow \min_{\theta} J(\theta) & \quad \frac{\partial}{\partial \theta_j} J(\theta) \stackrel{\text{set } 0}{=} 0 \quad \text{for } j = 0, 1, \dots, n \\ \rightarrow \theta &= (X^T X + \lambda \underbrace{\begin{bmatrix} 0 & 1 & \cdots & 0 \\ 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix}}_{(n+1) \times (n+1)})^{-1} X^T y \end{aligned}$$

### d. Regularized Logistic Regression:

## Regularized logistic regression.



Cost function:

$$\rightarrow J(\theta) = - \left[ \frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad | \boxed{\theta_0, \theta_1, \dots, \theta_n}$$

## Gradient descent

Repeat {

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

$$\rightarrow \theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad \begin{matrix} (j = 1, 2, 3, \dots, n) \\ \theta_0, \dots, \theta_n \end{matrix}$$

}

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

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

treat  $\theta_0$  separately.

## Advanced optimization

```

→ function [jVal, gradient] = costFunction(theta)
    jVal = [ code to compute J(theta) ];
    →  $J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$ 
    → gradient(1) = [ code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$  ];
    →  $\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$ 
    → gradient(2) = [ code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$  ];
    →  $\left( \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} \right) + \frac{\lambda}{m} \theta_1$ 
    → gradient(3) = [ code to compute  $\frac{\partial}{\partial \theta_2} J(\theta)$  ];
    :    $\left( \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} \right) + \frac{\lambda}{m} \theta_2$ 
    gradient(n+1) = [ code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$  ];
  
```